

Package ‘tidyclust’

September 25, 2023

Title A Common API to Clustering

Version 0.2.0

Description A common interface to specifying clustering models, in the same style as 'parsnip'. Creates unified interface across different functions and computational engines.

License MIT + file LICENSE

URL <https://github.com/tidymodels/tidyclust>,
<https://tidyclust.tidymodels.org/>

BugReports <https://github.com/tidymodels/tidyclust/issues>

Depends R (>= 3.6)

Imports cli (>= 3.0.0), dials (>= 1.1.0), dplyr (>= 1.0.9), flexclust (>= 1.3-6), foreach, generics (>= 0.1.2), glue (>= 1.6.2), hardhat (>= 1.0.0), modelenv (>= 0.1.0), parsnip (>= 1.0.2), prettyunits (>= 1.1.0), Rfast (>= 2.0.6), rlang (>= 1.0.6), rsample (>= 1.0.0), stats, tibble (>= 3.1.0), tidyr (>= 1.2.0), tune (>= 1.0.0), utils, vctrs (>= 0.5.0)

Suggests cluster, ClusterR, clustMixType (>= 0.3-5), covr, klaR, knitr, modeldata (>= 1.0.0), RcppHungarian, recipes (>= 1.0.0), rmarkdown, testthat (>= 3.0.0), workflows (>= 1.1.2)

Config/Needs/website pkgdown, tidymodels, tidyverse, palmerpenguins, patchwork, ggforce, tidyverse/tidytemplate

Config/testthat/edition 3

Encoding UTF-8

RoxygenNote 7.2.3

NeedsCompilation no

Author Emil Hvitfeldt [aut, cre] (<<https://orcid.org/0000-0002-0679-1945>>),
Kelly Bodwin [aut],
Posit Software, PBC [cph, fnd]

Maintainer Emil Hvitfeldt <emil.hvitfeldt@posit.co>

Repository CRAN

Date/Publication 2023-09-25 18:20:06 UTC

R topics documented:

augment.cluster_fit	2
cluster_fit	3
cluster_metric_set	4
cluster_spec	5
control_cluster	7
cut_height	8
extract-tidyclust	8
extract_centroids	9
extract_cluster_assignment	11
extract_fit_summary	12
finalize_model_tidyclust	13
fit.cluster_spec	14
get_centroid_dists	16
glance.cluster_fit	16
hier_clust	17
k_means	18
linkage_method	19
min_grid.cluster_spec	20
new_cluster_metric	20
predict.cluster_fit	21
prep_data_dist	23
reconcile_clusterings_mapping	23
set_args.cluster_spec	24
set_engine.cluster_spec	25
set_mode.cluster_spec	25
silhouette	26
silhouette_avg	27
sse_ratio	28
sse_total	29
sse_within	30
sse_within_total	31
tidy.cluster_fit	32
translate_tidyclust	32
tune_cluster	33
update.hier_clust	35

Index	37
--------------	-----------

augment.cluster_fit	<i>Augment data with predictions</i>
---------------------	--------------------------------------

Description

augment() will add column(s) for predictions to the given data.

Usage

```
## S3 method for class 'cluster_fit'
augment(x, new_data, ...)
```

Arguments

x	A <code>cluster_fit</code> object produced by <code>fit.cluster_spec()</code> or <code>fit_xy.cluster_spec()</code> .
new_data	A data frame or matrix.
...	Not currently used.

Details

For partition models, a `.pred_cluster` column is added.

Value

A `tibble::tibble()` with containing `new_data` with columns added depending on the mode of the model.

Examples

```
kmeans_spec <- k_means(num_clusters = 5) %>%
  set_engine("stats")

kmeans_fit <- fit(kmeans_spec, ~., mtcars)

kmeans_fit %>%
  augment(new_data = mtcars)
```

cluster_fit

Model Fit Object Information

Description

An object with class "cluster_fit" is a container for information about a model that has been fit to the data.

Details

The following model types are implemented in `tidyclust`:

- K-Means in `k_means()`
- Hierarchical (Agglomerative) Clustering in `hier_clust()`

The main elements of the object are:

- `spec`: A `cluster_spec` object.

- `fit`: The object produced by the fitting function.
- `preproc`: This contains any data-specific information required to process new a sample point for prediction. For example, if the underlying model function requires arguments `x` and the user passed a formula to `fit`, the `preproc` object would contain items such as the `terms` object and so on. When no information is required, this is `NA`.

As discussed in the documentation for `cluster_spec`, the original arguments to the specification are saved as quosures. These are evaluated for the `cluster_fit` object prior to fitting. If the resulting model object prints its call, any user-defined options are shown in the call preceded by a tilde (see the example below). This is a result of the use of quosures in the specification.

This class and structure is the basis for how **tidyclust** stores model objects after seeing the data and applying a model.

`cluster_metric_set` *Combine metric functions*

Description

`cluster_metric_set()` allows you to combine multiple metric functions together into a new function that calculates all of them at once.

Usage

```
cluster_metric_set(...)
```

Arguments

... The bare names of the functions to be included in the metric set. These functions must be cluster metrics such as `sse_total()`, `sse_ratio()`, or `silhouette_avg()`.

Details

All functions must be:

- Only cluster metrics

Value

A `cluster_metric_set()` object, combining the use of all input metrics.

Description

An object with class "cluster_spec" is a container for information about a model that will be fit.

Details

The following model types are implemented in tidyclust:

- K-Means in `k_means()`
- Hierarchical (Agglomerative) Clustering in `hier_clust()`

The main elements of the object are:

- `args`: A vector of the main arguments for the model. The names of these arguments may be different from their counterparts in the underlying model function. For example, for a `k_means()` model, the argument name for the number of clusters are called "num_clusters" instead of "k" to make it more general and usable across different types of models (and to not be specific to a particular model function). The elements of `args` can be tuned with the use of `tune_cluster()`.

For more information see <https://www.tidymodels.org/start/tuning/>. If left to their defaults (NULL), the arguments will use the underlying model functions default value. As discussed below, the arguments in `args` are captured as quosures and are not immediately executed.

- `...`: Optional model-function-specific parameters. As with `args`, these will be quosures and can be tuned with `tune()`.
- `mode`: The type of model, such as "partition". Other modes will be added once the package adds more functionality.
- `method`: This is a slot that is filled in later by the model's constructor function. It generally contains lists of information that are used to create the fit and prediction code as well as required packages and similar data.
- `engine`: This character string declares exactly what software will be used. It can be a package name or a technology type.

This class and structure is the basis for how **tidyclust** stores model objects prior to seeing the data.

Argument Details

An important detail to understand when creating model specifications is that they are intended to be functionally independent of the data. While it is true that some tuning parameters are *data dependent*, the model specification does not interact with the data at all.

For example, most R functions immediately evaluate their arguments. For example, when calling `mean(dat_vec)`, the object `dat_vec` is immediately evaluated inside of the function.

`tidyclust` model functions do not do this. For example, using

```
k_means(num_clusters = ncol(mtcars) / 5)
```

does not execute `ncol(mtcars) / 5` when creating the specification. This can be seen in the output:

```
> k_means(num_clusters = ncol(mtcars) / 5)
K Means Cluster Specification (partition)
```

Main Arguments:

```
num_clusters = ncol(mtcars)/5
```

Computational engine: stats

The model functions save the argument *expressions* and their associated environments (a.k.a. a quosure) to be evaluated later when either `fit.cluster_spec()` or `fit_xy.cluster_spec()` are called with the actual data.

The consequence of this strategy is that any data required to get the parameter values must be available when the model is fit. The two main ways that this can fail is if:

1. The data have been modified between the creation of the model specification and when the model fit function is invoked.
2. If the model specification is saved and loaded into a new session where those same data objects do not exist.

The best way to avoid these issues is to not reference any data objects in the global environment but to use data descriptors such as `.cols()`. Another way of writing the previous specification is

```
k_means(num_clusters = .cols() / 5)
```

This is not dependent on any specific data object and is evaluated immediately before the model fitting process begins.

One less advantageous approach to solving this issue is to use quasiquote. This would insert the actual R object into the model specification and might be the best idea when the data object is small. For example, using

```
k_means(num_clusters = ncol(!mtcars) - 1)
```

would work (and be reproducible between sessions) but embeds the entire `mtcars` data set into the `num_clusters` expression:

```
> k_means(num_clusters = ncol(!mtcars) / 5)
K Means Cluster Specification (partition)
```

Main Arguments:

```
num_clusters = ncol(structure(list(mpg = c(21, 21, 22.8, 21.4, 18.7,<snip>
```

Computational engine: stats

However, if there were an object with the number of columns in it, this wouldn't be too bad:

```
> num_clusters_val <- ncol(mtcars) / 5
> num_clusters_val
[1] 10
> k_means(num_clusters = !!num_clusters_val)
K Means Cluster Specification (partition)
```

Main Arguments:
num_clusters = 2.2

More information on quosures and quasiquotation can be found at <https://adv-r.hadley.nz/quasiquotation.html>.

control_cluster *Control the fit function*

Description

Options can be passed to the `fit.cluster_spec()` function that control the output and computations.

Usage

```
control_cluster(verbosity = 1L, catch = FALSE)
```

Arguments

verbosity	An integer where a value of zero indicates that no messages or output should be shown when packages are loaded or when the model is fit. A value of 1 means that package loading is quiet but model fits can produce output to the screen (depending on if they contain their own verbose-type argument). A value of 2 or more indicates that any output should be seen.
catch	A logical where a value of TRUE will evaluate the model inside of <code>try(, silent = TRUE)</code> . If the model fails, an object is still returned (without an error) that inherits the class "try-error".

Value

An S3 object with class "control_cluster" that is a named list with the results of the function call

Examples

```
control_cluster()
control_cluster(catch = TRUE)
```

cut_height	<i>Cut Height</i>
------------	-------------------

Description

Used in most `tidyclust::hier_clust()` models.

Usage

```
cut_height(range = c(0, dials::unknown()), trans = NULL)
```

Arguments

range	A two-element vector holding the <i>defaults</i> for the smallest and largest possible values, respectively. If a transformation is specified, these values should be in the <i>transformed units</i> .
trans	A trans object from the scales package, such as <code>scales::log10_trans()</code> or <code>scales::reciprocal_trans()</code> . If not provided, the default is used which matches the units used in range. If no transformation, NULL.

Examples

```
cut_height()
```

extract-tidyclust	<i>Extract elements of a tidyclust model object</i>
-------------------	---

Description

These functions extract various elements from a clustering object. If they do not exist yet, an error is thrown.

- `extract_fit_engine()` returns the engine specific fit embedded within a tidyclust model fit. For example, when using `k_means()` with the "lm" engine, this returns the underlying kmeans object.
- `extract_parameter_set_dials()` returns a set of dials parameter objects.

Usage

```
## S3 method for class 'cluster_fit'
extract_fit_engine(x, ...)
```

```
## S3 method for class 'cluster_spec'
extract_parameter_set_dials(x, ...)
```


Arguments

x A `cluster_fit` object or a `cluster_spec` object.
 ... Not currently used.

Details

Extracting the underlying engine fit can be helpful for describing the model (via `print()`, `summary()`, `plot()`, etc.) or for variable importance/explainers.

However, users should not invoke the `predict()` method on an extracted model. There may be preprocessing operations that `tidyclust` has executed on the data prior to giving it to the model. Bypassing these can lead to errors or silently generating incorrect predictions.

Good:

```
tidyclust_fit %>% predict(new_data)
```

Bad:

```
tidyclust_fit %>% extract_fit_engine() %>% predict(new_data)
```

Value

The extracted value from the `tidyclust` object, `x`, as described in the description section.

Examples

```
kmeans_spec <- k_means(num_clusters = 2)
kmeans_fit <- fit(kmeans_spec, ~., data = mtcars)

extract_fit_engine(kmeans_fit)
```

`extract_centroids` *Extract clusters from model*

Description

When applied to a fitted cluster specification, returns a tibble with cluster location. When such locations doesn't make sense for the model, a mean location is used.

Usage

```
extract_centroids(object, ...)
```

Arguments

object An fitted `cluster_spec` object.
 ... Other arguments passed to methods. Using the `prefix` allows you to change the prefix in the levels of the factor levels.

Details

Some model types such as K-means as seen in `k_means()` stores the centroid in the object itself. leading the use of this function to act as a simple extract. Other model types such as Hierarchical (Agglomerative) Clustering as seen in `hier_clust()`, are fit in such a way that the number of clusters can be determined at any time after the fit. Setting the `num_clusters` or `cut_height` in this function will be used to determine the clustering when reported.

Further more, some models like `hier_clust()`, doesn't have a notion of "centroids". The mean of the observation within each cluster assignment is returned as the centroid.

The ordering of the clusters is such that the first observation in the training data set will be in cluster 1, the next observation that doesn't belong to cluster 1 will be in cluster 2, and so on and forth. As the ordering of clustering doesn't matter, this is done to avoid identical sets of clustering having different labels if fit multiple times.

Related functions:

`extract_centroids()` is a part of a trio of functions doing similar things:

- `extract_cluster_assignment()` returns the cluster assignments of the training observations
- `extract_centroids()` returns the location of the centroids
- `predict()` returns the cluster a new observation belongs to

Value

A `tibble::tibble()` with 1 row for each centroid and their position. `.cluster` denotes the cluster name for the centroid. The remaining variables match variables passed into model.

See Also

`extract_cluster_assignment()` `predict.cluster_fit()`

Examples

```
set.seed(1234)
kmeans_spec <- k_means(num_clusters = 5) %>%
  set_engine("stats")

kmeans_fit <- fit(kmeans_spec, ~., mtcars)

kmeans_fit %>%
  extract_centroids()

# Some models such as `hier_clust()` fits in such a way that you can specify
# the number of clusters after the model is fit.
# A Hierarchical (Agglomerative) Clustering method doesn't technically have
# clusters, so the center of the observation within each cluster is returned
# instead.
hclust_spec <- hier_clust() %>%
  set_engine("stats")

hclust_fit <- fit(hclust_spec, ~., mtcars)
```

```
hclust_fit %>%  
  extract_centroids(num_clusters = 2)
```

```
hclust_fit %>%  
  extract_centroids(cut_height = 250)
```

extract_cluster_assignment

Extract cluster assignments from model

Description

When applied to a fitted cluster specification, returns a tibble with cluster assignments of the data used to train the model.

Usage

```
extract_cluster_assignment(object, ...)
```

Arguments

object	An fitted <code>cluster_spec</code> object.
...	Other arguments passed to methods. Using the prefix allows you to change the prefix in the levels of the factor levels.

Details

Some model types such as K-means as seen in `k_means()` stores the cluster assignments in the object itself. leading the use of this function to act as a simple extract. Other model types such as Hierarchical (Agglomerative) Clustering as seen in `hier_clust()`, are fit in such a way that the number of clusters can be determined at any time after the fit. Setting the `num_clusters` or `cut_height` in this function will be used to determine the clustering when reported.

The ordering of the clusters is such that the first observation in the training data set will be in cluster 1, the next observation that doesn't belong to cluster 1 will be in cluster 2, and so on and forth. As the ordering of clustering doesn't matter, this is done to avoid identical sets of clustering having different labels if fit multiple times.

Related functions:

`extract_cluster_assignment()` is a part of a trio of functions doing similar things:

- `extract_cluster_assignment()` returns the cluster assignments of the training observations
- `extract_centroids()` returns the location of the centroids
- `predict()` returns the cluster a new observation belongs to

Value

A `tibble::tibble()` with 1 column named `.cluster`. This tibble will correspond to the training data set.

See Also

[extract_centroids\(\)](#) [predict.cluster_fit\(\)](#)

Examples

```
kmeans_spec <- k_means(num_clusters = 5) %>%
  set_engine("stats")

kmeans_fit <- fit(kmeans_spec, ~., mtcars)

kmeans_fit %>%
  extract_cluster_assignment()

kmeans_fit %>%
  extract_cluster_assignment(prefix = "C_")

# Some models such as `hier_clust()` fits in such a way that you can specify
# the number of clusters after the model is fit
hclust_spec <- hier_clust() %>%
  set_engine("stats")

hclust_fit <- fit(hclust_spec, ~., mtcars)

hclust_fit %>%
  extract_cluster_assignment(num_clusters = 2)

hclust_fit %>%
  extract_cluster_assignment(cut_height = 250)
```

`extract_fit_summary` *S3 method to get fitted model summary info depending on engine*

Description

S3 method to get fitted model summary info depending on engine

Usage

```
extract_fit_summary(object, ...)
```

Arguments

<code>object</code>	a fitted cluster_spec object
<code>...</code>	other arguments passed to methods

Details

The elements `cluster_names` and `cluster_assignments` will be factors.

Value

A list with various summary elements

Examples

```
kmeans_spec <- k_means(num_clusters = 5) %>%
  set_engine("stats")

kmeans_fit <- fit(kmeans_spec, ~., mtcars)

kmeans_fit %>%
  extract_fit_summary()
```

`finalize_model_tidyclust`

Splice final parameters into objects

Description

The `finalize_*` functions take a list or tibble of tuning parameter values and update objects with those values.

Usage

```
finalize_model_tidyclust(x, parameters)

finalize_workflow_tidyclust(x, parameters)
```

Arguments

<code>x</code>	A recipe, parsnip model specification, or workflow.
<code>parameters</code>	A list or 1-row tibble of parameter values. Note that the column names of the tibble should be the id fields attached to <code>tune()</code> . For example, in the Examples section below, the model has <code>tune("K")</code> . In this case, the parameter tibble should be "K" and not "neighbors".

Value

An updated version of `x`.

Examples

```
kmeans_spec <- k_means(num_clusters = tune())
kmeans_spec

best_params <- data.frame(num_clusters = 5)
best_params

finalize_model_tidyclust(kmeans_spec, best_params)
```

<code>fit.cluster_spec</code>	<i>Fit a Model Specification to a Data Set</i>
-------------------------------	--

Description

`fit()` and `fit_xy()` take a model specification, `translate_tidyclust` the required code by substituting arguments, and execute the model fit routine.

Usage

```
## S3 method for class 'cluster_spec'
fit(object, formula, data, control = control_cluster(), ...)

## S3 method for class 'cluster_spec'
fit_xy(object, x, case_weights = NULL, control = control_cluster(), ...)
```

Arguments

<code>object</code>	An object of class <code>cluster_spec</code> that has a chosen engine (via <code>set_engine()</code>).
<code>formula</code>	An object of class <code>formula</code> (or one that can be coerced to that class): a symbolic description of the model to be fitted.
<code>data</code>	Optional, depending on the interface (see Details below). A data frame containing all relevant variables (e.g. predictors, case weights, etc). Note: when needed, a <i>named argument</i> should be used.
<code>control</code>	A named list with elements <code>verbosity</code> and <code>catch</code> . See <code>control_cluster()</code> .
<code>...</code>	Not currently used; values passed here will be ignored. Other options required to fit the model should be passed using <code>set_engine()</code> .
<code>x</code>	A matrix, sparse matrix, or data frame of predictors. Only some models have support for sparse matrix input. See <code>modelenv::get_encoding()</code> for details. <code>x</code> should have column names.
<code>case_weights</code>	An optional classed vector of numeric case weights. This must return <code>TRUE</code> when <code>hardhat::is_case_weights()</code> is run on it. See <code>hardhat::frequency_weights()</code> and <code>hardhat::importance_weights()</code> for examples.

Details

`fit()` and `fit_xy()` substitute the current arguments in the model specification into the computational engine's code, check them for validity, then fit the model using the data and the engine-specific code. Different model functions have different interfaces (e.g. formula or x/y) and these functions translate `tidyclust` between the interface used when `fit()` or `fit_xy()` was invoked and the one required by the underlying model.

When possible, these functions attempt to avoid making copies of the data. For example, if the underlying model uses a formula and `fit()` is invoked, the original data are references when the model is fit. However, if the underlying model uses something else, such as x/y , the formula is evaluated and the data are converted to the required format. In this case, any calls in the resulting model objects reference the temporary objects used to fit the model.

If the model engine has not been set, the model's default engine will be used (as discussed on each model page). If the verbosity option of `control_cluster()` is greater than zero, a warning will be produced.

If you would like to use an alternative method for generating contrasts when supplying a formula to `fit()`, set the global option `contrasts` to your preferred method. For example, you might set it to: `options(contrasts = c(unordered = "contr.helmert", ordered = "contr.poly"))`. See the help page for `stats::contr.treatment()` for more possible contrast types.

Value

A `cluster_fit` object that contains several elements:

- `spec`: The model specification object (object in the call to `fit`)
- `fit`: when the model is executed without error, this is the model object. Otherwise, it is a try-error object with the error message.
- `preproc`: any objects needed to convert between a formula and non-formula interface (such as the terms object)

The return value will also have a class related to the fitted model (e.g. `"_kmeans"`) before the base class of `"cluster_fit"`.

A fitted `cluster_fit` object.

See Also

[set_engine\(\)](#), [control_cluster\(\)](#), [cluster_spec](#), [cluster_fit](#)

Examples

```
library(dplyr)

kmeans_mod <- k_means(num_clusters = 5)

using_formula <-
  kmeans_mod %>%
  set_engine("stats") %>%
  fit(., data = mtcars)
```

```

using_x <-
  kmeans_mod %>%
  set_engine("stats") %>%
  fit_xy(x = mtcars)

using_formula
using_x

```

get_centroid_dists	<i>Computes distance from observations to centroids</i>
--------------------	---

Description

Computes distance from observations to centroids

Usage

```
get_centroid_dists(new_data, centroids, dist_fun = Rfast::dista)
```

Arguments

new_data	A data frame
centroids	A data frame where each row is a centroid.
dist_fun	A function for computing matrix-to-matrix distances. Defaults to <code>Rfast::dista()</code>

glance.cluster_fit	<i>Construct a single row summary "glance" of a model, fit, or other object</i>
--------------------	---

Description

This method glances the model in a tidyclust model object, if it exists.

Usage

```
## S3 method for class 'cluster_fit'
glance(x, ...)
```

Arguments

x	model or other R object to convert to single-row data frame
...	other arguments passed to methods

Value

a tibble

 hier_clust

Hierarchical (Agglomerative) Clustering

Description

hier_clust() defines a model that fits clusters based on a distance-based dendrogram

There are different ways to fit this model, and the method of estimation is chosen by setting the model engine. The engine-specific pages for this model are listed below.

- [stats](#)

Usage

```
hier_clust(
  mode = "partition",
  engine = "stats",
  num_clusters = NULL,
  cut_height = NULL,
  linkage_method = "complete"
)
```

Arguments

mode	A single character string for the type of model. The only possible value for this model is "partition".
engine	A single character string specifying what computational engine to use for fitting. Possible engines are listed below. The default for this model is "stats".
num_clusters	Positive integer, number of clusters in model (optional).
cut_height	Positive double, height at which to cut dendrogram to obtain cluster assignments (only used if num_clusters is NULL)
linkage_method	the agglomeration method to be used. This should be (an unambiguous abbreviation of) one of "ward.D", "ward.D2", "single", "complete", "average" (= UPGMA), "mcquitty" (= WPGMA), "median" (= WPGMC) or "centroid" (= UPGMC).

Details

What does it mean to predict?:

To predict the cluster assignment for a new observation, we find the closest cluster. How we measure "closeness" is dependent on the specified type of linkage in the model:

- *single linkage*: The new observation is assigned to the same cluster as its nearest observation from the training data.
- *complete linkage*: The new observation is assigned to the cluster with the smallest maximum distances between training observations and the new observation.

- *average linkage*: The new observation is assigned to the cluster with the smallest average distances between training observations and the new observation.
- *centroid method*: The new observation is assigned to the cluster with the closest centroid, as in prediction for `k_means`.
- *Ward's method*: The new observation is assigned to the cluster with the smallest increase in **error sum of squares (ESS)** due to the new addition. The ESS is computed as the sum of squared distances between observations in a cluster, and the centroid of the cluster.

Value

A `hier_clust` cluster specification.

Examples

```
# Show all engines
modelenv::get_from_env("hier_clust")

hier_clust()
```

k_means	<i>K-Means</i>
---------	----------------

Description

`k_means()` defines a model that fits clusters based on distances to a number of centers. This definition doesn't just include K-means, but includes models like K-prototypes.

There are different ways to fit this model, and the method of estimation is chosen by setting the model engine. The engine-specific pages for this model are listed below.

- [stats](#): Classical K-means
- [ClusterR](#): Classical K-means
- [klaR](#): K-Modes
- [clustMixType](#): K-prototypes

Usage

```
k_means(mode = "partition", engine = "stats", num_clusters = NULL)
```

Arguments

mode	A single character string for the type of model. The only possible value for this model is "partition".
engine	A single character string specifying what computational engine to use for fitting. Possible engines are listed below. The default for this model is "stats".
num_clusters	Positive integer, number of clusters in model.

Details**What does it mean to predict?:**

For a K-means model, each cluster is defined by a location in the predictor space. Therefore, prediction in tidyclust is defined by calculating which cluster centroid an observation is closest too.

Value

A k_means cluster specification.

Examples

```
# Show all engines
modelenv::get_from_env("k_means")

k_means()
```

linkage_method	<i>The agglomeration Linkage method</i>
----------------	---

Description

The agglomeration Linkage method

Usage

```
linkage_method(values = values_linkage_method)

values_linkage_method
```

Arguments

values A character string of possible values. See linkage_methods in examples below.

Format

An object of class character of length 8.

Details

This parameter is used in tidyclust models for hier_clust().

Examples

```
values_linkage_method
linkage_method()
```

`min_grid.cluster_spec` *Determine the minimum set of model fits*

Description

Determine the minimum set of model fits

Usage

```
## S3 method for class 'cluster_spec'
min_grid(x, grid, ...)
```

Arguments

<code>x</code>	A cluster specification.
<code>grid</code>	A tibble with tuning parameter combinations.
<code>...</code>	Not currently used.

Value

A tibble with the minimum tuning parameters to fit and an additional list column with the parameter combinations used for prediction.

`new_cluster_metric` *Construct a new clustering metric function*

Description

These functions provide convenient wrappers to create the one type of metric functions in celr: clustering metrics. They add a metric-specific class to `fn`. These features are used by `cluster_metric_set()` and by `tune_cluster()` when tuning.

Usage

```
new_cluster_metric(fn, direction)
```

Arguments

<code>fn</code>	A function.
<code>direction</code>	A string. One of: <ul style="list-style-type: none"> • "maximize" • "minimize" • "zero"

Value

A `cluster_metric` object.

predict.cluster_fit *Model predictions*

Description

Apply to a model to create different types of predictions. `predict()` can be used for all types of models and uses the "type" argument for more specificity.

Usage

```
## S3 method for class 'cluster_fit'  
predict(object, new_data, type = NULL, opts = list(), ...)  
  
## S3 method for class 'cluster_fit'  
predict_raw(object, new_data, opts = list(), ...)
```

Arguments

object	An object of class <code>cluster_fit</code> .
new_data	A rectangular data object, such as a data frame.
type	A single character value or NULL. Possible values are "cluster", or "raw". When NULL, <code>predict()</code> will choose an appropriate value based on the model's mode.
opts	A list of optional arguments to the underlying predict function that will be used when type = "raw". The list should not include options for the model object or the new data being predicted.
...	Arguments to the underlying model's prediction function cannot be passed here (see opts).

Details

If "type" is not supplied to `predict()`, then a choice is made:

- type = "cluster" for clustering models

`predict()` is designed to provide a tidy result (see "Value" section below) in a tibble output format.

The ordering of the clusters is such that the first observation in the training data set will be in cluster 1, the next observation that doesn't belong to cluster 1 will be in cluster 2, and so on and forth. As the ordering of clustering doesn't matter, this is done to avoid identical sets of clustering having different labels if fit multiple times.

What does it mean to predict?:

Prediction is not always formally defined for clustering models. Therefore, each `cluster_spec` method will have their own section on how "prediction" is interpreted, and done if implemented.

Related functions:

`predict()` when used with tidyclus objects is a part of a trio of functions doing similar things:

- `extract_cluster_assignment()` returns the cluster assignments of the training observations
- `extract_centroids()` returns the location of the centroids
- `predict()` returns the cluster a new observation belongs to

Value

With the exception of `type = "raw"`, the results of `predict.cluster_fit()` will be a tibble as many rows in the output as there are rows in `new_data` and the column names will be predictable.

For clustering results the tibble will have a `.pred_cluster` column.

Using `type = "raw"` with `predict.cluster_fit()` will return the unadulterated results of the prediction function.

When the model fit failed and the error was captured, the `predict()` function will return the same structure as above but filled with missing values. This does not currently work for multivariate models.

See Also

`extract_cluster_assignment()` `extract_centroids()`

Examples

```
kmeans_spec <- k_means(num_clusters = 5) %>%
  set_engine("stats")

kmeans_fit <- fit(kmeans_spec, ~., mtcars)

kmeans_fit %>%
  predict(new_data = mtcars)

# Some models such as `hier_clust()` fits in such a way that you can specify
# the number of clusters after the model is fit
hclust_spec <- hier_clust() %>%
  set_engine("stats")

hclust_fit <- fit(hclust_spec, ~., mtcars)

hclust_fit %>%
  predict(new_data = mtcars[4:6, ], num_clusters = 2)

hclust_fit %>%
  predict(new_data = mtcars[4:6, ], cut_height = 250)
```

```
prep_data_dist      Prepares data and distance matrices for metric calculation
```

Description

Prepares data and distance matrices for metric calculation

Usage

```
prep_data_dist(object, new_data = NULL, dists = NULL, dist_fun = Rfast::Dist)
```

Arguments

object	A fitted <code>cluster_spec</code> object.
new_data	A dataset to calculate predictions on. If NULL, the trained cluster assignments from the fitted object are used.
dists	A distance matrix for the data. If NULL, distance is computed on <code>new_data</code> using the <code>stats::dist()</code> function.
dist_fun	A custom distance functions.

Value

A list

```
reconcile_clusterings_mapping
      Relabels clusters to match another cluster assignment
```

Description

When forcing one-to-one, the user needs to decide what to prioritize:

- "accuracy": optimize raw count of all observations with the same label across the two assignments
- "precision": optimize the average percent of each alt cluster that matches the corresponding primary cluster

Usage

```
reconcile_clusterings_mapping(
  primary,
  alternative,
  one_to_one = TRUE,
  optimize = "accuracy"
)
```

Arguments

primary	A vector containing cluster labels, to be matched
alternative	Another vector containing cluster labels, to be changed
one_to_one	Boolean; should each alt cluster match only one primary cluster?
optimize	One of "accuracy" or "precision"; see description.

Details

Retains the cluster labels of the primary assignment, and relabel the alternate assignment to match as closely as possible. The user must decide whether clusters are forced to be "one-to-one"; that is, are we allowed to assign multiple labels from the alternate assignment to the same primary label?

Value

A tibble with 3 columns; primary, alt, alt_recoded

Examples

```
factor1 <- c("Apple", "Apple", "Carrot", "Carrot", "Banana", "Banana")
factor2 <- c("Dog", "Dog", "Cat", "Dog", "Fish", "Fish")
reconcile_clusterings_mapping(factor1, factor2)

factor1 <- c("Apple", "Apple", "Carrot", "Carrot", "Banana", "Banana")
factor2 <- c("Dog", "Dog", "Cat", "Dog", "Fish", "Parrot")
reconcile_clusterings_mapping(factor1, factor2, one_to_one = FALSE)
```

set_args.cluster_spec *Change arguments of a cluster specification*

Description

Change arguments of a cluster specification

Usage

```
## S3 method for class 'cluster_spec'
set_args(object, ...)
```

Arguments

object	A model specification.
...	One or more named model arguments.

Value

An updated `cluster_spec` object.

`set_engine.cluster_spec`*Change engine of a cluster specification*

Description

Change engine of a cluster specification

Usage

```
## S3 method for class 'cluster_spec'  
set_engine(object, engine, ...)
```

Arguments

<code>object</code>	A model specification.
<code>engine</code>	A character string for the software that should be used to fit the model. This is highly dependent on the type of model (e.g. linear regression, random forest, etc.).
<code>...</code>	Any optional arguments associated with the chosen computational engine. These are captured as quosures and can be tuned with <code>tune()</code> .

Value

An updated `cluster_spec` object.

`set_mode.cluster_spec` *Change mode of a cluster specification*

Description

Change mode of a cluster specification

Usage

```
## S3 method for class 'cluster_spec'  
set_mode(object, mode)
```

Arguments

<code>object</code>	A model specification.
<code>mode</code>	A character string for the model type (e.g. "classification" or "regression")

Value

An updated `cluster_spec` object.

silhouette	<i>Measures silhouette between clusters</i>
------------	---

Description

Measures silhouette between clusters

Usage

```
silhouette(object, new_data = NULL, dists = NULL, dist_fun = Rfast::Dist)
```

Arguments

object	A fitted tidyclust model
new_data	A dataset to predict on. If NULL, uses trained clustering.
dists	A distance matrix. Used if new_data is NULL.
dist_fun	A function for calculating distances between observations. Defaults to Euclidean distance on processed data.

Details

[silhouette_avg\(\)](#) is the corresponding cluster metric function that returns the average of the values given by `silhouette()`.

Value

A tibble giving the silhouette for each observation.

Examples

```
kmeans_spec <- k_means(num_clusters = 5) %>%  
  set_engine("stats")  
  
kmeans_fit <- fit(kmeans_spec, ~., mtcars)  
  
dists <- mtcars %>%  
  as.matrix() %>%  
  dist()  
  
silhouette(kmeans_fit, dists = dists)
```

silhouette_avg	<i>Measures average silhouette across all observations</i>
----------------	--

Description

Measures average silhouette across all observations

Usage

```
silhouette_avg(object, ...)  
  
## S3 method for class 'cluster_spec'  
silhouette_avg(object, ...)  
  
## S3 method for class 'cluster_fit'  
silhouette_avg(object, new_data = NULL, dists = NULL, dist_fun = NULL, ...)  
  
## S3 method for class 'workflow'  
silhouette_avg(object, new_data = NULL, dists = NULL, dist_fun = NULL, ...)  
  
silhouette_avg_vec(  
  object,  
  new_data = NULL,  
  dists = NULL,  
  dist_fun = Rfast::Dist,  
  ...  
)
```

Arguments

object	A fitted kmeans tidyclust model
...	Other arguments passed to methods.
new_data	A dataset to predict on. If NULL, uses trained clustering.
dists	A distance matrix. Used if new_data is NULL.
dist_fun	A function for calculating distances between observations. Defaults to Euclidean distance on processed data.

Details

Not to be confused with [silhouette\(\)](#) that returns a tibble with silhouette for each observation.

Value

A double; the average silhouette.

See Also

Other cluster metric: [sse_ratio\(\)](#), [sse_total\(\)](#), [sse_within_total\(\)](#)

Examples

```
kmeans_spec <- k_means(num_clusters = 5) %>%
  set_engine("stats")

kmeans_fit <- fit(kmeans_spec, ~., mtcars)

dists <- mtcars %>%
  as.matrix() %>%
  dist()

silhouette_avg(kmeans_fit, dists = dists)

silhouette_avg_vec(kmeans_fit, dists = dists)
```

sse_ratio

Compute the ratio of the WSS to the total SSE

Description

Compute the ratio of the WSS to the total SSE

Usage

```
sse_ratio(object, ...)

## S3 method for class 'cluster_spec'
sse_ratio(object, ...)

## S3 method for class 'cluster_fit'
sse_ratio(object, new_data = NULL, dist_fun = NULL, ...)

## S3 method for class 'workflow'
sse_ratio(object, new_data = NULL, dist_fun = NULL, ...)

sse_ratio_vec(object, new_data = NULL, dist_fun = Rfast::dista, ...)
```

Arguments

object	A fitted kmeans tidyclust model
...	Other arguments passed to methods.
new_data	A dataset to predict on. If NULL, uses trained clustering.
dist_fun	A function for calculating distances to centroids. Defaults to Euclidean distance on processed data.

Value

A tibble with 3 columns: `.metric`, `.estimator`, and `.estimate`.

See Also

Other cluster metric: `silhouette_avg()`, `sse_total()`, `sse_within_total()`

Examples

```
kmeans_spec <- k_means(num_clusters = 5) %>%
  set_engine("stats")

kmeans_fit <- fit(kmeans_spec, ~., mtcars)

sse_ratio(kmeans_fit)

sse_ratio_vec(kmeans_fit)
```

sse_total	<i>Compute the total sum of squares</i>
-----------	---

Description

Compute the total sum of squares

Usage

```
sse_total(object, ...)

## S3 method for class 'cluster_spec'
sse_total(object, ...)

## S3 method for class 'cluster_fit'
sse_total(object, new_data = NULL, dist_fun = NULL, ...)

## S3 method for class 'workflow'
sse_total(object, new_data = NULL, dist_fun = NULL, ...)

sse_total_vec(object, new_data = NULL, dist_fun = Rfast::dista, ...)
```

Arguments

<code>object</code>	A fitted kmeans tidyclust model
<code>...</code>	Other arguments passed to methods.
<code>new_data</code>	A dataset to predict on. If NULL, uses trained clustering.
<code>dist_fun</code>	A function for calculating distances to centroids. Defaults to Euclidean distance on processed data.

Value

A tibble with 3 columns: `.metric`, `.estimator`, and `.estimate`.

See Also

Other cluster metric: [silhouette_avg\(\)](#), [sse_ratio\(\)](#), [sse_within_total\(\)](#)

Examples

```
kmeans_spec <- k_means(num_clusters = 5) %>%
  set_engine("stats")

kmeans_fit <- fit(kmeans_spec, ~., mtcars)

sse_total(kmeans_fit)

sse_total_vec(kmeans_fit)
```

sse_within

Calculates Sum of Squared Error in each cluster

Description

Calculates Sum of Squared Error in each cluster

Usage

```
sse_within(object, new_data = NULL, dist_fun = Rfast::dista)
```

Arguments

<code>object</code>	A fitted kmeans tidyclust model
<code>new_data</code>	A dataset to predict on. If NULL, uses trained clustering.
<code>dist_fun</code>	A function for calculating distances to centroids. Defaults to Euclidean distance on processed data.

Details

[sse_within_total\(\)](#) is the corresponding cluster metric function that returns the sum of the values given by `sse_within()`.

Value

A tibble with two columns, the cluster name and the SSE within that cluster.

Examples

```
kmeans_spec <- k_means(num_clusters = 5) %>%
  set_engine("stats")

kmeans_fit <- fit(kmeans_spec, ~., mtcars)

sse_within(kmeans_fit)
```

sse_within_total	<i>Compute the sum of within-cluster SSE</i>
------------------	--

Description

Compute the sum of within-cluster SSE

Usage

```
sse_within_total(object, ...)

## S3 method for class 'cluster_spec'
sse_within_total(object, ...)

## S3 method for class 'cluster_fit'
sse_within_total(object, new_data = NULL, dist_fun = NULL, ...)

## S3 method for class 'workflow'
sse_within_total(object, new_data = NULL, dist_fun = NULL, ...)

sse_within_total_vec(object, new_data = NULL, dist_fun = Rfast::dista, ...)
```

Arguments

object	A fitted kmeans tidyclust model
...	Other arguments passed to methods.
new_data	A dataset to predict on. If NULL, uses trained clustering.
dist_fun	A function for calculating distances to centroids. Defaults to Euclidean distance on processed data.

Details

Not to be confused with [sse_within\(\)](#) that returns a tibble with within-cluster SSE, one row for each cluster.

Value

A tibble with 3 columns: `.metric`, `.estimator`, and `.estimate`.

See Also

Other cluster metric: [silhouette_avg\(\)](#), [sse_ratio\(\)](#), [sse_total\(\)](#)

Examples

```
kmeans_spec <- k_means(num_clusters = 5) %>%
  set_engine("stats")

kmeans_fit <- fit(kmeans_spec, ~., mtcars)

sse_within_total(kmeans_fit)

sse_within_total_vec(kmeans_fit)
```

tidy.cluster_fit	<i>Turn a tidyclust model object into a tidy tibble</i>
------------------	---

Description

This method tidies the model in a tidyclust model object, if it exists.

Usage

```
## S3 method for class 'cluster_fit'
tidy(x, ...)
```

Arguments

`x` An object to be converted into a tidy [tibble::tibble\(\)](#).

`...` Additional arguments to tidying method.

Value

a tibble

translate_tidyclust	<i>Resolve a Model Specification for a Computational Engine</i>
---------------------	---

Description

`translate_tidyclust()` will `translate_tidyclust` a model specification into a code object that is specific to a particular engine (e.g. R package). It translate tidyclust generic parameters to their counterparts.

Usage

```
translate_tidyclust(x, ...)

## Default S3 method:
translate_tidyclust(x, engine = x$engine, ...)
```

Arguments

x	A model specification.
...	Not currently used.
engine	The computational engine for the model (see ?set_engine).

Details

translate_tidyclust() produces a *template* call that lacks the specific argument values (such as data, etc). These are filled in once fit() is called with the specifics of the data for the model. The call may also include tune() arguments if these are in the specification. To handle the tune() arguments, you need to use the [tune package](https://www.tidymodels.org/start/tuning/). For more information see <https://www.tidymodels.org/start/tuning/>

It does contain the resolved argument names that are specific to the model fitting function/engine.

This function can be useful when you need to understand how tidyclust goes from a generic model specific to a model fitting function.

Note: this function is used internally and users should only use it to understand what the underlying syntax would be. It should not be used to modify the cluster specification.

Value

Prints translated code.

tune_cluster	<i>Model tuning via grid search</i>
--------------	-------------------------------------

Description

[tune_cluster\(\)](#) computes a set of performance metrics (e.g. accuracy or RMSE) for a pre-defined set of tuning parameters that correspond to a model or recipe across one or more resamples of the data.

Usage

```
tune_cluster(object, ...)

## S3 method for class 'cluster_spec'
tune_cluster(
  object,
```

```

preprocessor,
resamples,
...,
param_info = NULL,
grid = 10,
metrics = NULL,
control = tune::control_grid()
)

## S3 method for class 'workflow'
tune_cluster(
  object,
  resamples,
  ...,
  param_info = NULL,
  grid = 10,
  metrics = NULL,
  control = tune::control_grid()
)

```

Arguments

object	A tidyclust model specification or a <code>workflows::workflow()</code> .
...	Not currently used.
preprocessor	A traditional model formula or a recipe created using <code>recipes::recipe()</code> .
resamples	An <code>rset()</code> object.
param_info	A <code>dials::parameters()</code> object or NULL. If none is given, a parameters set is derived from other arguments. Passing this argument can be useful when parameter ranges need to be customized.
grid	A data frame of tuning combinations or a positive integer. The data frame should have columns for each parameter being tuned and rows for tuning parameter candidates. An integer denotes the number of candidate parameter sets to be created automatically.
metrics	A <code>cluster_metric_set()</code> or NULL.
control	An object used to modify the tuning process. Defaults to <code>tune::control_grid()</code> .

Value

An updated version of `resamples` with extra list columns for `.metrics` and `.notes` (optional columns are `.predictions` and `.extracts`). `.notes` contains warnings and errors that occur during execution.

Examples

```

library(recipes)
library(rsample)
library(workflows)

```

```

library(tune)

rec_spec <- recipe(~., data = mtcars) %>%
  step_normalize(all_numeric_predictors()) %>%
  step_pca(all_numeric_predictors())

kmeans_spec <- k_means(num_clusters = tune())

wflow <- workflow() %>%
  add_recipe(rec_spec) %>%
  add_model(kmeans_spec)

grid <- tibble(num_clusters = 1:3)

set.seed(4400)
folds <- vfold_cv(mtcars, v = 2)

res <- tune_cluster(
  wflow,
  resamples = folds,
  grid = grid
)
res

collect_metrics(res)

```

update.hier_clust *Update a cluster specification*

Description

If parameters of a cluster specification need to be modified, `update()` can be used in lieu of recreating the object from scratch.

Usage

```

## S3 method for class 'hier_clust'
update(
  object,
  parameters = NULL,
  num_clusters = NULL,
  cut_height = NULL,
  linkage_method = NULL,
  fresh = FALSE,
  ...
)

## S3 method for class 'k_means'
update(object, parameters = NULL, num_clusters = NULL, fresh = FALSE, ...)

```

Arguments

object	A cluster specification.
parameters	A 1-row tibble or named list with <i>main</i> parameters to update. Use either parameters or the main arguments directly when updating. If the main arguments are used, these will supersede the values in parameters. Also, using engine arguments in this object will result in an error.
num_clusters	Positive integer, number of clusters in model.
cut_height	Positive double, height at which to cut dendrogram to obtain cluster assignments (only used if num_clusters is NULL)
linkage_method	the agglomeration method to be used. This should be (an unambiguous abbreviation of) one of "ward.D", "ward.D2", "single", "complete", "average" (= UPGMA), "mcquitty" (= WPGMA), "median" (= WPGMC) or "centroid" (= UPGMC).
fresh	A logical for whether the arguments should be modified in-place or replaced wholesale.
...	Not used for update().

Value

An updated cluster specification.

Examples

```
kmeans_spec <- k_means(num_clusters = 5)
kmeans_spec
update(kmeans_spec, num_clusters = 1)
update(kmeans_spec, num_clusters = 1, fresh = TRUE)

param_values <- tibble::tibble(num_clusters = 10)

kmeans_spec %>% update(param_values)
```

Index

- * **cluster metric**
 - silhouette_avg, 27
 - sse_ratio, 28
 - sse_total, 29
 - sse_within_total, 31
- * **datasets**
 - linkage_method, 19
- augment.cluster_fit, 2
- cluster_fit, 3, 3, 9, 15, 21
- cluster_metric_set, 4
- cluster_metric_set(), 20, 34
- cluster_spec, 3, 4, 5, 9, 11, 12, 14, 15, 21, 23–25
- ClusterR, 18
- clustMixType, 18
- control_cluster, 7
- control_cluster(), 14, 15
- cut_height, 8
- dials::parameters(), 34
- extract-tidyclust, 8
- extract_centroids, 9
- extract_centroids(), 10–12, 22
- extract_cluster_assignment, 11
- extract_cluster_assignment(), 10, 11, 22
- extract_fit_engine.cluster_fit
 - (extract-tidyclust), 8
- extract_fit_summary, 12
- extract_parameter_set_dials.cluster_spec
 - (extract-tidyclust), 8
- finalize_model_tidyclust, 13
- finalize_workflow_tidyclust
 - (finalize_model_tidyclust), 13
- fit.cluster_spec, 14
- fit.cluster_spec(), 3, 6, 7
- fit_xy.cluster_spec(fit.cluster_spec), 14
- fit_xy.cluster_spec(), 3, 6
- get_centroid_dists, 16
- glance.cluster_fit, 16
- hardhat::frequency_weights(), 14
- hardhat::importance_weights(), 14
- hardhat::is_case_weights(), 14
- hier_clust, 17
- hier_clust(), 3, 5, 10, 11
- k_means, 18
- k_means(), 3, 5, 8, 10, 11
- klaR, 18
- linkage_method, 19
- min_grid.cluster_spec, 20
- new_cluster_metric, 20
- predict(), 9–11, 22
- predict.cluster_fit, 21
- predict.cluster_fit(), 10, 12
- predict_raw.cluster_fit
 - (predict.cluster_fit), 21
- prep_data_dist, 23
- recipes::recipe(), 34
- reconcile_clusterings_mapping, 23
- set_args.cluster_spec, 24
- set_engine(), 14, 15
- set_engine.cluster_spec, 25
- set_mode.cluster_spec, 25
- silhouette, 26
- silhouette(), 27
- silhouette_avg, 27, 29, 30, 32
- silhouette_avg(), 4, 26
- silhouette_avg_vec(silhouette_avg), 27
- sse_ratio, 28, 28, 30, 32

`sse_ratio()`, 4
`sse_ratio_vec (sse_ratio)`, 28
`sse_total`, 28, 29, 29, 32
`sse_total()`, 4
`sse_total_vec (sse_total)`, 29
`sse_within`, 30
`sse_within()`, 31
`sse_within_total`, 28–30, 31
`sse_within_total()`, 30
`sse_within_total_vec`
 (`sse_within_total`), 31
`stats`, 17, 18
`stats::contr.treatment()`, 15

`tibble::tibble()`, 32
`tidy.cluster_fit`, 32
`tidyclust_update (update.hier_clust)`, 35
`translate_tidyclust`, 32
`tune_cluster`, 33
`tune_cluster()`, 5, 20, 33

`update.hier_clust`, 35
`update.k_means (update.hier_clust)`, 35

`values_linkage_method (linkage_method)`,
 19

`workflows::workflow()`, 34