

# Package ‘lightgbm’

January 18, 2024

**Type** Package

**Title** Light Gradient Boosting Machine

**Version** 4.3.0

**Date** 2024-01-17

**Description** Tree based algorithms can be improved by introducing boosting frameworks.

'LightGBM' is one such framework, based on Ke, Guolin et al. (2017) <<https://papers.nips.cc/paper/6907-lightgbm-a-highly-efficient-gradient-boosting-decision>>.

This package offers an R interface to work with it.

It is designed to be distributed and efficient with the following advantages:

1. Faster training speed and higher efficiency.
2. Lower memory usage.
3. Better accuracy.
4. Parallel learning supported.
5. Capable of handling large-scale data.

In recognition of these advantages, 'LightGBM' has been widely-used in many winning solutions of machine learning competitions.

Comparison experiments on public datasets suggest that 'LightGBM' can outperform existing boosting frameworks on both efficiency and accuracy, with significantly lower memory consumption. In addition, parallel experiments suggest that in certain circumstances, 'LightGBM' can achieve a linear speed-up in training time by using multiple machines.

**Encoding** UTF-8

**License** MIT + file LICENSE

**URL** <https://github.com/Microsoft/LightGBM>

**BugReports** <https://github.com/Microsoft/LightGBM/issues>

**NeedsCompilation** yes

**Biarch** true

**VignetteBuilder** knitr

**Suggests** knitr, markdown, RhpCBLASctl, testthat

**Depends** R (>= 3.5)

**Imports** R6 (>= 2.0), data.table (>= 1.9.6), graphics, jsonlite (>= 1.0), Matrix (>= 1.1-0), methods, parallel, utils

**SystemRequirements** C++17

**RoxygenNote** 7.2.3

**Author** Yu Shi [aut],

Guolin Ke [aut],

Damien Soukhavong [aut],

James Lamb [aut, cre],

Qi Meng [aut],

Thomas Finley [aut],

Taifeng Wang [aut],

Wei Chen [aut],

Weidong Ma [aut],

Qiwei Ye [aut],

Tie-Yan Liu [aut],

Nikita Titov [aut],

Yachen Yan [ctb],

Microsoft Corporation [cph],

Dropbox, Inc. [cph],

Alberto Ferreira [ctb],

Daniel Lemire [ctb],

Victor Zverovich [cph],

IBM Corporation [ctb],

David Cortes [aut],

Michael Mayer [ctb]

**Maintainer** James Lamb <jaylamb20@gmail.com>

**Repository** CRAN

**Date/Publication** 2024-01-18 11:50:09 UTC

## R topics documented:

agaricus.test . . . . .	3
agaricus.train . . . . .	4
bank . . . . .	4
dim.lgb.Dataset . . . . .	5
dimnames.lgb.Dataset . . . . .	6
getLGBMThreads . . . . .	7
get_field . . . . .	7
lgb.configure_fast_predict . . . . .	8
lgb.convert_with_rules . . . . .	11
lgb.cv . . . . .	12
lgb.Dataset . . . . .	15
lgb.Dataset.construct . . . . .	17
lgb.Dataset.create.valid . . . . .	18
lgb.Dataset.save . . . . .	19
lgb.Dataset.set.categorical . . . . .	20
lgb.Dataset.set.reference . . . . .	21
lgb.drop_serialized . . . . .	22

lgb.dump . . . . . 22  
 lgb.get.eval.result . . . . . 23  
 lgb.importance . . . . . 25  
 lgb.interprete . . . . . 26  
 lgb.load . . . . . 27  
 lgb.make\_serializable . . . . . 28  
 lgb.model.dt.tree . . . . . 29  
 lgb.plot.importance . . . . . 30  
 lgb.plot.interpretation . . . . . 32  
 lgb.restore\_handle . . . . . 33  
 lgb.save . . . . . 34  
 lgb.train . . . . . 35  
 lightgbm . . . . . 38  
 predict.lgb.Booster . . . . . 41  
 print.lgb.Booster . . . . . 44  
 setLGBMThreads . . . . . 44  
 set\_field . . . . . 45  
 slice . . . . . 46  
 summary.lgb.Booster . . . . . 47

**Index** **48**

agaricus.test *Test part from Mushroom Data Set*

**Description**

This data set is originally from the Mushroom data set, UCI Machine Learning Repository. This data set includes the following fields:

- label: the label for each record
- data: a sparse Matrix of dgCMatrix class, with 126 columns.

**Usage**

```
data(agaricus.test)
```

**Format**

A list containing a label vector, and a dgCMatrix object with 1611 rows and 126 variables

**References**

<https://archive.ics.uci.edu/ml/datasets/Mushroom>  
 Bache, K. & Lichman, M. (2013). UCI Machine Learning Repository [<http://archive.ics.uci.edu/ml>].  
 Irvine, CA: University of California, School of Information and Computer Science.

---

agaricus.train      *Training part from Mushroom Data Set*

---

**Description**

This data set is originally from the Mushroom data set, UCI Machine Learning Repository. This data set includes the following fields:

- label: the label for each record
- data: a sparse Matrix of dgCMatix class, with 126 columns.

**Usage**

```
data(agaricus.train)
```

**Format**

A list containing a label vector, and a dgCMatix object with 6513 rows and 127 variables

**References**

<https://archive.ics.uci.edu/ml/datasets/Mushroom>

Bache, K. & Lichman, M. (2013). UCI Machine Learning Repository [<http://archive.ics.uci.edu/ml>]. Irvine, CA: University of California, School of Information and Computer Science.

---

bank      *Bank Marketing Data Set*

---

**Description**

This data set is originally from the Bank Marketing data set, UCI Machine Learning Repository.

It contains only the following: bank.csv with 10 randomly selected from 3 (older version of this dataset with less inputs).

**Usage**

```
data(bank)
```

**Format**

A data.table with 4521 rows and 17 variables

**References**

<http://archive.ics.uci.edu/ml/datasets/Bank+Marketing>

S. Moro, P. Cortez and P. Rita. (2014) A Data-Driven Approach to Predict the Success of Bank Telemarketing. Decision Support Systems

---

dim.lgb.Dataset	<i>Dimensions of an lgb.Dataset</i>
-----------------	-------------------------------------

---

## Description

Returns a vector of numbers of rows and of columns in an lgb.Dataset.

## Usage

```
## S3 method for class 'lgb.Dataset'  
dim(x)
```

## Arguments

x                    Object of class lgb.Dataset

## Details

Note: since nrow and ncol internally use dim, they can also be directly used with an lgb.Dataset object.

## Value

a vector of numbers of rows and of columns

## Examples

```
data(agaricus.train, package = "lightgbm")  
train <- agaricus.train  
dtrain <- lgb.Dataset(train$data, label = train$label)  
  
stopifnot(nrow(dtrain) == nrow(train$data))  
stopifnot(ncol(dtrain) == ncol(train$data))  
stopifnot(all(dim(dtrain) == dim(train$data)))
```

---

dimnames.lgb.Dataset *Handling of column names of lgb.Dataset*

---

### Description

Only column names are supported for `lgb.Dataset`, thus setting of row names would have no effect and returned row names would be `NULL`.

### Usage

```
## S3 method for class 'lgb.Dataset'  
dimnames(x)  
  
## S3 replacement method for class 'lgb.Dataset'  
dimnames(x) <- value
```

### Arguments

<code>x</code>	object of class <code>lgb.Dataset</code>
<code>value</code>	a list of two elements: the first one is ignored and the second one is column names

### Details

Generic `dimnames` methods are used by `colnames`. Since row names are irrelevant, it is recommended to use `colnames` directly.

### Value

A list with the dimension names of the dataset

### Examples

```
data(agaricus.train, package = "lightgbm")  
train <- agaricus.train  
dtrain <- lgb.Dataset(train$data, label = train$label)  
lgb.Dataset.construct(dtrain)  
dimnames(dtrain)  
colnames(dtrain)  
colnames(dtrain) <- make.names(seq_len(ncol(train$data)))  
print(dtrain, verbose = TRUE)
```

---

getLGBMThreads	<i>Get default number of threads used by LightGBM</i>
----------------	---

---

### Description

LightGBM attempts to speed up many operations by using multi-threading. The number of threads used in those operations can be controlled via the `num_threads` parameter passed through `params` to functions like `lgb.train` and `lgb.Dataset`. However, some operations (like materializing a model from a text file) are done via code paths that don't explicitly accept thread-control configuration.

Use this function to see the default number of threads LightGBM will use for such operations.

### Usage

```
getLGBMthreads()
```

### Value

number of threads as an integer. -1 means that in situations where parameter `num_threads` is not explicitly supplied, LightGBM will choose a number of threads to use automatically.

### See Also

[setLGBMthreads](#)

---

get_field	<i>Get one attribute of a lgb.Dataset</i>
-----------	---

---

### Description

Get one attribute of a `lgb.Dataset`

### Usage

```
get_field(dataset, field_name)
```

```
## S3 method for class 'lgb.Dataset'
get_field(dataset, field_name)
```

### Arguments

<code>dataset</code>	Object of class <code>lgb.Dataset</code>
<code>field_name</code>	String with the name of the attribute to get. One of the following. <ul style="list-style-type: none"> <li>• <code>label</code>: label lightgbm learns from ;</li> <li>• <code>weight</code>: to do a weight rescale ;</li> </ul>

- `group`: used for learning-to-rank tasks. An integer vector describing how to group rows together as ordered results from the same set of candidate results to be ranked. For example, if you have a 100-document dataset with `group = c(10, 20, 40, 10, 10, 10)`, that means that you have 6 groups, where the first 10 records are in the first group, records 11-30 are in the second group, etc.
- `init_score`: initial score is the base prediction lightgbm will boost from.

**Value**

requested attribute

**Examples**

```
data(agaricus.train, package = "lightgbm")
train <- agaricus.train
dtrain <- lgb.Dataset(train$data, label = train$label)
lgb.Dataset.construct(dtrain)

labels <- lightgbm::get_field(dtrain, "label")
lightgbm::set_field(dtrain, "label", 1 - labels)

labels2 <- lightgbm::get_field(dtrain, "label")
stopifnot(all(labels2 == 1 - labels))
```

---

`lgb.configure_fast_predict`

*Configure Fast Single-Row Predictions*

---

**Description**

Pre-configures a LightGBM model object to produce fast single-row predictions for a given input data type, prediction type, and parameters.

**Usage**

```
lgb.configure_fast_predict(
  model,
  csr = FALSE,
  start_iteration = NULL,
  num_iteration = NULL,
  type = "response",
  params = list()
)
```



**Arguments**

model	LighGBM model object (class <code>lgb.Booster</code> ). <b>The object will be modified in-place.</b>
csr	Whether the prediction function is going to be called on sparse CSR inputs. If FALSE, will be assumed that predictions are going to be called on single-row regular R matrices.
start_iteration	int or None, optional (default=None) Start index of the iteration to predict. If None or $\leq 0$ , starts from the first iteration.
num_iteration	int or None, optional (default=None) Limit number of iterations in the prediction. If None, if the best iteration exists and start_iteration is None or $\leq 0$ , the best iteration is used; otherwise, all iterations from start_iteration are used. If $\leq 0$ , all iterations from start_iteration are used (no limits).
type	Type of prediction to output. Allowed types are: <ul style="list-style-type: none"> <li>• "response": will output the predicted score according to the objective function being optimized (depending on the link function that the objective uses), after applying any necessary transformations - for example, for objective="binary", it will output class probabilities.</li> <li>• "class": for classification objectives, will output the class with the highest predicted probability. For other objectives, will output the same as "response". Note that "class" is not a supported type for <a href="#">lgb.configure_fast_predict</a> (see the documentation of that function for more details).</li> <li>• "raw": will output the non-transformed numbers (sum of predictions from boosting iterations' results) from which the "response" number is produced for a given objective function - for example, for objective="binary", this corresponds to log-odds. For many objectives such as "regression", since no transformation is applied, the output will be the same as for "response".</li> <li>• "leaf": will output the index of the terminal node / leaf at which each observations falls in each tree in the model, outputted as integers, with one column per tree.</li> <li>• "contrib": will return the per-feature contributions for each prediction, including an intercept (each feature will produce one column).</li> </ul> <p>Note that, if using custom objectives, types "class" and "response" will not be available and will default towards using "raw" instead.</p> <p>If the model was fit through function <a href="#">lightgbm</a> and it was passed a factor as labels, passing the prediction type through <code>params</code> instead of through this argument might result in factor levels for classification objectives not being applied correctly to the resulting output.</p> <p><i>New in version 4.0.0</i></p>
params	a list of additional named parameters. See <a href="#">the "Predict Parameters" section of the documentation</a> for a list of parameters and valid values. Where these conflict with the values of keyword arguments to this function, the values in <code>params</code> take precedence.

## Details

Calling this function multiple times with different parameters might not override the previous configuration and might trigger undefined behavior.

Any saved configuration for fast predictions might be lost after making a single-row prediction of a different type than what was configured (except for types "response" and "class", which can be switched between each other at any time without losing the configuration).

In some situations, setting a fast prediction configuration for one type of prediction might cause the prediction function to keep using that configuration for single-row predictions even if the requested type of prediction is different from what was configured.

Note that this function will not accept argument `type="class"` - for such cases, one can pass `type="response"` to this function and then `type="class"` to the `predict` function - the fast configuration will not be lost or altered if the switch is between "response" and "class".

The configuration does not survive de-serializations, so it has to be generated anew in every R process that is going to use it (e.g. if loading a model object through `readRDS`, whatever configuration was there previously will be lost).

Requesting a different prediction type or passing parameters to `predict.lgb.Booster` will cause it to ignore the fast-predict configuration and take the slow route instead (but be aware that an existing configuration might not always be overridden by supplying different parameters or prediction type, so make sure to check that the output is what was expected when a prediction is to be made on a single row for something different than what is configured).

Note that, if configuring a non-default prediction type (such as leaf indices), then that type must also be passed in the call to `predict.lgb.Booster` in order for it to use the configuration. This also applies for `start_iteration` and `num_iteration`, but **the params list must be empty** in the call to `predict`.

Predictions about feature contributions do not allow a fast route for CSR inputs, and as such, this function will produce an error if passing `csr=TRUE` and `type="contrib"` together.

## Value

The same model that was passed as input, invisibly, with the desired configuration stored inside it and available to be used in future calls to `predict.lgb.Booster`.

## Examples

```
library(lightgbm)
data(mtcars)
X <- as.matrix(mtcars[, -1L])
y <- mtcars[, 1L]
dtrain <- lgb.Dataset(X, label = y, params = list(max_bin = 5L))
params <- list(
  min_data_in_leaf = 2L
  , num_threads = 2L
)
model <- lgb.train(
```

```

    params = params
  , data = dtrain
  , obj = "regression"
  , nrounds = 5L
  , verbose = -1L
)
lgb.configure_fast_predict(model)

x_single <- X[11L, , drop = FALSE]
predict(model, x_single)

# Will not use it if the prediction to be made
# is different from what was configured
predict(model, x_single, type = "leaf")

```

---

lgb.convert\_with\_rules

*Data preparator for LightGBM datasets with rules (integer)*

---

## Description

Attempts to prepare a clean dataset to prepare to put in a `lgb.Dataset`. Factor, character, and logical columns are converted to integer. Missing values in factors and characters will be filled with 0L. Missing values in logicals will be filled with -1L.

This function returns and optionally takes in "rules" the describe exactly how to convert values in columns.

Columns that contain only NA values will be converted by this function but will not show up in the returned rules.

NOTE: In previous releases of LightGBM, this function was called `lgb.prepare_rules2`.

## Usage

```
lgb.convert_with_rules(data, rules = NULL)
```

## Arguments

<code>data</code>	A <code>data.frame</code> or <code>data.table</code> to prepare.
<code>rules</code>	A set of rules from the data preparator, if already used. This should be an R list, where names are column names in <code>data</code> and values are named character vectors whose names are column values and whose values are new values to replace them with.

## Value

A list with the cleaned dataset (`data`) and the rules (`rules`). Note that the data must be converted to a matrix format (`as.matrix`) for input in `lgb.Dataset`.

## Examples

```
data(iris)

str(iris)

new_iris <- lgb.convert_with_rules(data = iris)
str(new_iris$data)

data(iris) # Erase iris dataset
iris$Species[1L] <- "NEW FACTOR" # Introduce junk factor (NA)

# Use conversion using known rules
# Unknown factors become 0, excellent for sparse datasets
newer_iris <- lgb.convert_with_rules(data = iris, rules = new_iris$rules)

# Unknown factor is now zero, perfect for sparse datasets
newer_iris$data[1L, ] # Species became 0 as it is an unknown factor

newer_iris$data[1L, 5L] <- 1.0 # Put back real initial value

# Is the newly created dataset equal? YES!
all.equal(new_iris$data, newer_iris$data)

# Can we test our own rules?
data(iris) # Erase iris dataset

# We remapped values differently
personal_rules <- list(
  Species = c(
    "setosa" = 3L
    , "versicolor" = 2L
    , "virginica" = 1L
  )
)
newest_iris <- lgb.convert_with_rules(data = iris, rules = personal_rules)
str(newest_iris$data) # SUCCESS!
```

---

lgb.cv

*Main CV logic for LightGBM*

---

## Description

Cross validation logic used by LightGBM

## Usage

```
lgb.cv(
```

```

params = list(),
data,
nrounds = 100L,
nfold = 3L,
label = NULL,
weight = NULL,
obj = NULL,
eval = NULL,
verbose = 1L,
record = TRUE,
eval_freq = 1L,
showsd = TRUE,
stratified = TRUE,
folds = NULL,
init_model = NULL,
colnames = NULL,
categorical_feature = NULL,
early_stopping_rounds = NULL,
callbacks = list(),
reset_data = FALSE,
serializable = TRUE,
eval_train_metric = FALSE
)

```

## Arguments

params	a list of parameters. See <a href="#">the "Parameters" section of the documentation</a> for a list of parameters and valid values.
data	a <code>lgb.Dataset</code> object, used for training. Some functions, such as <code>lgb.cv</code> , may allow you to pass other types of data like matrix and then separately supply label as a keyword argument.
nrounds	number of training rounds
nfold	the original dataset is randomly partitioned into nfold equal size subsamples.
label	Vector of labels, used if data is not an <code>lgb.Dataset</code>
weight	vector of response values. If not NULL, will set to dataset
obj	objective function, can be character or custom objective function. Examples include regression, regression_l1, huber, binary, lambdarank, multiclass, multiclass
eval	evaluation function(s). This can be a character vector, function, or list with a mixture of strings and functions. <ul style="list-style-type: none"> <li><b>a. character vector:</b> If you provide a character vector to this argument, it should contain strings with valid evaluation metrics. See <a href="#">The "metric" section of the documentation</a> for a list of valid metrics.</li> <li><b>b. function:</b> You can provide a custom evaluation function. This should accept the keyword arguments <code>preds</code> and <code>dtrain</code> and should return a named list with three elements:</li> </ul>

- name: A string with the name of the metric, used for printing and storing results.
- value: A single number indicating the value of the metric for the given predictions and true values
- higher\_better: A boolean indicating whether higher values indicate a better fit. For example, this would be FALSE for metrics like MAE or RMSE.

- **c. list:** If a list is given, it should only contain character vectors and functions. These should follow the requirements from the descriptions above.

verbose	verbosity for output, if $\leq 0$ and <code>valids</code> has been provided, also will disable the printing of evaluation during training
record	Boolean, TRUE will record iteration message to <code>booster\$record_evals</code>
eval_freq	evaluation output frequency, only effective when <code>verbose &gt; 0</code> and <code>valids</code> has been provided
showsd	boolean, whether to show standard deviation of cross validation. This parameter defaults to TRUE. Setting it to FALSE can lead to a slight speedup by avoiding unnecessary computation.
stratified	a boolean indicating whether sampling of folds should be stratified by the values of outcome labels.
folds	list provides a possibility to use a list of pre-defined CV folds (each element must be a vector of test fold's indices). When folds are supplied, the <code>nfold</code> and <code>stratified</code> parameters are ignored.
init_model	path of model file or <code>lgb.Booster</code> object, will continue training from this model
colnames	feature names, if not null, will use this to overwrite the names in dataset
categorical_feature	categorical features. This can either be a character vector of feature names or an integer vector with the indices of the features (e.g. <code>c(1L, 10L)</code> to say "the first and tenth columns").
early_stopping_rounds	int. Activates early stopping. When this parameter is non-null, training will stop if the evaluation of any metric on any validation set fails to improve for <code>early_stopping_rounds</code> consecutive boosting rounds. If training stops early, the returned model will have attribute <code>best_iter</code> set to the iteration number of the best iteration.
callbacks	List of callback functions that are applied at each iteration.
reset_data	Boolean, setting it to TRUE (not the default value) will transform the booster model into a predictor model which frees up memory and the original datasets
serializable	whether to make the resulting objects serializable through functions such as <code>save</code> or <code>saveRDS</code> (see section "Model serialization").
eval_train_metric	boolean, whether to add the cross validation results on the training data. This parameter defaults to FALSE. Setting it to TRUE will increase run time.

### Value

a trained model `lgb.CVBooster`.

## Early Stopping

"early stopping" refers to stopping the training process if the model's performance on a given validation set does not improve for several consecutive iterations.

If multiple arguments are given to `eval`, their order will be preserved. If you enable early stopping by setting `early_stopping_rounds` in `params`, by default all metrics will be considered for early stopping.

If you want to only consider the first metric for early stopping, pass `first_metric_only = TRUE` in `params`. Note that if you also specify `metric` in `params`, that metric will be considered the "first" one. If you omit `metric`, a default metric will be used based on your choice for the parameter `obj` (keyword argument) or `objective` (passed into `params`).

## Examples

```
data(agaricus.train, package = "lightgbm")
train <- agaricus.train
dtrain <- lgb.Dataset(train$data, label = train$label)
params <- list(
  objective = "regression"
  , metric = "l2"
  , min_data = 1L
  , learning_rate = 1.0
  , num_threads = 2L
)
model <- lgb.cv(
  params = params
  , data = dtrain
  , nrounds = 5L
  , nfold = 3L
)
```

---

`lgb.Dataset`

*Construct lgb.Dataset object*

---

## Description

Construct `lgb.Dataset` object from dense matrix, sparse matrix or local file (that was created previously by saving an `lgb.Dataset`).

## Usage

```
lgb.Dataset(
  data,
  params = list(),
  reference = NULL,
```

```

    colnames = NULL,
    categorical_feature = NULL,
    free_raw_data = TRUE,
    label = NULL,
    weight = NULL,
    group = NULL,
    init_score = NULL
  )

```

### Arguments

<code>data</code>	a matrix object, a <code>dgCMatrx</code> object, a character representing a path to a text file (CSV, TSV, or LibSVM), or a character representing a path to a binary <code>lgb.Dataset</code> file
<code>params</code>	a list of parameters. See <a href="#">The "Dataset Parameters" section of the documentation</a> for a list of parameters and valid values.
<code>reference</code>	reference dataset. When LightGBM creates a Dataset, it does some preprocessing like binning continuous features into histograms. If you want to apply the same bin boundaries from an existing dataset to new data, pass that existing Dataset to this argument.
<code>colnames</code>	names of columns
<code>categorical_feature</code>	categorical features. This can either be a character vector of feature names or an integer vector with the indices of the features (e.g. <code>c(1L, 10L)</code> to say "the first and tenth columns").
<code>free_raw_data</code>	LightGBM constructs its data format, called a "Dataset", from tabular data. By default, that Dataset object on the R side does not keep a copy of the raw data. This reduces LightGBM's memory consumption, but it means that the Dataset object cannot be changed after it has been constructed. If you'd prefer to be able to change the Dataset object after construction, set <code>free_raw_data = FALSE</code> .
<code>label</code>	vector of labels to use as the target variable
<code>weight</code>	numeric vector of sample weights
<code>group</code>	used for learning-to-rank tasks. An integer vector describing how to group rows together as ordered results from the same set of candidate results to be ranked. For example, if you have a 100-document dataset with <code>group = c(10, 20, 40, 10, 10, 10)</code> , that means that you have 6 groups, where the first 10 records are in the first group, records 11-30 are in the second group, etc.
<code>init_score</code>	initial score is the base prediction lightgbm will boost from

### Value

constructed dataset

### Examples



```
data(agaricus.train, package = "lightgbm")
train <- agaricus.train
dtrain <- lgb.Dataset(train$data, label = train$label)
data_file <- tempfile(fileext = ".data")
lgb.Dataset.save(dtrain, data_file)
dtrain <- lgb.Dataset(data_file)
lgb.Dataset.construct(dtrain)
```

---

`lgb.Dataset.construct` *Construct Dataset explicitly*

---

### **Description**

Construct Dataset explicitly

### **Usage**

```
lgb.Dataset.construct(dataset)
```

### **Arguments**

dataset            Object of class `lgb.Dataset`

### **Value**

constructed dataset

### **Examples**

```
data(agaricus.train, package = "lightgbm")
train <- agaricus.train
dtrain <- lgb.Dataset(train$data, label = train$label)
lgb.Dataset.construct(dtrain)
```

---

```
lgb.Dataset.create.valid
```

*Construct validation data*

---

## Description

Construct validation data according to training data

## Usage

```
lgb.Dataset.create.valid(
  dataset,
  data,
  label = NULL,
  weight = NULL,
  group = NULL,
  init_score = NULL,
  params = list()
)
```

## Arguments

dataset	lgb.Dataset object, training data
data	a matrix object, a dgCMatrx object, a character representing a path to a text file (CSV, TSV, or LibSVM), or a character representing a path to a binary Dataset file
label	vector of labels to use as the target variable
weight	numeric vector of sample weights
group	used for learning-to-rank tasks. An integer vector describing how to group rows together as ordered results from the same set of candidate results to be ranked. For example, if you have a 100-document dataset with <code>group = c(10, 20, 40, 10, 10, 10)</code> , that means that you have 6 groups, where the first 10 records are in the first group, records 11-30 are in the second group, etc.
init_score	initial score is the base prediction lightgbm will boost from
params	a list of parameters. See <a href="#">The "Dataset Parameters" section of the documentation</a> for a list of parameters and valid values. If this is an empty list (the default), the validation Dataset will have the same parameters as the Dataset passed to argument dataset.

## Value

constructed dataset

**Examples**

```
data(agaricus.train, package = "lightgbm")
train <- agaricus.train
dtrain <- lgb.Dataset(train$data, label = train$label)
data(agaricus.test, package = "lightgbm")
test <- agaricus.test
dtest <- lgb.Dataset.create.valid(dtrain, test$data, label = test$label)

# parameters can be changed between the training data and validation set,
# for example to account for training data in a text file with a header row
# and validation data in a text file without it
train_file <- tempfile(pattern = "train_", fileext = ".csv")
write.table(
  data.frame(y = rnorm(100L), x1 = rnorm(100L), x2 = rnorm(100L))
  , file = train_file
  , sep = ","
  , col.names = TRUE
  , row.names = FALSE
  , quote = FALSE
)

valid_file <- tempfile(pattern = "valid_", fileext = ".csv")
write.table(
  data.frame(y = rnorm(100L), x1 = rnorm(100L), x2 = rnorm(100L))
  , file = valid_file
  , sep = ","
  , col.names = FALSE
  , row.names = FALSE
  , quote = FALSE
)

dtrain <- lgb.Dataset(
  data = train_file
  , params = list(has_header = TRUE)
)
dtrain$construct()

dvalid <- lgb.Dataset(
  data = valid_file
  , params = list(has_header = FALSE)
)
dvalid$construct()
```

**Description**

Please note that `init_score` is not saved in binary file. If you need it, please set it again after loading Dataset.

**Usage**

```
lgb.Dataset.save(dataset, fname)
```

**Arguments**

<code>dataset</code>	object of class <code>lgb.Dataset</code>
<code>fname</code>	object filename of output file

**Value**

the dataset you passed in

**Examples**

```
data(agaricus.train, package = "lightgbm")
train <- agaricus.train
dtrain <- lgb.Dataset(train$data, label = train$label)
lgb.Dataset.save(dtrain, tempfile(fileext = ".bin"))
```

---

```
lgb.Dataset.set.categorical
      Set categorical feature of lgb.Dataset
```

---

**Description**

Set the categorical features of an `lgb.Dataset` object. Use this function to tell LightGBM which features should be treated as categorical.

**Usage**

```
lgb.Dataset.set.categorical(dataset, categorical_feature)
```

**Arguments**

<code>dataset</code>	object of class <code>lgb.Dataset</code>
<code>categorical_feature</code>	categorical features. This can either be a character vector of feature names or an integer vector with the indices of the features (e.g. <code>c(1L, 10L)</code> to say "the first and tenth columns").

**Value**

the dataset you passed in

**Examples**

```
data(agaricus.train, package = "lightgbm")
train <- agaricus.train
dtrain <- lgb.Dataset(train$data, label = train$label)
data_file <- tempfile(fileext = ".data")
lgb.Dataset.save(dtrain, data_file)
dtrain <- lgb.Dataset(data_file)
lgb.Dataset.set.categorical(dtrain, 1L:2L)
```

---

lgb.Dataset.set.reference

*Set reference of lgb.Dataset*

---

**Description**

If you want to use validation data, you should set reference to training data

**Usage**

```
lgb.Dataset.set.reference(dataset, reference)
```

**Arguments**

dataset	object of class lgb.Dataset
reference	object of class lgb.Dataset

**Value**

the dataset you passed in

**Examples**

```
# create training Dataset
data(agaricus.train, package = "lightgbm")
train <- agaricus.train
dtrain <- lgb.Dataset(train$data, label = train$label)
```

```
# create a validation Dataset, using dtrain as a reference
data(agaricus.test, package = "lightgbm")
test <- agaricus.test
dtest <- lgb.Dataset(test$data, label = test$label)
lgb.Dataset.set.reference(dtest, dtrain)
```

---

`lgb.drop_serialized`     *Drop serialized raw bytes in a LightGBM model object*

---

### Description

If a LightGBM model object was produced with argument ‘serializable=TRUE’, the R object will keep a copy of the underlying C++ object as raw bytes, which can be used to reconstruct such object after getting serialized and de-serialized, but at the cost of extra memory usage. If these raw bytes are not needed anymore, they can be dropped through this function in order to save memory. Note that the object will be modified in-place.

*New in version 4.0.0*

### Usage

```
lgb.drop_serialized(model)
```

### Arguments

`model`             `lgb.Booster` object which was produced with ‘serializable=TRUE’.

### Value

`lgb.Booster` (the same ‘model’ object that was passed as input, as invisible).

### See Also

[lgb.restore\\_handle](#), [lgb.make\\_serializable](#).

---

`lgb.dump`             *Dump LightGBM model to json*

---

### Description

Dump LightGBM model to json

### Usage

```
lgb.dump(booster, num_iteration = NULL)
```

**Arguments**

`booster`            Object of class `lgb.Booster`  
`num_iteration`    number of iteration want to predict with, NULL or  $\leq 0$  means use best iteration

**Value**

json format of model

**Examples**

```
library(lightgbm)

data(agaricus.train, package = "lightgbm")
train <- agaricus.train
dtrain <- lgb.Dataset(train$data, label = train$label)
data(agaricus.test, package = "lightgbm")
test <- agaricus.test
dtest <- lgb.Dataset.create.valid(dtrain, test$data, label = test$label)
params <- list(
  objective = "regression"
  , metric = "l2"
  , min_data = 1L
  , learning_rate = 1.0
  , num_threads = 2L
)
valids <- list(test = dtest)
model <- lgb.train(
  params = params
  , data = dtrain
  , nrounds = 10L
  , valids = valids
  , early_stopping_rounds = 5L
)
json_model <- lgb.dump(model)
```

---

`lgb.get.eval.result`    *Get record evaluation result from booster*

---

**Description**

Given a `lgb.Booster`, return evaluation results for a particular metric on a particular dataset.

**Usage**

```
lgb.get.eval.result(  
  booster,  
  data_name,  
  eval_name,  
  iters = NULL,  
  is_err = FALSE  
)
```

**Arguments**

<code>booster</code>	Object of class <code>lgb.Booster</code>
<code>data_name</code>	Name of the dataset to return evaluation results for.
<code>eval_name</code>	Name of the evaluation metric to return results for.
<code>iters</code>	An integer vector of iterations you want to get evaluation results for. If <code>NULL</code> (the default), evaluation results for all iterations will be returned.
<code>is_err</code>	<code>TRUE</code> will return evaluation error instead

**Value**

numeric vector of evaluation result

**Examples**

```
# train a regression model  
data(agaricus.train, package = "lightgbm")  
train <- agaricus.train  
dtrain <- lgb.Dataset(train$data, label = train$label)  
data(agaricus.test, package = "lightgbm")  
test <- agaricus.test  
dtest <- lgb.Dataset.create.valid(dtrain, test$data, label = test$label)  
params <- list(  
  objective = "regression"  
  , metric = "l2"  
  , min_data = 1L  
  , learning_rate = 1.0  
  , num_threads = 2L  
)  
valids <- list(test = dtest)  
model <- lgb.train(  
  params = params  
  , data = dtrain  
  , nrounds = 5L  
  , valids = valids  
)  
  
# Examine valid data_name values
```



```
print(setdiff(names(model$record_evals), "start_iter"))

# Examine valid eval_name values for dataset "test"
print(names(model$record_evals[["test"]]))

# Get L2 values for "test" dataset
lgb.get.eval.result(model, "test", "l2")
```

---

lgb.importance	<i>Compute feature importance in a model</i>
----------------	--

---

### Description

Creates a data.table of feature importances in a model.

### Usage

```
lgb.importance(model, percentage = TRUE)
```

### Arguments

model	object of class lgb.Booster.
percentage	whether to show importance in relative percentage.

### Value

For a tree model, a data.table with the following columns:

- Feature: Feature names in the model.
- Gain: The total gain of this feature's splits.
- Cover: The number of observation related to this feature.
- Frequency: The number of times a feature splitted in trees.

### Examples

```
data(agaricus.train, package = "lightgbm")
train <- agaricus.train
dtrain <- lgb.Dataset(train$data, label = train$label)

params <- list(
  objective = "binary"
  , learning_rate = 0.1
  , max_depth = -1L
  , min_data_in_leaf = 1L
```

```

    , min_sum_hessian_in_leaf = 1.0
    , num_threads = 2L
  )
model <- lgb.train(
  params = params
  , data = dtrain
  , nrounds = 5L
)

tree_imp1 <- lgb.importance(model, percentage = TRUE)
tree_imp2 <- lgb.importance(model, percentage = FALSE)

```

---

lgb.interprete	<i>Compute feature contribution of prediction</i>
----------------	---

---

### Description

Computes feature contribution components of rawscore prediction.

### Usage

```
lgb.interprete(model, data, idxset, num_iteration = NULL)
```

### Arguments

model	object of class lgb.Booster.
data	a matrix object or a dgCMatrix object.
idxset	an integer vector of indices of rows needed.
num_iteration	number of iteration want to predict with, NULL or <= 0 means use best iteration.

### Value

For regression, binary classification and lambdarank model, a list of data.table with the following columns:

- Feature: Feature names in the model.
- Contribution: The total contribution of this feature's splits.

For multiclass classification, a list of data.table with the Feature column and Contribution columns to each class.

## Examples

```
Logit <- function(x) log(x / (1.0 - x))
data(agaricus.train, package = "lightgbm")
train <- agaricus.train
dtrain <- lgb.Dataset(train$data, label = train$label)
set_field(
  dataset = dtrain
  , field_name = "init_score"
  , data = rep(Logit(mean(train$label)), length(train$label))
)
data(agaricus.test, package = "lightgbm")
test <- agaricus.test

params <- list(
  objective = "binary"
  , learning_rate = 0.1
  , max_depth = -1L
  , min_data_in_leaf = 1L
  , min_sum_hessian_in_leaf = 1.0
  , num_threads = 2L
)
model <- lgb.train(
  params = params
  , data = dtrain
  , nrounds = 3L
)

tree_interpretation <- lgb.interprete(model, test$data, 1L:5L)
```

---

lgb.load

*Load LightGBM model*

---

## Description

Load LightGBM takes in either a file path or model string. If both are provided, Load will default to loading from file

## Usage

```
lgb.load(filename = NULL, model_str = NULL)
```

## Arguments

filename	path of model file
model_str	a str containing the model (as a character or raw vector)

**Value**

lgb.Booster

**Examples**

```
data(agaricus.train, package = "lightgbm")
train <- agaricus.train
dtrain <- lgb.Dataset(train$data, label = train$label)
data(agaricus.test, package = "lightgbm")
test <- agaricus.test
dtest <- lgb.Dataset.create.valid(dtrain, test$data, label = test$label)
params <- list(
  objective = "regression"
  , metric = "l2"
  , min_data = 1L
  , learning_rate = 1.0
  , num_threads = 2L
)
valids <- list(test = dtest)
model <- lgb.train(
  params = params
  , data = dtrain
  , nrounds = 5L
  , valids = valids
  , early_stopping_rounds = 3L
)
model_file <- tempfile(fileext = ".txt")
lgb.save(model, model_file)
load_booster <- lgb.load(filename = model_file)
model_string <- model$save_model_to_string(NULL) # saves best iteration
load_booster_from_str <- lgb.load(model_str = model_string)
```

---

`lgb.make_serializable` *Make a LightGBM object serializable by keeping raw bytes*

---

**Description**

If a LightGBM model object was produced with argument ‘serializable=FALSE’, the R object will not be serializable (e.g. cannot save and load with `saveRDS` and `readRDS`) as it will lack the raw bytes needed to reconstruct its underlying C++ object. This function can be used to forcibly produce those serialized raw bytes and make the object serializable. Note that the object will be modified in-place.

*New in version 4.0.0*

**Usage**

```
lgb.make_serializable(model)
```

**Arguments**

model                    lgb.Booster object which was produced with ‘serializable=FALSE’.

**Value**

lgb.Booster (the same ‘model’ object that was passed as input, as invisible).

**See Also**

[lgb.restore\\_handle](#), [lgb.drop\\_serialized](#).

---

lgb.model.dt.tree            *Parse a LightGBM model json dump*

---

**Description**

Parse a LightGBM model json dump into a data.table structure.

**Usage**

```
lgb.model.dt.tree(model, num_iteration = NULL)
```

**Arguments**

model                    object of class lgb.Booster  
num\_iteration            number of iterations you want to predict with. NULL or <= 0 means use best iteration

**Value**

A data.table with detailed information about model trees’ nodes and leaves.

The columns of the data.table are:

- tree\_index: ID of a tree in a model (integer)
- split\_index: ID of a node in a tree (integer)
- split\_feature: for a node, it’s a feature name (character); for a leaf, it simply labels it as "NA"
- node\_parent: ID of the parent node for current node (integer)
- leaf\_index: ID of a leaf in a tree (integer)
- leaf\_parent: ID of the parent node for current leaf (integer)
- split\_gain: Split gain of a node

- `threshold`: Splitting threshold value of a node
- `decision_type`: Decision type of a node
- `default_left`: Determine how to handle NA value, TRUE -> Left, FALSE -> Right
- `internal_value`: Node value
- `internal_count`: The number of observation collected by a node
- `leaf_value`: Leaf value
- `leaf_count`: The number of observation collected by a leaf

## Examples

```
data(agaricus.train, package = "lightgbm")
train <- agaricus.train
dtrain <- lgb.Dataset(train$data, label = train$label)

params <- list(
  objective = "binary"
  , learning_rate = 0.01
  , num_leaves = 63L
  , max_depth = -1L
  , min_data_in_leaf = 1L
  , min_sum_hessian_in_leaf = 1.0
  , num_threads = 2L
)
model <- lgb.train(params, dtrain, 10L)

tree_dt <- lgb.model.dt.tree(model)
```

---

`lgb.plot.importance`     *Plot feature importance as a bar graph*

---

## Description

Plot previously calculated feature importance: Gain, Cover and Frequency, as a bar graph.

## Usage

```
lgb.plot.importance(
  tree_imp,
  top_n = 10L,
  measure = "Gain",
  left_margin = 10L,
  cex = NULL
)
```

## Arguments

tree_imp	a data.table returned by <code>lgb.importance</code> .
top_n	maximal number of top features to include into the plot.
measure	the name of importance measure to plot, can be "Gain", "Cover" or "Frequency".
left_margin	(base R barplot) allows to adjust the left margin size to fit feature names.
cex	(base R barplot) passed as <code>cex.names</code> parameter to <code>barplot</code> . Set a number smaller than 1.0 to make the bar labels smaller than R's default and values greater than 1.0 to make them larger.

## Details

The graph represents each feature as a horizontal bar of length proportional to the defined importance of a feature. Features are shown ranked in a decreasing importance order.

## Value

The `lgb.plot.importance` function creates a barplot and silently returns a processed data.table with `top_n` features sorted by defined importance.

## Examples

```
data(agaricus.train, package = "lightgbm")
train <- agaricus.train
dtrain <- lgb.Dataset(train$data, label = train$label)

params <- list(
  objective = "binary"
  , learning_rate = 0.1
  , min_data_in_leaf = 1L
  , min_sum_hessian_in_leaf = 1.0
  , num_threads = 2L
)

model <- lgb.train(
  params = params
  , data = dtrain
  , nrounds = 5L
)

tree_imp <- lgb.importance(model, percentage = TRUE)
lgb.plot.importance(tree_imp, top_n = 5L, measure = "Gain")
```

---

`lgb.plot.interpretation`*Plot feature contribution as a bar graph*

---

### Description

Plot previously calculated feature contribution as a bar graph.

### Usage

```
lgb.plot.interpretation(  
  tree_interpretation_dt,  
  top_n = 10L,  
  cols = 1L,  
  left_margin = 10L,  
  cex = NULL  
)
```

### Arguments

<code>tree_interpretation_dt</code>	a <code>data.table</code> returned by <a href="#">lgb.interprete</a> .
<code>top_n</code>	maximal number of top features to include into the plot.
<code>cols</code>	the column numbers of layout, will be used only for multiclass classification feature contribution.
<code>left_margin</code>	(base R barplot) allows to adjust the left margin size to fit feature names.
<code>cex</code>	(base R barplot) passed as <code>cex.names</code> parameter to <code>barplot</code> .

### Details

The graph represents each feature as a horizontal bar of length proportional to the defined contribution of a feature. Features are shown ranked in a decreasing contribution order.

### Value

The `lgb.plot.interpretation` function creates a barplot.

### Examples

```
Logit <- function(x) {  
  log(x / (1.0 - x))  
}  
data(agaricus.train, package = "lightgbm")  
labels <- agaricus.train$label
```



```
dtrain <- lgb.Dataset(  
  agaricus.train$data  
  , label = labels  
)  
set_field(  
  dataset = dtrain  
  , field_name = "init_score"  
  , data = rep(Logit(mean(labels)), length(labels))  
)  
  
data(agaricus.test, package = "lightgbm")  
  
params <- list(  
  objective = "binary"  
  , learning_rate = 0.1  
  , max_depth = -1L  
  , min_data_in_leaf = 1L  
  , min_sum_hessian_in_leaf = 1.0  
  , num_threads = 2L  
)  
model <- lgb.train(  
  params = params  
  , data = dtrain  
  , nrounds = 5L  
)  
  
tree_interpretation <- lgb.interprete(  
  model = model  
  , data = agaricus.test$data  
  , idxset = 1L:5L  
)  
lgb.plot.interpretation(  
  tree_interpretation_dt = tree_interpretation[[1L]]  
  , top_n = 3L  
)
```

---

`lgb.restore_handle`      *Restore the C++ component of a de-serialized LightGBM model*

---

### Description

After a LightGBM model object is de-serialized through functions such as `save` or `saveRDS`, its underlying C++ object will be blank and needs to be restored to be able to use it. Such object is restored automatically when calling functions such as `predict`, but this function can be used to forcibly restore it beforehand. Note that the object will be modified in-place.

*New in version 4.0.0*

### Usage

```
lgb.restore_handle(model)
```

## Arguments

`model` `lgb.Booster` object which was de-serialized and whose underlying C++ object and R handle need to be restored.

## Details

Be aware that fast single-row prediction configurations are not restored through this function. If you wish to make fast single-row predictions using a `lgb.Booster` loaded this way, call [lgb.configure\\_fast\\_predict](#) on the loaded `lgb.Booster` object.

## Value

`lgb.Booster` (the same ‘`model`’ object that was passed as input, invisibly).

## See Also

[lgb.make\\_serializable](#), [lgb.drop\\_serialized](#).

## Examples

```
library(lightgbm)

data("agaricus.train")
model <- lightgbm(
  agaricus.train$data
  , agaricus.train$label
  , params = list(objective = "binary")
  , nrounds = 5L
  , verbose = 0
  , num_threads = 2L
)
fname <- tempfile(fileext="rds")
saveRDS(model, fname)

model_new <- readRDS(fname)
model_new$check_null_handle()
lgb.restore_handle(model_new)
model_new$check_null_handle()
```

---

`lgb.save`

*Save LightGBM model*

---

## Description

Save LightGBM model

**Usage**

```
lgb.save(booster, filename, num_iteration = NULL)
```

**Arguments**

booster	Object of class lgb.Booster
filename	saved filename
num_iteration	number of iteration want to predict with, NULL or <= 0 means use best iteration

**Value**

lgb.Booster

**Examples**

```
library(lightgbm)
data(agaricus.train, package = "lightgbm")
train <- agaricus.train
dtrain <- lgb.Dataset(train$data, label = train$label)
data(agaricus.test, package = "lightgbm")
test <- agaricus.test
dtest <- lgb.Dataset.create.valid(dtrain, test$data, label = test$label)
params <- list(
  objective = "regression"
  , metric = "l2"
  , min_data = 1L
  , learning_rate = 1.0
  , num_threads = 2L
)
valids <- list(test = dtest)
model <- lgb.train(
  params = params
  , data = dtrain
  , nrounds = 10L
  , valids = valids
  , early_stopping_rounds = 5L
)
lgb.save(model, tempfile(fileext = ".txt"))
```

## Description

Low-level R interface to train a LightGBM model. Unlike `lightgbm`, this function is focused on performance (e.g. speed, memory efficiency). It is also less likely to have breaking API changes in new releases than `lightgbm`.

## Usage

```
lgb.train(
  params = list(),
  data,
  nrounds = 100L,
  valids = list(),
  obj = NULL,
  eval = NULL,
  verbose = 1L,
  record = TRUE,
  eval_freq = 1L,
  init_model = NULL,
  colnames = NULL,
  categorical_feature = NULL,
  early_stopping_rounds = NULL,
  callbacks = list(),
  reset_data = FALSE,
  serializable = TRUE
)
```

## Arguments

<code>params</code>	a list of parameters. See <a href="#">the "Parameters" section of the documentation</a> for a list of parameters and valid values.
<code>data</code>	a <code>lgb.Dataset</code> object, used for training. Some functions, such as <code>lgb.cv</code> , may allow you to pass other types of data like <code>matrix</code> and then separately supply <code>label</code> as a keyword argument.
<code>nrounds</code>	number of training rounds
<code>valids</code>	a list of <code>lgb.Dataset</code> objects, used for validation
<code>obj</code>	objective function, can be character or custom objective function. Examples include <code>regression</code> , <code>regression_l1</code> , <code>huber</code> , <code>binary</code> , <code>lambda_rank</code> , <code>multiclass</code> , <code>multiclass</code>
<code>eval</code>	evaluation function(s). This can be a character vector, function, or list with a mixture of strings and functions. <ul style="list-style-type: none"> <li><b>a. character vector:</b> If you provide a character vector to this argument, it should contain strings with valid evaluation metrics. See <a href="#">The "metric" section of the documentation</a> for a list of valid metrics.</li> <li><b>b. function:</b> You can provide a custom evaluation function. This should accept the keyword arguments <code>preds</code> and <code>dtrain</code> and should return a named list with three elements:</li> </ul>

	<ul style="list-style-type: none"> <li>- name: A string with the name of the metric, used for printing and storing results.</li> <li>- value: A single number indicating the value of the metric for the given predictions and true values</li> <li>- higher_better: A boolean indicating whether higher values indicate a better fit. For example, this would be FALSE for metrics like MAE or RMSE.</li> </ul> <ul style="list-style-type: none"> <li>• <b>c. list:</b> If a list is given, it should only contain character vectors and functions. These should follow the requirements from the descriptions above.</li> </ul>
verbose	verbosity for output, if $\leq 0$ and <code>valids</code> has been provided, also will disable the printing of evaluation during training
record	Boolean, TRUE will record iteration message to <code>booster\$record_evals</code>
eval_freq	evaluation output frequency, only effective when <code>verbose &gt; 0</code> and <code>valids</code> has been provided
init_model	path of model file or <code>lgb.Booster</code> object, will continue training from this model
colnames	feature names, if not null, will use this to overwrite the names in dataset
categorical_feature	categorical features. This can either be a character vector of feature names or an integer vector with the indices of the features (e.g. <code>c(1L, 10L)</code> to say "the first and tenth columns").
early_stopping_rounds	int. Activates early stopping. When this parameter is non-null, training will stop if the evaluation of any metric on any validation set fails to improve for <code>early_stopping_rounds</code> consecutive boosting rounds. If training stops early, the returned model will have attribute <code>best_iter</code> set to the iteration number of the best iteration.
callbacks	List of callback functions that are applied at each iteration.
reset_data	Boolean, setting it to TRUE (not the default value) will transform the booster model into a predictor model which frees up memory and the original datasets
serializable	whether to make the resulting objects serializable through functions such as <code>save</code> or <code>saveRDS</code> (see section "Model serialization").

**Value**

a trained booster model `lgb.Booster`.

**Early Stopping**

"early stopping" refers to stopping the training process if the model's performance on a given validation set does not improve for several consecutive iterations.

If multiple arguments are given to `eval`, their order will be preserved. If you enable early stopping by setting `early_stopping_rounds` in `params`, by default all metrics will be considered for early stopping.

If you want to only consider the first metric for early stopping, pass `first_metric_only = TRUE` in `params`. Note that if you also specify `metric` in `params`, that metric will be considered the "first" one. If you omit `metric`, a default metric will be used based on your choice for the parameter `obj` (keyword argument) or `objective` (passed into `params`).

## Examples

```
data(agaricus.train, package = "lightgbm")
train <- agaricus.train
dtrain <- lgb.Dataset(train$data, label = train$label)
data(agaricus.test, package = "lightgbm")
test <- agaricus.test
dtest <- lgb.Dataset.create.valid(dtrain, test$data, label = test$label)
params <- list(
  objective = "regression"
  , metric = "l2"
  , min_data = 1L
  , learning_rate = 1.0
  , num_threads = 2L
)
valids <- list(test = dtest)
model <- lgb.train(
  params = params
  , data = dtrain
  , nrounds = 5L
  , valids = valids
  , early_stopping_rounds = 3L
)
```

---

lightgbm

*Train a LightGBM model*

---

## Description

High-level R interface to train a LightGBM model. Unlike `lgb.train`, this function is focused on compatibility with other statistics and machine learning interfaces in R. This focus on compatibility means that this interface may experience more frequent breaking API changes than `lgb.train`. For efficiency-sensitive applications, or for applications where breaking API changes across releases is very expensive, use `lgb.train`.

## Usage

```
lightgbm(
  data,
  label = NULL,
  weights = NULL,
  params = list(),
  nrounds = 100L,
  verbose = 1L,
  eval_freq = 1L,
```

```

    early_stopping_rounds = NULL,
    init_model = NULL,
    callbacks = list(),
    serializable = TRUE,
    objective = "auto",
    init_score = NULL,
    num_threads = NULL,
    ...
)

```

## Arguments

data	a <code>lgb.Dataset</code> object, used for training. Some functions, such as <code>lgb.cv</code> , may allow you to pass other types of data like matrix and then separately supply label as a keyword argument.
label	Vector of labels, used if data is not an <code>lgb.Dataset</code>
weights	Sample / observation weights for rows in the input data. If <code>NULL</code> , will assume that all observations / rows have the same importance / weight. <i>Changed from 'weight', in version 4.0.0</i>
params	a list of parameters. See the <a href="#">"Parameters" section of the documentation</a> for a list of parameters and valid values.
nrounds	number of training rounds
verbose	verbosity for output, if $\leq 0$ and <code>valids</code> has been provided, also will disable the printing of evaluation during training
eval_freq	evaluation output frequency, only effective when <code>verbose &gt; 0</code> and <code>valids</code> has been provided
early_stopping_rounds	int. Activates early stopping. When this parameter is non-null, training will stop if the evaluation of any metric on any validation set fails to improve for <code>early_stopping_rounds</code> consecutive boosting rounds. If training stops early, the returned model will have attribute <code>best_iter</code> set to the iteration number of the best iteration.
init_model	path of model file or <code>lgb.Booster</code> object, will continue training from this model
callbacks	List of callback functions that are applied at each iteration.
serializable	whether to make the resulting objects serializable through functions such as <code>save</code> or <code>saveRDS</code> (see section "Model serialization").
objective	Optimization objective (e.g. "regression", "binary", etc.). For a list of accepted objectives, see the <a href="#">"objective" item of the "Parameters" section of the documentation</a> . If passing "auto" and data is not of type <code>lgb.Dataset</code> , the objective will be determined according to what is passed for label: <ul style="list-style-type: none"> <li>• If passing a factor with two variables, will use objective "binary".</li> <li>• If passing a factor with more than two variables, will use objective "multiclass" (note that parameter <code>num_class</code> in this case will also be determined automatically from label).</li> </ul>

- Otherwise (or if passing `lgb.Dataset` as input), will use objective "regression".

*New in version 4.0.0*

`init_score` initial score is the base prediction lightgbm will boost from  
*New in version 4.0.0*

`num_threads` Number of parallel threads to use. For best speed, this should be set to the number of physical cores in the CPU - in a typical x86-64 machine, this corresponds to half the number of maximum threads.  
Be aware that using too many threads can result in speed degradation in smaller datasets (see the parameters documentation for more details).  
If passing zero, will use the default number of threads configured for OpenMP (typically controlled through an environment variable `OMP_NUM_THREADS`).  
If passing NULL (the default), will try to use the number of physical cores in the system, but be aware that getting the number of cores detected correctly requires package `RhpcBLASctl` to be installed.  
This parameter gets overridden by `num_threads` and its aliases under `params` if passed there.  
*New in version 4.0.0*

... Additional arguments passed to `lgb.train`. For example

- `valids`: a list of `lgb.Dataset` objects, used for validation
- `obj`: objective function, can be character or custom objective function. Examples include `regression`, `regression_l1`, `huber`, `binary`, `lambdarank`, `multiclass`, `multiclass`
- `eval`: evaluation function, can be (a list of) character or custom eval function
- `record`: Boolean, TRUE will record iteration message to `booster$record_evals`
- `colnames`: feature names, if not null, will use this to overwrite the names in dataset
- `categorical_feature`: categorical features. This can either be a character vector of feature names or an integer vector with the indices of the features (e.g. `c(1L, 10L)` to say "the first and tenth columns").
- `reset_data`: Boolean, setting it to TRUE (not the default value) will transform the booster model into a predictor model which frees up memory and the original datasets

**Value**

a trained `lgb.Booster`

**Early Stopping**

"early stopping" refers to stopping the training process if the model's performance on a given validation set does not improve for several consecutive iterations.

If multiple arguments are given to `eval`, their order will be preserved. If you enable early stopping by setting `early_stopping_rounds` in `params`, by default all metrics will be considered for early stopping.



If you want to only consider the first metric for early stopping, pass `first_metric_only = TRUE` in `params`. Note that if you also specify `metric` in `params`, that `metric` will be considered the "first" one. If you omit `metric`, a default metric will be used based on your choice for the parameter `obj` (keyword argument) or `objective` (passed into `params`).

---

predict.lgb.Booster     *Predict method for LightGBM model*

---

## Description

Predicted values based on class `lgb.Booster`

*New in version 4.0.0*

## Usage

```
## S3 method for class 'lgb.Booster'
predict(
  object,
  newdata,
  type = "response",
  start_iteration = NULL,
  num_iteration = NULL,
  header = FALSE,
  params = list(),
  ...
)
```

## Arguments

- |                      |  |
|----------------------|--|
| <code>object</code>  | Object of class <code>lgb.Booster</code>   |
| <code>newdata</code> | a matrix object, a <code>dgCMatrix</code> , a <code>dgRMatrix</code> object, a <code>dsparseVector</code> object, or a character representing a path to a text file (CSV, TSV, or LibSVM).<br>For sparse inputs, if predictions are only going to be made for a single row, it will be faster to use CSR format, in which case the data may be passed as either a single-row CSR matrix (class <code>dgRMatrix</code> from package <code>Matrix</code> ) or as a sparse numeric vector (class <code>dsparseVector</code> from package <code>Matrix</code> ).<br>If single-row predictions are going to be performed frequently, it is recommended to pre-configure the model object for fast single-row sparse predictions through function <a href="#">lgb.configure_fast_predict</a> .<br><i>Changed from 'data', in version 4.0.0</i> |
| <code>type</code>    | Type of prediction to output. Allowed types are: <ul style="list-style-type: none"> <li>"response": will output the predicted score according to the objective function being optimized (depending on the link function that the objective uses), after applying any necessary transformations - for example, for <code>objective="binary"</code>, it will output class probabilities.</li> </ul>  |

- "class": for classification objectives, will output the class with the highest predicted probability. For other objectives, will output the same as "response". Note that "class" is not a supported type for [lgb.configure\\_fast\\_predict](#) (see the documentation of that function for more details).
- "raw": will output the non-transformed numbers (sum of predictions from boosting iterations' results) from which the "response" number is produced for a given objective function - for example, for objective="binary", this corresponds to log-odds. For many objectives such as "regression", since no transformation is applied, the output will be the same as for "response".
- "leaf": will output the index of the terminal node / leaf at which each observations falls in each tree in the model, outputted as integers, with one column per tree.
- "contrib": will return the per-feature contributions for each prediction, including an intercept (each feature will produce one column).

Note that, if using custom objectives, types "class" and "response" will not be available and will default towards using "raw" instead.

If the model was fit through function [lightgbm](#) and it was passed a factor as labels, passing the prediction type through `params` instead of through this argument might result in factor levels for classification objectives not being applied correctly to the resulting output.

*New in version 4.0.0*

<code>start_iteration</code>	int or None, optional (default=None) Start index of the iteration to predict. If None or $\leq 0$ , starts from the first iteration.
<code>num_iteration</code>	int or None, optional (default=None) Limit number of iterations in the prediction. If None, if the best iteration exists and <code>start_iteration</code> is None or $\leq 0$ , the best iteration is used; otherwise, all iterations from <code>start_iteration</code> are used. If $\leq 0$ , all iterations from <code>start_iteration</code> are used (no limits).
<code>header</code>	only used for prediction for text file. True if text file has header
<code>params</code>	a list of additional named parameters. See the <a href="#">"Predict Parameters" section of the documentation</a> for a list of parameters and valid values. Where these conflict with the values of keyword arguments to this function, the values in <code>params</code> take precedence.
<code>...</code>	ignored

### Details

If the model object has been configured for fast single-row predictions through [lgb.configure\\_fast\\_predict](#), this function will use the prediction parameters that were configured for it - as such, extra prediction parameters should not be passed here, otherwise the configuration will be ignored and the slow route will be taken.

### Value

For prediction types that are meant to always return one output per observation (e.g. when predicting `type="response"` or `type="raw"` on a binary classification or regression objective), will return a vector with one element per row in `newdata`.

For prediction types that are meant to return more than one output per observation (e.g. when predicting `type="response"` or `type="raw"` on a multi-class objective, or when predicting `type="leaf"`, regardless of objective), will return a matrix with one row per observation in `newdata` and one column per output.

For `type="leaf"` predictions, will return a matrix with one row per observation in `newdata` and one column per tree. Note that for multiclass objectives, LightGBM trains one tree per class at each boosting iteration. That means that, for example, for a multiclass model with 3 classes, the leaf predictions for the first class can be found in columns 1, 4, 7, 10, etc.

For `type="contrib"`, will return a matrix of SHAP values with one row per observation in `newdata` and columns corresponding to features. For regression, ranking, cross-entropy, and binary classification objectives, this matrix contains one column per feature plus a final column containing the Shapley base value. For multiclass objectives, this matrix will represent `num_classes` such matrices, in the order "feature contributions for first class, feature contributions for second class, feature contributions for third class, etc."

If the model was fit through function `lightgbm` and it was passed a factor as labels, predictions returned from this function will retain the factor levels (either as values for `type="class"`, or as column names for `type="response"` and `type="raw"` for multi-class objectives). Note that passing the requested prediction type under `params` instead of through `type` might result in the factor levels not being present in the output.

## Examples

```
data(agaricus.train, package = "lightgbm")
train <- agaricus.train
dtrain <- lgb.Dataset(train$data, label = train$label)
data(agaricus.test, package = "lightgbm")
test <- agaricus.test
dtest <- lgb.Dataset.create.valid(dtrain, test$data, label = test$label)
params <- list(
  objective = "regression"
  , metric = "l2"
  , min_data = 1L
  , learning_rate = 1.0
  , num_threads = 2L
)
valids <- list(test = dtest)
model <- lgb.train(
  params = params
  , data = dtrain
  , nrounds = 5L
  , valids = valids
)
preds <- predict(model, test$data)

# pass other prediction parameters
preds <- predict(
  model,
```

```

    test$data,
    params = list(
      predict_disable_shape_check = TRUE
    )
  )
)

```

---

```
print.lgb.Booster
```

*Print method for LightGBM model*

---

### Description

Show summary information about a LightGBM model object (same as `summary`).

*New in version 4.0.0*

### Usage

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

### Arguments

<code>x</code>	Object of class <code>lgb.Booster</code>
<code>...</code>	Not used

### Value

The same input `x`, returned as invisible.

---

```
setLGBMThreads
```

*Set maximum number of threads used by LightGBM*

---

### Description

LightGBM attempts to speed up many operations by using multi-threading. The number of threads used in those operations can be controlled via the `num_threads` parameter passed through `params` to functions like [lgb.train](#) and [lgb.Dataset](#). However, some operations (like materializing a model from a text file) are done via code paths that don't explicitly accept thread-control configuration.

Use this function to set the maximum number of threads LightGBM will use for such operations.

This function affects all LightGBM operations in the same process.

So, for example, if you call `setLGBMthreads(4)`, no other multi-threaded LightGBM operation in the same process will use more than 4 threads.

Call `setLGBMthreads(-1)` to remove this limitation.

**Usage**

```
setLGBMthreads(num_threads)
```

**Arguments**

num\_threads      maximum number of threads to be used by LightGBM in multi-threaded operations

**See Also**

[getLGBMthreads](#)

---

set_field	<i>Set one attribute of a lgb.Dataset object</i>
-----------	--

---

**Description**

Set one attribute of a lgb.Dataset

**Usage**

```
set_field(dataset, field_name, data)
```

```
## S3 method for class 'lgb.Dataset'
set_field(dataset, field_name, data)
```

**Arguments**

dataset	Object of class lgb.Dataset
field_name	String with the name of the attribute to set. One of the following. <ul style="list-style-type: none"> <li>• label: label lightgbm learns from ;</li> <li>• weight: to do a weight rescale ;</li> <li>• group: used for learning-to-rank tasks. An integer vector describing how to group rows together as ordered results from the same set of candidate results to be ranked. For example, if you have a 100-document dataset with <code>group = c(10, 20, 40, 10, 10, 10)</code>, that means that you have 6 groups, where the first 10 records are in the first group, records 11-30 are in the second group, etc.</li> <li>• init_score: initial score is the base prediction lightgbm will boost from.</li> </ul>
data	The data for the field. See examples.

**Value**

The lgb.Dataset you passed in.

## Examples

```
data(agaricus.train, package = "lightgbm")
train <- agaricus.train
dtrain <- lgb.Dataset(train$data, label = train$label)
lgb.Dataset.construct(dtrain)

labels <- lightgbm::get_field(dtrain, "label")
lightgbm::set_field(dtrain, "label", 1 - labels)

labels2 <- lightgbm::get_field(dtrain, "label")
stopifnot(all.equal(labels2, 1 - labels))
```

---

slice

*Slice a dataset*

---

## Description

Get a new `lgb.Dataset` containing the specified rows of original `lgb.Dataset` object

## Usage

```
slice(dataset, idxset)

## S3 method for class 'lgb.Dataset'
slice(dataset, idxset)
```

## Arguments

<code>dataset</code>	Object of class <code>lgb.Dataset</code>
<code>idxset</code>	an integer vector of indices of rows needed

## Value

constructed sub dataset

## Examples

```
data(agaricus.train, package = "lightgbm")
train <- agaricus.train
dtrain <- lgb.Dataset(train$data, label = train$label)

dsub <- lightgbm::slice(dtrain, seq_len(42L))
```

```
lgb.Dataset.construct(dsub)
labels <- lightgbm::get_field(dsub, "label")
```

---

summary.lgb.Booster      *Summary method for LightGBM model*

---

### **Description**

Show summary information about a LightGBM model object (same as print).

*New in version 4.0.0*

### **Usage**

```
## S3 method for class 'lgb.Booster'
summary(object, ...)
```

### **Arguments**

object	Object of class lgb.Booster
...	Not used

### **Value**

The same input object, returned as invisible.

# Index

- \* **datasets**
  - agaricus.test, 3
  - agaricus.train, 4
  - bank, 4
- agaricus.test, 3
- agaricus.train, 4
- bank, 4
- barplot, 31
- dim.lgb.Dataset, 5
- dimnames.lgb.Dataset, 6
- dimnames<-lgb.Dataset  
(dimnames.lgb.Dataset), 6
- get\_field, 7
- getLGBMThreads, 7
- getLGBMthreads, 45
- getLGBMthreads(getLGBMThreads), 7
- lgb.configure\_fast\_predict, 8, 9, 34, 41, 42
- lgb.convert\_with\_rules, 11
- lgb.cv, 12, 13, 36, 39
- lgb.Dataset, 7, 13, 15, 39, 44
- lgb.Dataset.construct, 17
- lgb.Dataset.create.valid, 18
- lgb.Dataset.save, 19
- lgb.Dataset.set.categorical, 20
- lgb.Dataset.set.reference, 21
- lgb.drop\_serialized, 22, 29, 34
- lgb.dump, 22
- lgb.get.eval.result, 23
- lgb.importance, 25, 31
- lgb.interprete, 26, 32
- lgb.load, 27
- lgb.make\_serializable, 22, 28, 34
- lgb.model.dt.tree, 29
- lgb.plot.importance, 30
- lgb.plot.interpretation, 32
- lgb.restore\_handle, 22, 29, 33
- lgb.save, 34
- lgb.train, 7, 35, 38, 40, 44
- lightgbm, 9, 36, 38, 42, 43
- predict.lgb.Booster, 10, 41
- print.lgb.Booster, 44
- set\_field, 45
- setLGBMThreads, 44
- setLGBMthreads, 7
- setLGBMthreads(setLGBMThreads), 44
- slice, 46
- summary.lgb.Booster, 47