

Package ‘fastml’

July 24, 2025

Type Package

Title Fast Machine Learning Model Training and Evaluation

Version 0.6.2

Description Streamlines the training, evaluation, and comparison of multiple machine learning models with minimal code by providing comprehensive data preprocessing and support for a wide range of algorithms with hyperparameter tuning.
It offers performance metrics and visualization tools to facilitate efficient and effective machine learning workflows.

Encoding UTF-8

License MIT + file LICENSE

URL <https://github.com/selcukorkmaz/fastml>

BugReports <https://github.com/selcukorkmaz/fastml/issues>

LazyData true

Imports recipes, dplyr, ggplot2, reshape2, rsample, parsnip, tune, workflows, yardstick, tibble, rlang, dials, RColorBrewer, baguette, bonsai, discrim, doFuture, finetune, future, plsmod, probably, viridisLite, DALEX, magrittr, patchwork, pROC, janitor, stringr, DT, GGally, UpSetR, VIM, broom, dbscan, ggpubr, gridExtra, htmlwidgets, kableExtra, moments, naniar, plotly, scales, skimr, tidyr, knitr, rmarkdown, purrr, mice, missForest

Suggests testthat (>= 3.0.0), C50, glmnet, xgboost, ranger, crayon, kernlab, klaR, kkn, keras, lightgbm, rstanarm, mixOmics,

RoxygenNote 7.3.2

Config/testthat/edition 3

NeedsCompilation no

Author Selcuk Korkmaz [aut, cre] (ORCID:

<<https://orcid.org/0000-0003-4632-6850>>),

Dincer Goksuluk [aut] (ORCID: <<https://orcid.org/0000-0002-2752-7668>>),

Eda Karaismailoglu [aut] (ORCID:

<<https://orcid.org/0000-0003-3085-7809>>)

Maintainer Selcuk Korkmaz <selcukorkmaz@gmail.com>
Depends R (>= 3.5.0)
Repository CRAN
Date/Publication 2025-07-24 14:30:08 UTC

Contents

availableMethods	2
evaluate_models	3
fastexplain	4
fastexplore	5
fastml	8
flatten_and_rename_models	12
framingham	13
get_best_model_idx	13
get_best_model_names	14
get_best_workflows	15
get_default_engine	15
get_default_params	16
get_default_tune_params	17
get_engine_names	18
get_model_engine_names	18
load_model	19
plot.fastml	20
predict.fastml	21
process_model	22
sanitize	24
save.fastml	24
summary.fastml	25
train_models	26
Index	28

availableMethods	<i>Get Available Methods</i>
------------------	------------------------------

Description

Returns a character vector of algorithm names available for either classification or regression tasks.

Usage

```
availableMethods(type = c("classification", "regression"), ...)
```

Arguments

type	A character string specifying the type of task. Must be either "classification" or "regression". Defaults to c("classification", "regression") and uses match.arg to select one.
...	Additional arguments (currently not used).

Details

Depending on the specified type, the function returns a different set of algorithm names:

- For "classification", it returns algorithms such as "logistic_reg", "multinom_reg", "decision_tree", "C5_rules", "rand_forest", "xgboost", "lightgbm", "svm_linear", "svm_rbf", "nearest_neighbor", "naive_Bayes", "mlp", "discrim_linear", "discrim_quad", and "bag_tree".
- For "regression", it returns algorithms such as "linear_reg", "ridge_regression", "lasso_regression", "elastic_net", "decision_tree", "rand_forest", "xgboost", "lightgbm", "svm_linear", "svm_rbf", "nearest_neighbor", "mlp", "pls", and "bayes_glm".

Value

A character vector containing the names of the available algorithms for the specified task type.

evaluate_models	<i>Evaluate Models Function</i>
-----------------	---------------------------------

Description

Evaluates the trained models on the test data and computes performance metrics.

Usage

```
evaluate_models(  
  models,  
  train_data,  
  test_data,  
  label,  
  task,  
  metric = NULL,  
  event_class  
)
```

Arguments

<code>models</code>	A list of trained model objects.
<code>train_data</code>	Preprocessed training data frame.
<code>test_data</code>	Preprocessed test data frame.
<code>label</code>	Name of the target variable.
<code>task</code>	Type of task: "classification" or "regression".
<code>metric</code>	The performance metric to optimize (e.g., "accuracy", "rmse").
<code>event_class</code>	A single string. Either "first" or "second" to specify which level of truth to consider as the "event".

Value

A list with two elements:

performance A named list of performance metric tibbles for each model.

predictions A named list of data frames with columns including truth, predictions, and probabilities per model.

fastexplain

FastExplain the fastml (DALEX + SHAP + Permutation-based VI)

Description

Provides model explainability using DALEX. This function:

- Creates a DALEX explainer.
- Computes permutation-based variable importance with boxplots showing variability, displays the table and plot.
- Computes partial dependence-like model profiles if 'features' are provided.
- Computes Shapley values (SHAP) for a sample of the training observations, displays the SHAP table, and plots a summary bar chart of mean(|SHAP value|) per feature. For classification, it shows separate bars for each class.

Usage

```
fastexplain(
  object,
  method = "dalex",
  features = NULL,
  grid_size = 20,
  shap_sample = 5,
  vi_iterations = 10,
  seed = 123,
  loss_function = NULL,
  ...
)
```

Arguments

object	A fastml object.
method	Currently only "dalex" is supported.
features	Character vector of feature names for partial dependence (model profiles). Default NULL.
grid_size	Number of grid points for partial dependence. Default 20.
shap_sample	Integer number of observations from processed training data to compute SHAP values for. Default 5.
vi_iterations	Integer. Number of permutations for variable importance (B). Default 10.
seed	Integer. A value specifying the random seed.
loss_function	Function. The loss function for model_parts. <ul style="list-style-type: none"> • If NULL and task = 'classification', defaults to DALEX: :loss_cross_entropy. • If NULL and task = 'regression', defaults to DALEX: :loss_root_mean_square.
...	Additional arguments (not currently used).

Details**1. Custom number of permutations for VI (vi_iterations):**

You can now specify how many permutations (B) to use for permutation-based variable importance. More permutations yield more stable estimates but take longer.

2. Better error messages and checks:

Improved checks and messages if certain packages or conditions are not met.

3. Loss Function:

A loss_function argument has been added to let you pick a different performance measure (e.g., loss_cross_entropy for classification, loss_root_mean_square for regression).

4. Parallelization Suggestion:**Value**

Prints DALEX explanations: variable importance table & plot, model profiles (if any), and SHAP table & summary plot.

fastexplore

Explore and Summarize a Dataset Quickly

Description

fastexplore provides a fast and comprehensive exploratory data analysis (EDA) workflow. It automatically detects variable types, checks for missing and duplicated data, suggests potential ID columns, and provides a variety of plots (histograms, boxplots, scatterplots, correlation heatmaps, etc.). It also includes optional outlier detection, normality testing, and feature engineering.

Usage

```

fastexplore(
  data,
  label = NULL,
  visualize = c("histogram", "boxplot", "barplot", "heatmap", "scatterplot"),
  save_results = TRUE,
  output_dir = NULL,
  sample_size = NULL,
  interactive = FALSE,
  corr_threshold = 0.9,
  auto_convert_numeric = TRUE,
  visualize_missing = TRUE,
  imputation_suggestions = FALSE,
  report_duplicate_details = TRUE,
  detect_near_duplicates = TRUE,
  auto_convert_dates = FALSE,
  feature_engineering = FALSE,
  outlier_method = c("iqr", "zscore", "dbscan", "lof"),
  run_distribution_checks = TRUE,
  normality_tests = c("shapiro"),
  pairwise_matrix = TRUE,
  max_scatter_cols = 5,
  grouped_plots = TRUE,
  use_upset_missing = TRUE
)

```

Arguments

<code>data</code>	A <code>data.frame</code> . The dataset to analyze.
<code>label</code>	A character string specifying the name of the target or label column (optional). If provided, certain grouped plots and class imbalance checks will be performed.
<code>visualize</code>	A character vector specifying which visualizations to produce. Possible values: <code>c("histogram", "boxplot", "barplot", "heatmap", "scatterplot")</code> .
<code>save_results</code>	Logical. If <code>TRUE</code> , saves plots and a rendered report (HTML) into a timestamped <code>EDA_Results_</code> folder inside <code>output_dir</code> .
<code>output_dir</code>	A character string specifying the output directory for saving results (if <code>save_results = TRUE</code>). Defaults to current working directory.
<code>sample_size</code>	An integer specifying a random sample size for the data to be used in visualizations. If <code>NULL</code> , uses the entire dataset.
<code>interactive</code>	Logical. If <code>TRUE</code> , attempts to produce interactive Plotly heatmaps and other interactive elements. If required packages are not installed, falls back to static plots.
<code>corr_threshold</code>	Numeric. Threshold above which correlations (in absolute value) are flagged as high. Defaults to 0.9.
<code>auto_convert_numeric</code>	Logical. If <code>TRUE</code> , automatically converts factor/character columns that look numeric (only digits, minus sign, or decimal point) to numeric.

<code>visualize_missing</code>	Logical. If TRUE, attempts to visualize missingness patterns (e.g., via an UpSet plot, if UpSetR is available, or VIM , naniar).
<code>imputation_suggestions</code>	Logical. If TRUE, prints simple text suggestions for imputation strategies.
<code>report_duplicate_details</code>	Logical. If TRUE, shows top duplicated rows and their frequency.
<code>detect_near_duplicates</code>	Logical. Placeholder for near-duplicate (fuzzy) detection. Currently not implemented.
<code>auto_convert_dates</code>	Logical. If TRUE, attempts to detect and convert date-like strings (YYYY-MM-DD) to Date format.
<code>feature_engineering</code>	Logical. If TRUE, automatically engineers derived features (day, month, year) from any date/time columns, and identifies potential ID columns.
<code>outlier_method</code>	A character string indicating which outlier detection method(s) to apply. One of <code>c("iqr", "zscore", "dbscan", "lof")</code> . Only the first match will be used in the code (though the function is designed to handle multiple).
<code>run_distribution_checks</code>	Logical. If TRUE, runs normality tests (e.g., Shapiro-Wilk) on numeric columns.
<code>normality_tests</code>	A character vector specifying which normality tests to run. Possible values include "shapiro" or "ks" (Kolmogorov-Smirnov). Only used if <code>run_distribution_checks = TRUE</code> .
<code>pairwise_matrix</code>	Logical. If TRUE, produces a scatterplot matrix (using GGally) for numeric columns.
<code>max_scatter_cols</code>	Integer. Maximum number of numeric columns to include in the pairwise matrix.
<code>grouped_plots</code>	Logical. If TRUE, produce grouped histograms, violin plots, and density plots by label (if the label is a factor).
<code>use_upset_missing</code>	Logical. If TRUE, attempts to produce an UpSet plot for missing data if UpSetR is available.

Details

This function automates many steps of EDA:

1. Automatically detects numeric vs. categorical variables.
2. Auto-converts columns that look numeric (and optionally date-like).
3. Summarizes data structure, missingness, duplication, and potential ID columns.
4. Computes correlation matrix and flags highly correlated pairs.
5. (Optional) Outlier detection using IQR, Z-score, DBSCAN, or LOF methods.
6. (Optional) Normality tests on numeric columns.
7. Saves all results and an R Markdown report if `save_results = TRUE`.

Value

A (silent) list containing:

- `data_overview` - A basic overview (head, unique values, skim summary).
- `summary_stats` - Summary statistics for numeric columns.
- `freq_tables` - Frequency tables for factor columns.
- `missing_data` - Missing data overview (count, percentage).
- `duplicated_rows` - Count of duplicated rows.
- `class_imbalance` - Class distribution if label is provided and is categorical.
- `correlation_matrix` - The correlation matrix for numeric variables.
- `zero_variance_cols` - Columns with near-zero variance.
- `potential_id_cols` - Columns with unique values in every row.
- `date_time_cols` - Columns recognized as date/time.
- `high_corr_pairs` - Pairs of variables with correlation above `corr_threshold`.
- `outlier_method` - The chosen method for outlier detection.
- `outlier_summary` - Outlier proportions or metrics (if computed).

If `save_results = TRUE`, additional side effects include saving figures, a correlation heatmap, and an R Markdown report in the specified directory.

fastml

Fast Machine Learning Function

Description

Trains and evaluates multiple classification or regression models automatically detecting the task based on the target variable type.

Usage

```
fastml(
  data = NULL,
  train_data = NULL,
  test_data = NULL,
  label,
  algorithms = "all",
  task = "auto",
  test_size = 0.2,
  resampling_method = "cv",
  folds = ifelse(grepl("cv", resampling_method), 10, 25),
  repeats = ifelse(resampling_method == "repeatedcv", 1, NA),
  event_class = "first",
  exclude = NULL,
```



```

recipe = NULL,
tune_params = NULL,
metric = NULL,
algorithm_engines = NULL,
n_cores = 1,
stratify = TRUE,
impute_method = "error",
impute_custom_function = NULL,
encode_categoricals = TRUE,
scaling_methods = c("center", "scale"),
balance_method = c("none", "upsample", "downsample"),
resamples = NULL,
summaryFunction = NULL,
use_default_tuning = FALSE,
tuning_strategy = "grid",
tuning_iterations = 10,
early_stopping = FALSE,
adaptive = FALSE,
learning_curve = FALSE,
seed = 123,
verbose = FALSE
)

```

Arguments

data	A data frame containing the complete dataset. If both 'train_data' and 'test_data' are 'NULL', 'fastml()' will split this into training and testing sets according to 'test_size' and 'stratify'. Defaults to 'NULL'.
train_data	A data frame pre-split for model training. If provided, 'test_data' must also be supplied, and no internal splitting will occur. Defaults to 'NULL'.
test_data	A data frame pre-split for model evaluation. If provided, 'train_data' must also be supplied, and no internal splitting will occur. Defaults to 'NULL'.
label	A string specifying the name of the target variable.
algorithms	A vector of algorithm names to use. Default is "all" to run all supported algorithms.
task	Character string specifying model type selection. Use "auto" to let the function detect whether the target is for classification or regression based on the data, or explicitly set to "classification" or "regression".
test_size	A numeric value between 0 and 1 indicating the proportion of the data to use for testing. Default is 0.2.
resampling_method	A string specifying the resampling method for model evaluation. Default is "cv" (cross-validation). Other options include "none", "boot", "repeatedcv", etc.
folds	An integer specifying the number of folds for cross-validation. Default is 10 for methods containing "cv" and 25 otherwise.
repeats	Number of times to repeat cross-validation (only applicable for methods like "repeatedcv").

<code>event_class</code>	A single string. Either "first" or "second" to specify which level of truth to consider as the "event". Default is "first".
<code>exclude</code>	A character vector specifying the names of the columns to be excluded from the training process.
<code>recipe</code>	A user-defined recipe object for custom preprocessing. If provided, internal recipe steps (imputation, encoding, scaling) are skipped.
<code>tune_params</code>	A named list of tuning ranges for each algorithm and engine pair. Example: <code>list(rand_forest = list(ranger = list(mtry = c(1, 3))))</code> will override the defaults for the ranger engine. Default is NULL.
<code>metric</code>	The performance metric to optimize during training.
<code>algorithm_engines</code>	A named list specifying the engine to use for each algorithm.
<code>n_cores</code>	An integer specifying the number of CPU cores to use for parallel processing. Default is 1.
<code>stratify</code>	Logical indicating whether to use stratified sampling when splitting the data. Default is TRUE for classification and FALSE for regression.
<code>impute_method</code>	Method for handling missing values. Options include: "medianImpute" Impute missing values using median imputation (recipe-based). "knnImpute" Impute missing values using k-nearest neighbors (recipe-based). "bagImpute" Impute missing values using bagging (recipe-based). "remove" Remove rows with missing values from the data (recipe-based). "mice" Impute missing values using MICE (Multiple Imputation by Chained Equations). "missForest" Impute missing values using the missForest algorithm. "custom" Use a user-provided imputation function (see 'impute_custom_function'). "error" Do not perform imputation; if missing values are detected, stop execution with an error. NULL Equivalent to "error". No imputation is performed, and the function will stop if missing values are present. Default is "error".
<code>impute_custom_function</code>	A function that takes a data.frame as input and returns an imputed data.frame. Used only if <code>impute_method = "custom"</code> .
<code>encode_categoricals</code>	Logical indicating whether to encode categorical variables. Default is TRUE.
<code>scaling_methods</code>	Vector of scaling methods to apply. Default is <code>c("center", "scale")</code> .
<code>balance_method</code>	Method to handle class imbalance. One of "none", "upsample", or "downsample". Applied to the training set for classification tasks. Default is "none".
<code>resamples</code>	Optional rsample object providing custom resampling splits. If supplied, <code>resampling_method</code> , <code>folds</code> , and <code>repeats</code> are ignored.
<code>summaryFunction</code>	A custom summary function for model evaluation. Default is NULL.

<code>use_default_tuning</code>	Logical; if TRUE and <code>tune_params</code> is NULL, tuning is performed using default grids. Tuning also occurs when custom <code>tune_params</code> are supplied. When FALSE and no custom parameters are given, models are fitted once with default settings. Default is FALSE.
<code>tuning_strategy</code>	A string specifying the tuning strategy. Must be one of "grid", "bayes", or "none". Default is "grid". If custom <code>tune_params</code> are provided while <code>tuning_strategy = "none"</code> , they will be ignored with a warning.
<code>tuning_iterations</code>	Number of iterations for Bayesian tuning. Ignored when <code>tuning_strategy</code> is not "bayes". Validation of this argument only occurs for the Bayesian strategy. Default is 10.
<code>early_stopping</code>	Logical indicating whether to use early stopping in Bayesian tuning methods (if supported). Default is FALSE.
<code>adaptive</code>	Logical indicating whether to use adaptive/racing methods for tuning. Default is FALSE.
<code>learning_curve</code>	Logical. If TRUE, generate learning curves (performance vs. training size).
<code>seed</code>	An integer value specifying the random seed for reproducibility.
<code>verbose</code>	Logical; if TRUE, prints progress messages during the training and evaluation process.

Details

Fast Machine Learning Function

Trains and evaluates multiple classification or regression models. The function automatically detects the task based on the target variable type and can perform advanced hyperparameter tuning using various tuning strategies.

Value

An object of class `fastml` containing the best model, performance metrics, and other information.

Examples

```
# Example 1: Using the iris dataset for binary classification (excluding 'setosa')
data(iris)
iris <- iris[iris$Species != "setosa", ] # Binary classification
iris$Species <- factor(iris$Species)

# Define a custom tuning grid for the ranger engine
tune <- list(
  rand_forest = list(
    ranger = list(mtry = c(1, 3))
  )
)

# Train models with custom tuning
```

```
model <- fastml(  
  data = iris,  
  label = "Species",  
  algorithms = "rand_forest",  
  tune_params = tune,  
  use_default_tuning = TRUE  
)  
  
# View model summary  
summary(model)
```

`flatten_and_rename_models`*Flatten and Rename Models*

Description

Flattens a nested list of models and renames the elements by combining the outer and inner list names.

Usage

```
flatten_and_rename_models(models)
```

Arguments

<code>models</code>	A nested list of models. The outer list should have names. If an inner element is a named list, the names will be combined with the outer name in the format "outer_name (inner_name)".
---------------------	---

Details

The function iterates over each element of the outer list. For each element, if it is a list with names, the function concatenates the outer list name and the inner names using `paste0` and `setNames`. If an element is not a list or does not have names, it is included in the result without modification.

Value

A flattened list with each element renamed according to its original outer and inner list names.

framingham

*Framingham Heart Study Data***Description**

This dataset is derived from the Framingham Heart Study and contains various clinical and demographic variables used to predict coronary heart disease risk over a ten-year period.

Format

male Integer indicator for male sex.
age Participant age in years.
education Education level.
currentSmoker Whether the participant currently smokes.
cigsPerDay Number of cigarettes smoked per day.
BPMeds Whether blood pressure medication is used.
prevalentStroke History of stroke at baseline.
prevalentHyp History of hypertension at baseline.
diabetes Diabetes diagnosis.
totChol Total cholesterol.
sysBP Systolic blood pressure.
diaBP Diastolic blood pressure.
BMI Body mass index.
heartRate Heart rate.
glucose Glucose level.
TenYearCHD Ten year risk of coronary heart disease.

get_best_model_idx

*Get Best Model Indices by Metric and Group***Description**

Identifies and returns the indices of rows in a data frame where the specified metric reaches the overall maximum within groups defined by one or more columns.

Usage

```
get_best_model_idx(df, metric, group_cols = c("Model", "Engine"))
```

Arguments

df	A data frame containing model performance metrics and grouping columns.
metric	A character string specifying the name of the metric column in df. The metric values are converted to numeric for comparison.
group_cols	A character vector of column names used for grouping. Defaults to c("Model", "Engine").

Details

The function converts the metric values to numeric and creates a combined grouping factor using the specified group_cols. It then computes the maximum metric value within each group and determines the overall best metric value across the entire data frame. Finally, it returns the indices of rows belonging to groups that achieve this overall maximum.

Value

A numeric vector of row indices in df corresponding to groups whose maximum metric equals the overall best metric value.

get_best_model_names *Get Best Model Names*

Description

Extracts and returns the best engine names from a named list of model workflows.

Usage

```
get_best_model_names(models)
```

Arguments

models	A named list where each element corresponds to an algorithm and contains a list of model workflows. Each workflow should be compatible with <code>tune::extract_fit_parsnip</code> .
--------	--

Details

For each algorithm, the function extracts the engine names from the model workflows using `tune::extract_fit_parsnip`. It then chooses "randomForest" if it is available; otherwise, it selects the first non-NA engine. If no engine names can be extracted for an algorithm, `NA_character_` is returned.

Value

A named character vector. The names of the vector correspond to the algorithm names, and the values represent the chosen best engine name for that algorithm.

get_best_workflows	<i>Get Best Workflows</i>
--------------------	---------------------------

Description

Extracts the best workflows from a nested list of model workflows based on the provided best model names.

Usage

```
get_best_workflows(models, best_model_name)
```

Arguments

models	A nested list of model workflows. Each element should correspond to an algorithm and contain sublists keyed by engine names.
best_model_name	A named character vector where the names represent algorithm names and the values represent the chosen best engine for each algorithm.

Details

The function iterates over each element in `best_model_name` and attempts to extract the corresponding workflow from `models` using the specified engine. If the workflow for an algorithm-engine pair is not found, a warning is issued and NULL is returned for that entry.

Value

A named list of workflows corresponding to the best engine for each algorithm. Each list element is named in the format "algorithm (engine)".

get_default_engine	<i>Get Default Engine</i>
--------------------	---------------------------

Description

Returns the default engine corresponding to the specified algorithm.

Usage

```
get_default_engine(algo)
```

Arguments

algo	A character string specifying the name of the algorithm. The value should match one of the supported algorithm names.
------	---

Details

The function uses a switch statement to select the default engine based on the given algorithm. If the provided algorithm does not have a defined default engine, the function terminates with an error.

Value

A character string containing the default engine name associated with the provided algorithm.

get_default_params	<i>Get Default Parameters for an Algorithm</i>
--------------------	--

Description

Returns a list of default tuning parameters for the specified algorithm based on the task type, number of predictors, and engine.

Usage

```
get_default_params(algo, task, num_predictors = NULL, engine = NULL)
```

Arguments

algo	A character string specifying the algorithm name. Supported values include: "rand_forest", "C5_rules", "xgboost", "lightgbm", "logistic_reg", "multinom_reg", "decision_tree", "svm_linear", "svm_rbf", "nearest_neighbor", "naive_Bayes", "mlp", "deep_learning", "discrim_linear", "discrim_quad", "bag_tree", "elastic_net", "bayes_glm", "pls", "linear_reg", "ridge_regression", and "lasso_regression".
task	A character string specifying the task type, typically "classification" or "regression".
num_predictors	An optional numeric value indicating the number of predictors. This value is used to compute default values for parameters such as mtry. Defaults to NULL.
engine	An optional character string specifying the engine to use. If not provided, a default engine is chosen where applicable.

Details

The function employs a switch statement to select and return a list of default parameters tailored for the given algorithm, task, and engine. The defaults vary by algorithm and, in some cases, by engine. For example:

- For "rand_forest", if engine is not provided, it defaults to "ranger". The parameters such as mtry, trees, and min_n are computed based on the task and the number of predictors.
- For "C5_rules", the defaults include trees, min_n, and sample_size.
- For "xgboost" and "lightgbm", default values are provided for parameters like tree depth, learning rate, and sample size.

- For "logistic_reg" and "multinom_reg", the function returns defaults for regularization parameters (penalty and mixture) that vary with the specified engine.
- For "decision_tree", the parameters (such as tree_depth, min_n, and cost_complexity) are set based on the engine (e.g., "rpart", "C5.0", "partykit", "spark").
- Other algorithms, including "svm_linear", "svm_rbf", "nearest_neighbor", "naive_Bayes", "mlp", "deep_learning", "elastic_net", "bayes_glm", "pls", "linear_reg", "ridge_regression", and "lasso_regression", have their respective default parameter lists.

Value

A list of default parameter settings for the specified algorithm. If the algorithm is not recognized, the function returns NULL.

```
get_default_tune_params
```

Get Default Tuning Parameters

Description

Returns a list of default tuning parameter ranges for a specified algorithm based on the provided training data, outcome label, and engine.

Usage

```
get_default_tune_params(algo, train_data, label, engine)
```

Arguments

algo	A character string specifying the algorithm name. Supported values include: "rand_forest", "C5_rules", "xgboost", "lightgbm", "logistic_reg", "multinom_reg", "decision_tree", "svm_linear", "svm_rbf", "nearest_neighbor", "naive_Bayes", "mlp", "deep_learning", "discrim_linear", "discrim_quad", "bag_tree", "elastic_net", "bayes_glm", "pls", "linear_reg", "ridge_regression", and "lasso_regression".
train_data	A data frame containing the training data.
label	A character string specifying the name of the outcome variable in train_data. This column is excluded when calculating the number of predictors.
engine	A character string specifying the engine to be used for the algorithm. Different engines may have different tuning parameter ranges.

Details

The function first determines the number of predictors by removing the outcome variable (specified by label) from train_data. It then uses a switch statement to select a list of default tuning parameter ranges tailored for the specified algorithm and engine. The tuning ranges have been adjusted for efficiency and may include parameters such as mtry, trees, min_n, and others depending on the algorithm.

Value

A list of tuning parameter ranges for the specified algorithm. If no tuning parameters are defined for the given algorithm, the function returns NULL.

get_engine_names	<i>Get Engine Names from Model Workflows</i>
------------------	--

Description

Extracts and returns a list of unique engine names from a list of model workflows.

Usage

```
get_engine_names(models)
```

Arguments

models	A list where each element is a list of model workflows. Each workflow is expected to contain a fitted model that can be processed with <code>tune::extract_fit_parsnip</code> .
--------	---

Details

The function applies `tune::extract_fit_parsnip` to each model workflow to extract the fitted model object. It then retrieves the engine name from the model specification (`spec$engine`). If the extraction fails, `NA_character_` is returned for that workflow. Finally, the function removes any duplicate engine names using `unique`.

Value

A list of character vectors. Each vector contains the unique engine names extracted from the corresponding element of `models`.

get_model_engine_names	<i>Get Model Engine Names</i>
------------------------	-------------------------------

Description

Extracts and returns a named vector mapping algorithm names to engine names from a nested list of model workflows.

Usage

```
get_model_engine_names(models)
```

Arguments

`models` A nested list of model workflows. Each inner list should contain model objects from which a fitted model can be extracted using `tune::extract_fit_parsnip`.

Details

The function iterates over a nested list of model workflows and, for each workflow, attempts to extract the fitted model object using `tune::extract_fit_parsnip`. If successful, it retrieves the algorithm name from the first element of the class attribute of the model specification and the engine name from the specification. The results are combined into a named vector.

Value

A named character vector where the names correspond to algorithm names (e.g., "rand_forest", "logistic_reg") and the values correspond to the associated engine names (e.g., "ranger", "glm").

`load_model`*Load Model Function*

Description

Loads a trained model object from a file.

Usage

```
load_model(filepath)
```

Arguments

`filepath` A string specifying the file path to load the model from.

Value

An object of class `fastml`.

plot.fastml

*Plot Methods for fastml Objects***Description**

plot.fastml produces visual diagnostics for a trained fastml object.

Usage

```
## S3 method for class 'fastml'
plot(
  x,
  algorithm = "best",
  type = c("all", "bar", "roc", "calibration", "residual"),
  ...
)
```

Arguments

x	A fastml object (output of <code>fastml()</code>).
algorithm	Character vector specifying which algorithm(s) to include when generating certain plots (e.g., ROC curves). Defaults to "best".
type	Character vector indicating which plot(s) to produce. Options are: "bar" Bar plot of performance metrics across all models/engines. "roc" ROC curve(s) for binary classification models. "calibration" Calibration plot for the best model(s). "residual" Residual diagnostics for the best model. "all" Produce all available plots.
...	Additional arguments (currently unused).

Details

When type = "all", plot.fastml will produce a bar plot of metrics, ROC curves (classification), calibration plot, and residual diagnostics (regression). If you specify a subset of types, only those will be drawn.

Examples

```
## Create a binary classification dataset from iris
data(iris)
iris <- iris[iris$Species != "setosa",]
iris$Species <- factor(iris$Species)

## Fit fastml model on binary classification task
model <- fastml(data = iris, label = "Species", algorithms = c("rand_forest", "svm_rbf"))
```

```

## 1. Plot all available diagnostics
plot(model, type = "all")

## 2. Bar plot of performance metrics
plot(model, type = "bar")

## 3. ROC curves (only for classification models)
plot(model, type = "roc")

## 4. Calibration plot (requires 'probably' package)
plot(model, type = "calibration")

## 5. ROC curves for specific algorithm(s) only
plot(model, type = "roc", algorithm = "rand_forest")

## 6. Residual diagnostics (only available for regression tasks)
model <- fastml(data = mtcars, label = "mpg", algorithms = c("linear_reg", "xgboost"))
plot(model, type = "residual")

```

predict.fastml

Predict method for fastml objects

Description

Generates predictions from a trained ‘fastml’ object on new data. Supports both single-model and multi-model workflows, and handles classification and regression tasks with optional post-processing and verbosity.

Usage

```

## S3 method for class 'fastml'
predict(
  object,
  newdata,
  type = "auto",
  model_name = NULL,
  verbose = FALSE,
  postprocess_fn = NULL,
  ...
)

```

Arguments

object	A fitted ‘fastml’ object created by the ‘fastml()’ function.
newdata	A data frame or tibble containing new predictor data for which to generate predictions.

type	Type of prediction to return. One of "auto" (default), "class", "prob", or "numeric". - "auto": chooses "class" for classification and "numeric" for regression. - "prob": returns class probabilities (only for classification). - "class": returns predicted class labels. - "numeric": returns predicted numeric values (for regression).
model_name	(Optional) Name of a specific model to use when 'object\$best_model' contains multiple models.
verbose	Logical; if 'TRUE', prints progress messages showing which models are used during prediction.
postprocess_fn	(Optional) A function to apply to the final predictions (e.g., inverse transforms, thresholding).
...	Additional arguments (currently unused).

Value

A vector of predictions, or a named list of predictions (if multiple models are used). If 'postprocess_fn' is supplied, its output will be returned instead.

Examples

```
## Not run:
set.seed(123)
model <- fastml(iris, label = "Species")
test_data <- iris[sample(1:150, 20),-5]

## Best model(s) predictions
preds <- predict(model, newdata = test_data)

## Predicted class probabilities using best model(s)
probs <- predict(model, newdata = test_data, type = "prob")

## Prediction from a specific model by name
single_model_preds <- predict(model, newdata = test_data, model_name = "rand_forest (ranger)")

## End(Not run)
```

Description

This function processes a fitted model or a tuning result, finalizes the model if tuning was used, makes predictions on the test set, and computes performance metrics depending on the task type (classification or regression). It supports binary and multiclass classification, and handles probabilistic outputs when supported by the modeling engine.

Usage

```
process_model(
  model_obj,
  model_id,
  task,
  test_data,
  label,
  event_class,
  engine,
  train_data,
  metric
)
```

Arguments

<code>model_obj</code>	A fitted model or a tuning result ('tune_results' object).
<code>model_id</code>	A character identifier for the model (used in warnings).
<code>task</code>	Type of task, either "classification" or "regression".
<code>test_data</code>	A data frame containing the test data.
<code>label</code>	The name of the outcome variable (as a character string).
<code>event_class</code>	For binary classification, specifies which class is considered the positive class: "first" or "second".
<code>engine</code>	A character string indicating the model engine (e.g., "xgboost", "randomForest"). Used to determine if class probabilities are supported. If 'NULL', probabilities are skipped.
<code>train_data</code>	A data frame containing the training data, required to refit finalized workflows.
<code>metric</code>	The name of the metric (e.g., "roc_auc", "accuracy", "rmse") used for selecting the best tuning result.

Details

- If the input 'model_obj' is a 'tune_results' object, the function finalizes the model using the best hyperparameters according to the specified 'metric', and refits the model on the full training data.
- For classification tasks, performance metrics include accuracy, kappa, sensitivity, specificity, precision, F1-score, and ROC AUC (if probabilities are available).
- For regression tasks, RMSE, R-squared, and MAE are returned.
- For models with missing prediction lengths, a helpful imputation error is thrown to guide data preprocessing.

Value

A list with two elements:

performance A tibble with computed performance metrics.

predictions A tibble with predicted values and corresponding truth values, and probabilities (if applicable).

sanitize	<i>Clean Column Names or Character Vectors by Removing Special Characters</i>
----------	---

Description

This function can operate on either a data frame or a character vector:

- **Data frame:** Detects columns whose names contain any character that is not a letter, number, or underscore, removes colons, replaces slashes with underscores, and spaces with underscores.
- **Character vector:** Applies the same cleaning rules to every element of the vector.

Usage

```
sanitize(x)
```

Arguments

x	A data frame or character vector to be cleaned.
---	---

Value

- If x is a data frame: returns a data frame with cleaned column names.
- If x is a character vector: returns a character vector with cleaned elements.

save.fastml	<i>Save Model Function</i>
-------------	----------------------------

Description

Saves the trained model object to a file.

Usage

```
save.fastml(model, filepath)
```

Arguments

model	An object of class fastml.
filepath	A string specifying the file path to save the model.

Value

No return value, called for its side effect of saving the model object to a file.

summary.fastml	<i>Summary Function for fastml (Using yardstick for ROC Curves)</i>
----------------	---

Description

Summarizes the results of machine learning models trained using the ‘fastml’ package. Depending on the task type (classification or regression), it provides customized output such as performance metrics, best hyperparameter settings, and confusion matrices. It is designed to be informative and readable, helping users quickly interpret model results.

Usage

```
## S3 method for class 'fastml'
summary(
  object,
  algorithm = "best",
  type = c("all", "metrics", "params", "conf_mat"),
  sort_metric = NULL,
  ...
)
```

Arguments

object	An object of class fastml.
algorithm	A vector of algorithm names to display summary. Default is "best".
type	Character vector indicating which outputs to produce. Options are "all" (all available outputs), "metrics" (performance metrics), "params" (best hyperparameters), and "conf_mat" (confusion matrix). Default is "all".
sort_metric	The metric to sort by. Default uses optimized metric.
...	Additional arguments.

Details

For classification tasks, the summary includes metrics such as Accuracy, F1 Score, Kappa, Precision, ROC AUC, Sensitivity, and Specificity. A confusion matrix is also provided for the best model(s). For regression tasks, the summary reports RMSE, R-squared, and MAE.

Users can control the type of output with the ‘type’ argument: ‘metrics’ displays model performance metrics. ‘params’ shows the best hyperparameter settings. ‘conf_mat’ prints confusion matrices (only for classification). ‘all’ includes all of the above.

If multiple algorithms are trained, the summary highlights the best model based on the optimized metric.

Value

Prints summary of fastml models.

train_models

*Train Specified Machine Learning Algorithms on the Training Data***Description**

Trains specified machine learning algorithms on the preprocessed training data.

Usage

```
train_models(
  train_data,
  label,
  task,
  algorithms,
  resampling_method,
  folds,
  repeats,
  resamples = NULL,
  tune_params,
  metric,
  summaryFunction = NULL,
  seed = 123,
  recipe,
  use_default_tuning = FALSE,
  tuning_strategy = "grid",
  tuning_iterations = 10,
  early_stopping = FALSE,
  adaptive = FALSE,
  algorithm_engines = NULL
)
```

Arguments

train_data	Preprocessed training data frame.
label	Name of the target variable.
task	Type of task: "classification" or "regression".
algorithms	Vector of algorithm names to train.
resampling_method	Resampling method for cross-validation (e.g., "cv", "repeatedcv", "boot", "none").
folds	Number of folds for cross-validation.
repeats	Number of times to repeat cross-validation (only applicable for methods like "repeatedcv").
resamples	Optional rsample object. If provided, custom resampling splits will be used instead of those created internally.

tune_params	A named list of tuning ranges. For each algorithm, supply a list of engine-specific parameter values, e.g. <code>list(rand_forest = list(ranger = list(mtry = c(1, 3))))</code> .
metric	The performance metric to optimize.
summaryFunction	A custom summary function for model evaluation. Default is NULL.
seed	An integer value specifying the random seed for reproducibility.
recipe	A recipe object for preprocessing.
use_default_tuning	Logical; if TRUE and tune_params is NULL, tuning is performed using default grids. Tuning also occurs when custom tune_params are supplied. When FALSE and no custom parameters are given, the model is fitted once with default settings.
tuning_strategy	A string specifying the tuning strategy. Must be one of "grid", "bayes", or "none". Adaptive methods may be used with "grid". If "none" is selected, the workflow is fitted directly without tuning. If custom tune_params are supplied with tuning_strategy = "none", they will be ignored with a warning.
tuning_iterations	Number of iterations for Bayesian tuning. Ignored when tuning_strategy is not "bayes"; validation occurs only for the Bayesian strategy.
early_stopping	Logical for early stopping in Bayesian tuning.
adaptive	Logical indicating whether to use adaptive/racing methods.
algorithm_engines	A named list specifying the engine to use for each algorithm.

Value

A list of trained model objects.

Index

`availableMethods`, [2](#)

`evaluate_models`, [3](#)

`fastexplain`, [4](#)

`fastexplore`, [5](#)

`fastml`, [8](#), [20](#)

`flatten_and_rename_models`, [12](#)

`framingham`, [13](#)

`get_best_model_idx`, [13](#)

`get_best_model_names`, [14](#)

`get_best_workflows`, [15](#)

`get_default_engine`, [15](#)

`get_default_params`, [16](#)

`get_default_tune_params`, [17](#)

`get_engine_names`, [18](#)

`get_model_engine_names`, [18](#)

`load_model`, [19](#)

`match.arg`, [3](#)

`plot.fastml`, [20](#)

`predict.fastml`, [21](#)

`process_model`, [22](#)

`sanitize`, [24](#)

`save.fastml`, [24](#)

`summary.fastml`, [25](#)

`train_models`, [26](#)