

Package ‘sgboost’

May 19, 2024

Title Sparse-Group Boosting

Version 0.1.3

Description

Sparse-group boosting to be used in conjunction with the 'mboost' for modeling grouped data. Applicable to all sparse-group lasso type problems where within-group and between-group sparsity is desired. Interprets and visualizes individual variables and groups.

Imports dplyr, mboost, stringr, rlang, tibble, ggplot2, ggforce

License MIT + file LICENSE

Encoding UTF-8

RoxygenNote 7.3.1

URL <https://github.com/FabianObster/sgboost>

BugReports <https://github.com/FabianObster/sgboost/issues>

Suggests knitr, rmarkdown, spelling, testthat (>= 3.0.0)

Config/testthat/edition 3

VignetteBuilder knitr

Language en-US

NeedsCompilation no

Author Fabian Obster [aut, cre, cph] (<<https://orcid.org/0000-0002-6951-9869>>)

Maintainer Fabian Obster <fabian.obster@unibw.de>

Repository CRAN

Date/Publication 2024-05-19 16:50:02 UTC

R topics documented:

| | |
|--------------------------|---|
| create_formula | 2 |
| get_coef | 3 |
| get_coef_path | 4 |
| get_varimp | 5 |

| | |
|--------------|---|
| plot_effects | 6 |
| plot_path | 8 |
| plot_varimp | 9 |

Index 11

| | |
|----------------|---|
| create_formula | <i>Create a sparse-group boosting formula</i> |
|----------------|---|

Description

Creates a `mboost` formula that allows to fit a sparse-group boosting model based on boosted Ridge Regression with mixing parameter `alpha`. The formula consists of a group baselearner part with degrees of freedom $1-\alpha$ and individual baselearners with degrees of freedom `alpha`. Groups should be defined through `group_df`. The corresponding modeling data should not contain categorical variables with more than two categories, as they are then treated as a group only.

Usage

```
create_formula(
  alpha = 0.3,
  group_df = NULL,
  blearner = "bols",
  outcome_name = "y",
  group_name = "group_name",
  var_name = "var_name",
  intercept = FALSE
)
```

Arguments

| | |
|---------------------------|--|
| <code>alpha</code> | Numeric mixing parameter. For <code>alpha = 0</code> only group baselearners and for <code>alpha = 1</code> only individual baselearners are defined. |
| <code>group_df</code> | input <code>data.frame</code> containing variable names with group structure. |
| <code>blearner</code> | Type of baselearner. Default is 'bols'. |
| <code>outcome_name</code> | String indicating the name of dependent variable. Default is "y" |
| <code>group_name</code> | Name of column in <code>group_df</code> indicating the group structure of the variables. Default is "group_name". |
| <code>var_name</code> | Name of column in <code>group_df</code> containing the variable names to be used as predictors. Default is "var_name". should not contain categorical variables with more than two categories, as they are then treated as a group only. |
| <code>intercept</code> | Logical, should intercept be used? |

Value

Character containing the formula to be passed to `mboost::mboost()` yielding the sparse-group boosting for a given value mixing parameter `alpha`.

Examples

```

library(mboost)
library(dplyr)
set.seed(1)
df <- data.frame(
  x1 = rnorm(100), x2 = rnorm(100), x3 = rnorm(100),
  x4 = rnorm(100), x5 = runif(100)
)
df <- df %>%
  mutate_all(function(x) {
    as.numeric(scale(x))
  })
df$y <- df$x1 + df$x4 + df$x5
group_df <- data.frame(
  group_name = c(1, 1, 1, 2, 2),
  var_name = c("x1", "x2", "x3", "x4", "x5")
)

sgb_formula <- create_formula(alpha = 0.3, group_df = group_df)
sgb_model <- mboost(formula = sgb_formula, data = df)
summary(sgb_model)

```

get_coef

Aggregated and raw coefficients in a sparse group boosting model

Description

Computes the aggregated coefficients from group and individual baselearners. Also returns the raw coefficients associated with each baselearner.

Usage

```
get_coef(sgb_model)
```

Arguments

sgb_model Model of type mboost to compute the coefficients for.

Details

in a sparse group boosting models a variable in a dataset can be selected as an individual variable or as a group. Therefore there can be two associated effect sizes for the same variable. This function aggregates both and returns it in a data.frame.

Value

List of data.frames containing the a data.frame '\$raw' with the variable and the raw (Regression) coefficients and the data.frame '\$aggregated' with the aggregated (Regression) coefficients.

Examples

```

library(mboost)
library(dplyr)
set.seed(1)
df <- data.frame(
  x1 = rnorm(100), x2 = rnorm(100), x3 = rnorm(100),
  x4 = rnorm(100), x5 = runif(100)
)
df <- df %>%
  mutate_all(function(x) {
    as.numeric(scale(x))
  })
df$y <- df$x1 + df$x4 + df$x5
group_df <- data.frame(
  group_name = c(1, 1, 1, 2, 2),
  var_name = c("x1", "x2", "x3", "x4", "x5")
)

sgb_formula <- create_formula(alpha = 0.3, group_df = group_df)
sgb_model <- mboost(formula = sgb_formula, data = df)
sgb_coef <- get_coef(sgb_model)

```

| | |
|---------------|---|
| get_coef_path | <i>Path of aggregated and raw coefficients in a sparse-group boosting model</i> |
|---------------|---|

Description

Computes the aggregated coefficients from group and individual baselearners for each boosting iteration.

Usage

```
get_coef_path(sgb_model)
```

Arguments

sgb_model Model of type mboost to compute the coefficient path for .

Details

in a sparse-group boosting models a variable in a dataset can be selected as an individual variable or as a group. Therefore there can be two associated effect sizes for the same variable. This function aggregates both and returns it in a data.frame for each boosting iteration

Value

List of data.frames containing the a data.frame \$raw with the variable and the raw (Regression) coefficients and the data.frame \$aggregated with the aggregated (Regression) coefficients.

See Also[get_coef\(\)](#)**Examples**

```

library(mboost)
library(dplyr)
set.seed(1)
df <- data.frame(
  x1 = rnorm(100), x2 = rnorm(100), x3 = rnorm(100),
  x4 = rnorm(100), x5 = runif(100)
)
df <- df %>%
  mutate_all(function(x) {
    as.numeric(scale(x))
  })
df$y <- df$x1 + df$x4 + df$x5
group_df <- data.frame(
  group_name = c(1, 1, 1, 2, 2),
  var_name = c("x1", "x2", "x3", "x4", "x5")
)

sgb_formula <- create_formula(alpha = 0.3, group_df = group_df)
sgb_model <- mboost(formula = sgb_formula, data = df)
sgb_coef_path <- get_coef_path(sgb_model)

```

get_varimp

*Variable importance of a sparse-group boosting model***Description**

Variable importance is computed as relative reduction of loss-function attributed to each predictor (groups and individual variables). Returns a list of two data.frames. The first contains the variable importance of a sparse-group model in a data.frame for each predictor. The second one contains the aggregated relative importance of all groups vs. individual variables.

Usage

```
get_varimp(sgb_model)
```

Arguments

`sgb_model` Model of type `mboost` to compute the variable importance for.

Value

List of two data.frames. `$raw` contains the name of the variables, group structure and variable importance on both group and individual variable basis. `$group_importance` contains the aggregated relative importance of all group baselearners and of all individual variables.

See Also

`mboost::varimp()` which this function uses.

Examples

```
library(mboost)
library(dplyr)
set.seed(1)
df <- data.frame(
  x1 = rnorm(100), x2 = rnorm(100), x3 = rnorm(100),
  x4 = rnorm(100), x5 = runif(100)
)
df <- df %>%
  mutate_all(function(x) {
    as.numeric(scale(x))
  })
df$y <- df$x1 + df$x4 + df$x5
group_df <- data.frame(
  group_name = c(1, 1, 1, 2, 2),
  var_name = c("x1", "x2", "x3", "x4", "x5")
)

sgb_formula <- as.formula(create_formula(alpha = 0.3, group_df = group_df))
sgb_model <- mboost(formula = sgb_formula, data = df)
sgb_varimp <- get_varimp(sgb_model)
```

plot_effects

Visualizing a sparse-group boosting model

Description

Radar or scatter/lineplot visualizing the effects sizes relative to the variable importance in a sparse-group boosting model. Works also for a regular mboost model.

Usage

```
plot_effects(
  sgb_model,
  plot_type = "radar",
  prop = 0,
  n_predictors = 30,
  max_char_length = 5,
  base_size = 8
)
```

Arguments

| | |
|-----------------|--|
| sgb_model | Model of type mboost to be used. |
| plot_type | String indicating the type of visualization to use. 'radar' refers to a radar plot using polar coordinates. Here the angle is relative to the cumulative relative importance of predictors and the radius is proportional to the effect size. "clock" does the same as "radar" but uses clock coordinates instead of polar coordinates. "scatter" uses the effect size as y-coordinate and the cumulative relative importance as x-axis in a classical Scatter plot. |
| prop | Numeric value indicating the minimal importance a predictor/baselearner has to have to be plotted. Default value is zero, meaning all predictors are plotted. By increasing prop the number of plotted variables can be reduced. One can also use n_predictors for limiting the number of variables to be plotted directly. |
| n_predictors | The maximum number of predictors to be plotted. Default is 30. Alternative to prop. |
| max_char_length | The maximum character length of a predictor to be printed. Default is 5. For long variable names one may adjust this number. |
| base_size | The base_size argument to be passed to the ggplot2 theme <code>ggplot2::theme_classic</code> to be used to control the overall size of the figure. Default value is 8. |

Value

ggplot2 object mapping the effect sizes and variable importance.

See Also

[get_coef\(\)](#), [get_varimp\(\)](#) which this function uses.

Examples

```
library(mboost)
library(dplyr)
set.seed(1)
df <- data.frame(
  x1 = rnorm(100), x2 = rnorm(100), x3 = rnorm(100),
  x4 = rnorm(100), x5 = runif(100)
)
df <- df %>%
  mutate_all(function(x) {
    as.numeric(scale(x))
  })
df$y <- df$x1 + df$x4 + df$x5
group_df <- data.frame(
  group_name = c(1, 1, 1, 2, 2),
  var_name = c("x1", "x2", "x3", "x4", "x5")
)

sgb_formula <- as.formula(create_formula(alpha = 0.3, group_df = group_df))
sgb_model <- mboost(formula = sgb_formula, data = df)
```

```
plot_effects(sgb_model)
```

| | |
|-----------|--|
| plot_path | <i>Coefficient path of a sparse-group boosting model</i> |
|-----------|--|

Description

Shows how the effect sizes change throughout the boosting iterations in a sparse-group boosting model. Works also for a regular mboost models. Color indicates the selection of group or individual variables within a boosting iteration.

Usage

```
plot_path(sgb_model, max_char_length = 5, base_size = 8)
```

Arguments

| | |
|-----------------|--|
| sgb_model | Model of type mboost to be used. |
| max_char_length | The maximum character length of a predictor to be printed. Default is 5. For long variable names one may adjust this number. |
| base_size | The base_size argument to be passed to the ggplot2 theme ggplot2::theme_bw to be used to control the overall size of the figure. Default value is 8. |

Value

ggplot2 object mapping the effect sizes and variable importance.

See Also

[get_coef_path\(\)](#) which this function uses.

Examples

```
library(mboost)
library(dplyr)
set.seed(1)
df <- data.frame(
  x1 = rnorm(100), x2 = rnorm(100), x3 = rnorm(100),
  x4 = rnorm(100), x5 = runif(100)
)
df <- df %>%
  mutate_all(function(x) {
    as.numeric(scale(x))
  })
df$y <- df$x1 + df$x4 + df$x5
group_df <- data.frame(
  group_name = c(1, 1, 1, 2, 2),
  var_name = c("x1", "x2", "x3", "x4", "x5")
```



```

)

sgb_formula <- as.formula(create_formula(alpha = 0.4, group_df = group_df))
sgb_model <- mboost(formula = sgb_formula, data = df)
plot_path(sgb_model)

```

plot_varimp

Variable importance bar plot of a sparse group boosting model

Description

Visualizes the variable importance of a sparse-group boosting model. Color indicates if a predictor is an individual variable or a group.

Usage

```

plot_varimp(
  sgb_model,
  prop = 0,
  n_predictors = 30,
  max_char_length = 15,
  base_size = 8
)

```

Arguments

| | |
|-----------------|---|
| sgb_model | Model of type mboost to plot the variable importance. |
| prop | Numeric value indicating the minimal importance a predictor/baselearner has to have. Default value is zero, meaning all predictors are plotted. By increasing prop the number of plotted variables can be reduced. One can also use 'n_predictors' for limiting the number of variables to be plotted directly. |
| n_predictors | The maximum number of predictors to be plotted. Default is 30. Alternative to 'prop'. |
| max_char_length | The maximum character length of a predictor to be printed. Default is 15. For larger groups or long variable names one may adjust this number to differentiate variables from groups. |
| base_size | The base_size argument to be passed to the ggplot2 theme ggplot2::theme_bw to be used to control the overall size of the figure. Default value is 8. |

Details

Note that aggregated group and individual variable importance printed in the legend is based only on the plotted variables and not on all variables that were selected in the sparse-group boosting model.

Value

object of type ggplot2.

See Also

[get_varimp](#) which this function uses.

Examples

```
library(mboost)
library(dplyr)
set.seed(1)
df <- data.frame(
  x1 = rnorm(100), x2 = rnorm(100), x3 = rnorm(100),
  x4 = rnorm(100), x5 = runif(100)
)
df <- df %>%
  mutate_all(function(x) {
    as.numeric(scale(x))
  })
df$y <- df$x1 + df$x4 + df$x5
group_df <- data.frame(
  group_name = c(1, 1, 1, 2, 2),
  var_name = c("x1", "x2", "x3", "x4", "x5")
)

sgb_formula <- as.formula(create_formula(alpha = 0.3, group_df = group_df))
sgb_model <- mboost(formula = sgb_formula, data = df)
sgb_varimp <- plot_varimp(sgb_model)
```

Index

`create_formula`, 2

`get_coef`, 3
`get_coef()`, 5, 7
`get_coef_path`, 4
`get_coef_path()`, 8
`get_varimp`, 5, 10
`get_varimp()`, 7
`ggplot2::theme_bw`, 8, 9
`ggplot2::theme_classic`, 7

`mboost::mboost()`, 2
`mboost::varimp()`, 6

`plot_effects`, 6
`plot_path`, 8
`plot_varimp`, 9