# Package 'MTPS'

March 17, 2025

**Type** Package

**Title** Multi-Task Prediction using Stacking Algorithms

**Version** 1.1.1

**License** GPL (>= 2)

**URL** https://doi.org/10.1093/bioinformatics/btz531

**Encoding** UTF-8

**LazyData** true

**Depends** R (>= 3.5.0)

**Imports** caret, aftgee, class, dplyr, e1071, glmnet, MASS, rpart, survival

**Suggests** knitr, rmarkdown, ggplot2, reshape2

**VignetteBuilder** knitr

**NeedsCompilation** no

**RoxygenNote** 7.2.3

**Author** Li Xing [aut, cre],
Xiaowen Cao [aut],
Shuai You [aut],
Yuying Huang [aut],
Peijie Xie [ctb],
Mary Lesperance [aut],
Xuekui Zhang [aut]

**Repository** CRAN

**Date/Publication** 2023-5-20 04:30:03 UTC

**Maintainer** Li Xing <sfulxing@gmail.com>

## R topics documented:

1

check_match                     *Check Model-Family Match*

## Description

Check if the specified model function matches the given family type.

## Usage

```
check_match(family, FUN)
```

## Arguments

| | |
|---|---|
| family | A character string indicating the family (e.g., "gaussian", "binomial", "survival"). |
| FUN | The model function to check. |

## Value

A logical value indicating whether the function matches the family type.

## Examples

```
# Example usage of the check_match function
check_match("gaussian", lm)
check_match("binomial", glm)
check_match("survival", surv)
```

---

cv_MTPS                         *Evaluation using Cross-Validation*

---

## Description

Use cross-validation to evaluate model performance.

## Usage

```
cv_MTPS(
  xmat,
  ymat,
  family,
  nfolds = 5,
  cv = FALSE,
  residual = TRUE,
  cv_stacking_nfold = 5,
  method_step1,
  method_step2,
  resid_type = c("deviance", "pearson", "raw"),
  resid_std = FALSE,
  dist1 = NULL,
  weights = NULL
)
```

## Arguments

| | |
|---|---|
| xmat | Matrix of predictors, each row is an observation vector. |
| ymat | Matrix of outcomes. Quantitative for family = **"gaussian"**. A factor of two levels for family = **"binomial"**. A survival object for family = **"survival"**. |
| family | Response type for each response. If all response variables are within the same family, it can be **"gaussian"**, **"binomial"** or **"survival"**, otherwise it is a vector with elements **"gaussian"**, **"binomial"** and **"survival"** to indicate each response family. |
| nfolds | Integer, number of folds for Cross-Validation to evaluate the performance of stacking algorithms. |
| cv | Logical, indicate if use Cross-Validation Stacking algorithm. |
| residual | Logical, indicate if use Residual Stacking algorithm. |
| cv_stacking_nfold | |
| | Integer, number of folds for Cross-Validation Stacking algorithm. The default value is 5. |
| method_step1 | Base Learners for fitting models in Step 1 of Stacking Algorithm. It can be one base learner function for all outcomes or a list of base learner functions for each outcome. The list of all base learners can be obtained by list.learners(). |
| method_step2 | Base Learners for fitting models in Step 2 of Stacking Algorithm. |
| resid_type | The residual type for Residual Stacking. |
| resid_std | Logical, whether or not to use standardized residuals. |

| dist1 | Assumed distribution for survival response. If the argument is a character string, it is assumed to name an element from survreg.distributions. These include **"weibull"**, **"exponential"**, **"gaussian"**, **"logistic"**, **"lognormal"** and **"loglogistic"**. Otherwise, it is assumed to be a user-defined list conforming to the format described in "survreg.distributions". |
| --- | --- |
| weights | Weight for survival response. |

### Details

The most frequently used evaluation metric of survival models is the concordance index (C-index). It is a measure of rank correlation between predicted risk scores. Only uncensored observations are used.

### Value

It returns the mean squared error of continuous outcomes, and AUC, accuracy, recall, and precision for binary outcomes of predictions using cross-validation.

### Examples

```
data("HIV")
cv_MTPS(xmat = XX, ymat = YY, family = "gaussian", nfolds = 2, method_step1 = rpart1, method_step2 = lm1)
```

---

cv_multiFit                     *Cross-validation for Multi-Fit Model*

---

### Description

Use cross-validation to evaluate a multi-fit model's performance.

### Usage

```
cv_multiFit(
  xmat,
  ymat,
  xmat_list = NULL,
  nfold = 5,
  method,
  family = family,
  dist1 = NULL,
  weights = NULL
)
```

### Arguments

| xmat | Matrix of predictors, each row is an observation vector. |
| --- | --- |
| ymat | Matrix of outcomes (quantitative for family = "gaussian", factor of two levels for family = "binomial", survival object for family = "survival"). |
| xmat_list | A list of predictor matrices (optional). |
| nfold | Integer, number of folds for cross-validation. |

| | |
|---|---|
| `method` | A model fitting method or a list of methods for each outcome. |
| `family` | The response type for each response (can be "gaussian", "binomial", or "survival"). |
| `dist1` | Assumed distribution for survival response (optional). |
| `weights` | Weights for the survival response (optional). |

## Details

The function performs cross-validation using different methods for fitting models and evaluates them based on the family type (gaussian, binomial, survival).

## Value

A list containing the fitted models, predictions, and other model details.

## Examples

```
cv_multiFit(xmat = XX, ymat = YY, family = "gaussian", nfold = 2, method = rpart1)
```

---

| c_index | *Concordance Index (C-index)* |
|---|---|

---

## Description

Compute the C-index in survival analysis.

## Usage

```
c_index(pre, object)
```

## Arguments

| | |
|---|---|
| `pre` | A numeric vector of predicted survival times. |
| `object` | A dataframe or matrix, where the first column is observed survival time and the second column is survival status. |

## Details

The most frequently used evaluation metric of survival models is the concordance index (C-index). It is a measure of rank correlation between predicted risk scores. Only uncensored observations are considered.

## Value

The value of the C-index.

## Examples

```
set.seed(1)
# Simulate predicted survival times
x1 <- rnorm(100)
# Simulate observed survival times
x2 <- rnorm(100)
# Simulate observed survival status
x3 <- rep(c(1, 0), 50)
# Calculate C-index
c_index(x1, cbind(x2, x3))
```

---

find_km_weights_mat          *Calculate Kaplan-Meier Weights*

---

## Description

Produce the Kaplan-Meier weights

## Usage

```
find_km_weights_mat(multi_dat, num_outcome)
```

## Arguments

| | |
|---|---|
| multi_dat | A data frame containing the time and status of multi-variate survival outcomes for every subject. The name of each column corresponding to the time or status must contain the characters "time" or "status", respectively. |
| num_outcome | Number of the outcomes. |

## Details

The Kaplan-Meier weights are defined as:

$$w_{n1} = \frac{d_{(1)}}{n}, w_{ni} = \frac{d_{(i)}}{n-i+1} \prod_{j=1}^{i-1} (\frac{n-j}{n-j+1})^{d_{(j)}}, i = 2, ..., n.$$

## Value

A matrix of Kaplan-Meier weights for each survival outcome.

## References

Huang J, Ma S, Xie H. Regularized estimation in the accelerated failure time model with high-dimensional covariates. Biometrics. 2006;62(3):813-820.

## Examples

```
# multi-outcome
library(dplyr)
data(simdat_mtps)
find_km_weights_mat(ymat, 2)
```

| HIV | *HIV Drug Resistance Database* |

## Description

The data from HIV Drug Resistance Database used for demonstration. After processing, YY contains 5 response variables for 1246 observations and XX are 228 predictors of those 1246 observations.

## Format

Data objects used for demonstration

## Details

In the HIV database, the resistance of five Nucleoside RT Inhibitor (NRTI) drugs were used as multivariate outcomes, including Lamivudine (3TC), Abacavir(ABC), Zidovudine (AZT), Stavudine (D4T), Didanosine (DDI). The mutation variables are used as the predictors. Some mutation variables were removed as they do not contain enough variation. The final outcome data is a matrix of size 1246 × 5, and the predictor data is a matrix of 1246 × 228 values, which is provided in the package called "HIV". In the example data in the package, "YY" refers the outcome data and "XX" refers the predictor data.

## References

Rhee SY, Taylor J, Wadhera G, Ben-Hur A, Brutlag DL, Shafer RW. Genotypic predictors of human immunodeficiency virus type 1 drug resistance. Proceedings of the National Academy of Sciences. 2006 Nov 14;103(46):17355-60.

Rhee SY, Taylor J, Fessel WJ, Kaufman D, Towner W, Troia P, Ruane P, Hellinger J, Shirvani V, Zolopa A, Shafer RW. (2010). HIV-1 protease mutations and protease inhibitor cross-resistance. Antimicrobial Agents and Chemotherapy, 2010 Oct

Melikian GL, Rhee SY, Taylor J, Fessel WJ, Kaufman D, Towner W, Troia-Cancio PV, Zolopa A, Robbins GK, Kagan R, Israelski D, Shafer RW (2012). Standardized comparison of the relative impacts of HIV-1 reverse transcriptase (RT) mutations on nucleoside RT inhibitor susceptibility. Antimicrobial Agents and Chemother. 2012 May;56(5):2305-13.

Melikian GL, Rhee SY, Varghese V, Porter D, White K, Taylor J, Towner W, Troia P, Burack J, Dejesus E, Robbins GK, Razzeca K, Kagan R, Liu TF, Fessel WJ, Israelski D, Shafer RW (2013). Non-nucleoside reverse transcriptase inhibitor (NNRTI) cross-resistance: implications for preclinical evaluation of novel NNRTIs and clinical genotypic resistance testing. J Antimicrob Chemother, 2013 Aug 9.

## Examples

```
data(HIV)
```

---

list.learners                   *List Available Base Learners*

---

### Description

This function lists all base learners provided in the package.

### Usage

```
list_learners()
```

### Details

- lm1: Linear regression - glm1: Generalized linear models - glmnet1: Performs k-fold cross-validation to choose the best alpha and lambda for generalized linear models via penalized maximum likelihood. - glmnet_lasso: LASSO, lambda is chosen by k-fold cross-validation for glmnet. - glmnet_ridge: Ridge regression, lambda is chosen by k-fold cross-validation for glmnet. - rpart1: Regression tree - lda1: Linear discriminant analysis - qda1: Quadratic discriminant analysis - KNN1: K-nearest neighbor classification, k is chosen by cross-validation. - svm1: Support vector machine - surv: Parametric AFT model - ela1: Elastic Net AFT model, weights are the weight of survival time.

### Value

The name of all base learners provided in the package.

### Examples

```
list_learners()
```

---

modify_parameter             *Modify Default Parameters For Base Learner*

---

### Description

Modify default parameters for methods provided in the package.

### Usage

```
modify_parameter(FUN, ...)
```

### Arguments

| | |
|---|---|
| FUN | Method. |
| ... | Modified arguments. |

### Value

It returns a new function with modified parameters.

## Examples

```
glmnet_lasso <- modify_parameter(glmnet1, alpha = 1)
glmnet_ridge <- modify_parameter(glmnet1, alpha = 0)
```

---

mse                     *Mean Square Error Loss*

---

## Description

Compute the mean square error loss in survival analysis.

## Usage

```
mse(pre, object)
```

## Arguments

| | |
|---|---|
| pre | A numeric vector of predicted survival time. |
| object | A dataframe or matrix, the first column is observed survival time, the second column is survival status. |

## Details

The mean square error is calculated on log-scale, only considering uncensored observations.

## Value

The value of MSE.

## Examples

```
set.seed(1)
# simulate predicted survival time
x1 <- rnorm(100) + 5
# simulate observed survival time
x2 <- rnorm(100) + 10
# simulate observed survival status
x3 <- rep(c(1, 0), 50)
# calculate MSE
mse(x1, cbind(x2, x3))
```

---

MTPS                          *Fit Models using Revised Stacking Algorithm*

---

**Description**

Fit a model using standard stacking algorithm or revised stacking algorithms to simultaneously predict multiple outcomes.

**Usage**

```
MTPS(
  xmat,
  ymat,
  xmat_list = NULL,
  family,
  cv = FALSE,
  residual = TRUE,
  nfold = 5,
  method_step1,
  method_step2,
  resid_type = c("deviance", "pearson", "raw"),
  resid_std = FALSE,
  dist1 = NULL,
  weights = NULL,
  weight2 = FALSE
)
```

**Arguments**

| | |
|---|---|
| xmat | Predictor matrix, each row is an observation vector. |
| ymat | Responses matrix. Quantitative for family = "gaussian". A factor of two levels for family = "binomial". A survival object for family = "survival". |
| xmat_list | The user defines a list to specify the number of subset columns of the xmat for each outcome. |
| family | Response type for each response. If all response variable are within the same family, it can be "gaussian", "binomial" or "survival", otherwise it is a vector with elements "gaussian", "binomial", and "survival" to indicate each response family. |
| cv | Logical, indicate if use Cross-Validation Stacking algorithm. |
| residual | Logical, indicate if use Residual Stacking algorithm. |
| nfold | Integer, number of folds for Cross-Validation Stacking algorithm. The default value is 5. |
| method_step1 | Base Learners for fitting models in Step 1 of Stacking Algorithm. It can be one base learner function for all outcomes or a list of base learner functions for each outcome. The list of all base learners can be obtained by `list.learners()`. |
| method_step2 | Base Learners for fitting models in Step 2 of Stacking Algorithm. |
| resid_type | The residual type for Residual Stacking. |
| resid_std | Logical, whether or not use standardized residual. |

| dist1 | Assumed distribution for response variable. If the argument is a character string, then it is assumed to name an element from "survreg.distributions". These include "weibull", "exponential", "gaussian", "logistic", "lognormal", and "loglogistic". |
|---|---|
| weights | A logical indicating whether to use weight for response in step1. |
| weight2 | A logical indicating whether to use weight for response in step2. The default is FALSE. |

## Value

It returns an MTPS object. It is a list of 4 parameters containing information about step 1 and step 2 models and the revised stacking algorithm method.

## Examples

```
data("HIV")
set.seed(1)
xmat <- as.matrix(XX)
ymat <- as.matrix(YY)
id <- createFolds(rowMeans(XX), k=5, list=FALSE)
training_id <- id != 1
y_train <- ymat[training_id, ]
y_test  <- ymat[!training_id, ]
x_train <- xmat[training_id, ]
x_test  <- xmat[!training_id, ]

# Residual Stacking
fit_rs <- MTPS(xmat = x_train, ymat = y_train,
               family = "gaussian", cv = FALSE, residual = TRUE,
               method_step1 = rpart1, method_step2 = lm1)
pre1 <- predict(fit_rs, x_test)

# Using different base learners for different outcomes
fit_mix_out <- MTPS(xmat = x_train, ymat = y_train,
                    family = "gaussian", cv = FALSE, residual = TRUE,
                    method_step1 = c(rpart1, lm1, rpart1, lm1, lm1),
                    method_step2 = c(rpart1, lm1, lm1, lm1, lm1))
pre2 <- predict(fit_mix_out, x_test)

# Residual Stacking for Survival Analysis
set.seed(1)
data("simdat_mtps")
id_train <- sample(1:100, 80)
xmat_train <- xmat[id_train, ]
xmat_test <- xmat[-id_train, ]
ymat_train <- cbind(list(survival::Surv(ymat[id_train, "time01"], ymat[id_train, "status01"])),
            list(survival::Surv(ymat[id_train, "time02"], ymat[id_train, "status02"])))
weights <- find_km_weights_mat(ymat[id_train, ], num_outcome = 2)

xmat_list <- list(c(1), c(2))
fit <- MTPS(xmat_train, ymat_train, xmat_list = xmat_list, family = 'survival',
            cv = FALSE, residual = TRUE, nfold = 5, method_step1 = surv,
            method_step2 = lm1, dist1 = "lognormal", weights = weights)
pre3 <- predict.MTPS(fit, xmat_test)
```

---

multiFit                          *Fit models on multiple outcomes.*

---

**Description**

This function fits individual models to predict each outcome separately.

**Usage**

```
multiFit(
  xmat,
  ymat,
  xmat_list = NULL,
  method,
  family,
  dist1 = NULL,
  weights = NULL
)
```

**Arguments**

| | |
|---|---|
| xmat | Matrix of predictors, each row is an observation vector. |
| ymat | Matrix of outcomes. Quantitative for family = "gaussian". A factor of two levels for family = "binomial". A survival object for family = "survival". |
| xmat_list | The user defines a numeric vector to specify the number of subset columns of the xmat. |
| method | Method for fitting models. It can be one base learner function for all outcomes or a list of base learner functions for each outcome. The list of all base learners can be obtained by list.learners(). |
| family | Response type for each response. If all response variable are within the same family it can be "gaussian", "binomial" or "survival", otherwise it is a vector with elements "gaussian", "binomial" and "survival" to indicate each response family. |
| dist1 | Assumed distribution for response variable. If the argument is a character string, then it is assumed to name an element from ″survreg.distributions″. These include "weibull", "exponential", "gaussian", "logistic", "lognormal" and "loglogistic". |
| weights | Weight for response. |

**Value**

It returns a multiFit object. It is a list of 6 parameters containing information about the fitted models and fitted values for each outcome.

**Examples**

```
data("HIV")
set.seed(1)
xmat <- as.matrix(XX)
ymat <- as.matrix(YY)
```

```
id <- createFolds(rowMeans(XX), k = 5, list = FALSE)
training_id <- id != 1
y_train <- ymat[training_id, ]
y_test  <- ymat[!training_id, ]
x_train <- xmat[training_id, ]
x_test  <- xmat[!training_id, ]
fit <- multiFit(xmat = x_train, ymat = y_train,
                method = rpart1, family = "gaussian")
pre1 <- predict(fit, x_test)

# Using different base learners for different outcomes
fit_mix_out <- multiFit(xmat = x_train, ymat = y_train,
                        method = c(rpart1, rpart1, lm1, lm1, lm1),
                        family = "gaussian")
pre2 <- predict(fit_mix_out, x_test)
```

---

predict.MTPS                    *Make predictions from a "MTPS" model*

---

### Description

This function makes predictions from a revised stacking model.

### Usage

```
## S3 method for class 'MTPS'
predict(object, newdata, ...)
```

### Arguments

| | |
|---|---|
| object | A fitted object from "MTPS" |
| newdata | Matrix of new predictors at which predictions are to be made |
| ... | additional arguments affecting the predictions produced |

### Value

The predicted value from new predictors.

### Examples

```
data("HIV")
set.seed(1)
xmat <- as.matrix(XX)
ymat <- as.matrix(YY)
id <- createFolds(rowMeans(XX), k=5, list=FALSE)
training_id <- id != 1
y_train <- ymat[training_id, ]
y_test  <- ymat[!training_id, ]
x_train <- xmat[training_id, ]
x_test  <- xmat[!training_id, ]
# Cross-Validation Residual Stacking
fit_rs <- MTPS(xmat = x_train, ymat = y_train,
  family = "gaussian",cv = FALSE, residual = TRUE,
```

```
   method_step1 = rpart1, method_step2 = lm1)
pred_rs <- predict(fit_rs, x_test)

# Residual Stacking for Survival Analysis
  set.seed(1)
  data("simdat_mtps")

 # prepare training and test set
 id_train <- sample(1:100, 80)
 xmat_train <- xmat[id_train,]
 xmat_test <- xmat[-id_train,]
 ymat_train <- cbind(list(survival::Surv(ymat[id_train,"time01"],ymat[id_train,"status01"])),
 list(survival::Surv(ymat[id_train,"time02"],ymat[id_train,"status02"])))
 # Produce the Kaplan-Meier estimator
 weights <- find_km_weights_mat(ymat[id_train,],num_outcome = 2)
 # fit Residual Stacking Model for Survival Data
 fit <- MTPS(xmat_train, ymat_train, family = 'survival', cv=FALSE,
 residual = TRUE, nfold=5, method_step1 = surv,
 method_step2 = lm1, dist1 = "lognormal", weights = weights)
 # predict the survival time on test set
 pre <- predict.MTPS(fit, xmat_test)
```

---

| predict.multiFit | *Make predictions for multiple outcomes* |
|---|---|

---

### Description

This function makes predictions from a multiFit object.

### Usage

```
## S3 method for class 'multiFit'
predict(object, newdata, ...)
```

### Arguments

| | |
|---|---|
| object | A fitted object from "multiFit" |
| newdata | Matrix of new predictors at which predictions are to be made |
| ... | additional arguments affecting the predictions produced |

### Value

The predicted value from new predictors.

### Examples

```
data("HIV")
set.seed(1)
xmat <- as.matrix(XX)
ymat <- as.matrix(YY)
```

```
id <- createFolds(rowMeans(XX), k=5, list=FALSE)
training_id <- id != 1
y_train <- ymat[training_id, ]
y_test  <- ymat[!training_id, ]
x_train <- xmat[training_id, ]
x_test  <- xmat[!training_id, ]
fit <- multiFit(xmat = x_train, ymat = y_train,
                method = rpart1, family = "gaussian")
predict(fit, x_test)
```

---

simdat_mtps                 *A Simulated Database for Survival analysis*

---

## Description

The data is simulated for predicting multiple survival outcomes. ymat contains 2 survival response variables with 10% censoring for 100 observations and xmat contains 2 predictors of those observations.

## Details

In the simdat_mtps database, the simulated survival time and status are used as multivariate outcomes, including **time01**,**status01**,**time02**,**status02**. The predictors both follow a standard normal distribution. In this simulated data, "ymat" refers multiple survival outcomes and "xmat" refers the predictor data.

## Examples

```
data(simdat_mtps)
```