

# Package ‘highMLR’

May 23, 2026

**Title** Machine Learning Feature Selection for High Dimensional Survival Data

**Version** 1.0.1

**Date** 2026-05-23

**Description** A unified, flexible framework for high dimensional feature selection in the presence of a survival outcome. Provides multiple machine learning approaches (Cox elastic net, random survival forest, accelerated oblique random survival forest, gradient-boosted Cox, stability selection, classical univariate Cox screening, pseudo-observation bridging to arbitrary regression learners, and Fine-Gray competing risks selection) under a single interface. Adds causal survival forest estimation of heterogeneous treatment effects on survival (experimental), conformal survival prediction with finite-sample coverage guarantees, and time-dependent ‘SHAP’ explanations via ‘SurvSHAP(t)’. Methodology is based on regularised Cox regression (2011) <doi:10.18637/jss.v039.i05>, random survival forests (2008) <doi:10.1214/08-AOAS169>, oblique random survival forests (2024) <doi:10.1080/10618600.2023.2231048>, stability selection (2010) <doi:10.1111/j.1467-9868.2010.00740.x>, causal survival forests (2023) <doi:10.1111/rssb.12538>, time-dependent survival explanations (2023) <doi:10.1016/j.knosys.2022.110234>, conformal survival prediction (2023) <doi:10.1093/biomet/asad043>, the Fine-Gray model for competing risks (1999) <doi:10.1080/01621459.1999.10474144>, and pseudo-observation regression (2010) <doi:10.1177/0962280209105020>.

**Depends** R (>= 4.1.0)

**Imports** survival, glmnet, ranger, aorsf, xgboost, stabs, survex, grf, prodlim, cmprsk, future, future.apply, tibble, ggplot2, rlang, stats, utils

**Suggests** knitr, rmarkdown, testthat (>= 3.0.0), mice, riskRegression

**License** GPL-3

**Encoding** UTF-8

**Language** en-GB

**LazyData** true  
**LazyDataCompression** xz  
**RoxygenNote** 7.3.3  
**VignetteBuilder** knitr  
**Config/testthat/edition** 3  
**NeedsCompilation** no  
**Author** Atanu Bhattacharjee [aut, cre]  
**Maintainer** Atanu Bhattacharjee <atanustat@gmail.com>  
**Repository** CRAN  
**Date/Publication** 2026-05-23 12:30:02 UTC

## Contents

coef.highmlr_fit . . . . .	2
highmlr . . . . .	3
highmlr_causal . . . . .	4
highmlr_compare . . . . .	6
highmlr_conformal . . . . .	7
highmlr_explain . . . . .	8
highmlr_report . . . . .	9
highmlr_screen . . . . .	10
highmlr_stability . . . . .	11
hnsec . . . . .	11
plot.highmlr_conformal . . . . .	12
plot.highmlr_fit . . . . .	12
predict.highmlr_fit . . . . .	13
print.highmlr_conformal . . . . .	13
print.highmlr_fit . . . . .	14
srdata . . . . .	14
summary.highmlr_fit . . . . .	15
<b>Index</b>	<b>16</b>

---

coef.highmlr_fit	<i>Coefficients from a highmlr_fit</i>
------------------	--

---

## Description

Coefficients from a highmlr\_fit

## Usage

```
## S3 method for class 'highmlr_fit'
coef(object, ...)
```

**Arguments**

object	A 'highmlr_fit' object.
...	Unused.

**Value**

A named numeric vector of coefficients (where defined) or importance scores otherwise.

---

highmlr	<i>Machine learning feature selection for high dimensional survival data</i>
---------	--

---

**Description**

Fits one of several survival ML methods and returns a unified 'highmlr\_fit' object summarising the selected features, their importance/coefficients, and (optionally) out-of-sample performance.

**Usage**

```
highmlr(
  data,
  time,
  status,
  features = NULL,
  method = c("coxnet", "rsf", "aorsf", "xgboost", "stability", "univariate", "pseudo",
    "finegray"),
  engine = NULL,
  recipe = NULL,
  resampling = c("cv", "bootstrap", "holdout", "none"),
  folds = 5L,
  tune = FALSE,
  top_n = 50L,
  parallel = FALSE,
  seed = NULL,
  ...
)
```

**Arguments**

data	A data frame containing 'time', 'status', and the candidate features (or a superset). Rows with missing time/status are dropped.
time	Character scalar: name of the survival time column.
status	Character scalar: name of the event indicator column. For right-censored methods: 1 = event, 0 = censored. For Fine-Gray (method = "finegray"): 0 = censored, 1 = event of interest, 2+ = competing event(s).
features	Character vector of candidate feature column names. If 'NULL' (default), all columns except 'time' and 'status' are used.

method	One of "coxnet", "rsf", "aorsf", "xgboost", "stability", "univariate", "pseudo", "finegray".
engine	Optional engine override.
recipe	Optional preprocessing recipe object (currently accepted for forward compatibility; not yet applied).
resampling	One of "cv", "bootstrap", "holdout", "none".
folds	Integer, number of CV folds (default 5).
tune	Logical. Internal tuning (currently coxnet only).
top_n	Integer. For ranking-based methods, keep this many top features (default 50).
parallel	Logical. Use future-based parallelism for the embarrassingly parallel parts.
seed	Optional integer for reproducibility.
...	Additional arguments passed to the method-specific fitter.

**Value**

An object of class 'highmlr\_fit'. See [new\_highmlr\_fit()].

**Examples**

```
if (requireNamespace("glmnet", quietly = TRUE)) {
  data(hnsc)
  fit <- highmlr(hnsc, time = "OS", status = "Death",
                method = "coxnet", resampling = "cv", folds = 5)
  print(fit)
}
```

---

highmlr_causal	<i>Causal survival forest for heterogeneous treatment effects (experimental)</i>
----------------	--

---

**Description**

Estimates patient-level conditional average treatment effects (CATEs) on a survival outcome using 'grf::causal\_survival\_forest'. Unlike the rest of 'highMLR', this function answers a different question: not "which features predict survival?" but "for which patients does treatment T extend (or shorten) survival, and which features modify that effect?".

**Usage**

```

highmlr_causal(
  data,
  time,
  status,
  treatment,
  covariates = NULL,
  horizon = NULL,
  num.trees = 2000L,
  target = c("RMST", "survival.probability"),
  honesty = TRUE,
  seed = NULL,
  ...
)

## S3 method for class 'highmlr_causal'
print(x, n = 10, ...)

## S3 method for class 'highmlr_causal'
plot(x, ...)

```

**Arguments**

<code>data</code>	A data frame.
<code>time</code>	Character: name of the survival time column.
<code>status</code>	Character: name of the event indicator (0/1).
<code>treatment</code>	Character: name of the binary treatment column (0 = control, 1 = treated). Must be exactly two levels.
<code>covariates</code>	Character vector of covariate column names. If 'NULL', all columns other than 'time', 'status', 'treatment'.
<code>horizon</code>	Numeric. The time horizon at which the treatment effect on the survival probability is estimated. Defaults to the median observed time.
<code>num.trees</code>	Number of trees in the forest (default 2000).
<code>target</code>	One of "RMST" (restricted mean survival time difference up to 'horizon') or "survival.probability" (difference in survival probability at 'horizon').
<code>honesty</code>	Logical (default TRUE) – honest splitting per 'grf'.
<code>seed</code>	Optional integer seed.
<code>...</code>	Passed to 'grf::causal_survival_forest'.
<code>x</code>	A 'highmlr_causal' object.
<code>n</code>	Number of top covariates to print (default 10).

**Value**

An object of class 'highmlr\_causal' containing the fitted forest, per-patient CATE estimates with standard errors, and covariate importance.



**Value**

A list with two elements: 'fits' (named list of 'highmlr\_fit' objects) and 'summary' (a tibble of method, n\_selected, key metric).

**Examples**

```
## Not run:
data(hnsc)
cmp <- highmlr_compare(hnsc, "OS", "Death",
                      methods = c("coxnet", "rsf", "univariate"))
cmp$summary

## End(Not run)
```

---

highmlr\_conformal      *Conformal prediction intervals for survival times*

---

**Description**

Computes calibrated lower bounds on survival time for each new subject using a split-conformal procedure with inverse probability of censoring weights (Candes, Lei and Ren, 2023). The returned lower bound satisfies a marginal coverage guarantee approximately equal to one minus alpha under standard conformal assumptions and a consistent censoring model.

**Usage**

```
highmlr_conformal(
  fit,
  new_data,
  calibration_data = NULL,
  alpha = 0.1,
  calibration_split = 0.3,
  time = NULL,
  status = NULL,
  seed = NULL
)
```

**Arguments**

fit	A highmlr_fit object whose predict() method returns a linear predictor or risk score.
new_data	Data frame on which to compute prediction intervals.
calibration_data	Data frame on which to compute conformity scores. If NULL, a random calibration_split fraction of new_data is held out for calibration and the rest is used as the test set (split-conformal).

alpha	Miscoverage level; default 0.1 (so 90 percent coverage).
calibration_split	Fraction of new_data to use for calibration when calibration_data is NULL. Default 0.3.
time	Name of the survival time column in calibration data. Defaults to the column used in fit.
status	Name of the event column in calibration data.
seed	Optional integer seed for the split.

**Value**

An object of class `highmlr_conformal` containing per-subject point predictions and lower confidence bounds for survival time.

**Examples**

```
## Not run:
fit <- highmlr(d_train, "OS", "Death", method = "coxnet")
intv <- highmlr_conformal(fit, new_data = d_test, alpha = 0.1)
print(intv)
plot(intv)

## End(Not run)
```

---

highmlr\_explain

*Time-dependent SHAP explanations for a highmlr\_fit (SurvSHAP(t))*


---

**Description**

Computes `SurvSHAP(t)` attributions (Krzyzinski et al., 2023) – SHAP values that vary with follow-up time – for the top features in a fitted ‘`highmlr_fit`’. Returns the `survex` explainer, per-feature aggregated importance, and a plotting helper.

**Usage**

```
highmlr_explain(
  fit,
  new_data = NULL,
  top_n = 10L,
  times = NULL,
  method = c("survshap", "permutation", "break_down"),
  n_explain = 25L,
  seed = NULL,
  ...
)
```

```
## S3 method for class 'highmlr_explain'
print(x, n = 10, ...)
```

```
## S3 method for class 'highmlr_explain'
plot(x, top_n = 10, ...)
```

### Arguments

fit	A 'highmlr_fit' object with a stored model.
new_data	Data on which to compute explanations.
top_n	Number of top features to explain (default 10).
times	Optional numeric vector of time points at which SHAP values are computed. Defaults to a 20-point grid spanning the observed time range.
method	SHAP method passed through to 'survex'. Default "survshap" (time-dependent). Other options: "permutation", "break_down".
n_explain	How many test rows to compute SHAP for. Default 25 (SHAP is expensive; full-cohort computation is rarely needed).
seed	Optional integer for reproducibility of subsampling.
...	Passed to 'survex::model_survshap()' or 'survex::explain_survival()'.
x	A 'highmlr_explain' object.
n	Number of top features to print (default 10).

### Value

A list with class 'highmlr\_explain' containing: \* 'explainer' – the 'survex' explainer object \* 'survshap' – the time-dependent SHAP object (if applicable) \* 'top\_features' – the top features table from the fit \* 'aggregated' – tibble of mean absolute SHAP per feature, averaged across time and explained rows

---

highmlr_report	<i>Generate a Quarto/Rmd report skeleton for a highmlr_fit</i>
----------------	--

---

### Description

Writes a self-contained Rmd file that, when rendered, produces a standard biomarker report (selected features, hazard ratios where available, performance, forest plot).

### Usage

```
highmlr_report(fit, file = "highmlr_report.Rmd", render = FALSE)
```

### Arguments

fit	A 'highmlr_fit' object.
file	Output '.Rmd' path (default "highmlr_report.Rmd").
render	Logical: if 'TRUE', also render via 'rmarkdown::render()'.

**Value**

Invisibly, the path to the written file.

---

highmlr_screen	<i>Pre-screen features when <math>p</math> is very large</i>
----------------	--

---

**Description**

Lightweight filter before the main pipeline (e.g. to drop features with low variance or low marginal association).

**Usage**

```
highmlr_screen(
  data,
  time,
  status,
  features = NULL,
  filter = c("variance", "univariate_p", "none"),
  keep = 1000L
)
```

**Arguments**

data, time, status, features  
As in [highmlr()].

filter One of "variance", "univariate\_p", "none".

keep Integer, how many features to retain (default 1000).

**Value**

Character vector of retained feature names.

**Examples**

```
## Not run:
data(srdata)
keep <- highmlr_screen(srdata, "OS", "event",
  filter = "variance", keep = 500)
fit <- highmlr(srdata, "OS", "event", features = keep, method = "coxnet")

## End(Not run)
```

---

highmlr_stability	<i>Post-hoc stability analysis of a fitted highmlr_fit</i>
-------------------	--

---

**Description**

Runs stability selection on the data used in ‘fit’, returning a selection frequency per feature.

**Usage**

```
highmlr_stability(fit, B = 100L, cutoff = 0.75, PFER = 1, ...)
```

**Arguments**

fit	A ‘highmlr_fit’ object (used only for the data / call).
B	Number of subsamples (default 100).
cutoff	Selection probability threshold (default 0.75).
PFER	Per-family error rate bound (default 1).
...	Passed to [fit_stability()].

**Value**

A new ‘highmlr\_fit’ with ‘method = "stability"’.

---

hnscc	<i>High dimensional head and neck cancer survival and gene expression data</i>
-------	--

---

**Description**

Survival and gene expression measurements for head and neck squamous cell carcinoma patients, used to demonstrate high-dimensional feature selection.

**Usage**

```
hnscc
```

**Format**

A data frame with 565 rows (one per patient) and 104 columns. The first five columns are the identifier and outcome variables: ID (patient identifier), Death (overall survival event indicator, 1 = death, 0 = censored), OS (overall survival time), PFS (progression-free survival time), and Prog (progression event indicator, 1 = progression, 0 = none). The remaining 99 columns are numeric gene expression features named by gene symbol (for example GJB1, HPN, PROM1).

**Source**

Bundled with the package since highMLR v0.1.1.

---

`plot.highmlr_conformal`*Plot method for highmlr\_conformal objects*

---

**Description**

Plot method for highmlr\_conformal objects

**Usage**

```
## S3 method for class 'highmlr_conformal'  
plot(x, ...)
```

**Arguments**

x	A highmlr_conformal object.
...	Unused.

**Value**

A ggplot object.

---

`plot.highmlr_fit`*Forest / importance plot for a highmlr\_fit*

---

**Description**

Forest / importance plot for a highmlr\_fit

**Usage**

```
## S3 method for class 'highmlr_fit'  
plot(x, top_n = 20, ...)
```

**Arguments**

x	A 'highmlr_fit' object.
top_n	Number of top features to plot (default 20).
...	Unused.

**Value**

A 'ggplot' object.

---

predict.highmlr\_fit *Predict from a highmlr\_fit*

---

### Description

Predict from a highmlr\_fit

### Usage

```
## S3 method for class 'highmlr_fit'
predict(object, new_data, type = c("linear_pred", "survival", "risk"), ...)
```

### Arguments

object	A 'highmlr_fit' object.
new_data	A data frame containing the features used in fitting.
type	One of "linear_pred", "survival", or "risk". Availability depends on the underlying model.
...	Passed to the underlying model's predict method.

### Value

Predicted values (vector or tibble depending on 'type').

---

print.highmlr\_conformal  
*Print method for highmlr\_conformal objects*

---

### Description

Print method for highmlr\_conformal objects

### Usage

```
## S3 method for class 'highmlr_conformal'
print(x, n = 10, ...)
```

### Arguments

x	A highmlr_conformal object.
n	Number of rows to display in the preview table (default 10).
...	Unused.

### Value

Invisibly returns x.

`print.highmlr_fit`      *Print method for highmlr\_fit*

---

### Description

Print method for highmlr\_fit

### Usage

```
## S3 method for class 'highmlr_fit'
print(x, n = 10, ...)
```

### Arguments

<code>x</code>	A 'highmlr_fit' object.
<code>n</code>	Number of top features to display (default 10).
<code>...</code>	Unused.

### Value

Invisibly returns 'x'.

---

`srdata`      *High dimensional protein gene expression survival data*

---

### Description

Protein expression measurements with a survival outcome, used to demonstrate high-dimensional feature selection.

### Usage

```
srdata
```

### Format

A data frame with 288 rows and 250 columns. The first four columns are the identifier and outcome variables: `ID` (subject identifier), `Visit` (visit number), `OS` (overall survival time), and `event` (survival event indicator, 1 = event, 0 = censored). The remaining 246 columns are numeric protein expression features named by protein or marker (for example `C6kine`, `ActivinA`, `Adiponectin`).

### Source

Bundled with the package since highMLR v0.1.1.

---

summary.highmlr\_fit    *Summary method for highmlr\_fit*

---

### **Description**

Summary method for highmlr\_fit

### **Usage**

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

### **Arguments**

object	A 'highmlr_fit' object.
...	Unused.

### **Value**

A list with the full selected feature table and performance.

# Index

## \* datasets

hnscc, 11  
srdata, 14

coef.highmlr\_fit, 2

highmlr, 3  
highmlr\_causal, 4  
highmlr\_compare, 6  
highmlr\_conformal, 7  
highmlr\_explain, 8  
highmlr\_report, 9  
highmlr\_screen, 10  
highmlr\_stability, 11  
hnscc, 11

plot.highmlr\_causal (highmlr\_causal), 4  
plot.highmlr\_conformal, 12  
plot.highmlr\_explain (highmlr\_explain),  
8  
plot.highmlr\_fit, 12  
predict.highmlr\_fit, 13  
print.highmlr\_causal (highmlr\_causal), 4  
print.highmlr\_conformal, 13  
print.highmlr\_explain  
(highmlr\_explain), 8  
print.highmlr\_fit, 14

srdata, 14  
summary.highmlr\_fit, 15