

# Package ‘stratamatch’

January 8, 2021

**Type** Package

**Date** 2021-01-05

**Title** Stratification and Matching for Large Observational Data Sets

**Version** 0.1.6

**Maintainer** Rachael C. Aikens <rockyaikens@gmail.com>

**BugReports** <https://github.com/raikens1/stratamatch/issues>

**Description** A pilot matching design to automatically stratify and match large datasets. The `manual_stratify()` function allows users to manually stratify a dataset based on categorical variables of interest, while the `auto_stratify()` function does automatically by allocating a held-aside (pilot) data set, fitting a prognostic score (see Hansen (2008) <doi:10.1093/biomet/asn004>) on the pilot set, and stratifying the data set based on prognostic score quantiles. The `strata_match()` function then does optimal matching of the data set in parallel within strata.

**License** GPL-3

**Encoding** UTF-8

**LazyData** true

**Imports** dplyr (>= 0.8.3), Hmisc (>= 4.2-0), magrittr (>= 1.5), rlang (>= 0.4.0), survival (>= 2.44.1.1)

**Depends** R (>= 3.4.0)

**Suggests** knitr, optmatch (>= 0.9-11), rmarkdown, testthat (>= 2.1.0)

**URL** <https://github.com/raikens1/stratamatch>

**RoxygenNote** 7.1.0

**VignetteBuilder** knitr

**NeedsCompilation** no

**Author** Rachael C. Aikens [aut, cre],  
Joseph Rigdon [aut],  
Justin Lee [aut],  
Michael Baiocchi [aut],  
Jonathan Chen [aut]

**Repository** CRAN

**Date/Publication** 2021-01-08 20:00:02 UTC

## R topics documented:

auto_stratify . . . . .	2
extract_cut_points . . . . .	5
extract_cut_points.auto_strata . . . . .	6
ICU_data . . . . .	6
is.auto_strata . . . . .	8
is.manual_strata . . . . .	8
is.strata . . . . .	9
make_sample_data . . . . .	9
manual_stratify . . . . .	10
plot.strata . . . . .	11
print.auto_strata . . . . .	12
print.manual_strata . . . . .	13
split_pilot_set . . . . .	13
stratamatch . . . . .	14
strata_match . . . . .	15
summary.strata . . . . .	16
<b>Index</b>	<b>17</b>

---

auto_stratify	<i>Auto Stratify</i>
---------------	----------------------

---

## Description

Automatically creates strata for matching based on a prognostic score formula or a vector of prognostic scores already estimated by the user. Creates a `auto_strata` object, which can be passed to `strata_match` for stratified matching or unpacked by the user to be matched by some other means.

## Usage

```
auto_stratify(
  data,
  treat,
  prognosis,
  outcome = NULL,
  size = 2500,
  pilot_fraction = 0.1,
  pilot_size = NULL,
  pilot_sample = NULL,
  group_by_covariates = NULL
)
```

**Arguments**

<code>data</code>	<code>data.frame</code> with observations as rows, features as columns
<code>treat</code>	string giving the name of column designating treatment assignment
<code>prognosis</code>	information on how to build prognostic scores. Three different input types are allowed: <ol style="list-style-type: none"> <li>1. vector of prognostic scores for all individuals in the data set. Should be in the same order as the rows of <code>data</code>.</li> <li>2. a formula for fitting a prognostic model</li> <li>3. an already-fit prognostic score model</li> </ol>
<code>outcome</code>	string giving the name of column with outcome information. Required if <code>prognostic_scores</code> is specified. Otherwise it will be inferred from <code>prog_formula</code>
<code>size</code>	numeric, desired size of strata (default = 2500)
<code>pilot_fraction</code>	numeric between 0 and 1 giving the proportion of controls to be allotted for building the prognostic score (default = 0.1)
<code>pilot_size</code>	alternative to <code>pilot_fraction</code> . Approximate number of observations to be used in pilot set. Note that the actual pilot set size returned may not be exactly <code>pilot_size</code> if <code>group_by_covariates</code> is specified because balancing by covariates may result in deviations from desired size. If <code>pilot_size</code> is specified, <code>pilot_fraction</code> is ignored.
<code>pilot_sample</code>	a <code>data.frame</code> of held aside samples for building prognostic score model. If <code>pilot_sample</code> is specified, <code>pilot_size</code> and <code>pilot_fraction</code> are both ignored.
<code>group_by_covariates</code>	character vector giving the names of covariates to be grouped by (optional). If specified, the pilot set will be sampled in a stratified manner, so that the composition of the pilot set reflects the composition of the whole data set in terms of these covariates. The specified covariates must be categorical.

**Details**

Stratifying by prognostic score quantiles can be more effective than manually stratifying a data set because the prognostic score is continuous, thus the strata produced tend to be of equal size with similar prognosis.

Automatic stratification requires information on how the prognostic scores should be derived. This is primarily determined by the specification of the `prognosis` argument. Three main forms of input for `prognosis` are allowed:

1. A vector of prognostic scores. This vector should be the same length and order of the rows in the data set. If this method is used, the `outcome` argument must also be specified; this is simply a string giving the name of the column which contains outcome information.
2. A formula for prognosis (e.g. `outcome ~ X1 + X2`). If this method is used, `auto_stratify` will automatically split the data set into a `pilot_set` and an `analysis_set`. The pilot set will be used to fit a logistic regression model for outcome in the absence of treatment, and this model will be used to estimate prognostic scores on the analysis set. The analysis set will then be stratified based on the estimated prognostic scores. In this case the `outcome` argument need not be specified since it can be inferred from the input formula.

3. A model for prognosis (e.g. a glm object). If this method is used, the outcome argument must also be specified

## Value

Returns an auto\_strata object. This contains:

- `outcome` - a string giving the name of the column where outcome information is stored
- `treat` - a string giving the name of the column encoding treatment assignment
- `analysis_set` - the data set with strata assignments
- `call` - the call to `auto_stratify` used to generate this object
- `issue_table` - a table of each stratum and potential issues of size and treat:control balance
- `strata_table` - a table of each stratum and the prognostic score quantile bin to which it corresponds
- `prognostic_scores` - a vector of prognostic scores.
- `prognostic_model` - a model for prognosis fit on a pilot data set. Will be NULL if a vector of prognostic scores was provided as the prognosis argument to `auto_stratify` rather than a model or formula.
- `pilot_set` - the set of controls used to fit the prognostic model. These are excluded from subsequent analysis so that the prognostic score is not overfit to the data used to estimate the treatment effect. Will be NULL if a pre-fit model or a vector of prognostic scores was provided as the prognosis argument to `auto_stratify` rather than formula.

## Troubleshooting

This section suggests fixes for common errors that appear while fitting the prognostic score or using it to estimate prognostic scores on the analysis set.

- Encountered an error while fitting the prognostic model... numeric probabilities 0 or 1 produced. This error means that the prognostic model can perfectly separate positive from negative outcomes. Estimating a treatment effect in this case is unwise since an individual's baseline characteristics perfectly determine their outcome, regardless of whether they receive the treatment. This error may also appear on rare occasions when your pilot set is very small (number of observations approximately  $\leq$  number of covariates in the prognostic model), so that perfect separation happens by chance.
- Encountered an error while estimating prognostic scores ... factor X has new levels ... This may indicate that some value(s) of one or more categorical variables appear in the analysis set which were not seen in the pilot set. This means that when we try to obtain prognostic scores for our analysis set, we run into some new value that our prognostic model was not prepared to handle. There are a few options we have to troubleshoot this problem:
  - **Rejection sampling.** Run `auto_stratify` again with the same arguments until this error does not occur (i.e. until some observations with the missing value are randomly selected into the pilot set)
  - **Eliminate this covariate from the prognostic formula.**
  - **Remove observations with the rare covariate value from the entire data set.** Consider carefully how this exclusion might affect your results.

Other errors or warnings can occur if the pilot set is too small and the prognostic formula is too complicated. Always make sure that the number of observations in the pilot set is large enough that you can confidently fit a prognostic model with the number of covariates you want.

### See Also

[manual\\_stratify](#), [new\\_auto\\_strata](#)

### Examples

```
# make sample data set
set.seed(111)
dat <- make_sample_data(n = 75)

# construct a pilot set, build a prognostic score for `outcome` based on X2
# and stratify the data set based on the scores into sets of about 25
# observations
a.strat_formula <- auto_stratify(dat, "treat", outcome ~ X2, size = 25)

# stratify the data set based on a model for prognosis
pilot_data <- make_sample_data(n = 30)
prognostic_model <- glm(outcome ~ X2, pilot_data, family = "binomial")
a.strat_model <- auto_stratify(dat, "treat", prognostic_model,
  outcome = "outcome", size = 25
)

# stratify the data set based on a vector of prognostic scores
prognostic_scores <- predict(prognostic_model,
  newdata = dat,
  type = "response"
)
a.strat_scores <- auto_stratify(dat, "treat", prognostic_scores,
  outcome = "outcome", size = 25
)

# diagnostic plots
plot(a.strat_formula)
plot(a.strat_formula, type = "AC", propensity = treat ~ X1, stratum = 1)
plot(a.strat_formula, type = "hist", propensity = treat ~ X1, stratum = 1)
plot(a.strat_formula, type = "residual")
```

---

extract\_cut\_points      *Extract cutoffs between strata*

---

### Description

By default, returns only the internal cut points. Cutoffs at 0 and 1 are implied.

### Usage

```
extract_cut_points(x)
```

**Arguments**

x                    an autostrata object

**Value**

a vector of the score values delineating cutoffs between strata

**Examples**

```
dat <- make_sample_data()
a.strat <- auto_stratify(dat, "treat", outcome ~ X1 + X2)
cutoffs <- extract_cut_points(a.strat)
```

---

```
extract_cut_points.auto_strata
```

*Extract cutoffs between strata*

---

**Description**

Extract cutoffs between strata

**Usage**

```
## S3 method for class 'auto_strata'
extract_cut_points(x)
```

**Arguments**

x                    an autostrata object

**Value**

a vector of the score values delineating cutoffs between strata

---

ICU\_data                    *Demographics and comorbidities of 10,157 ICU patients*

---

**Description**

An deidentified data set containing the demographics, comorbidities, DNR code status, and surgical team assignment of 10,157 patients in the Stanford University Hospital Intensive Care Unit (ICU). This data was extracted from the electronic record system, deidentified, and made publically available by Chavez et al (2018) <doi:10.1371/journal.pone.0190569>. It was reprocessed for use in the stratamatch package as a sample data set. For more details on the data extraction and inclusion criteria, see Chavez et al.

**Usage**

ICU\_data

**Format**

A data frame with 10157 rows and 29 variables:

**patid** patient id, numeric

**Birth.preTimeDays** age of patient at time of admission to the ICU in days, numeric

**Female.pre** whether the patient was documented to be female prior to ICU visit, binary

**RaceAsian.pre** whether the patient's race/ethnicity was documented as Asian prior to ICU visit, binary

**RaceUnknown.pre** whether the patient's race/ethnicity was unknown prior to ICU visit, binary

**RaceOther.pre** whether the patient's race/ethnicity was documented as "Other" prior to ICU visit, binary

**RaceBlack.pre** whether the patient's race/ethnicity was documented as Black/African American prior to ICU visit, binary

**RacePacificIslander.pre** whether the patient's race/ethnicity was documented as PacificIslander prior to ICU visit, binary

**RaceNativeAmerican.pre** whether the patient's race/ethnicity was documented as Native American prior to ICU visit, binary

**self\_pay** whether the patient was "self pay" (i.e. uninsured), binary

**all\_latinos** whether the patient was documented to be latino prior to ICU visit, binary

**DNR** whether the patient had code status set to any DNR "Do not resuscitate" order at any point during their ICU stay, binary

**surgicalTeam** whether the patient was assigned to a surgical team at any point during their ICU stay, binary

**Details**

License information for this data is as follows:

Copyright (c) 2016, Stanford University

Permission is hereby granted, free of charge, to any person obtaining a copy of this software and associated documentation files (the "Software"), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions:

The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Software.

THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.

**Source**

[https://simtk.org/frs/download\\_confirm.php/latestzip/1969/ICUDNR-latest.zip?group\\_id=892](https://simtk.org/frs/download_confirm.php/latestzip/1969/ICUDNR-latest.zip?group_id=892)

---

is.auto\_strata      *Checks auto\_strata class*

---

**Description**

Checks if the target object is an auto\_strata object.

**Usage**

```
is.auto_strata(object)
```

**Arguments**

object            any R object

**Value**

Returns TRUE if its argument has auto\_strata among its classes and FALSE otherwise.

**Examples**

```
dat <- make_sample_data()
a.strat <- auto_stratify(dat, "treat", outcome ~ X1 + X2)
is.auto_strata(a.strat) # returns TRUE
```

---

is.manual\_strata      *Checks manual\_strata class*

---

**Description**

Checks if the target object is a manual\_strata object.

**Usage**

```
is.manual_strata(object)
```

**Arguments**

object            any R object

**Value**

Returns TRUE if its argument has manual\_strata among its classes and FALSE otherwise.



**Examples**

```
dat <- make_sample_data()
m.strat <- manual_stratify(dat, treat ~ C1)
is.manual_strata(m.strat) # returns TRUE
```

---

is.strata	<i>Checks strata class</i>
-----------	----------------------------

---

**Description**

Checks if the target object is a strata object.

**Usage**

```
is.strata(object)
```

**Arguments**

object            any R object

**Value**

Returns TRUE if its argument has strata among its classes and FALSE otherwise.

**Examples**

```
dat <- make_sample_data()
m.strat <- manual_stratify(dat, treat ~ C1)
is.strata(m.strat) # returns TRUE
```

---

make_sample_data	<i>Make sample data</i>
------------------	-------------------------

---

**Description**

Makes a simple data frame with treat (binary), outcome (binary), and five covariates: X1 (continuous), X2 (continuous), B1 (binary), B2 (binary), and C1 (categorical). Probability outcome = 1 is  $\text{sigmoid}(\text{treat} + X1)$ . Probability treatment = 1 is  $\text{sigmoid}(-0.2 * X1 + X2 - B1 + 2 * B2)$

**Usage**

```
make_sample_data(n = 100)
```

**Arguments**

n                    the size of the desired data set

**Examples**

```
# make sample data set of 30 observations
dat <- make_sample_data(n = 30)
```

---

manual_stratify	<i>Manual Stratify</i>
-----------------	------------------------

---

**Description**

Stratifies a data set based on a set of blocking covariates specified by the user. Creates a `manual_strata` object, which can be passed to `strata_match` for stratified matching or unpacked by the user to be matched by some other means.

**Usage**

```
manual_stratify(data, strata_formula, force = FALSE)
```

**Arguments**

<code>data</code>	data.frame with observations as rows, features as columns
<code>strata_formula</code>	the formula to be used for stratification. (e.g. <code>treat ~ X1</code> ) the variable on the left is taken to be the name of the treatment assignment column, and the variables on the left are taken to be the variables by which the data should be stratified
<code>force</code>	a boolean. If true, run even if a variable appears continuous. (default = FALSE)

**Value**

Returns a `manual_strata` object. This contains:

- `treat` - a string giving the name of the column encoding treatment assignment
- `covariates` - a character vector with the names of the categorical columns on which the data were stratified
- `analysis_set` - the data set with strata assignments
- `call` - the call to `manual_stratify` used to generate this object
- `issue_table` - a table of each stratum and potential issues of size and treat:control balance
- `strata_table` - a table of each stratum and the covariate bin to which it corresponds

**See Also**

[auto\\_stratify](#), [new\\_manual\\_strata](#)

## Examples

```
# make sample data set
dat <- make_sample_data(n = 75)

# stratify based on B1 and B2
m.strat <- manual_stratify(dat, treat ~ B1 + B2)

# diagnostic plot
plot(m.strat)
```

---

plot.strata

*Plot method for strata object*

---

## Description

Generates diagnostic plots for the product of a stratification by [auto\\_stratify](#) or [manual\\_stratify](#). There are three plot types:

1. "SR" (default) - produces a scatter plot of strata by size and treat:control ratio
2. "hist" - produces a histogram of propensity scores within a stratum
3. "AC" - produces a Assignment-Control plot of individuals within a stratum (not supported for manual strata objects)
4. "residual" - produces a residual plot for the prognostic model (not supported for manual strata objects)

## Usage

```
## S3 method for class 'strata'
plot(
  x,
  type = "SR",
  label = FALSE,
  stratum = "all",
  strata_lines = TRUE,
  jitter_prognosis,
  jitter_propensity,
  propensity,
  ...
)
```

## Arguments

x a strata object returned by [auto\\_stratify](#) or [manual\\_stratify](#)

type string giving the plot type (default = "SR"). Other options are "hist", "AC" and "residual"

label	ignored unless type = "SR". If TRUE, a clickable plot is produced. The user may click on any number of strata and press finish to have those strata labeled. Note: uses <code>identify</code> , which may not be supported on some devices
stratum	ignored unless type = "hist" or type = "AC". A number specifying which stratum to plot.
strata_lines	default = TRUE. Ignored unless type = "AC". If TRUE, lines on the plot indicate strata cut points.
jitter_prognosis	ignored unless type = "AC". Amount of uniform random noise to add to prognostic scores in plot.
jitter_propensity	ignored unless type = "AC". Amount of uniform random noise to add to propensity scores in plot.
propensity	ignored unless type = "hist" or type = "AC". Specifies propensity score information for plots where this is required. Accepts either a vector of propensity scores, a <code>glm</code> model for propensity scores, or a formula for fitting a propensity score model.
...	other arguments

**See Also**

Aikens et al. (preprint) <https://arxiv.org/abs/1908.09077> . Section 3.2 for an explanation of Assignment-Control plots

**Examples**

```
dat <- make_sample_data()
a.strat <- auto_stratify(dat, "treat", outcome ~ X1 + X2)
plot(a.strat) # makes size-ratio scatter plot
plot(a.strat, type = "hist", propensity = treat ~ X1, stratum = 1)
plot(a.strat, type = "AC", propensity = treat ~ X1, stratum = 1)
plot(a.strat, type = "residual")
```

---

`print.auto_strata`      *Print Auto Strata*

---

**Description**

Print method for `auto_strata` object

**Usage**

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

**Arguments**

x,                    an auto\_strata object  
...                   other arguments

**Examples**

```
dat <- make_sample_data()
a.strat <- auto_stratify(dat, "treat", outcome ~ X1 + X2)
print(a.strat) # prints information about a.strat
```

---

print.manual\_strata    *Print Manual Strata*

---

**Description**

Print method for manual\_strata object

**Usage**

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

**Arguments**

x,                    a manual\_strata object  
...                   other arguments

**Examples**

```
dat <- make_sample_data()
m.strat <- manual_stratify(dat, treat ~ C1)
print(m.strat) # prints information about m.strat
```

---

split\_pilot\_set        *Split data into pilot and analysis sets*

---

**Description**

Given a data set and some parameters about how to split the data, this function partitions the data accordingly and returns the partitioned data as a list containing the analysis\_set and pilot\_set.

**Usage**

```
split_pilot_set(  
  data,  
  treat,  
  pilot_fraction = 0.1,  
  pilot_size = NULL,  
  group_by_covariates = NULL  
)
```

**Arguments**

**data** `data.frame` with observations as rows, features as columns

**treat** string giving the name of column designating treatment assignment

**pilot\_fraction** numeric between 0 and 1 giving the proportion of controls to be allotted for building the prognostic score (default = 0.1)

**pilot\_size** alternative to `pilot_fraction`. Approximate number of observations to be used in pilot set. Note that the actual pilot set size returned may not be exactly `pilot_size` if `group_by_covariates` is specified because balancing by covariates may result in deviations from desired size. If `pilot_size` is specified, `pilot_fraction` is ignored.

**group\_by\_covariates** character vector giving the names of covariates to be grouped by (optional). If specified, the pilot set will be sampled in a stratified manner, so that the composition of the pilot set reflects the composition of the whole data set in terms of these covariates. The specified covariates must be categorical.

**Value**

a list with `analysis_set` and `pilot_set`

**Examples**

```
dat <- make_sample_data()  
splt <- split_pilot_set(dat, "treat", 0.2)  
# can be passed into auto_stratify if desired  
a.strat <- auto_stratify(splt$analysis_set, "treat", outcome ~ X1,  
  pilot_sample = splt$pilot_set  
)
```

**Description**

This package employs a pilot matching design to automatically stratify and match large datasets. The `manual_stratify` function allows users to manually stratify a dataset based on categorical variables of interest, while the `auto_stratify` function does automatically by allocating a held-aside (pilot) data set, fitting a prognostic score (see Hansen (2008) <doi:10.1093/biomet/asn004>) on the pilot set, and stratifying the data set based on prognostic score quantiles. The `strata_match` function then does optimal matching of the data set within strata.

**See Also**

1. <https://github.com/raikens1/stratamatch>

---

strata\_match

*Strata Match*


---

**Description**

Match within strata in series using optmatch. Requires optmatch package to be installed.

**Usage**

```
strata_match(object, propensity = NULL, k = 1)
```

**Arguments**

object	a strata object
propensity	(optional) formula for propensity score. If left blank, all columns of the dataset will be used as covariates (except outcome, treatment and stratum)
k	numeric, the number of control individuals to be matched to each treated individual

**Value**

a named factor with matching assignments

**Examples**

```
# make a sample data set
set.seed(1)
dat <- make_sample_data(n = 75)

# stratify with auto_stratify
a.strat <- auto_stratify(dat, "treat", outcome ~ X2, size = 25)

# 1:1 match based on propensity formula: treat ~ X1 + X2
strata_match(a.strat, propensity = treat ~ X1 + X2, k = 1)
```

---

summary.strata	<i>Summary for strata object</i>
----------------	----------------------------------

---

**Description**

Summarize number and sizes of strata in a strata object. Also prints number of strata with potential issues.

**Usage**

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

**Arguments**

object	a strata object
...	other arguments

**Details**

For more information, access the issue table for your strata object with `mystrata$issue_table`.

**Examples**

```
dat <- make_sample_data()  
m.strat <- manual_stratify(dat, treat ~ C1)  
summary(m.strat) # Summarizes strata in m.strat
```



# Index

## \* datasets

- ICU\_data, [6](#)
  
- auto\_stratify, [2](#), [10](#), [11](#), [15](#)
  
- extract\_cut\_points, [5](#)
- extract\_cut\_points.auto\_strata, [6](#)
  
- ICU\_data, [6](#)
- identify, [12](#)
- is.auto\_strata, [8](#)
- is.manual\_strata, [8](#)
- is.strata, [9](#)
  
- make\_sample\_data, [9](#)
- manual\_stratify, [5](#), [10](#), [11](#), [15](#)
  
- new\_auto\_strata, [5](#)
- new\_manual\_strata, [10](#)
  
- plot.strata, [11](#)
- print.auto\_strata, [12](#)
- print.manual\_strata, [13](#)
  
- split\_pilot\_set, [13](#)
- strata\_match, [2](#), [10](#), [15](#), [15](#)
- stratamatch, [14](#)
- summary.strata, [16](#)