

Package ‘modeldata’

December 6, 2019

Title Data Sets Used Useful for Modeling Packages

Version 0.0.1

Description

Data sets used for demonstrating or testing model-related packages are contained in this package.

Depends R (>= 2.10)

License MIT + file LICENSE

Encoding UTF-8

LazyData false

RoxygenNote 7.0.1.9000

URL <https://github.com/tidymodels/modeldata>

NeedsCompilation no

Author Max Kuhn [aut, cre],
RStudio [cph]

Maintainer Max Kuhn <max@rstudio.com>

Repository CRAN

Date/Publication 2019-12-06 21:20:06 UTC

R topics documented:

ad_data	2
attrition	3
biomass	4
bivariate	4
car_prices	5
cells	5
check_times	6
Chicago	7
concrete	8
covers	8
credit_data	9
drinks	9

hpc_cv	10
hpc_data	10
lending_club	11
meats	12
mlc_churn	12
oils	13
okc	13
okc_text	14
parabolic	14
pathology	15
pd_speech	15
Sacramento	16
scat	17
small_fine_foods	17
Smithsonian	18
solubility_test	18
stackoverflow	19
two_class_dat	19
two_class_example	20
wa_churn	20

Index	22
--------------	-----------

ad_data	<i>Alzheimer's disease data</i>
---------	---------------------------------

Description

Alzheimer's disease data

Details

Craig-Schapiro et al. (2011) describe a clinical study of 333 patients, including some with mild (but well-characterized) cognitive impairment as well as healthy individuals. CSF samples were taken from all subjects. The goal of the study was to determine if subjects in the early states of impairment could be differentiated from cognitively healthy individuals. Data collected on each subject included:

- Demographic characteristics such as age and gender
- Apolipoprotein E genotype
- Protein measurements of Abeta, Tau, and a phosphorylated version of Tau (called pTau)
- Protein measurements of 124 exploratory biomarkers, and
- Clinical dementia scores

For these analyses, we have converted the scores to two classes: impaired and healthy. The goal of this analysis is to create classification models using the demographic and assay data to predict which patients have early stages of disease.

Value

ad_data a tibble

Source

Kuhn, M., Johnson, K. (2013) *Applied Predictive Modeling*, Springer.

Craig-Schapiro R, Kuhn M, Xiong C, Pickering EH, Liu J, Misko TP, et al. (2011) Multiplexed Immunoassay Panel Identifies Novel CSF Biomarkers for Alzheimer's Disease Diagnosis and Prognosis. *PLoS ONE* 6(4): e18850.

Examples

```
data(ad_data)
str(ad_data)
```

attrition	<i>Job attrition</i>
-----------	----------------------

Description

Job attrition

Details

These data are from the IBM Watson Analytics Lab. The website describes the data with “Uncover the factors that lead to employee attrition and explore important questions such as ‘show me a breakdown of distance from home by job role and attrition’ or ‘compare average monthly income by education and attrition’. This is a fictional data set created by IBM data scientists.”. There are 1470 rows.

Value

attrition a data frame

Source

The IBM Watson Analytics Lab website <https://www.ibm.com/communities/analytics/watson-analytics-blog/hr-employee-attrition/>

Examples

```
data(attrition)
str(attrition)
```

`biomass`*Biomass data*

Description

Ghugare et al (2014) contains a data set where different biomass fuels are characterized by the amount of certain molecules (carbon, hydrogen, oxygen, nitrogen, and sulfur) and the corresponding higher heating value (HHV). These data are from their Table S.2 of the Supplementary Materials

Value

`biomass` a data frame

Source

Ghugare, S. B., Tiwary, S., Elangovan, V., and Tambe, S. S. (2013). Prediction of Higher Heating Value of Solid Biomass Fuels Using Artificial Intelligence Formalisms. *BioEnergy Research*, 1-12.

Examples

```
data(biomass)
str(biomass)
```

`bivariate`*Example bivariate classification data*

Description

Example bivariate classification data

Details

These data are a simplified version of the segmentation data contained in `caret`. There are three columns: A and B are predictors and the column `Class` is a factor with levels "One" and "Two". There are three data sets: one for training (n = 1009), validation (n = 300), and testing (n = 710).

Value

`bivariate_train`, `bivariate_test`, `bivariate_val`
tibbles

Examples

```
data(bivariate)
```

car_prices

Kelly Blue Book resale data for 2005 model year GM cars

Description

Kuiper (2008) collected data on Kelly Blue Book resale data for 804 GM cars (2005 model year).

Value

cars data frame of the suggested retail price (column Price) and various characteristics of each car (columns Mileage, Cylinder, Doors, Cruise, Sound, Leather, Buick, Cadillac, Chevy, Pontiac, Saab, Saturn, convertible, coupe, hatchback, sedan and wagon)

Source

Kuiper, S. (2008). Introduction to Multiple Regression: How Much Is Your Car Worth?, *Journal of Statistics Education*, Vol. 16 http://jse.amstat.org/jse_archive.htm#2008.

cells

Cell body segmentation

Description

Hill, LaPan, Li and Haney (2007) develop models to predict which cells in a high content screen were well segmented. The data consists of 119 imaging measurements on 2019. The original analysis used 1009 for training and 1010 as a test set (see the column called case).

Details

The outcome class is contained in a factor variable called class with levels "PS" for poorly segmented and "WS" for well segmented.

The raw data used in the paper can be found at the Biomedcentral website. The version contained in cells is modified. First, several discrete versions of some of the predictors (with the suffix "Status") were removed. Second, there are several skewed predictors with minimum values of zero (that would benefit from some transformation, such as the log). A constant value of 1 was added to these fields: avg_inten_ch_2, fiber_align_2_ch_3, fiber_align_2_ch_4, spot_fiber_count_ch_4 and total_inten_ch_2.

Value

cells a tibble

Source

Hill, LaPan, Li and Haney (2007). Impact of image segmentation on high-content screening data quality for SK-BR-3 cells, *BMC Bioinformatics*, Vol. 8, pg. 340, <http://www.biomedcentral.com/1471-2105/8/340>.

check_times

Execution time data

Description

These data were collected from the CRAN web page for 13,626 R packages. The time to complete the standard package checking routine was collected. In some cases, the package checking process is stopped due to errors and these data are treated as censored. It is less than 1 percent.

Details

As predictors, the associated package source code were downloaded and parsed to create predictors, including

- authors: The number of authors in the author field.
- imports: The number of imported packages.
- suggests: The number of packages suggested.
- depends: The number of hard dependencies.
- Roxygen: a binary indicator for whether Roxygen was used for documentation.
- gh: a binary indicator for whether the URL field contained a GitHub link.
- rforge: a binary indicator for whether the URL field contained a link to R-forge.
- descr: The number of characters (or, in some cases, bytes) in the description field.
- r_count: The number of R files in the R directory.
- r_size: The total disk size of the R files.
- ns_import: Estimated number of imported functions or methods.
- ns_export: Estimated number of exported functions or methods.
- s3_methods: Estimated number of S3 methods.
- s4_methods: Estimated number of S4 methods.
- doc_count: How many Rmd or Rnw files in the vignettes directory.
- doc_size: The disk size of the Rmd or Rnw files.
- src_count: The number of files in the src directory.
- src_size: The size on disk of files in the src directory.
- data_count: The number of files in the data directory.
- data_size: The size on disk of files in the data directory.
- testthat_count: The number of files in the testthat directory.

- `testthat_size`: The size on disk of files in the `testthat` directory.
- `check_time`: The time (in seconds) to run R CMD check using the "r-devel-windows-ix86+x86_64" flavor.
- `status`: An indicator for whether the tests completed.

Data were collected on 2019-01-20.

Value

`check_times` a data frame

Source

CRAN

Examples

```
data(check_times)
str(check_times)
```

Chicago

Chicago ridership data

Description

Chicago ridership data

Details

These data are from Kuhn and Johnson (2020) and contain an *abbreviated* training set for modeling the number of people (in thousands) who enter the Clark and Lake L station.

The date column corresponds to the current date. The columns with station names (Austin through California) are a *sample* of the columns used in the original analysis (for file size reasons). These are 14 day lag variables (i.e. `date - 14 days`). There are columns related to weather and sports team schedules.

The station at 35th and Archer is contained in the column `Archer_35th` to make it a valid R column name.

Value

`Chicago` a tibble
`stations` a vector of station names

Source

Kuhn and Johnson (2020), *Feature Engineering and Selection*, Chapman and Hall/CRC . <https://bookdown.org/max/FES/> and <https://github.com/topepo/FES>

Examples

```
data(Chicago)
str(Chicago)
stations
```

concrete

Compressive strength of concrete mixtures

Description

Yeh (2006) describes an aggregated data set for experimental designs used to test the compressive strength of concrete mixtures. The data are used by Kuhn and Johnson (2013).

Value

concrete a tibble

Source

Yeh I (2006). "Analysis of Strength of Concrete Using Design of Experiments and Neural Networks." *Journal of Materials in Civil Engineering*, 18, 597-604.

Kuhn, M., Johnson, K. (2013) *Applied Predictive Modeling*, Springer.

Examples

```
data(concrete)
```

covers

Raw cover type data

Description

These data are raw data describing different types of forest cover-types from the UCI Machine Learning Database (see link below). There is one column in the data that has a few difference pieces of textual information (of variable lengths).

Value

covers a data frame

Source

<https://archive.ics.uci.edu/ml/machine-learning-databases/covtype/covtype.info>

Examples

```
data(covers)
str(covers)
```

credit_data	<i>Credit data</i>
-------------	--------------------

Description

These data are from the website of Dr. Lluís A. Belanche Muñoz by way of a github repository of Dr. Gaston Sanchez. One data point is a missing outcome was removed from the original data.

Value

credit_data a data frame

Source

<https://github.com/gastonstat/CreditScoring>, <http://bit.ly/2kkBFrk>

Examples

```
data(credit_data)
str(credit_data)
```

drinks	<i>Sample time series data</i>
--------	--------------------------------

Description

Sample time series data

Details

Drink sales. The exact name of the series from FRED is: "Merchant Wholesalers, Except Manufacturers' Sales Branches and Offices Sales: Nondurable Goods: Beer, Wine, and Distilled Alcoholic Beverages Sales"

Value

drinks a tibble

Source

The Federal Reserve Bank of St. Louis website <https://fred.stlouisfed.org/series/S4248SM144NCEN>

Examples

```
data(drinks)
str(drinks)
```

hpc_cv	<i>Class probability predictions</i>
--------	--------------------------------------

Description

Class probability predictions

Details

This data frame contains the predicted classes and class probabilities for a linear discriminant analysis model fit to the HPC data set from Kuhn and Johnson (2013). These data are the assessment sets from a 10-fold cross-validation scheme. The data column columns for the true class (obs), the class prediction (pred) and columns for each class probability (columns VF, F, M, and L). Additionally, a column for the resample indicator is included.

Value

hpc_cv	a data frame
--------	--------------

Source

Kuhn, M., Johnson, K. (2013) *Applied Predictive Modeling*, Springer

Examples

```
data(hpc_cv)
str(hpc_cv)
```

hpc_data	<i>High-performance computing system data</i>
----------	---

Description

Kuhn and Johnson (2013) describe a data set where characteristics of unix jobs were used to classify their completion times as either very fast (1 min or less, VF), fast (1–50 min, F), moderate (5–30 min, M), or long (greater than 30 min, L).

Value

hpc_data	a tibble
----------	----------

Source

Kuhn, M., Johnson, K. (2013) *Applied Predictive Modeling*, Springer.

Examples

```
data(hpc_data)
```

lending_club	<i>Loan data</i>
--------------	------------------

Description

Loan data

Details

These data were downloaded from the Lending Club access site (see below) and are from the first quarter of 2016. A subset of the rows and variables are included here. The outcome is in the variable `Class` and is either "good" (meaning that the loan was fully paid back or currently on-time) or "bad" (charged off, defaulted, of 21-120 days late). A data dictionary can be found on the source website.

Value

lending_club a data frame

Source

Lending Club Statistics <https://www.lendingclub.com/info/download-data.action>

Examples

```
data(lending_club)
str(lending_club)
```

meats	<i>Fat, water and protein content of meat samples</i>
-------	---

Description

"These data are recorded on a Tecator Infratec Food and Feed Analyzer working in the wavelength range 850 - 1050 nm by the Near Infrared Transmission (NIT) principle. Each sample contains finely chopped pure meat with different moisture, fat and protein contents.

Details

If results from these data are used in a publication we want you to mention the instrument and company name (Tecator) in the publication. In addition, please send a preprint of your article to Karin Thente, Tecator AB, Box 70, S-263 21 Hoganas, Sweden

The data are available in the public domain with no responsibility from the original data source. The data can be redistributed as long as this permission note is attached."

"For each meat sample the data consists of a 100 channel spectrum of absorbances and the contents of moisture (water), fat and protein. The absorbance is $-\log_{10}$ of the transmittance measured by the spectrometer. The three contents, measured in percent, are determined by analytic chemistry."

Included here are the training, monitoring and test sets.

Value

meats	a tibble
-------	----------

Examples

```
data(meats)
```

mlc_churn	<i>Customer churn data</i>
-----------	----------------------------

Description

A data set from the MLC++ machine learning software for modeling customer churn. There are 19 predictors, mostly numeric: state (categorical), account_length area_code international_plan (yes/no), voice_mail_plan (yes/no), number_vmail_messages total_day_minutes total_day_calls total_day_charge total_eve_minutes total_eve_calls total_eve_charge total_night_minutes total_night_calls total_night_charge total_intl_minutes total_intl_calls total_intl_charge, and number_customer_service_calls.

Details

The outcome is contained in a column called churn (also yes/no). A note in one of the source files states that the data are "artificial based on claims similar to real world".

Value

mlc_churn a tibble

Source

<http://www.sgi.com/tech/mlc/>

oils

Fatty acid composition of commercial oils

Description

Fatty acid concentrations of commercial oils were measured using gas chromatography. The data is used to predict the type of oil. Note that only the known oils are in the data set. Also, the authors state that there are 95 samples of known oils. However, we count 96 in Table 1 (pgs. 33-35).

Value

oils a tibble

Source

Brodnjak-Voncina et al. (2005). Multivariate data analysis in classification of vegetable oils characterized by the content of fatty acids, *Chemometrics and Intelligent Laboratory Systems*, Vol. 75:31-45.

okc

OkCupid data

Description

These are a sample of columns of users of OkCupid dating website. The data are from Kim and Escobedo-Land (2015). Permission to use this data set was explicitly granted by OkCupid.

Value

okc a data frame

Source

Kim, A. Y., and A. Escobedo-Land. 2015. "OkCupid Data for Introductory Statistics and Data Science Courses." *Journal of Statistics Education: An International Journal on the Teaching and Learning of Statistics*.

Examples

```
data(okc)
str(okc)
```

okc_text	<i>OkCupid text data</i>
----------	--------------------------

Description

These are a sample of columns and users of OkCupid dating website. The data are from Kim and Escobedo-Land (2015). Permission to use this data set was explicitly granted by OkCupid. The data set contains 10 text fields filled out by users.

Value

okc_text a tibble

Source

Kim, A. Y., and A. Escobedo-Land. 2015. "OkCupid Data for Introductory Statistics and Data Science Courses." *Journal of Statistics Education: An International Journal on the Teaching and Learning of Statistics*.

Examples

```
data(okc_text)
str(okc_text)
```

parabolic	<i>Parabolic class boundary data</i>
-----------	--------------------------------------

Description

Parabolic class boundary data

Details

These data were simulated. There are two correlated predictors and two classes in the factor outcome.

Value

parabolic a data frame

Examples

```
data(parabolic)
```

pathology	<i>Liver pathology data</i>
-----------	-----------------------------

Description

Liver pathology data

Details

These data have the results of a *x*-ray examination to determine whether liver is abnormal or not (in the scan column) versus the more extensive pathology results that approximate the truth (in pathology).

Value

pathology a data frame

Source

Altman, D.G., Bland, J.M. (1994) "Diagnostic tests 1: sensitivity and specificity," *British Medical Journal*, vol 308, 1552.

Examples

```
data(pathology)
str(pathology)
```

pd_speech	<i>Parkinson's disease speech classification data set</i>
-----------	---

Description

Parkinson's disease speech classification data set

Details

From the UCI ML archive, the description is "The data used in this study were gathered from 188 patients with PD (107 men and 81 women) with ages ranging from 33 to 87 (65.1 p/m 10.9) at the Department of Neurology in Cerrahpaşa Faculty of Medicine, Istanbul University. The control group consists of 64 healthy individuals (23 men and 41 women) with ages varying between 41 and 82 (61.1 p/m 8.9). During the data collection process, the microphone is set to 44.1 KHz and following the physician's examination, the sustained phonation of the vowel /a/ was collected from each subject with three repetitions."

The data here are averaged over the replicates.

Value

pd_speech a data frame

Source

UCI ML repository (data) <https://archive.ics.uci.edu/ml/datasets/Parkinson%27s+Disease+Classification#>, Sakar et al (2019), "A comparative analysis of speech signal processing algorithms for Parkinson's disease classification and the use of the tunable Q-factor wavelet transform", *Applied Soft Computing*, V74, pg 255-263.

Examples

```
data(pd_speech)
str(pd_speech)
```

Sacramento

Sacramento CA home prices

Description

This data frame contains house and sale price data for 932 homes in Sacramento CA. The original data were obtained from the website for the SpatialKey software. From their website: "The Sacramento real estate transactions file is a list of 985 real estate transactions in the Sacramento area reported over a five-day period, as reported by the Sacramento Bee." Google was used to fill in missing/incorrect data.

Value

Sacramento a tibble

Source

SpatialKey website: <https://support.spatialkey.com/spatialkey-sample-csv-data>

Examples

```
data(Sacramento)
```

scat	<i>Morphometric data on scat</i>
------	----------------------------------

Description

Reid (2015) collected data on animal feces in coastal California. The data consist of DNA verified species designations as well as fields related to the time and place of the collection and the scat itself. The data are on the three main species.

Value

scat	a tibble
------	----------

Source

Reid, R. E. B. (2015). A morphometric modeling approach to distinguishing among bobcat, coyote and gray fox scats. *Wildlife Biology*, 21(5), 254-262

small_fine_foods	<i>Fine foods example data</i>
------------------	--------------------------------

Description

Fine foods example data

Details

These data are from Amazon, who describe it as "This dataset consists of reviews of fine foods from amazon. The data span a period of more than 10 years, including all ~500,000 reviews up to October 2012. Reviews include product and user information, ratings, and a plaintext review."

A subset of the data are contained here and are split into a training and test set. The training set sampled 10 products and retained all of their individual reviews. Since the reviews within these products are correlated, we recommend resampling the data using a leave-one-product-out approach. The test set sampled 500 products that were not included in the training set and selected a single review at random for each.

There is a column for the product, a column for the text of the review, and a factor column for a class variable. The outcome is whether the reviewer gave the product a 5-star rating or not.

Value

training_data, testing_data
tibbles

Source

<https://snap.stanford.edu/data/web-FineFoods.html>

Examples

```
data(small_fine_foods)
str(training_data)
```

Smithsonian	<i>Smithsonian museums</i>
-------------	----------------------------

Description

Geocodes for the Smithsonian museums (circa 2018).

Value

Smithsonian a tibble

Source

https://en.wikipedia.org/wiki/List_of_Smithsonian_museums

Examples

```
data(Smithsonian)
Smithsonian
```

solubility_test	<i>Solubility predictions from MARS model</i>
-----------------	---

Description

Solubility predictions from MARS model

Details

For the solubility data in Kuhn and Johnson (2013), these data are the test set results for the MARS model. The observed solubility (in column solubility) and the model results (prediction) are contained in the data.

Value

solubility_test
a data frame

Source

Kuhn, M., Johnson, K. (2013) *Applied Predictive Modeling*, Springer

Examples

```
data(solubility_test)
str(solubility_test)
```

stackoverflow	<i>Annual Stack Overflow Developer Survey Data</i>
---------------	--

Description

Annual Stack Overflow Developer Survey Data

Details

These data are a collection of 5,594 data points collected on developers. These data could be used to try to predict who works remotely (as used in the source listed below).

Value

stackoverflow a tibble

Source

Julia Silge, *Supervised Machine Learning Case Studies in R*
<https://supervised-ml-course.netlify.com/chapter2>
Raw data: <https://insights.stackoverflow.com/survey/>

Examples

```
data(stackoverflow)
```

two_class_dat	<i>Two class data</i>
---------------	-----------------------

Description

Two class data

Details

There are artificial data with two predictors (A and B) and a factor outcome variable (Class).

Value

two_class_dat a data frame

Examples

```
data(two_class_dat)
str(two_class_dat)
```

two_class_example	<i>Two class predictions</i>
-------------------	------------------------------

Description

Two class predictions

Details

These data are a test set form a model built for two classes ("Class1" and "Class2"). There are columns for the true and predicted classes and column for the probabilities for each class.

Value

```
two_class_example
      a data frame
```

Examples

```
data(two_class_example)
str(two_class_example)
```

wa_churn	<i>Watson churn data</i>
----------	--------------------------

Description

Watson churn data

Details

These data were downloaded from the IBM Watson site (see below) in September 2018. The data contain a factor for whether a customer churned or not. Alternatively, the tenure column presumably contains information on how long the customer has had an account. A survival analysis can be done on this column using the churn outcome as the censoring information. A data dictionary can be found on the source website.

Value

```
wa_churn      a data frame
```

wa_churn

21

Source

IBM Watson Analytics <https://ibm.co/2sOvyvy>

Examples

```
data(wa_churn)
str(wa_churn)
```

Index

*Topic **datasets**

- ad_data, 2
- attrition, 3
- biomass, 4
- bivariate, 4
- car_prices, 5
- cells, 5
- check_times, 6
- Chicago, 7
- concrete, 8
- covers, 8
- credit_data, 9
- drinks, 9
- hpc_cv, 10
- hpc_data, 10
- lending_club, 11
- meats, 12
- mlc_churn, 12
- oils, 13
- okc, 13
- okc_text, 14
- parabolic, 14
- pathology, 15
- pd_speech, 15
- Sacramento, 16
- scat, 17
- small_fine_foods, 17
- Smithsonian, 18
- solubility_test, 18
- stackoverflow, 19
- two_class_dat, 19
- two_class_example, 20
- wa_churn, 20

- ad_data, 2
- attrition, 3

- biomass, 4
- bivariate, 4
- bivariate_test (bivariate), 4

- bivariate_train (bivariate), 4
- bivariate_val (bivariate), 4

- car_prices, 5
- cells, 5
- check_times, 6
- Chicago, 7
- concrete, 8
- covers, 8
- credit_data, 9

- drinks, 9

- hpc_cv, 10
- hpc_data, 10

- lending_club, 11

- meats, 12
- mlc_churn, 12

- oils, 13
- okc, 13
- okc_text, 14

- parabolic, 14
- pathology, 15
- pd_speech, 15

- Sacramento, 16
- scat, 17
- small_fine_foods, 17
- Smithsonian, 18
- solubility_test, 18
- stackoverflow, 19
- stations (Chicago), 7

- testing_data (small_fine_foods), 17
- training_data (small_fine_foods), 17
- two_class_dat, 19
- two_class_example, 20

- wa_churn, 20