

Package ‘dslabs’

May 8, 2026

Title Data Science Labs

Version 0.9.1

Description Datasets and functions that can be used for data analysis practice, homework and projects in data science courses and workshops. 26 datasets are available for case studies in data visualization, statistical inference, modeling, linear regression, data wrangling and machine learning.

Depends R (>= 4.1.0)

Imports ggplot2

License Artistic-2.0

Encoding UTF-8

LazyData true

RoxygenNote 7.3.3

NeedsCompilation no

Author Rafael A. Irizarry [aut, cre],
Amy Gill [aut]

Maintainer Rafael A. Irizarry <rafael_irizarry@dfci.harvard.edu>

Repository CRAN

Date/Publication 2025-12-02 06:10:32 UTC

Contents

admissions	2
brca	3
brexit_polls	4
death_prob	5
divorce_margarine	5
ds_theme_set	6
fit_recommender_model	7
gapminder	9
greenhouse_gases	10
heights	11

historic_co2	11
mice_weights	12
mnist_127	13
mnist_27	14
movielens	15
murders	16
na_example	16
nyc_regents_scores	17
olive	18
outlier_example	19
polls_2008	19
polls_us_election_2016	20
pr_death_counts	21
read_mnist	22
reported_heights	23
research_funding_rates	24
results_us_election_2012	25
rfalling_object	26
stars	27
take_poll	27
temp_carbon	28
tissue_gene_expression	29
trump_tweets	29
us_contagious_diseases	30

Index 32

admissions	<i>Gender bias among graduate school admissions to UC Berkeley.</i>
------------	---

Description

The admission data for six majors for the fall of 1973; often used as an example of Simpson's paradox

Usage

admissions

Format

An object of class "data.frame".

Details

- major. The major or university department.
- gender. Men or women.
- admitted. Percent of students admitted.
- applicants. Total number of applicants.

Source

PJ Bickel, EA Hammel, and JW O'Connell. Science (1975)

Examples

admissions

brca	<i>Breast Cancer Wisconsin Diagnostic Dataset from UCI Machine Learning Repository</i>
------	--

Description

Biopsy features for classification of 569 malignant (cancer) and benign (not cancer) breast masses.

Usage

brca

Format

An object of class `list`.

Details

Features were computationally extracted from digital images of fine needle aspirate biopsy slides. Features correspond to properties of cell nuclei, such as size, shape and regularity. The mean, standard error, and worst value of each of 10 nuclear parameters is reported for a total of 30 features.

This is a classic dataset for training and benchmarking machine learning algorithms.

- `y`. The outcomes. A factor with two levels denoting whether a mass is malignant ("M") or benign ("B").
- `x`. The predictors. A matrix with the mean, standard error and worst value of each of 10 nuclear measurements on the slide, for 30 total features per biopsy:
 - `radius`. Nucleus radius (mean of distances from center to points on perimeter).
 - `texture`. Nucleus texture (standard deviation of grayscale values).
 - `perimeter`. Nucleus perimeter.
 - `area`. Nucleus area.
 - `smoothness`. Nucleus smoothness (local variation in radius lengths).
 - `compactness`. Nucleus compactness ($\text{perimeter}^2/\text{area} - 1$).
 - `concavity`. Nucleus concavity (severity of concave portions of the contour).
 - `concave_pts`. Number of concave portions of the nucleus contour.
 - `symmetry`. Nucleus symmetry.
 - `fractal_dim`. Nucleus fractal dimension ("coastline approximation" -1).

Source

[UCI Machine Learning Repository](#)

Examples

```
table(brca$y)
dim(brca$x)
head(brca$x)
```

brexit_polls

Brexit Poll Data

Description

Brexit (EU referendum) poll outcomes for 127 polls from January 2016 to the referendum date on June 23, 2016.

Usage

```
brexit_polls
```

Format

An object of class "data.frame".

Details

- startdate. Start date of poll.
- enddate. End date of poll.
- pollster. Pollster conducting the poll.
- poll_type. Online or telephone poll.
- samplesize. Sample size of poll.
- remain. Proportion voting Remain.
- leave. Proportion voting Leave.
- undecided. Proportion of undecided voters.
- spread. Spread calculated as remain - leave.

Source

[Wikipedia](#)

Examples

```
head(brexit_polls)
```

death_prob	<i>2015 US Period Life Table</i>
------------	----------------------------------

Description

Probability of death within 1 year by age and sex in the United States in 2015.

Usage

```
death_prob
```

Format

An object of class "data.frame".

Details

- age. Age strata, with each year a different stratum.
- sex. Male or Female.
- prob. Probability of death within 1 year given exact age and sex.

Source

[Social Security Administraton](#)

Examples

```
head(death_prob)
```

divorce_margarine	<i>Divorce rate and margarine consumption data</i>
-------------------	--

Description

Divorce rates in Maine and per capita consumption of margarine in US data

Usage

```
divorce_margarine
```

Format

An object of class "data.frame".

Details

- `divorce_rate_maine`. Divorce per 1000 in Maine.
- `margarine_consumption_per_capita`. US per capita consumption of margarine in pounds.
- `year`. Year.

Source

Spurious Correlations

Examples

```
with(divorce_margarine, plot(margarine_consumption_per_capita, divorce_rate_maine))
```

`ds_theme_set`

dslabs theme set

Description

This function sets a ggplot2 theme used throughout the data science labs. It can be called without arguments.

Usage

```
ds_theme_set(
  new = "theme_bw",
  args = NULL,
  base_size = 11,
  bold_title = TRUE,
  ...
)
```

Arguments

<code>new</code>	a prebuilt ggplot2 theme. Defaults to "theme_minimal"
<code>args</code>	the arguments to be passed along to the ggplot2 theme function. Defaults to "NULL".
<code>base_size</code>	if "args" is "NULL", <code>base_size</code> is one of the arguments passed to the theme function. It defaults to 11.
<code>bold_title</code>	if TRUE, sets titles to be bold
<code>...</code>	additional arguments to be used by theme

Value

None

Examples

```
library(ggplot2)
ds_theme_set()
qplot(hp, mpg, data=mtcars, color=am, facets=gear~cyl,
main="Scatterplots of MPG vs. Horsepower",
xlab="Horsepower", ylab="Miles per Gallon")
```

fit_recommender_model *Fit a Latent Factor Recommender Model via Alternating Least Squares*

Description

This function fits a penalized latent factor model for recommendation systems using an alternating least squares (ALS) algorithm. It estimates user effects, item effects, and latent factors simultaneously, with regularization to prevent overfitting. The implementation supports filtering out items with too few ratings.

Usage

```
fit_recommender_model(
  rating,
  user_id,
  item_id,
  K = 8,
  lambda_1 = 5e-05,
  lambda_2 = 1e-04,
  min_ratings = 20,
  maxit = 500,
  reltol = 1e-08,
  damping = 0.75,
  verbose = FALSE
)
```

Arguments

rating	A numeric vector of observed ratings.
user_id	A character vector identifying the user for each rating. Must be the same length as 'rating'.
item_id	A character vector identifying the item for each rating. Must be the same length as 'rating'.
K	Integer. The number of latent factors to estimate.
lambda_1	Numeric. Regularization parameter for user and item effects.
lambda_2	Numeric. Regularization parameter for latent factors.

min_ratings	Integer. Minimum number of ratings required for an item to be included in the estimation of latent factors.
maxit	Integer. Maximum number of iterations.
reltol	Numeric. Relative reletolerance for convergence, based on change in the objective function.
damping	Numeric between 0 and 1. Damping factor used to blend updates with the previous iteration for convergence stability.
verbose	Logical. If 'TRUE', prints progress messages during optimization.

Details

The model being fit is:

$$Y_{ij} = \mu + \alpha_i + \beta_j + \sum_{k=1}^K p_{ik}q_{jk} + \varepsilon_{ij}$$

where μ is the global mean, α_i are user effects, β_j are item effects, and the p_{ik} , q_{jk} are latent factors for users and items respectively. The estimation minimizes a penalized least squares criterion with separate penalties for user/item effects and latent factors.

Items with less than min_ratings observations are excluded from the estimation of 'p' and 'q'.

Value

A list with the following components:

mu	Global mean rating.
a	Named numeric vector of user-specific effects.
b	Named numeric vector of item-specific effects.
p	Matrix of user latent factors, with one row per user. The rownames of this matrix match the names of 'a'.
q	Matrix of item latent factors, with one row per item. The rownames of this matrix match the names of 'b'.
min_ratings	The threshold value used to filter items: only items with at least this many ratings are included in the estimation of latent factors.
n_item	Named integer vector of number of ratings per item.
n_user	Named integer vector of number of retained ratings per user.
fitted	Fitted values.

Examples

```
set.seed(2010)
## Simulation settings
n_users <- 200      # number of users
n_items <- 300     # number of items
K_true <- 4        # true number of latent factors
sparsity <- 0.25   # ~5% of user-item pairs are observed
```

```

## True parameters
mu <- 3.5
a_true <- rnorm(n_users, 0, 0.3)    # user effects
b_true <- rnorm(n_items, 0, 0.4)    # item effects
p_true <- matrix(rnorm(n_users * K_true, 0, 0.5), n_users, K_true)
q_true <- matrix(rnorm(n_items * K_true, 0, 0.5), n_items, K_true)

names(a_true) <- 1:n_users
names(b_true) <- 1:n_items
rownames(p_true) <- 1:n_users
rownames(q_true) <- 1:n_items
## Generate observed ratings matrix with sparsity
user_id <- rep(as.character(1:n_users), each = n_items)
item_id <- rep(as.character(1:n_items), times = n_users)

## Which entries are observed?
obs <- runif(length(user_id)) < sparsity

## Ratings with noise
rating_full <- mu + a_true[user_id] + b_true[item_id] +
  rowSums(p_true[user_id, ] * q_true[item_id, ]) +
  rnorm(length(user_id), 0, 0.25)

rating <- rating_full[obs]
user_id <- user_id[obs]
item_id <- item_id[obs]

## Call your recommender function
fit <- fit_recommender_model(rating, user_id, item_id, K = 4, reltol = 1e-5,
                             min_ratings = 5, verbose = TRUE)
plot(fit$fitted, rating)

```

gapminder

Gapminder Data

Description

Health and income outcomes for 184 countries from 1960 to 2016. Also includes two character vectors, `oecd` and `opec`, with the names of OECD and OPEC countries from 2016.

Usage

```
gapminder
```

Format

An object of class "data.frame".

Details

- country.
- year.
- infant_mortality. Infant deaths per 1000.
- life_expectancy. Life expectancy in years.
- fertility. Average number of children per woman.
- population. Country population.
- gpd. GDP according to World Bankdev.
- continent.
- region. Geographical region.

Examples

```
head(gapminder)
print(oecd)
print(opeca)
```

greenhouse_gases

Greenhouse gas concentrations over 2000 years

Description

Concentrations of the three main greenhouse gases carbon dioxide, methane and nitrous oxide. Measurements are from the Law Dome Ice Core in Antarctica. Selected measurements are provided every 20 years from 1-2000 CE.

Usage

```
greenhouse_gases
```

Format

An object of class "data.frame".

Details

- year. Year (CE).
- gas. Gas being measured: carbon dioxide ('CO2'), methane ('CH4') or nitrous oxide ('N2O').
- concentration. Gas concentration in ppm by volume ('CO2') or ppb by volume ('CH4', 'N2O').

Source

MacFarling Meure et al. 2006 via [NOAA](#).

Examples

```
head(greenhouse_gases)
```

heights	<i>Self-reported Heights in Inches</i>
---------	--

Description

Self-reported heights and sex. The heights were converted to inches from the original data included in [reported_heights](#).

Usage

```
heights
```

Format

An object of class "data.frame".

Details

- sex. A factor with the self-reported sex.
- height. A numeric vector with self-reported heights in inches.

See Also

[reported_heights](#) for the original data source.

Examples

```
head(heights)
```

historic_co2	<i>Atmospheric carbon dioxide concentration over 800,000 years</i>
--------------	--

Description

Concentration of carbon dioxide in ppm by volume from direct measurements at Mauna Loa (1959-2018 CE) and indirect measurements from a series of Antarctic ice cores (approx. -800,000-2001 CE).

Usage

```
historic_co2
```

Format

An object of class "data.frame".

Details

- year. Year (CE).
- co2. Carbon dioxide concentration in ppm by volume.
- source. Source of carbon dioxide measurement: direct CO2 annual mean concentrations from Mauna Loa ('Mauna Loa') or indirect CO2 concentrations from air trapped in ice cores ('Ice Cores').

Source

Mauna Loa data from [NOAA](#). Ice core data from Bereiter et al. 2015 via [NOAA](#).

Examples

```
head(historic_co2)
```

mice_weights

Mice weights

Description

Body weights, bone density, and percent fat for mice under two diets: chow and high fat. Data provided by Karen Svenson from Jackson Laboratories. Funding to generate these data came from NIH grant P50 GM070683 awarded to Gary Churchill.

Usage

```
mice_weights
```

Format

An object of class "data.frame".

Details

- body_weight. Body weight in grams at 19 weeks.
- bone_density. Body density.
- percent_fat. Percent fat.
- sex. The sex of the mice.
- diet. The diet of the mice: chow or high fat.
- gen. These are outbred mice. This variable denotes the generation.
- litter. Which of two litters mice belong to.

Source

Karen Svenson, Daniel M. Gatti, and Gary Churchill from Jackson Laboratories.

References

Daniel M. Gatti, Petr Simecek, Lisa Somes, Clifton T. Jeffrey, Matthew J. Vincent, Kwangbom Choi, Xingyao Chen, Gary A. Churchill, and Karen L. Svenson. "The Effects of Sex and Diet on Physiology and Liver Gene Expression in Diversity Outbred Mice". bioRxiv 098657; doi:[10.1101/098657](https://doi.org/10.1101/098657)

Examples

```
mice_weights |> head()
with(mice_weights, table(sex, diet))
```

mnist_127

Useful example for illustrating generative models based on MNIST data

Description

A randomly selected set of 1s, 2s and 7s along with the two predictors based on the proportion of dark pixels in the upper left and lower right quadrants respectively. The dataset is divided into training and test sets.

Usage

```
mnist_127
```

Format

An object of class `list`.

Details

- `train`. A data frame containing training data: labels and predictors.
- `test`. A data frame containing test data: labels and predictors.

References

Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. "Gradient-based learning applied to document recognition." Proceedings of the IEEE, 86(11):2278-2324, November 1998.

See Also

[`read_mnist()`, `mnist_27`]

Examples

```
with(mnist_127$train, plot(x_1, x_2, col = as.numeric(y)))
```

mnist_27

Useful example for illustrating machine learning algorithms based on MNIST data

Description

A randomly selected set of 2s and 7s along with the two predictors based on the proportion of dark pixels in the upper left and lower right quadrants respectively. The dataset is divided into training and test sets.

Usage

```
mnist_27
```

Format

An object of class `list`.

Details

- `train`. A data frame containing training data: labels and predictors.
- `test`. A data frame containing test data: labels and predictors.
- `index_train`. The index of the original mnist training data used for the training set.
- `index_test`. The index of the original mnist test data used for the test set.
- `true_p`. A data frame containing the two predictors `x_1` and `x_2` and the conditional probability of being a 7 for `x_1`, `x_2`.

References

Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. "Gradient-based learning applied to document recognition." *Proceedings of the IEEE*, 86(11):2278-2324, November 1998.

See Also

[`read_mnist()`]

Examples

```
with(mnist_27$train, plot(x_1, x_2, col = as.numeric(y)))
```

`movielens`*Movie ratings*

Description

MovieLens Latest Dataset (Small)

Usage

```
movielens
```

Format

Two object of class `data.frame`.

Details

- `movieId`. Unique ID for the movie.
- `title`. Movie title (not unique).
- `year`. Year the movie was released.
- `genres`. Genres associated with the movie.
- `userId`. Unique ID for the user.
- `rating`. A rating between 0 and 5 for the movie.
- `timestamp`. Date and time the rating was given.

Source

<https://files.grouplens.org/datasets/movielens/ml-latest-small.zip>

References

F. Maxwell Harper and Joseph A. Konstan. 2015. The MovieLens Datasets: History and Context. *ACM Transactions on Interactive Intelligent Systems (TiiS)* 5, 4, Article 19 (December 2015), 19 pages. DOI=<https://dx.doi.org/10.1145/2827872>

Examples

```
head(movielens)
```

murders	<i>US gun murders by state for 2010</i>
---------	---

Description

Gun murder data from FBI reports. Also contains the population of each state.

Usage

```
murders
```

Format

An object of class "data.frame".

Details

- state. US state
- abb. Abbreviation of US state
- region. Geographical US region
- population. State population (2010)
- total. Number of gun murders in state (2010)

Source

[Wikipedia](#)

Examples

```
print(murders)
```

na_example	<i>Count data with some missing values</i>
------------	--

Description

This dataset was randomly generated.

Usage

```
na_example
```

Format

An object of class "integer".

Examples

```
print(sum(is.na(na_example)))
```

nyc_regents_scores *NYC Regents exams scores 2010*

Description

Distribution of scores for New York City Regents algebra, global history, biology, English, and U.S. history exams. These data were used to make [this New York Times plot](#).

Usage

```
nyc_regents_scores
```

Format

An object of class "data.frame".

Details

- score. Test score from 0 to 100.
- integrated_algebra. Score frequency on Algebra exam.
- global_history. Score frequency on global history exam.
- living_environment. Score frequency on biology exam.
- english. Score frequency on English exam.
- us_history. Score frequency on U.S. history exam.

Source

New York City Department of Education via Amanda Cox.

Examples

```
print(nyc_regents_scores)
```

olive

Italian olive

Description

Composition in percentage of eight fatty acids found in the lipid fraction of 572 Italian olive oils

Usage

olive

Format

An object of class "data.frame".

Details

- region. General region of Italy.
- area. Area of Italy.
- palmitic. Percent palmitic acid of sample.
- palmitoleic. Percent palmitoleic of sample.
- stearic. Percent stearic acid of sample.
- oleic. Percent oleic acid of sample.
- linoleic. Percent linoleic acid of sample.
- linolenic. Percent linolenic acid of sample.
- arachidic. Percent arachidic acid of sample.
- eicosenoic. Percent eicosenoic acid of sample.

Source

J. Zupan, and J. Gasteiger. Neural Networks in Chemistry and Drug Design.

Examples

```
head(olive)
```

outlier_example	<i>Adult male heights in feet with outliers</i>
-----------------	---

Description

This dataset was randomly generated with a normal distribution (average: 5 feet 9 inches, standard deviation: 3 inches). One value was changed to be mistakenly reported in centimeters rather than feet.

Usage

```
outlier_example
```

Format

An object of class "numeric".

Examples

```
mean(outlier_example)
median(outlier_example)
```

polls_2008	<i>Poll data for popular vote in 2008 presidential election</i>
------------	---

Description

Data from different pollsters for the popular vote between Obama and McCain in the 2008 presidential election.

Usage

```
polls_2008
```

Format

An object of class `data.frame`.

Details

- `day`. Days until election day. Negative numbers are reported so that days can increase up to 0, which is election day.
- `margin`. Average difference between Obama and McCain for that day.

Source

<https://web.archive.org/web/20161108190914/http://www.pollster.com/08USPresGEMvO-2.html>

Examples

```
with(polls_2008, plot(day, margin))
```

```
polls_us_election_2016
```

Fivethirtyeight 2016 Poll Data

Description

Poll results from US 2016 presidential elections aggregated from HuffPost Pollster, RealClearPolitics, polling firms, and news reports. The dataset also includes election results (popular vote) and electoral college votes in `results_us_election_2016`.

Usage

```
polls_us_election_2016
```

Format

An object of class "data.frame".

Details

- `state`. State in which poll was taken. 'U.S' is for national polls.
- `startdate`. Poll's start date.
- `enddate`. Poll's end date.
- `pollster`. Pollster conducting the poll.
- `grade`. Grade assigned by fivethirtyeight to pollster.
- `samplesize`. Sample size.
- `population`. Type of population being polled.
- `rawpoll_clinton`. Percentage for Hillary Clinton.
- `rawpoll_trump`. Percentage for Donald Trump
- `rawpoll_johnson`. Percentage for Gary Johnson
- `rawpoll_mcmullin`. Percentage for Evan McMullin.
- `adjpoll_clinton`. Fivethirtyeight adjusted percentage for Hillary Clinton.
- `ajdpoll_trump`. Fivethirtyeight adjusted percentage for Donald Trump
- `adjpoll_johnson`. Fivethirtyeight adjusted percentage for Gary Johnson
- `adjpoll_mcmullin`. Fivethirtyeight adjusted percentage for Evan McMullin.

Source

The original csv file used to create `polls_us_election_2016` is here: https://projects.fivethirtyeight.com/general-model/president_general_polls_2016.csv

The data for `results_us_election_2016` is from Ballotpedia and can be found here: https://docs.google.com/spreadsheets/d/1zxyOQDjN0JS_UkzerorUCf20AdcMcIQEwRciKuYBIZ4/pubhtml?widget=true&headers=false#gid=658726802/

Examples

```
head(polls_us_election_2016)
```

pr_death_counts	<i>Puerto Rico daily mortality</i>
-----------------	------------------------------------

Description

A data frame with Puerto Rico daily mortality counts 2015 to May 2018. This includes the day hurricanes Maria made 2017-09-20.

Usage

```
pr_death_counts
```

Format

An object of class `data.frame`.

Details

- `date`. Date of the count.
- `deaths`. Number of deaths reported that day.

Source

Puerto Rico Demographic Registry. Data was extracted from PDF provided in `'system.file("extdata", "RD-Mortality-Report_2015-18-180531.pdf", package = "dslabs")'`

Examples

```
with(pr_death_counts, plot(date, deaths))
```

`read_mnist`*Download and read the mnist dataset*

Description

This function downloads the mnist training and test data available here <http://yann.lecun.com/exdb/mnist/>

Usage

```
read_mnist(  
  path = NULL,  
  download = FALSE,  
  destdir = tempdir(),  
  url = "https://www2.harvardx.harvard.edu/courses/IDS_08_v2_03/",  
  keep.files = TRUE  
)
```

Arguments

<code>path</code>	A character giving the full path of the directory to look for files. It assumes the filenames are the same as the originals. If <code>path</code> is <code>NULL</code> a download or direct read of the files is attempted.
<code>download</code>	If <code>TRUE</code> the files will be downloaded and saved in <code>destdir</code> .
<code>destdir</code>	A character giving the full path of the directory in which to save the downloaded files. The default is to use a temporary directory.
<code>url</code>	A character giving the URL from which to download files. Currently a copy of the data is available at https://www2.harvardx.harvard.edu/courses/IDS_08_v2_03/ , the current default URL.
<code>keep.files</code>	A logical. If <code>TRUE</code> the downloaded files will be saved in <code>destdir</code> . If <code>FALSE</code> the entire directory is erased. This argument is ignored if <code>download</code> is <code>FALSE</code> .

Value

A list with two components: `train` and `test`. Each of these is a list with two components: `images` and `labels`. The `images` component is a matrix with each column representing one of the $28 \times 28 = 784$ pixels. The values are integers between 0 and 255 representing grey scale. The `labels` components is a vector representing the digit shown in the image.

Note that the data is over 10MB, so the download may take several seconds depending on internet speed. If you plan to load the data more than once we recommend you download the data once and read it from disk in the future. See examples.

Author(s)

Samuela Pollack

Rafael A. Irizarry, <rafael_irizarry@dfci.harvard.edu>

References

Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. "Gradient-based learning applied to document recognition." Proceedings of the IEEE, 86(11):2278-2324, November 1998.

Examples

```
# this can take several seconds, depending on internet speed.

## Not run:
mnist <- read_mnist()
i <- 5
image(1:28, 1:28, matrix(mnist$test$images[i,], nrow=28)[ , 28:1],
      col = gray(seq(0, 1, 0.05)), xlab = "", ylab="")
## the labels for this image is:
mnist$test$labels[i]

## End(Not run)

# You can download and save the data to a directory like this:
## Not run:
mnist <- read_mnist(download = TRUE, destdir = "~/Downloads")

# and then, going forward, read from disk
mnist <- read_mnist("~/Downloads")

## End(Not run)
```

reported_heights

Self-reported Heights

Description

Students were asked to report their height (in inches) and sex in an anonymous online form. This table includes the results from combining data from four courses.

Usage

```
reported_heights
```

Format

An object of class "data.frame".

Details

- time_stamp. Time and date of the entry.
- sex. Sex as reported by the students.
- height. Height as reported by student by filling in a text free box.

Examples

```
head(reported_heights)
```

```
research_funding_rates
```

Gender bias in research funding in the Netherlands

Description

Table S1 from paper title "Gender contributes to personal research funding success in The Netherlands"

Usage

```
research_funding_rates
```

Format

An object of class "data.frame".

Details

- discipline. Research area discipline.
- applications_total. Total applications.
- applications_men. Total applications by men.
- applications_women. Total applications by women.
- awards_total. Total awards.
- awards_men. Total awards received by men.
- awards_women. Total awards received by women.
- success_rates_total. Overall success rate.
- success_rates_men. Success rate for men.
- success_rates_women. Success rate for women.

References

van der Lee R, Ellemers N. Gender contributes to personal research funding success in The Netherlands. Proc Natl Acad Sci U S A. 2015 Oct 6;112(40):12349-53. doi: 10.1073/pnas.1510159112. Epub 2015 Sep 21. PMID: 26392544; PMCID: PMC4603485.

Examples

```
research_funding_rates
# The raw data for this table is available from
invisible(raw_data_research_funding_rates)
```

```
results_us_election_2012
```

```
2012 US presidential election results
```

Description

Percentages for the four major candidates by state in the US 2012 presidential elections. It includes congressional districts for Maine and Nebraska.

Usage

```
results_us_election_2012
```

Format

An object of class "data.frame".

Details

- state. State in which poll was taken. 'U.S' is for national polls.
- electoral_votes. Electoral votes for that state.
- obama. Percent obtained by Barack Obama.
- romney. Percent obtained by Mitt Romney.
- johnson. Percent obtained by Gary Johnson.
- stein. Percent obtained by Jill Stein.

Source

Wikipedia: https://en.wikipedia.org/w/index.php?title=2012_United_States_presidential_election&oldid=1264588444

Examples

```
head(results_us_election_2012)
```

rfalling_object *Simulate falling object data*

Description

The function simulates a falling object's position. Default parameters are for dropping a weight from the tower of Pisa.

Usage

```
rfalling_object(  
  n = 14,  
  d_0 = 55.86,  
  v_0 = 0,  
  g = -9.8,  
  scale = 1,  
  time = seq(0, 3.25, length.out = n),  
  error_distribution = c("rnorm", "rt"),  
  df = 3  
)
```

Arguments

n	Sample size
d_0	Height from which object will fall in meters.
v_0	Initial velocity with which object will fall in meters per second.
g	Gravitational constant, 9.8 meters per second per second
scale	The measurement errors will be multiplied by this constant.
time	Numeric vector of times, in seconds, at which measurements were taken.
error_distribution	Character. Either rnorm for normal or rt for t-distribution.
df	If using t-distribution, the degrees of freedom.

Value

A data.frame with the time, the distance travelled, and the observed distance.

Examples

```
dat <- rfalling_object()  
with(dat, plot(time, observed_distance))  
with(dat, lines(time, distance, col = "blue"))
```

`stars`*Physical Properties of Stars*

Description

Physical properties of selected stars, including luminosity, temperature, and spectral class.

Usage

```
stars
```

Format

An object of class "data.frame".

Details

- `star`. Name of star.
- `magnitude`. Absolute magnitude of the star, which is a function of the star's luminosity and distance to the star.
- `temp`. Surface temperature in degrees Kelvin (K).
- `type`. Spectral class of star in the OBAFGKM system.

Source

Compiled from multiple open-access references on [VizieR](#).

Examples

```
head(stars)
```

`take_poll`*Models results from taking a poll*

Description

The function shows a plot of a random sample drawn from an urn with blue and red beads. The sample is taken with replacement. The proportion of blue beads is not shown so that students can try to estimate it.

Usage

```
take_poll(n, ...)
```

Arguments

n Sample size
... additional arguments to be used by the function sample.

Value

None

Examples

```
take_poll(25)
```

temp_carbon

Global temperature anomaly and carbon emissions, 1751-2018

Description

Annual mean global temperature anomaly on land, sea and combined, 1880-2018. Annual global carbon emissions, 1751-2014.

Usage

```
temp_carbon
```

Format

An object of class "data.frame".

Details

- year. Year (CE).
- temp_anomaly. Global annual mean temperature anomaly in degrees Celsius relative to the 20th century mean temperature. 1880-2018.
- land_anomaly. Annual mean temperature anomaly on land in degrees Celsius relative to the 20th century mean temperature. 1880-2018.
- ocean_anomaly. Annual mean temperature anomaly over ocean in degrees Celsius relative to the 20th century mean temperature. 1880-2018.
- carbon_emissions. Annual carbon emissions in millions of metric tons of carbon. 1751-2014.

Source

NOAA and Boden, T.A., G. Marland, and R.J. Andres (2017) via CDIAC

Examples

```
head(temp_carbon)
```

tissue_gene_expression

Gene expression profiles for 189 biological samples taken from seven different tissue types.

Description

This is a subset of the data provided by the `tissuesGeneExpression` package available from the `genomicsclass` GitHub repository. The predictors are gene expression measurements from 500 genes that are a random subset of the original 22,215.

Usage

```
tissue_gene_expression
```

Format

An object of class `list`.

Details

The example dataset is recommended for illustrating clustering and machine learning techniques.

- `x`. The predictors composed of 500 genes. Each row is a gene expression profile and each column is different gene. The column names are the gene symbols.
- `y`. The outcomes. A character vector representing the tissue. One of seven tissue types.

Source

<https://github.com/genomicsclass/tissuesGeneExpression>

Examples

```
table(tissue_gene_expression$y)
dim(tissue_gene_expression$x)
```

trump_tweets

Trump Tweets from 2009 to 2017

Description

This dataset contains all tweets from Donald Trump's Twitter account from 2009 to 2017. Additionally, the results of a sentiment analysis, conducted on tweets from the campaign period (2015-06-17 to 2016-11-08), are included in `sentiment_counts`.

Usage

```
trump_tweets
```

Format

An object of class "data.frame".

Details

- source. Device or service used to compose tweet.
- id_str. Tweet ID.
- text. Tweet.
- created_at. Data and time tweet was tweeted.
- retweet_count. How many times tweet had been retweeted at time dataset was created.
- in_reply_to_user_id_str. If a reply, the user id of person being replied to.
- favorite_count. Number of times tweet had been favored at time dataset was created.
- is_retweet. A logical telling us if it is a retweet or not.

Source

The Trump Twitter Archive: <https://www.thetrumparchive.com/>

Examples

```
head(trump_tweets)
```

us_contagious_diseases

Contagious disease data for US states

Description

Yearly counts for Hepatitis A, Measles, Mumps, Pertussis, Polio, Rubella, and Smallpox for US states. Original data courtesy of Tycho Project (<http://www.tycho.pitt.edu/>).

Usage

```
us_contagious_diseases
```

Format

An object of class "data.frame".

Details

- `disease`. A factor containing disease names.
- `state`. A factor containing state names.
- `year`.
- `weeks_reporting`. Number of weeks counts were reported that year.
- `count`. Total number of reported cases.
- `population`. State population, interpolated for non-census years.

Source

[Tycho Project](#)

References

Willem G. van Panhuis, John Grefenstette, Su Yon Jung, Nian Shong Chok, Anne Cross, Heather Eng, Bruce Y Lee, Vladimir Zadorozhny, Shawn Brown, Derek Cummings, Donald S. Burke. Contagious Diseases in the United States from 1888 to the present. *NEJM* 2013; 369(22): 2152-2158.

Examples

```
head(us_contagious_diseases)
```

Index

* datasets

- admissions, 2
 - brca, 3
 - brexit_polls, 4
 - death_prob, 5
 - divorce_margarine, 5
 - gapminder, 9
 - greenhouse_gases, 10
 - heights, 11
 - historic_co2, 11
 - mice_weights, 12
 - mnist_127, 13
 - mnist_27, 14
 - movielens, 15
 - murders, 16
 - na_example, 16
 - nyc_regents_scores, 17
 - olive, 18
 - outlier_example, 19
 - polls_2008, 19
 - polls_us_election_2016, 20
 - pr_death_counts, 21
 - reported_heights, 23
 - research_funding_rates, 24
 - results_us_election_2012, 25
 - stars, 27
 - temp_carbon, 28
 - tissue_gene_expression, 29
 - trump_tweets, 29
 - us_contagious_diseases, 30
-
- admissions, 2
 - brca, 3
 - brexit_polls, 4
 - death_prob, 5
 - divorce_margarine, 5
 - ds_theme_set, 6
 - fit_recommender_model, 7
 - gapminder, 9
 - greenhouse_gases, 10
 - heights, 11
 - historic_co2, 11
 - mice_weights, 12
 - mnist_127, 13
 - mnist_27, 14
 - movielens, 15
 - murders, 16
 - na_example, 16
 - nyc_regents_scores, 17
 - oecd (gapminder), 9
 - olive, 18
 - opec (gapminder), 9
 - outlier_example, 19
 - polls_2008, 19
 - polls_us_election_2016, 20
 - pr_death_counts, 21
 - raw_data_research_funding_rates (research_funding_rates), 24
 - read_mnist, 22
 - reported_heights, 11, 23
 - research_funding_rates, 24
 - results_us_election_2012, 25
 - results_us_election_2016 (polls_us_election_2016), 20
 - rf_falling_object, 26
 - sentiment_counts (trump_tweets), 29
 - stars, 27
 - take_poll, 27
 - temp_carbon, 28
 - tissue_gene_expression, 29
 - trump_tweets, 29
 - us_contagious_diseases, 30