

datawizard: An R Package for Easy Data Preparation and Statistical Transformations

Indrajeet Patil¹, Dominique Makowski², Mattan S. Ben-Shachar³,
Brenton M. Wiernik⁴, Etienne Bacher⁵, and Daniel Lüdtke⁶

¹ esqLABS GmbH, Germany ² Nanyang Technological University, Singapore ³ Ben-Gurion University of the Negev, Israel ⁴ Facebook ⁵ Luxembourg Institute of Socio-Economic Research, Luxembourg ⁶ University Medical Center Hamburg-Eppendorf, Germany

DOI:

Software

- [Review](#) ↗
- [Repository](#) ↗
- [Archive](#) ↗

Submitted:

Published:

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC-BY](#)).

Summary

The `{datawizard}` package for the R programming language ([R Core Team, 2021](#)) provides a lightweight toolbox to assist in key steps involved in any data analysis workflow: (1) wrangling the raw data to get it in the needed form, (2) applying preprocessing steps and statistical transformations, and (3) carrying out sanity checks for transformed data. Therefore, it can be a valuable tool for R users and developers looking for a lightweight option for data preparation.

Statement of Need

The `{datawizard}` package is part of `{easystats}`, a collection of R packages designed to make statistical analysis easier (Ben-Shachar et al. (2020), Lüdtke et al. (2020), Lüdtke, Ben-Shachar, et al. (2021), Lüdtke, Patil, et al. (2021), Lüdtke et al. (2019), Makowski et al. (2019), Makowski et al. (2020)). As this ecosystem follows a “0-external-hard-dependency” policy, a base R data manipulation package that relies only on base R needed to be created. In effect, `{datawizard}` provides data processing backend for this entire ecosystem. In addition to its usefulness to the `{easystats}` ecosystem, it also provides *an* option for R users and package developers if they wish to keep their (recursive) dependency weight to a minimum (for other options, see Dowle & Srinivasan (2021), Eastwood (2021), etc.).

Because `{datawizard}` is also meant to be used and adopted easily by a wide range of users, its workflow and syntax are designed to be similar to `{tidyverse}` (Wickham et al. (2019)), a widely used ecosystem of R packages. Thus, users familiar with the `{tidyverse}` can easily translate their knowledge and make full usage of `{datawizard}`.

In addition to being a lightweight solution to clean messy data, `{datawizard}` also provides helpers for the other important step of data analysis: applying statistical transformations to the cleaned data while setting up statistical models. This includes various types of data standardization, normalization, rank-transformation, and adjustment.

Lastly, `{datawizard}` also provides a toolbox to create a detailed profile of data properties.

Features

Data Preparation

The raw data is rarely in a state that it can be directly fed into a statistical model. It often needs to be modified in various ways. For example, columns need to be renamed, certain portions of the data need to be filtered out, reshape data, data scattered across multiple tables needs to be joined, etc.

{datawizard} provides various functions for cleaning and preparing data (see Table 1).

Table 1: The table below lists a few key functions offered by *datawizard* for data wrangling. To see the full list, see the package website: <https://easystats.github.io/datawizard/>

Function	Operation
<code>data_filter()</code>	to select only certain <i>observations</i>
<code>data_select()</code>	to select only a few <i>variables</i>
<code>data_extract()</code>	to extract a single <i>variable</i>
<code>data_rename()</code>	to rename variables
<code>reshape_longer()</code>	to convert data from wide to long
<code>reshape_wider()</code>	to convert data from long to wide
<code>data_join()</code>	to join two data frames
...	...

We will look at one example function that converts data in wide format to tidy/long format:

```
stocks <- data.frame(
  time = as.Date('2009-01-01') + 0:4,
  X = rnorm(5, 0, 1),
  Y = rnorm(5, 0, 2)
)

stocks
#>      time      X      Y
#> 1 2009-01-01 -0.5286321 -3.3399106
#> 2 2009-01-02 -2.5970577 -0.4214296
#> 3 2009-01-03  0.7577631  2.0283686
#> 4 2009-01-04  0.3187828 -2.4260454
#> 5 2009-01-05  2.2927804  0.6690921

data_to_long(
  stocks,
  select = -c("time"),
  colnames_to = "stock",
  values_to = "price"
)

#>      time stock      price
#> 1 2009-01-01      X -0.5286321
#> 2 2009-01-01      Y -3.3399106
#> 3 2009-01-02      X -2.5970577
#> 4 2009-01-02      Y -0.4214296
#> 5 2009-01-03      X  0.7577631
```

```
#> 6 2009-01-03 Y 2.0283686
#> 7 2009-01-04 X 0.3187828
#> 8 2009-01-04 Y -2.4260454
#> 9 2009-01-05 X 2.2927804
#> 10 2009-01-05 Y 0.6690921
```

Statistical Transformations

Even after getting the raw data in the needed format, we may need to transform certain variables further to meet requirements imposed by a statistical test.

{datawizard} provides a rich collection of such functions for transforming variables (see Table 2).

Table 2: The table below lists a few key functions offered by *datawizard* for data transformations. To see the full list, see the package website: <https://easystats.github.io/datawizard/>

Function	Operation
<code>standardize()</code>	to center and scale data
<code>normalize()</code>	to scale variables to 0-1 range
<code>adjust()</code>	to adjust data for effect of other variables
<code>data_shift()</code>	to shift numeric value range
<code>ranktransform()</code>	to convert numeric values to integer ranks
...	...

We will look at one example function that standardizes (i.e. centers and scales) data so that it can be expressed in terms of standard deviation:

```
d <- data.frame(
  a = c(-2, -1, 0, 1, 2),
  b = c(3, 4, 5, 6, 7)
)

standardize(d, center = c(3, 4), scale = c(2, 4))
#>      a      b
#> 1 -2.5 -0.25
#> 2 -2.0  0.00
#> 3 -1.5  0.25
#> 4 -1.0  0.50
#> 5 -0.5  0.75
```

Data Properties

The workhorse function to get a comprehensive summary of data properties is `describe_distribution()`, which combines a set of indices (e.g., measures of centrality, dispersion, range, skewness, kurtosis, etc.) computed by other functions in {datawizard}.

```
describe_distribution(mtcars)
```

Variable	Mean	SD	IQR	Min	Max	Skewness	Kurtosis	n	n_Missing
mpg	20.091	6.027	7.53	10.40	33.90	0.672	-0.022	32	0
cyl	6.188	1.786	4.00	4.00	8.00	-0.192	-1.763	32	0
disp	230.722	123.939	221.53	71.10	472.00	0.420	-1.068	32	0
hp	146.688	68.563	84.50	52.00	335.00	0.799	0.275	32	0
drat	3.597	0.535	0.84	2.76	4.93	0.293	-0.450	32	0
wt	3.217	0.978	1.19	1.51	5.42	0.466	0.417	32	0
qsec	17.849	1.787	2.02	14.50	22.90	0.406	0.865	32	0
vs	0.438	0.504	1.00	0.00	1.00	0.265	-2.063	32	0
am	0.406	0.499	1.00	0.00	1.00	0.401	-1.967	32	0
gear	3.688	0.738	1.00	3.00	5.00	0.582	-0.895	32	0
carb	2.812	1.615	2.00	1.00	8.00	1.157	2.020	32	0

Licensing and Availability

{datawizard} is licensed under the GNU General Public License (v3.0), with all source code openly developed and stored at GitHub (<https://github.com/easystats/datawizard>), along with a corresponding issue tracker for bug reporting and feature enhancements. In the spirit of honest and open science, we encourage requests, tips for fixes, feature updates, as well as general questions and concerns via direct interaction with contributors and developers.

Acknowledgments

{datawizard} is part of the collaborative *easystats* ecosystem. Thus, we thank the [members of easystats](#) as well as the users.

References

- Ben-Shachar, M. S., Lüdtke, D., & Makowski, D. (2020). effectsize: Estimation of effect size indices and standardized parameters. *Journal of Open Source Software*, 5(56), 2815. <https://doi.org/10.21105/joss.02815>
- Dowle, M., & Srinivasan, A. (2021). *Data.table: Extension of 'data.frame'*. <https://CRAN.R-project.org/package=data.table>
- Eastwood, N. (2021). *Poorman: A poor man's dependency free recreation of 'dplyr'*. <https://CRAN.R-project.org/package=poorman>
- Lüdtke, D., Ben-Shachar, M. S., Patil, I., & Makowski, D. (2020). Extracting, computing and exploring the parameters of statistical models using R. *Journal of Open Source Software*, 5(53), 2445. <https://doi.org/10.21105/joss.02445>
- Lüdtke, D., Ben-Shachar, M. S., Patil, I., Waggoner, P., & Makowski, D. (2021). performance: An R package for assessment, comparison and testing of statistical models. *Journal of Open Source Software*, 6(60), 3139. <https://doi.org/10.21105/joss.03139>
- Lüdtke, D., Patil, I., Ben-Shachar, M. S., Wiernik, B. M., Waggoner, P., & Makowski, D. (2021). see: An R package for visualizing statistical models. *Journal of Open Source Software*, 6(64), 3393. <https://doi.org/10.21105/joss.03393>
- Lüdtke, D., Waggoner, P., & Makowski, D. (2019). insight: A unified interface to access information from model objects in R. *Journal of Open Source Software*, 4(38), 1412. <https://doi.org/10.21105/joss.01412>
- Makowski, D., Ben-Shachar, M. S., & Lüdtke, D. (2019). bayestestR: Describing effects and their uncertainty, existence and significance within the Bayesian framework. *Journal of Open Source Software*, 4(40), 1541. <https://doi.org/10.21105/joss.01541>

- Makowski, D., Ben-Shachar, M. S., Patil, I., & Lüdtke, D. (2020). Methods and algorithms for correlation analysis in R. *Journal of Open Source Software*, 5(51), 2306. <https://doi.org/10.21105/joss.02306>
- R Core Team. (2021). *R: A language and environment for statistical computing*. R Foundation for Statistical Computing. <https://www.R-project.org/>
- Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L. D., François, R., Grolemund, G., Hayes, A., Henry, L., Hester, J., Kuhn, M., Pedersen, T. L., Miller, E., Bache, S. M., Müller, K., Ooms, J., Robinson, D., Seidel, D. P., Spinu, V., ... Yutani, H. (2019). Welcome to the tidyverse. *Journal of Open Source Software*, 4(43), 1686. <https://doi.org/10.21105/joss.01686>