



ggESDA: An R Package for Exploratory Symbolic Data Analysis using ggplot2

Bo-Syue Jiang

National Taipei University

Han-Ming Wu

National Chengchi University

Abstract

This paper presents the **ggESDA** package, which we developed for exploratory symbolic data analysis in R. Based on **ggplot2** Wickham (2009), the **ggESDA** package which is familiar programming structure with its parent provides a wide variety of graphical techniques such as histogram, 3D-scatterplot and radar plot. In addition, a general and customized transformation function `classic2sym()` is implemented for generating a symbolic data table from classical data frame by clustering algorithm, **RSDA** Rojas (2015) function and user-defined method. wait for edit.....

Keywords: data visualization, symbolic data analysis, exploratory data analysis, **ggplot2** extensions, interval-valued data, R.

1. Introduction

"In Data Science the aim is to extract new knowledge from Standard, Big, and complex data. Often these data are unstructured with variables defined on different kinds of units. They can also be multi-sources (as mixtures of numerical and textual data, with images and networks)." Diday and Edwin (2018). The statement indicates that not only conventional data but the unstructured data, also known as symbolic data, is vital for data science. Rather than the classical data represented by a single value, symbolic data with measurements on p random variables can be p -dimensional statistical units such as hypercubes or histograms. The field of symbolic data analysis (SDA) Billard and Diday (2007) is to broaden the application aspects of statistical methodologies, extend traditional cognition of a form of data unit and build a brand-new analysis system of data science. Recent developments in the field of big data analytics have led to a renewed interest in complex structure data such as symbolic data. As shown in Figure 1, the number of researches in SDA represents an increasing trend from 1998 to 2020, which outstands the importance of it during the years. Among ScienceDirect,

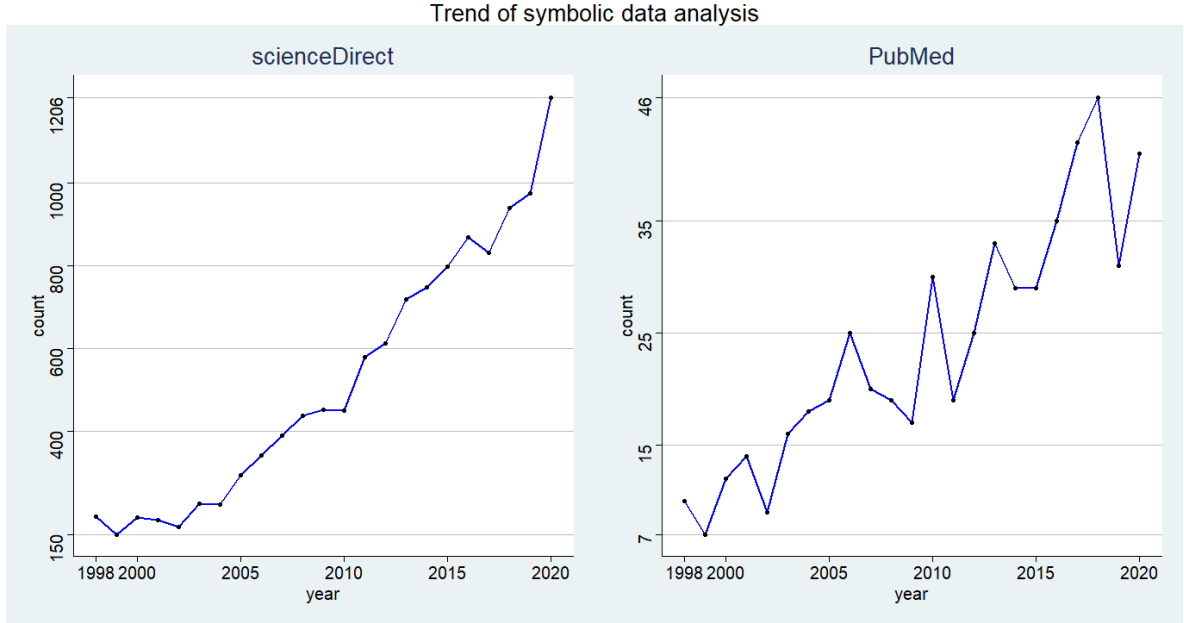


Figure 1: The number of "symbolic data analysis" or "interval-valued data" related articles in researches and applications according to PubMed and ScienceDirect online database over time from 1998 to 2020.

Engineering and Computer Science lead the subject areas obviously, shown in Figure 2.

In practice, the symbolic data is often generated by aggregating massive datasets into intervals in order to make the management easy and appropriate. An interval-valued symbolic random variable X , taking values in interval, can be denoted such as $X = [a, b] \subset R^1$, where $a \leq b$, and $a, b \in R^1$. Let the random variable X , for instance, be the weight, then $X = [50, 100]$ represents the interval covering the weight of people. With the advent of big data analytic, interval-valued data is becoming more common and accessible than ever. The researches for interval-valued data such as the sign test for COVID-19 data [Sherwani, Shakeel, Saleem, Awan, Aslam, and Farooq \(2021\)](#), the prediction via regularized artificial neural network [Yang, Lin, and Zhang \(2019\)](#), a bivariate Bayesian method for regression models [Xu and Qin \(2021\)](#), etc.

Exploratory Data Analysis (EDA) [Tukey \(1977\)](#) is primarily used to see what data can reveal beyond the formal modeling or hypothesis testing task, provides an overview of raw datasets and obtains a general understanding about the variables and their relationships.

References

- Billard L, Diday E (2007). *Symbolic Data Analysis: Conceptual Statistics and Data Mining*. Wiley, New Jersey.
- Diday, Edwin (2018). "New Advances on the Symbolic Data Analysis Framework: Basic

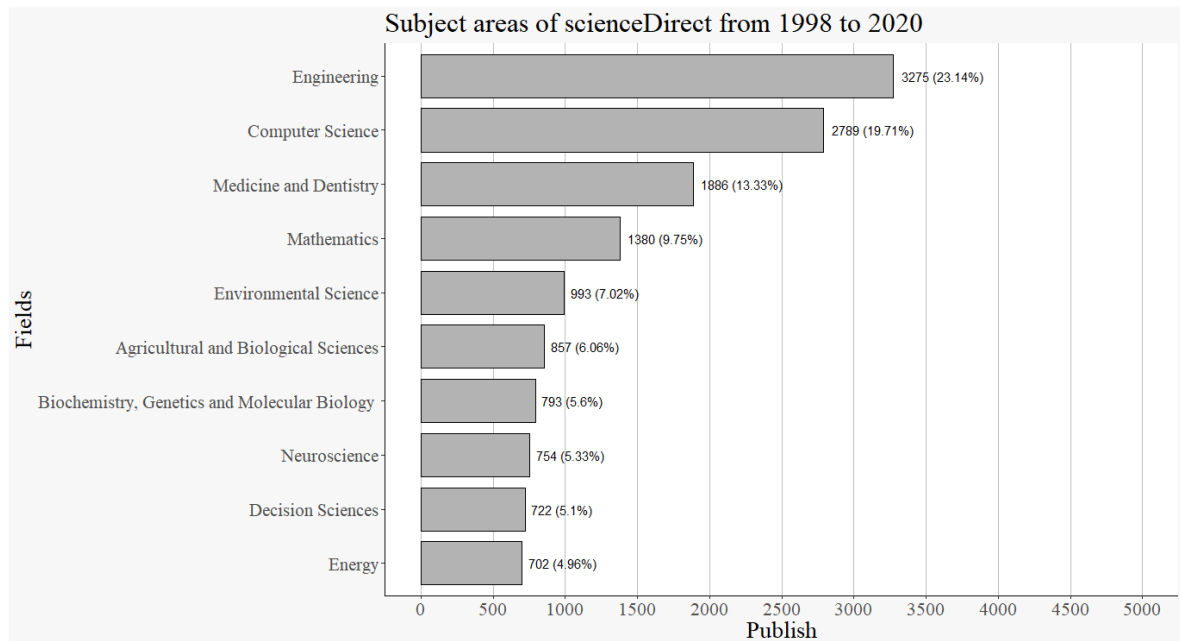


Figure 2: Top 10 researches and applications domains for SDA or interval-valued data (ScienceDirect) from 1998 to 2020

Theory, Explanatory Criteria, Improving Machine Learning, New Directions of Research.” p. 17.

Rojas OR (2015). “R to Symbolic Data Analysis.” URL https://www.imsbio.co.jp/RGM/R_rdfile?f=RSDA/man/RSDA-package.Rd&d=R_CC.

Sherwani RAK, Shakeel H, Saleem M, Awan WB, Aslam M, Farooq M (2021). “A new neutrosophic sign test: An application to COVID-19 data.” *Plos one*, **16**(8), e0255671.

Tukey JW (1977). *Exploratory data analysis*. Addison-Wesley series in behavioral science : quantitative methods. Addison-Wesley. ISBN 0201076160. URL <https://www.worldcat.org/oclc/03058187>.

Wickham H (2009). “**ggplot2**: Elegant Graphics for Data Analysis.” *Media*, **35**(211), 10–1007. doi:10.1007/978-0-387-98141-3.

Xu M, Qin Z (2021). “A bivariate Bayesian method for interval-valued regression models.” *Knowledge-Based Systems*, p. 107396.

Yang Z, Lin DK, Zhang A (2019). “Interval-valued data prediction via regularized artificial neural network.” *Neurocomputing*, **331**, 336–345.

Affiliation:

Firstname Lastname

Affiliation

Address, Country

E-mail: **name@address**URL: **http://link/to/webpage/**