



Databases and ontologies

RTCGA - The Family of R Packages Integrating Data from The Cancer Genome Atlas Study

Marcin Kosiński^{1,2,*}, Witold Chodor² and Przemysław Biecek^{1,2}

¹Faculty of Mathematics and Information Science, Warsaw University of Technology, 00-662 Warsaw, Poland and

²Faculty of Mathematics, Informatics, and Mechanics, University of Warsaw, City, 02-097 Warsaw, Poland.

*To whom correspondence should be addressed.

Associate Editor: XXXXXXX

Received on XXXXX; revised on XXXXX; accepted on XXXXX

Abstract

Summary: The following article presents **RTCGA** software package and a family of **R** (R Core Team, 2016) packages with data from The Cancer Genome Atlas Project (TCGA) study (Broad Institute of MIT and Harvard, 2014). TCGA is a comprehensive and coordinated effort to accelerate our understanding of the molecular basis of cancer through the application of genome analysis technologies, including large-scale genome sequencing [1]. We converted selected datasets from this study into few separate packages that are hosted on (?). These R packages make selected datasets easier to access and manage. Datasets in **RTCGA** packages are large and cover complex relations between clinical outcomes and genetic background. These packages will be useful for at least three audiences: biostatisticians that work with cancer data; researchers that are working on large scale algorithms; teachers that are presenting data analysis method on real data problems

Availability: **RTCGA** family of R packages is freely available at <http://rtcga.github.io/RTCGA/> and from the Bioconductor project at <http://bioconductor.org/packages/RTCGA/>.

Contact: m.p.kosinski@gmail.com

Motivation

The Cancer Genome Atlas (TCGA) Data Portal provides a platform for researchers to search, download, and analyze data sets generated by TCGA. It contains clinical information, genomic characterization data, and high level sequence analysis of the tumor genomes [1].

TCGA data is available through Firehose Broad GDAC portal [1]. One can select cancer type (cohort) and data type (e.g. clinical, RNA expression, methylation, ..) and download a 'tar.gz' file with compressed data.

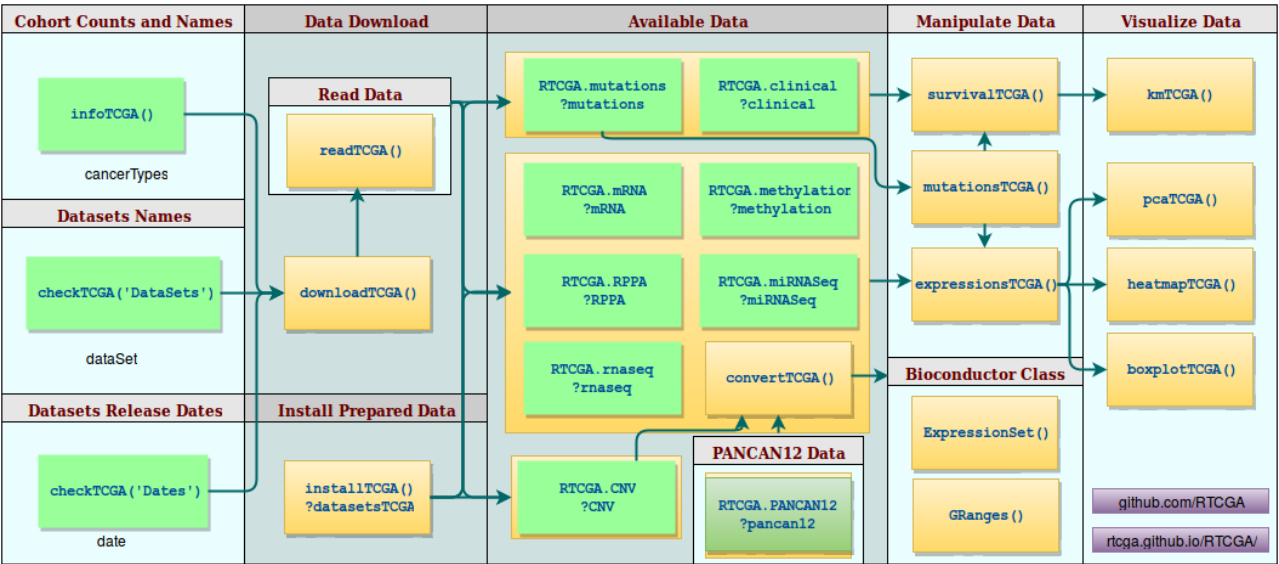
When working with many cancer types we find this approach burdensome:

- If one requires to download datasets containing e.g. information about genes' expressions for all available cohorts types (TCGA collected data for more than 30 various cancer types) one would have to go through the click-to-download process many times, which is inconvenient and time-consuming.
- Clinical datasets from TCGA project are not in a standard tidy data format, which is: one row for one observation and one column for

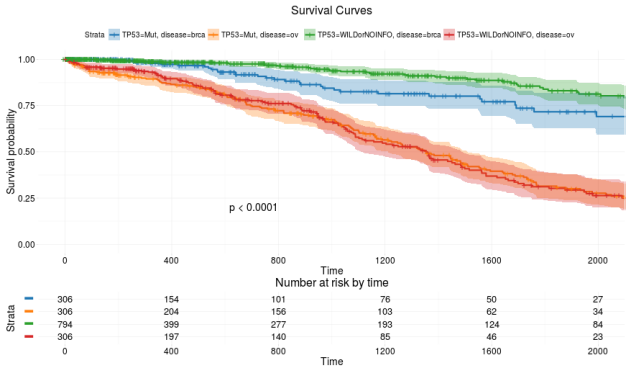
one variable. They are transposed which makes work with that data burdensome. That becomes more onerous when one would like to investigate many clinical datasets.

- Datasets containing information on some data types (e.g. gene's mutations) are not in one easy-to-handle file. Every patient has it's own file, what for many potential researchers may be an impassable barrier.
- Data governance for many datasets for various cohorts saved in different folders with strange (default after untarring) names may be exhausting and uncomfortable for researchers that are not very skilled in data management or data processing.

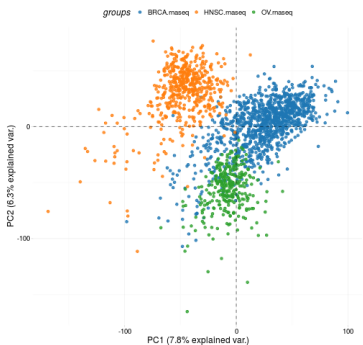
For these reasons we prepared selected datasets from the TCGA project in an easy to handle and process way and embed them in 9 separate R packages. All packages can be installed from BioConductor.



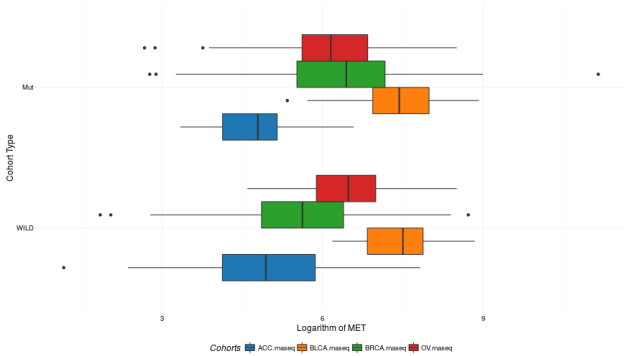
(a) The workflow of **RTCGA** family of software and data R packages. **RTCGA** consists of auxiliary functions: `infoTCGA()` and `checkTCGA` that enable to find out metadata about datasets and dates of their release provided by TCGA. Data download section consists of `downloadTCGA()` function that allows to download every dataset from TCGA study and `readTCGA()` function that enables to read (most popular) data into the tidy format. It is possible to use `installTCGA()` function to install, prepared in tidy format, most popular data types from TCGA that are included in data packages. Datasets in **RTCGA** data packages can be converted to Bioconductor format (`ExpressionSet`, `GRanges`) with `convertTCGA()` function. Effects of functions designed to manipulate and visualize **RTCGA** data are presented on subfigures (b)-(e).



(b) Kaplan-Meier (Kaplan and Meier, 1958) estimates of survival curves and risk set table for **Breast invasive carcinoma (BRCA)** and **Ovarian serous cystadenocarcinoma (OV)** cohorts divided on **TP53** mutations. The effect of `survivalTCGA()`, `mutationsTCGA()` and `kmTCGA()` functions.



(c) Plot of Principal Component Analysis (Krzanowski, 2000) performed for genes expressions (RNASeq) for **Breast invasive carcinoma (BRCA)**, **Head and Neck squamous cell carcinoma (HNSC)** and **Ovarian serous cystadenocarcinoma (OV)** cohorts. The effect of `expressionsTCGA()` and `pcaTCGA()` functions.



(d) Boxplots (Robert McGill, 1978) of logarithm of **MET** gene expression (RNASeq) for **Adrenocortical carcinoma (ACC)**, **Bladder urothelial carcinoma (BLCA)**, **Breast invasive carcinoma (BRCA)** and **Ovarian serous cystadenocarcinoma (OV)** divided on mutations in gene **TP53**. The effect of `expressionsTCGA()` and `boxplotTCGA()` functions.



(e) Heatmap (Friendly, 1994) presenting medians of **ZNF501** gene for **Adrenocortical carcinoma (ACC)**, **Bladder urothelial carcinoma (BLCA)**, **Breast invasive carcinoma (BRCA)** and **Ovarian serous cystadenocarcinoma (OV)** divided on **MET** gene quantiles. The effect of `expressionsTCGA()` and `heatmapTCGA()` functions.

Fig. 1: The workflow of **RTCGA** family of software and data R packages and effects of functions designed to manipulate and visualize **RTCGA** data.

Examples

Acknowledgements

Text Text Text Text Text Text Text. ? might want to know about text text text text

Funding

This work has been supported by the... Text Text Text Text.

References

Friendly, M. (1994). Mosaic displays for multi-way contingency tables. *Journal of the American Statistical Association*, **89**(425), 190–200.

Kaplan, E. L. and Meier, P. (1958). Nonparametric estimation from incomplete observations. *Journal of the American Statistical Association*, **53**(282), 457–481.

Krzanowski, W. (2000). *Principles of Multivariate Analysis*. Oxford Statistical Science Series. OUP Oxford.

R Core Team (2016). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.

Robert McGill, John W. Tukey, W. A. L. (1978). Variations of box plots. *The American Statistician*, **32**(1), 12–16.

Broad Institute of MIT and Harvard (2014). *Broad Institute TCGA Genome Data Analysis Center: Firehose stddata 2015 06 01 run*. DOI:10.7908/C1251HBG.