

RTCGA.data - The Family of R Packages with Data from The Cancer Genome Atlas Study

by Marcin Kosinski, Przemysław Biecek

Abstract The following article presents RTCGA.data: a family of R packages with data from The Cancer Genome Atlas Project (TCGA) study. 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 one GitHub repository. These R packages make selected datasets easier to access and manage. Data sets in RTCGA.data 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, for them RTCGA data will be a perfect blasting site; teachers that are presenting data analysis method on real data problems.

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 are 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 i.e. 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 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 what makes work with those 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 TCGA project in an easy to handle and process way and embed them in 5 separate R packages. All packages can be installed from GitHub by evaluating the following code:

```
if (!require(devtools)) {  
  install.packages("devtools")  
  require(devtools)  
}  
install_github(paste0("mi2-warsaw/RTCGA.data/",  
  subdir = paste0("RTCGA.",  
    c("clinical", "rnaseq", "mutations", "cnv", "PANCAN12"))))
```

One package, i.e. RTCGA.clinical can be installed with the command

```
if (!require(devtools)) {  
  install.packages("devtools")  
  require(devtools)  
}
```

```
install_github("mi2-warsaw/RTCGA.data",
              subdir = "RTCGA.clinical")
```

If you are using Windows, make sure you have rtools [3] installed on your computer, before evaluating aboved commands.

RTCGA.data family contains 5 packages:

- **RTCGA.clinical** package containing clinical datasets from TCGA. Each cohort contains one dataset prepared in a tidy format. Each row, marked with patients' barcode, corresponds to one patient. Clinical data format is explained here <https://wiki.nci.nih.gov/display/TCGA/Clinical+Data+Overview>
- **RTCGA.rnaseq** package containing genes' expressions datasets from TCGA. Each cohort contains one dataset with over 20 thousand of columns corresponding to genes' expression. Rows correspond to patients, that can be matched with patient's barcode. Genes' expressions data format is explained here <https://wiki.nci.nih.gov/display/TCGA/RNASEq+Version+2>
- **RTCGA.mutations** package containint genes' mutations datsets from TCGA. Each cohort contains one dataset with extra column specifying patient's barcode which enables to distinguish which rows correspond to which patient. Mutations' data format is explained here [https://wiki.nci.nih.gov/display/TCGA/Mutation+Annotation+Format+\(MAF\)+Specification](https://wiki.nci.nih.gov/display/TCGA/Mutation+Annotation+Format+(MAF)+Specification).
- **RTCGA.cnv** package **explanation needed**.
- **RTCGA.PANCAN12** package **explanation needed**.

More detailed information about datasets included in **RTCGA.data** family are shown in Table 1

How to work with RTCGA.data family

After installation, one can load any package from **RTCGA.data** family with commands

```
library(RTCGA.clinical)
library(RTCGA.rnaseq)
library(RTCGA.mutations)
library(RTCGA.PANCAN12)
library(RTCGA.cnv)
```

and one can check what datasets are available (Table 1) with commands

```
?clinical
?rnaseq
?mutations
?pancan12
?cnv
```

The data loading proceeds in a regular way. Simply type

```
data(cohort.package)
```

Where `cohort` corresponds to a specific Cohort of patients and `package` corresponds to the one of five packages from **RTCGA.data** family.

Examples

```
data("BRCA.cnv")
data("COAD.rnaseq")
data("GBMLGG.mutations")
```

Patient's barcode as a key to merge data

A TCGA barcode is composed of a collection of identifiers. Each specifically identifies a TCGA data element. Refer to the following figure for an illustration of how metadata identifiers comprise a barcode. An aliquot barcode, an example of which shows in the illustration, contains the highest number of identifiers [4]. An illustration on what each part of the patient's barcode stands for is shown on Figure 1.

More detailed information about patients' barcodes can be found on [4].

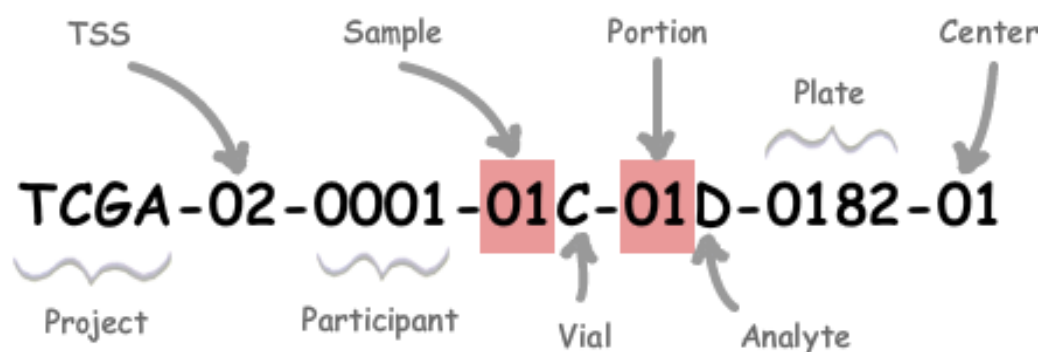


Figure 1: This figure of an aliquot barcode shows how it can be broken down into its components and translated into its metadata. Source [4]

Examples of applications

The Kaplan-Meier estimate of the survival curves with the clinical data

RTCGA.data family is excellent when one researches in a field of survival analysis and genomics. Survival times for patients are included in clinical datasets. The following example plots Kaplan-Meier [5] estimates of the survival functions for patients suffering from LUAD cancer, divided into stage of the cancer.

```
library(dplyr)
library(archivist)
library(RTCGA.clinical)
library(survival)
library(survMisc)
library(ggplot2)

createEmptyRepo("RTCGA.family")
setLocalRepo("RTCGA.family")
mergeStages <- function( patient.stage_event.pathologic_stage ){
  levels(patient.stage_event.pathologic_stage) <- c(rep("1",3), rep("2",3),
                                                    rep("3",2), "4")

  patient.stage_event.pathologic_stage %>%
    as.character %>% as.numeric()
}
LUAD.clinical %a%
  select( patient.days_to_last_followup,
          patient.stage_event.pathologic_stage,
          patient.drugs.drug.therapy_types.therapy_type,
          patient.vital_status,
          patient.days_to_death,
          patient.bcr_patient_barcode
        ) %a%
  mutate(
    patient.days_to_death = patient.days_to_death%>% as.character() %>% as.numeric(),
    patient.days_to_last_followup = patient.days_to_last_followup%>% as.character() %>% as.numeric(),
    patient.vital_status = ifelse(LUAD.clinical$patient.vital_status %>% as.character() == "dead", 1, 0),
    barcode = patient.bcr_patient_barcode %>% as.character()
  ) %a%
  mutate(
    times = ifelse( !is.na(patient.days_to_last_followup),
                    patient.days_to_last_followup,
                    patient.days_to_death),
    stage = mergeStages(LUAD.clinical$patient.stage_event.pathologic_stage)
  ) %a%
```

```

rename(
  therapy = patient.drugs.drug.therapy_types.therapy_type
) %>%
filter( !is.na(times) ) -> LUAD.clinical.selected
LUAD.clinical.selected %>%
survfit( Surv(times, patient.vital_status)~ stage, data = .) %>%
survMisc::autoplot.survfit( titleSize=12, type="CI") %>%
. [[2]] -> km_plot_luad

pdf(file = "km_plot_luad.pdf")
km_plot_luad
dev.off()

```

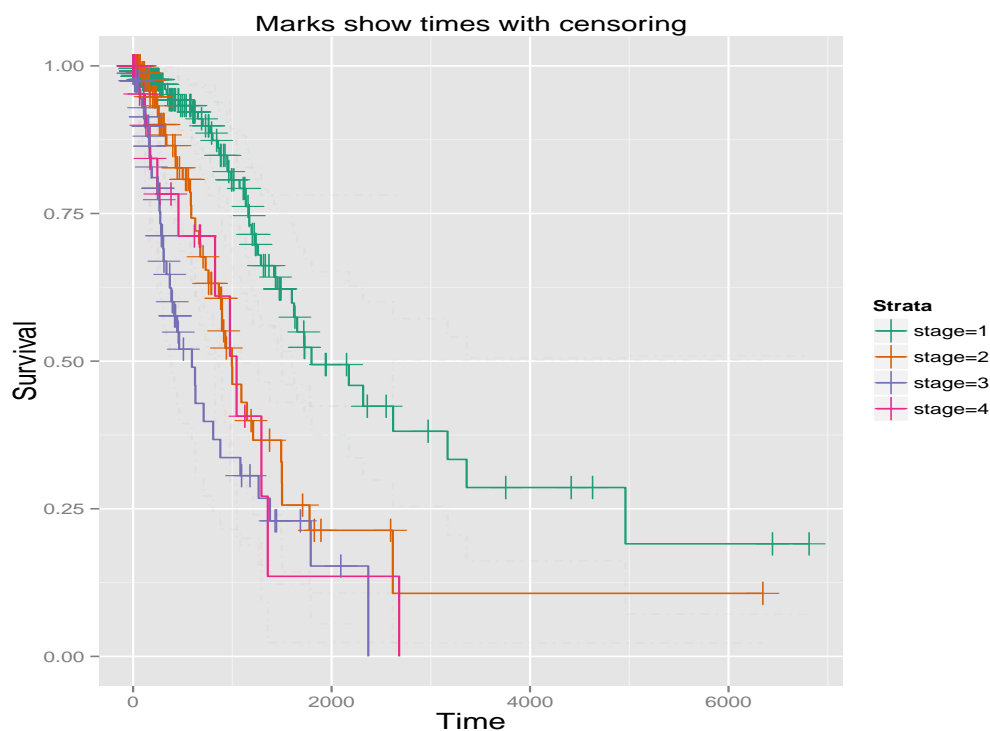


Figure 2: The Kaplan-Meier estimate of the survival curve for the LUAD cancer. The plot is available via code `aread(mi2-warsaw/rticle/RTCGA.family/857415c30aa6f9d1fac345c54c447b46)`

The Cox proportional hazards model with the genes' mutations data

In a simple way one can use previously selected data to merge them with genes' mutations data and to compute Cox proportional hazards model.

```
divideTP53 <- function( TP53 ){
  TP53 <- as.character(TP53)
  TP53 <- ifelse( (TP53 %in% c("")) | is.na(TP53),
    "WILD",
    TP53 )
  TP53 <- ifelse( grepl( "Misse", TP53 ),
    "Missense",
    TP53 )
  TP53 <- ifelse( !( TP53 %in% c("WILD", "Missense", NA) ),
    "Other",
    TP53 )
  return(TP53)
}
library(RTCGA.mutations)
LUAD.clinical.selected %a%
  left_join( y = LUAD.mutations %a%
    filter( Hugo_Symbol == "TP53") %>%
    mutate( barcode = barcode %>% as.character %>% tolower %>% substr(1,12) ) %>%
    select( barcode, Variant_Classification),
    by = "barcode") %a%
    mutate( Variant_Classification = divideTP53(Variant_Classification) ) ->
LUAD.clinical.mutations.selected

coxph(Surv(times, patient.vital_status)~ as.factor(stage)+Variant_Classification,
  data = LUAD.clinical.mutations.selected)
```

Call:

```
coxph(formula = Surv(times, patient.vital_status) ~ as.factor(stage) +
  Variant_Classification, data = LUAD.clinical.mutations.selected)
```

	coef	exp(coef)	se(coef)	z	p
as.factor(stage)2	0.8072	2.2417	0.2328	3.47	0.00053
as.factor(stage)3	1.3804	3.9764	0.2339	5.90	3.6e-09
as.factor(stage)4	1.1555	3.1756	0.3414	3.38	0.00071
Variant_ClassificationOther	0.4397	1.5523	0.3284	1.34	0.18058
Variant_ClassificationWILD	-0.0365	0.9642	0.2396	-0.15	0.87890

Likelihood ratio test=45.1 on 5 df, p=1.36e-08

n= 508, number of events= 126

(2 observations deleted due to missingness)

The model is available via code `aread("mi2-warsaw...")`.

The Principal Components Analysis for the rnaseq data

Enhancement text needed.

```
data(package = "RTCGA.rnaseq")\$results[1:6,3] %a%
  sapply(function(element){
    data(list=element,
          package = "RTCGA.rnaseq",
          envir = .GlobalEnv)
    get(element, envir = .GlobalEnv) %>%
      t() %>%
      .[-1,-1]
  }) %a%
do.call(rbind, .) %a%
apply( 2, function(x) as.numeric(as.character(x))) ->
rnaseq_sample_joined_numeric
      "8443b0fb0e5292f403260296458d7a97"
rnaseq_sample_joined_numeric %a%
  colSums() -> rnaseq_col_sums

(rnaseq_col_sums == 0 ) %a%
  which -> columns_with_only0

rnaseq_sample_joined_numeric[,-columns_with_only0] %a%
  prcomp( scale = TRUE ) -> PCA

data(package = "RTCGA.rnaseq")\$results[1:6,3] %a%
  sapply(function(element){
    get(element, envir = .GlobalEnv) %>%
      ncol()
  }) -> rnaseq_ncol

mapply(rep,
        data(package = "RTCGA.rnaseq")\$results[1:6,3],
        rnaseq_ncol-1) %a%
  unlist -> rnaseq_pca_labels

library(ggbiplot)
rownames(PCA$rotation) <- 1:nrow(PCA$rotation)
ggbiplot(PCA, obs.scale = 1, var.scale = 1,
  groups = rnaseq_pca_labels, ellipse = TRUE, circle = TRUE, var.axes=FALSE) +
  theme(legend.direction = 'horizontal', legend.position = 'top') ->x
```

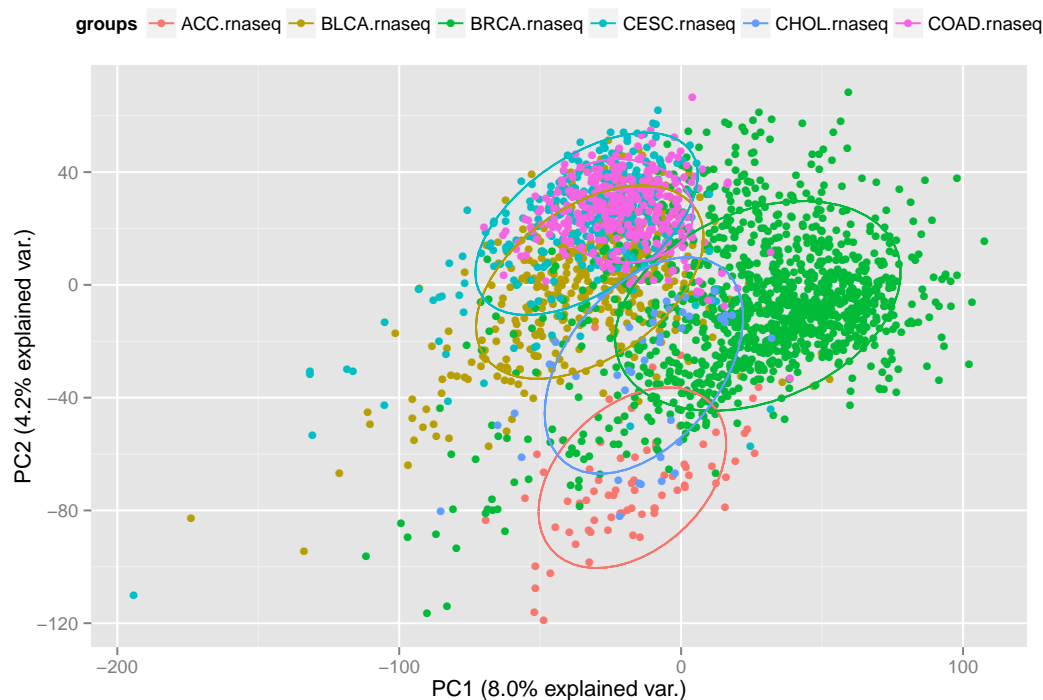


Figure 3: The biplot for 2 main components of principal component analysis of genes' expressions data for 6 various cancer types. The plot is available via code `aread(mi2-warsaw/rticle/TCGA.family/md5hash)`

[1] <http://cancergenome.nih.gov/>

[2] <http://gdac.broadinstitute.org/>

[3] <http://cran.r-project.org/bin/windows/Rtools/>

[4] <https://wiki.nci.nih.gov/display/TCGA/TCGA+barcode>

[5] Kaplan, E. L.; Meier, P. (1958). "Nonparametric estimation from incomplete observations". J. Amer. Statist. Assoc. 53(282): 457-481.

`\bibliography{RJreferences}`

Marcin Kosinski
Warsaw University of Technology
Faculty of Mathematics and Information Science
Koszykowa 75, 00-662 Warsaw, Poland
M.P.Kosinski@gmail.com

Przemysław Biecek
University of Warsaw
Faculty of Mathematics, Informatics, and Mechanics
Banacha 2, 02-097 Warsaw, Poland
Przemyslaw.Biecek@gmail.com

Table 1: Dimensions of available datasets in **RTCGA.family**.

	Disease Name	Cohort	Cases	clinical	cnv ^a	mutations	rnaseq ^b
1	Adrenocortical carcinoma	ACC	92	92 x 1046	21052	20255 x 53	80
2	Bladder urothelial carcinoma	BLCA	412	393 x 1978	105795	39441 x 96	428
3	Breast invasive carcinoma	BRCA	1098	1080 x 3464	284510	91471 x 68	1213
4	Cervical and endocervical cancers	CESC	307	304 x 1556	59450	46740 x 58	310
5	Cholangiocarcinoma	CHOL	36	36 x 794	7570	6789 x 49	46
6	Colon adenocarcinoma	COAD	460	453 x 2935	91166	62683 x 40	329
7	Colorectal adenocarcinoma	COADREAD	631	624 x 3241	126931		434
8	Lymphoid Neoplasm Diffuse ... ^c	DLBC	58	47 x 693	9343		29
9	Esophageal carcinoma	ESCA	185	174 x 1115	60803		
10	FFPE Pilot Phase II	FPPP	38	38 x 3277			
11	Glioblastoma multiforme	GBM	613	592 x 4935	146852	22362 x 80	167
12	Glioma	GBMLGG	1129	1067 x 5158	226643		697
13	Head and Neck squamous cell carcinoma	HNSC	528	522 x 1625	110289	52077 x 90	567
14	Kidney Chromophobe	KICH	113	110 x 854	10164	7624 x 37	92
15	Pan-kidney cohort (KICH+KIRC+KIRP)	KIPAN	973	914 x 2554	142122	73527 x 36	1021
16	Kidney renal clear cell carcinoma	KIRC	537	533 x 2474	85044	26785 x 36	607
17	Kidney renal papillary cell carcinoma	KIRP	323	271 x 1746	46914	15745 x 53	324
18	Acute Myeloid Leukemia	LAML	200	200 x 1148	28324	2781 x 65	174
19	Brain Lower Grade Glioma	LGG	516	475 x 1940	79791	10170 x 39	531
20	Liver hepatocellular carcinoma	LIHC	377	363 x 1441	93328	28089 x 49	424
21	Lung adenocarcinoma	LUAD	585	521 x 2765	122927	72770 x 92	577
22	Lung squamous cell carcinoma	LUSC	504	495 x 2487	134864	65482 x 87	553
23	Mesothelioma	MESO	87	77 x 855	18335		87
24	Ovarian serous cystadenocarcinoma	OV	602	591 x 3305	261680	20534 x 44	266
25	Pancreatic adenocarcinoma	PAAD	185	174 x 1148	34808	37850 x 85	184
26	Pheochromocytoma and Paraganglioma	PCPG	179	179 x 1102	31256	4784 x 91	188
27	Prostate adenocarcinoma	PRAD	499		117345	27687 x 91	551
28	Rectum adenocarcinoma	READ	171	171 x 2492	35765	22143 x 40	106
29	Sarcoma	SARC	260		106617		
30	Skin Cutaneous Melanoma	SKCM	470	447 x 1750	108084	290666 x 91	473
31	Stomach adenocarcinoma	STAD	443	438 x 1583	118389	148808 x 80	
32	Stomach and Esophageal carcinoma	STES	628	612 x 1719	179192	148808 x 80	
33	Testicular Germ Cell Tumors	TGCT	150	133 x 924	24952	14826 x 58	157
34	Thyroid carcinoma	THCA	503	501 x 1556	55377	7862 x 91	569
35	Thymoma	THYM	124	122 x 771	15571		123
36	Uterine Corpus Endometrial Carcinoma	UCEC	560	537 x 2038	127430	185108 x 50	202
37	Uterine Carcinosarcoma	UCS	57	57 x 876	19298	11395 x 91	58
38	Uveal Melanoma	UVM	80	80 x 541	12973	2607 x 91	81

^aThe second dimension is always equal to 6.^bThe first dimension is always equal to 20532.^cLymphoid Neoplasm Diffuse Large B-cell Lymphoma