

# Working with TCGAbiolinks package

*Antonio Colaprico, Tiago Chedraoui Silva, Luciano Garofano, Catharina Olsen, Davide Carolini, Claudia Cava, Isabella Castiglioni, Thais, Tathiane, Stefano Pagnotta, Michele Ceccarelli Houtan Nourshahre, Gianluca Bontempi*

2015-08-06

## Contents

Introduction	2
<b>TCGAquery: Searching TCGA open access data</b>	<b>2</b>
TCGAquery: some filtering examples	2
TCGAquery: variation: Multiple versions of the data in TCGA	5
TCGAquery_clinic & TCGAquery_clinicFilter: Working with clinical data	5
TCGAquery_subtypes: Working with molecular subtypes data	7
TCGAquery_integrate: Summary of the common numbers of patient samples in different platforms	8
TCGAquery: some examples	8
<b>TCGAdownload: Downloading open access data</b>	<b>8</b>
TCGAdownload: Example of use	9
TCGAdownload: Table of types available for downloading	9
<b>TCGAprepare: Preparing the data</b>	<b>9</b>
TCGAprepare: Example of use	10
TCGAprepare: Table of types available for the TCGAprepare	11
TCGAprepare: Preparing the data with parameter: toPackage	11
TCGAprepare: Preparing the data with CNV data (Genome_Wide_SNP_6)	12
<b>TCGAanalyze: Analyze data from TCGA</b>	<b>12</b>
TCGAanalyze_Preprocessing: Preprocessing of Gene Expression data (IlluminaHiSeq_RNASeqV2)	12
TCGAanalyze_DEA & TCGAanalyze_Lavaan: Differential expression analysis (DEA)	14
TCGAanalyze_KEGGpath & TCGAanalyze_KEGGpath: Enrichment Analysis	15
TCGAanalyze_survival: Survival Analysis: Cox Regression and dnet package	16
<b>TCGAvisualize: Visualize results from analysis functions with TCGA's data</b>	<b>18</b>
TCGAvisualize_PCA: Principal Component Analysis plot for differentially expressed genes	18
TCGAvisualize_SurvivalCoxNET: Survival Analysis: Cox Regression and dnet package	19
<b>TCGA Downstream Analysis some workflows and pipelines</b>	<b>20</b>
Downstream Analysis n.1	21
Downstream Analysis n.2 IlluminaHiSeq_RNASeq data	21
Downstream Analysis n.3 CG and GRM Integration (Heatmap and Cluster)	22
Downstream Analysis n.4 DNA methylation analysis	25
TCGAvisualize_momMethylation: Sample Mean DNA Methylation Analysis	26
TCGAanalyze_DMR: Differentially methylated regions Analysis	27
TCGAvisualize_starburst: Analyzing expression and methylation together	28
<b>TCGAinvestigator: Searching questions, answers and literature</b>	<b>29</b>
TCGAinvestigator: Find most studied TFs in PubMed	29
<b>TCGAassoc1: Searching questions, answers and literature</b>	<b>30</b>
TCGAassoc1: with BioConductor	30

Working with TCGAbiolinks package	2
TCGAquery with Bioconductor	31
References	34

## Introduction

**Motivation:** The Cancer Genome Atlas (TCGA) provides us with an enormous collection of data sets, not only spanning a large number of tumors but also a large number of experimental platforms. Even though the data can be accessed and downloaded from the database, the possibility to analyze these downloaded data directly in one single R package has not yet been available.

TCGAbiolinks consists of three parts or levels. Firstly, we provide different options to query and download from TCGA relevant data from all currently platforms and their subsequent pre-processing for commonly used bio-informatics (tools) packages in Bioconductor or CRAN. Secondly, the package allows to integrate different data types and it can be used for different types of analyses dealing with all platforms such as diff-expression, network-inference or survival analysis, etc. and thus it allows to visualize the obtained results. Thirdly we added a social level where a researcher can found a similar interest in a bioinformatic community, and allows both to find a validation of results in literature is published and also to raise questions and answers from site such as support.bioconductor.org, biocStars.org, stackoverflow.com, etc.

This document describes how to search, download and analyze TCGA data using the TCGAbiolinks package.

## TCGAquery: Searching TCGA open-access data

You can easily search TCGA samples using the TCGAquery function. Using a summary of filters as used in the TCGA portal, the function works with the following parameters:

- **tumor** Tumor or list of tumors. The list of tumor is shown in the examples.
- **platform** Platform or list of tumors. The list of platforms is shown in the examples.
- **samples** List of TCGA barcodes.
- **level** Options: 1,2,3 "mega-tab"
- **center**
- **version** List of Platform/Tumor/Version to be changed

### TCGAquery: some filtering examples

#### TCGAquery: Searching by tumor

You can filter the search by tumor using the tumor parameter.

```
query <- TCGAquery(tumor = "gliob")
```

If you don't remember the tumor name, or if you have incorrectly typed it, it will provide you with all the tumor names in TCGA. Also the names can be seen in the help pages TCGAquery.

```
query <- TCGAquery(tumor = "?")
```

```
##
##
## Table: TCGA tumors
##
## -----
## AOC   CML   GBV   LAML   LUSC   PDP   STAD   UCS
## BLCA   COAD   HNSC   LGG   MESO   PAAD   TCGT   UCE
## BRCA   BLGG   KICH   LGG   MISC   KIRC   HNSC   ACC
## ESCA   ESCA   KIRC   LIHC   OV   SARC   THCA   BLCA
```

```

44 CNCL  FPM  KIP  LMO  MLD  SKM  HSB  BSC
45 -----
46
47 ERROR: Discard not found. Select from the table above.
48 =====

```

#### TCGAquery: Searching by level

You can filter the search by level "1", "2", "3" or "mega-tab".

```

query <- TCGAquery(tumor = "gli", level = 1)
query <- TCGAquery(tumor = "gli", level = 2)
query <- TCGAquery(tumor = "gli", level = 3)
query <- TCGAquery(tumor = "gli", level = "mega-tab")

```

#### TCGAquery: Searching by platform

You can filter the search by platform using the platform parameter.

```
query <- TCGAquery(tumor = "gli", platform = "IlluminaHiSeq_RNASeqV2")
```

If you don't remember the platform, or if you have incorrectly typed it, it will provide you with all the platforms names in TCGA. Also the names can be used in the help pages TCGAquery.

```
query <- TCGAquery(tumor = "gli", platform = "")
```

```

49
50
51 Table: TCGA Platforms
52
53 -----
54
55 454      HumanMethylation27      IlluminaHiSeq 450S
56 ABI      HumanMethylation450    Kapping250K_Mop
57 Agilent04502A_07      IlluminaDNA Methylation_OXA002_CFI      Kapping250K_Sty
58 Agilent04502A_07_1    IlluminaDNA Methylation_OXA003_CFI      KDA_RPPA_Core
59 Agilent04502A_07_2    IlluminaGA_DNASeq      microarr_1
60 Agilent04502A_07_3    IlluminaGA_DNASeq_autocated      sinbio
61 bio      IlluminaGA_DNASeq_Count      sinbiotab
62 biotab    IlluminaGA_DNASeq_Count_autocated      Fixed_DNASeq
63 CCHM-1x1M_G4447A      IlluminaGA_DNASeq_Count_curated      Fixed_DNASeq_autocated
64 diagnostic_images      IlluminaGA_DNASeq_curated      Fixed_DNASeq_Count
65 fh_analyses      IlluminaGA_mRNASeq      Fixed_DNASeq_Count_autocated
66 fh_reports      IlluminaGA_mRNA_DCE      Fixed_DNASeq_Count_curated
67 fh_stddata      IlluminaGA_RNASeq      Fixed_DNASeq_curated
68 Genomic_Wide_SNP_5    IlluminaGA_RNASeqV2      Multicenter_mutation_calling_HCC
69 GenomicWideSNP_5      IlluminaGA      Multicenter_mutation_calling_HCC_Count
70 H-mRNA_8x10K      IlluminaHiSeq_150Seq      pathology_reports
71 H-mRNA_8x10Kv2      IlluminaHiSeq_DNASeq_autocated      SOLiD_DNASeq
72 H-mRNA_EarlyAccess    IlluminaHiSeq_DNASeq_Count      SOLiD_DNASeq_autocated
73 H-mRNA_1M4VGA      IlluminaHiSeq_150Seq_Count_autocated      SOLiD_DNASeq_Count
74 HG-CGR-244A      IlluminaHiSeq_150Seq_Count_curated      SOLiD_DNASeq_Count_autocated
75 HG-CGR-415K_G4124A    IlluminaHiSeq_150Seq_curated      SOLiD_DNASeq_Count_curated
76 HG-C188_P1us_2      IlluminaHiSeq_DNASeqC      SOLiD_DNASeq_curated
77 HG-C188A_2      IlluminaHiSeq_mRNASeq      supplemental_clinical
78 HT_50-V133A      IlluminaHiSeq_mRNA_DGE      tissue_images
79 HuEx-1_0-st-v2      IlluminaHiSeq_RNASeq      WHG-1x1M_G4112A

```

```

## Human1KDup      ILLUMINAHiSeq_RNASeqV2      NRG-4s44K_34x15F
## HumanRap550      ILLUMINAHiSeq_TotalRNASeqV2      NRG-02E_4x64S
## -----
##
## ERROR: Platform not found. Select from the table above.
## -----

```

#### TCGAquery: Searching by center

You can filter the search by center using the center parameter

```
query <- TCGAquery(tumor = "gli", center = "mskcc.org")
```

If you don't remember the center or if you have incorrectly typed it, it will provide you with all the center names in TCGA

```
query <- TCGAquery(tumor = "gli", center = "")
```

```

##
##
## Table: TCGA Centers
##
## -----
## bcga.ca      intgen.org      rubicongenomics.com
## brhad.mh.edu  jhu.usc.wcu      xenger.wu.ac
## broadinstitute.org  jhu.edu      xysterabiology.org
## cshbimed.CSCa  lbl.gov      ucsc.edu
## genosc.puhtl.edu  adanderson.org  unc.edu
## hgsc.bcm.edu     mskcc.org      usc.edu
## hsc.harvard.edu  nsticnquidechildren.org  vanderbilt.edu
## hudsonalpha.org  pml.gov      hogan.ca
## -----
##
## ERROR: Center not found. Select from the table above.
## -----

```

#### TCGAquery: Searching by samples

You can filter the search by samples using the samples parameter. You can give a list of barcodes or only one barcode. These barcode can be partial barcodes.

# You can define a list of samples to query and download providing relative TCGA barcodes.

```
listSamples <- c("TCGA-E9-A18G-11A-S2R-A14P-07", "TCGA-BE-A1FD-11A-S2R-A18Q-07",
  "TCGA-A7-A15C-11A-S1R-A15C-07", "TCGA-BE-A0DK-11A-13R-A0E9-07",
  "TCGA-B9-A15H-11A-S4R-A15G-07", "TCGA-BE-A0NU-01A-11R-A12P-07",
  "TCGA-CU-A1EJ-01A-11R-A13C-07", "TCGA-A7-A13D-01A-13R-A12P-07",
  "TCGA-B2-A0C9-01A-S3R-A11D-07", "TCGA-A9-A0Y5-01A-11R-A19R-07")
```

# Query all available platforms with a list of barcode

```
query <- TCGAquery(samples = listSamples)
```

# Query with a partial barcode

```
query <- TCGAquery(samples = "TCGA-B1-174S-01A")
```



### TCGAquery\_version: Retrieve versions of the data in TCGA

Query version for a specific platform for example IlluminaHiSeq\_RNASeqV2

```
library(TCGAdata.links)
```

```
BRCA_RNASeqV2_version <- TCGAquery_Version(tumor = "brca",
                                             platform = "IlluminaHiSeq_RNASeqV2")
```

The result is shown below:

Table 1: Table with version, number of samples and size (Mbyte) of BRCA IlluminaHiSeq\_RNASeqV2 Level 3

Version	Date	Samples	Size(Mbyte)
unc.edu_BRCA.IlluminaHiSeq_RNASeqV2_Level_3.1.11.0/	2015-01-26 08:15	1215	1740.6
unc.edu_BRCA.IlluminaHiSeq_RNASeqV2_Level_3.1.10.0/	2014-10-15 18:00	1215	1736.4
unc.edu_BRCA.IlluminaHiSeq_RNASeqV2_Level_3.1.9.0/	2014-07-14 18:13	1182	1689.6
unc.edu_BRCA.IlluminaHiSeq_RNASeqV2_Level_3.1.8.0/	2014-05-05 23:14	1172	1675.2
unc.edu_BRCA.IlluminaHiSeq_RNASeqV2_Level_3.1.7.0/	2014-02-13 20:07	1180	1677.9
unc.edu_BRCA.IlluminaHiSeq_RNASeqV2_Level_3.1.6.0/	2014-01-13 09:53	1140	1629.1
unc.edu_BRCA.IlluminaHiSeq_RNASeqV2_Level_3.1.5.0/	2013-08-22 18:05	1105	1580.8
unc.edu_BRCA.IlluminaHiSeq_RNASeqV2_Level_3.1.4.0/	2013-04-25 16:35	1032	1476.5
unc.edu_BRCA.IlluminaHiSeq_RNASeqV2_Level_3.1.3.0/	2013-04-12 15:29	958	1369.3
unc.edu_BRCA.IlluminaHiSeq_RNASeqV2_Level_3.1.2.0/	2012-12-17 18:23	956	1366.5
unc.edu_BRCA.IlluminaHiSeq_RNASeqV2_Level_3.1.1.0/	2012-07-27 17:52	919	1312.9
unc.edu_BRCA.IlluminaHiSeq_RNASeqV2_Level_3.1.0.0/	2012-05-18 12:21	858	1226.1

### TCGAquery: Searching old versions

The results from TCGAquery are always the last one from the TCGA data portal. As we have a preprocessing table you should always update TCGAdata.links package. We intent to update the database constantly.

In case you want an old version of the files we have the version parameter that should be a list of tuple values(platform, tumor, version). For example the code below will get the LGG and GBM tumor for platform HumanMethylation450 but for the LGG/HumanMethylation450, we want the version 5 of the files instead of the latest. This might take some seconds

```
query <- TCGAquery(tumor = c("LGG", "GBM"), platform = c("HumanMethylation450"),
                  level = 3,
                  version = list(c("HumanMethylation450", "LGG", 5)))
```

### TCGAquery\_clinic & TCGAquery\_clinicFilter: Working with clinical data

You can retrieve clinical data using the clinic function. The parameters of this function are:

- tumor ("OV", "BRCA", "GBM", etc)
- clinical\_data\_type ("clinical\_patient", "clinical\_drug", etc)

A full list of cancer and clinical data type can be found in the help of the function.

#### 4 Get clinical data

```
clinical_brca_data <- TCGAquery_clinic("brca", "clinical_patient")
clinical_ov_data_bis <- TCGAquery_clinic("ov", "biopsysection_normal_control")
clinical_brca_data_bis <- TCGAquery_clinic("brca", "biopsysection_normal_control")
clinical_brca_data <- TCGAquery_clinic("brca", "clinical_patient")
```

Also, some functions to work with clinical data are provided. For example the function `TCGAquery_clinicFilter` will filter your data, returning the list of barcodes that matches all the filter.

The parameters of `TCGAquery_clinicFilter` are:

- **barcode** List of barcodes
- **clinical\_patient\_data** clinical patient data obtained with `clinc` function. Ex: `clinical_patient_data <- TCGAquery_clinic("LGG","clinical_patient")`
- **HER** her2 neu immunohistochemistry receptor status: "Positive" or "Negative"
- **gender** "MALE" or "FEMALE"
- **PR** Progesterone receptor status: "Positive" or "Negative"
- **stage** Pathologic Stage: "stage\_IX", "stage\_I", "stage\_IIA", "stage\_II", "stage\_III", "stage\_IIIA", "stage\_IIIB", "stage\_IIIC", "stage\_IV"
- **ER** Estrogen receptor status: "Positive" or "Negative"

```
bar <- c("TCGA-BR-8338-02A-11R-1789-07", "TCGA-BR-8338-04A-11A-1789-07",
        "TCGA-BR-8338-06A-11R-1789-07", "TCGA-BR-8338-05A-11R-1789-07",
        "TCGA-BR-8338-11A-11R-1789-07", "TCGA-BR-7836-11A-11R-1789-07",
        "TCGA-BR-7836-04A-11R-1789-07", "TCGA-BR-7836-14A-11R-1789-07",
        "TCGA-BR-7836-02A-11R-1789-07", "TCGA-BR-7836-02A-11R-1789-07",
        "TCGA-BR-7836-11A-11R-1789-07", "TCGA-BR-7836-05A-11R-1789-07",
        "TCGA-BR-7836-10A-11R-1789-07", "TCGA-BR-A1E8-10A-11R-1789-07",
        "TCGA-BR-A1F0-10A-11R-1789-07", "TCGA-BR-A0E2-02A-11R-1789-07",
        "TCGA-BR-A0WY-04A-11R-1789-07", "TCGA-BR-A1F0-04A-11R-1789-08",
        "TCGA-BR-A1F5-04A-11R-2009-00", "TCGA-BR-A0FH-11A-11R-0789-08",
        "TCGA-BR-A2U9-10A-11R-3079-00", "TCGA-BR-A2U8-05A-11R-1789-07",
        "TCGA-BR-A1F8-04A-11R-5789-07", "TCGA-BR-A241-04A-55R-1789-07",
        "TCGA-BR-A0F5-05A-11R-1789-07", "TCGA-BR-A0B4-11A-12R-1789-07",
        "TCGA-BR-A1X5-06A-15R-1789-07", "TCGA-BR-A0J5-01A-11R-1789-07",
        "TCGA-BR-A0F5-01A-11R-1789-07", "TCGA-BR-6336-11A-11R-1789-07",
        "TCGA-BR-6336-11A-11R-1789-07", "TCGA-BR-6336-01A-11R-1789-07",
        "TCGA-BR-6340-01A-11R-1789-07", "TCGA-BR-6340-11A-11R-1789-07")

S <- TCGAquery_SampleTypes(bar,"TP")
S2 <- TCGAquery_SampleTypes(bar,"XS")

# Retrieve multiple tissues types FROM THE SAME PATIENTS
S8 <- TCGAquery_SampleTypes(bar,c("TP","XS"))

# Retrieve multiple tissues types FROM THE SAME PATIENTS
S8S <- TCGAquery_MatchedSampleTypes(bar,c("MT","TP"))

# Get clinical data
clinical_bre_data <- TCGAquery_clinic("breca","clinical_patient")
female_bre_data <- TCGAquery_clinicFilter(bar,clin, HER="Positive", gender="FEMALE", ER="Positive")

The result is shown below:
## ER Positive Samples:
##
##
## HER Positive Samples:
##
##
## GENDER FEMALE Samples:
## TCGA-BR-81F5
```

```

44 TCGA-BH-KLFB
44 TCGA-BH-ACE2
44 TCGA-BG-MONY
44 TCGA-BH-KLFB
44 TCGA-BH-KLFB
44 TCGA-AH-KOPN
44 TCGA-AR-KSLQ
44 TCGA-AR-KSLH
44 TCGA-BH-KLFB
44 TCGA-AR-KGAT
44 TCGA-AJ-KOJ5
44 TCGA-BG-KLKH
44 character(0)

```

TCGAquery\_subtypes: Working with molecular subtypes data.

```

# Check with subtypes from TCGAquery
require(xlsx)
GSK_path_subtypes <- TCGAquery_subtypes(tumor = "gsk", path = "../dataGSK")
GSK_subtypes <- as.data.frame(read.xlsx2(GSK_path_subtypes, 1, stringsAsFactors = FALSE))
GSK_subtypes <- GSK_subtypes[,c(1:10)]
GSK_subtypes <- GSK_subtypes[-1,]
colnames(GSK_subtypes) <- paste(" ", "", as.numeric(GSK_subtypes[,1]))
GSK_subtypes <- GSK_subtypes[-1,]
GSK_subtypes <- GSK_subtypes[!duplicate(GSK_subtypes$sample_id),]
rownames(GSK_subtypes) <- GSK_subtypes$sample_id
dim(GSK_subtypes)
# [1] 208 10

# starting find difference in subtypes
tableSubtypes_filt_GSK <- tableSubtypes_filt[tableSubtypes_filt$tumor_type == "GSK",]
write.table(tableSubtypes_filt_GSK$id, GSK_subtypes$sample_id)
# [1] "TCGA-10-0085"

require(xlsx)
LGG_path_subtypes <- TCGAquery_subtypes(tumor = "lgg", path = "../dataLGG")
LGG_subtypes <- as.data.frame(read.xlsx2(LGG_path_subtypes, 1, stringsAsFactors = FALSE))
rownames(LGG_subtypes) <- LGG_subtypes$tumor
dim(LGG_subtypes)
# [1] 228 9

tableSubtypes_filt_LGG <- tableSubtypes_filt[tableSubtypes_filt$tumor_type == "LGG",]
write.table(tableSubtypes_filt_LGG$id, LGG_subtypes$tumor)
# [228] "TCGA-08-A82S"

LGG_clinic <- TCGAquery_clinic(cancer = "LGG", clinical_data_type = "clinical_patient")
# table(LGG_clinic$id[,mutation_found])

STAD_path_subtypes <- TCGAquery_subtypes(tumor = "stad", path = "../dataSTAD")
STAD_subtypes <- as.data.frame(read.xlsx2(STAD_path_subtypes, 1, stringsAsFactors = FALSE))

```

### TCGAquery\_integrate: Summary of the common numbers of patient samples in different platforms

Some times researches would like to use samples from different platforms from the same patient. In order to help the user to have an overview of the number of samples in common we created the function `TCGAquery_integrate` that will receive the data frame returned from `TCGAquery` and produce a matrix  $n \text{ platforms} \times n \text{ platforms}$  with the values of samples in common.

Some search examples are shown below:

```
query <- TCGAquery(tumor = "breast", level = 3)
matSamples <- TCGAquery_integrate(query)
```

The result of the 3 platforms of `TCGAquery_integrate` result is shown below:

Table 2: Table common samples among platforms from TCGAquery

	AgilentG4502A_07_3	HumanMethylation450	HumanHiSeq_RNASeqV2
AgilentG4502A_07_3	694	224	530
HumanMethylation450	224	930	790
HumanHiSeq_RNASeqV2	530	790	1218

### TCGAquery: some examples

Some search examples are shown below:

```
query <- TCGAquery(tumor = c("gbml","lgg"),
  platform = c("HumanMethylation450","HumanMethylation27"))

query <- TCGAquery(tumor = "gbml", platform = "HumanMethylation450", level = "3")

query <- TCGAquery(samples = "TCGA-61-1743-G1A-00C")

query <- TCGAquery(samples = "TCGA-61-1743-G1A-050-0840-C4", level = 3)

query <- TCGAquery(samples = "TCGA-61-1743-G1A-030-0858-C4",
  tumor = "OV", platform = "TCGA-1x1M_G44426A")
```

### TCGAdownload: Downloading open-access data

You can easily download data using the `TCGAdownload` function.

The arguments are:

- **data** The `TCGAquery` output
- **path** Location to save the files. Default: "."
- **type** Filter the files to download by type
- **samples** List of samples to download
- **force** Download again if file already exists? Default: FALSE



### TCGAdownload: Example of use

```
# get all samples from the query and save them to the TCGA folder.
# samples from IlluminaHiSeq_RNASeq2 with type rsem_genes.results
# samples to normalize later.

TCGAdownload(query, path = "data", type = "rsem_genes.results")

TCGAdownload(query, path = "data", type = "rsem_isoforms.normalized.results")

TCGAdownload(query, path = "dataBcr", type = "rsem_genes.results",
  samples = c("TCGA-EB-A383-11X-R2B-614R-01",
    "TCGA-IB-K110-11X-02B-410Q-01")
)
```

Comment: The function will structure the folders to save the data as: Path given by the user/Experiment folder

### TCGAdownload: Table of types available for downloading

- **RNASeqV2:** junction\_quantification, rsem\_genes.results, rsem\_isoforms.results, rsem\_genes.normalized\_results, rsem\_isoforms.normalized\_results, bt exon\_quantification
- **RNASeq:** rsem\_quantification, spliced\_quantification, gene\_quantification
- **genome\_wide\_snp\_6:** hg18.ssg, hg19.ssg, nochr\_hg18.ssg, nochr\_hg19.ssg

### TCGAprepare: Preparing the data

You can easily read the downloaded data using the TCGAprepare function. This function will prepare the data into a [SummarizedExperiment](#) (Huber, Wolfgang and Carey, Vincent J and Gentleman, Robert and Anders, Simon and Carlson, Marc and Carvalho, Benjamin S and Bravo, Hector Corrado and Davis, Sean and Gatto, Laurent and Giese, Thomas and others 2015) object for downstream analysis. For the moment this function is working only with data level 3.

The arguments are:

- **query** Data frame as the one returned from TCGAquery
- **dir** Directory with the files
- **type** File to prepare.
- **samples** List of samples to prepare.
- **save** Save a rda object with the prepared object? Default: FALSE
- **filename** Name of the rda object that will be saved if save is TRUE
- **toPackage** Name of the package to prepare the data specific to that package.
- **summarizedExperiment** Should the output be a SummarizedExperiment object? Default: TRUE

In order to add useful information to researchers we added in the colData of the summarizedExperiment the subtypes classification for the LGG and GBM samples that can be found in the TCGA publication section. We intend to add more tumor types in the future.

Also in the metadata of the object we added the parameters used in TCGAprepare, the query matrix used for preparing, and file information (name, creation time and modification time) in order to help the user know which samples, versions, and parameters they used.

## TCGAbiolinks: Example of use

```
# get all samples from the query and save them to the TCGA folder.
# samples from IlluminaHiSeq_RNASeqV2 with type raw.genes.normalized
# samples for normalization later.
data <- TCGAprepare(query, dir = "data", save = TRUE, filename = "myfile.rda")
```

As an example, for the platform IlluminaHiSeq\_RNASeqV2 we prepared two samples (TCGA-DY-A105-01A-11R-A155-07 and TCGA-DY-A08A-01A-11R-A155-07) for the raw.genes.normalized\_results type. In order to create the object mapped the gene\_id to the hg18. The genes\_id not found are then removed from the final matrix. The default output is a SummarizedExperiment is shown below.

```
library(TCGAbiolinks)
library(SummarizedExperiment)

## Loading required package: GenomicRanges
## Loading required package: BiocGenerics
## Loading required package: parallel
##
## Attaching package: 'BiocGenerics'
##
## The following objects are masked from 'package:parallel':
##
##   clusterApply, clusterApplyLB, clusterCall, clusterCallB,
##   clusterExport, clusterMap, parApply, parCapply, parLapply,
##   parLapplyLB, parRapply, parSapply, parSapplyLB
##
## The following objects are masked from 'package:stats':
##
##   IQR, mad, rle
##
## The following objects are masked from 'package:base':
##
##   anyDuplicated, append, as.data.frame, as.vector, cbind,
##   colnames, do.call, duplicated, eval, evalq, Filter, Find, get,
##   grep, grepl, intersect, is.unsorted, lapply, Map, mapply,
##   match, mget, order, paste, pmax, pmax.int, pmin, pmin.int,
##   position, rank, rbind, rbindx, rep.int, rownames, sapply,
##   setdiff, sort, table, tapply, union, unique, unlist, unsplit
##
## Loading required package: R6
## Loading required package: stats4
## Loading required package: IRanges
## Loading required package: GenomicData
## Loading required package: Biobase
## Welcome to Biobase
##
## signature contain introductory material; view with
## browseSignatures(). To cite Biobase, see
## citation("Biobase"), and for packages 'citation("pkgname")'.
head(assay(dataREAD, "normalized_counts"))
```

	TCGA-DY-A105-01A-11R-A155-07	TCGA-DY-A08A-01A-11R-A155-07
A1B2.1	13.6782	13.0282
A1CF.26974	53.4379	140.6466
A2Y.2	5630.4792	1061.9360

```
## A2YL11144568      0.0000      18.0001
## AGOAL1153047      170.1189      30.0695
## AGOJ1161146       0.0000      0.0000
```

In order to create the `SummarizedExperiment` object we mapped the rows of the experiments into `GRanges`. In order to map miRNA we used the miRNA from the annotation database `TxDb.Hsapiens.HGSC.hg19.knownGene`, this will exclude the miRNA from viruses and bacteria. In order to map genes, `gene` alias, we used the biomaRt `hg19` database (`hsapiens_gene_ensembl` from `genefilter` package).

In case you prefer to have the raw data. You can get a data frame without any modification setting the `summarizedExperiment` to `false`.

```
library(TCGAbiolinks)
data <- getQuery(query, dataFrame = TRUE)

## [1] "data.frame"
dim(data)
## [1] 20631      2
head(data)
##           TCGA-CY-A1DE-G1A-11R-A155-C7 TCGA-CY-A02A-G1A-11R-A155-C7
## Y100130426                0.0000                0.0000
## Y100133144                11.5308                32.0677
## Y100134360                 4.1574                12.5126
## Y100887                   222.1458                102.8808
## Y10431                    1269.9776                774.5168
## Y106892                   0.0000                0.0000
```

#### TCGAbiolinks: Table of types available for the TCGAprepare

- **RNASeqV2:** `junction_quantification`, `rsam_genes_results`, `rsam_isoforms_results`, `rsam_genes_normalized_results`, `rsam_isoforms_normalized_results`, `txsam_quantification`
- **RNASeq:** `count_quantification`, `splnx_quantification`, `gene_quantification`
- **genome\_wide\_snp\_5:** `hg18seg`, `hg19seg`, `ncore_hg18seg`, `ncore_hg19seg`

#### TCGAbiolinks: Preparing the data with parameter - toPackage

This section will show how to integrate TCGAbiolinks with other packages. Our intention is to provide as many integrations as possible.

The example below shows how to use TCGAbiolinks with ELNEX package (expression/methylation analysis). The TCGAprepare for the DNA methylation data will remove probes with NA values in more than 0.80% samples and remove the annotation data, for the expression data it will take the  $\log_2(\text{expression} + 1)$  of the expression matrix is used. To linearize the relation between DNA methylation and expression also it will present the names as the specified by the package.

```
##### Get tumor samples with TCGAbiolinks
library(TCGAbiolinks)
query <- TCGAquery(tumor = "GBML", level = 3, platform = "HumanMethylation450")
# This function will take a lot of time depends on internet connection
TCGAdownload(query, path = "~/TGA/450k")
res <- TCGAprepare(query, dir = "TCGA/450k",
  name = "GBML",
  filename = "test.rda",
  toPackage = "ELNEX")
```

```

query.rna <- TCGAquery(tumor="GBM",level=3, platform="IlluminaHiSeq_RNASeqV2")
TCGAdownload(query.rna,path="TCGA/rna",type = "rsem.genes.normalized_results")
exp <- TCGAprepare(query.rna, dir="TCGA/rna", save = TRUE,
  filename = "exp.rna",toPackage = "X.MIM")

##### To ENRICH
library(ENRICH)

##### gene annotation
geneAnnot <- txi()
geneAnnot$GENEID <- paste0("ID",geneAnnot$GENEID)
geneInfo <- promoters(geneAnnot,upstream = 0, downstream = 0)
##### probe
probe <- get.feature(probe())

res.glm.glm.with.exp <- fetch.res(glm = glm.glm.m,
  exp = exp,
  probeInfo = probe,
  TCGA = "RNA",
  geneInfo = geneInfo)

```

### TCGAprepare: Preparing the data with CNV data (Genome\_Wide\_SNP\_6)

You can easily search TCGA samples, download and prepare a matrix of gene expression.

# Define a list of samples to query and download providing relative TCGA barcodes.

```

samplesList <- c("TCGA-02-0046-10A-01B-0182-01",
  "TCGA-02-0052-01A-01B-0182-01",
  "TCGA-02-0053-10A-01B-0182-01",
  "TCGA-02-0084-01A-01B-0182-01",
  "TCGA-02-0007-01A-01B-0182-01")

```

# Query platform Genome\_Wide\_SNP\_6 with a list of barcodes

```
query <- TCGAquery(tumor = "gbm", level = 3, platform = "Genome_Wide_SNP_6")
```

# Download a list of barcodes with platform Genome\_Wide\_SNP\_6

```
TCGAdownload(query, path = "samples")
```

# Prepare matrix

```
GBM_CNV <- TCGAprepare(query, dir = "samples", type = ".hg19.exp.txt")
```

### TCGAanalyze: Analyze data from TCGA.

---

You can easily analyze data using following functions:

#### TCGAanalyze\_Preprocessing Preprocessing of Gene Expression data (IlluminaHiSeq\_RNASeqV2).

You can easily search TCGA samples, download and prepare a matrix of gene expression.

# You can define a list of samples to query and download providing relative TCGA barcodes.



```

listSamples <- c("TCGA-B9-A18G-11A-52R-A14P-07", "TCGA-B9-A1FC-11A-52R-A18Q-07",
               "TCGA-A7-A13G-11A-51R-A15G-07", "TCGA-B9-A0DK-11A-13R-A0E9-07",
               "TCGA-B9-A13H-11A-54R-A15S-07", "TCGA-B9-A0AU-01A-11R-A12F-07",
               "TCGA-C9-A1EJ-01A-11R-A15G-07", "TCGA-A7-A13B-01A-13R-A12F-07",
               "TCGA-A2-A0CA-01A-53R-A11S-07", "TCGA-A2-A0YS-01A-11R-A10R-07")

# Query platform IlluminaHiSeq_RNASeqV2 with a list of barcodes
query <- TCGAquery(criteria = "barcode", samples = listSamples,
                  platform = "IlluminaHiSeq_RNASeqV2", level = "S")

# dont run
#TCGAdownload(query, path = "dataBrcs", type = "gene.quantification", samples = listSamples)

# Download a list of barcodes with platform IlluminaHiSeq_RNASeqV2
TCGAdownload(query, path = "../dataBrcs", type = "raw.gene.results", samples = listSamples)

# Prepare expression matrix with gene id in rows and samples (barcodes) in columns
# raw.gene.results as values
BRCAHiSeq assay <- TCGAprepare(query, "../dataBrcs", type = "raw.gene.results")

BRCAMatrix <- assay(BRCAHiSeq assay, "raw_counts")

# For gene expression if you need to see a boxplot correlation and AIC plot
# to define outliers you can run

BRCAHiSeq_correlation <- TCGAanalyze_Preprocessing(BRCAHiSeq assay)

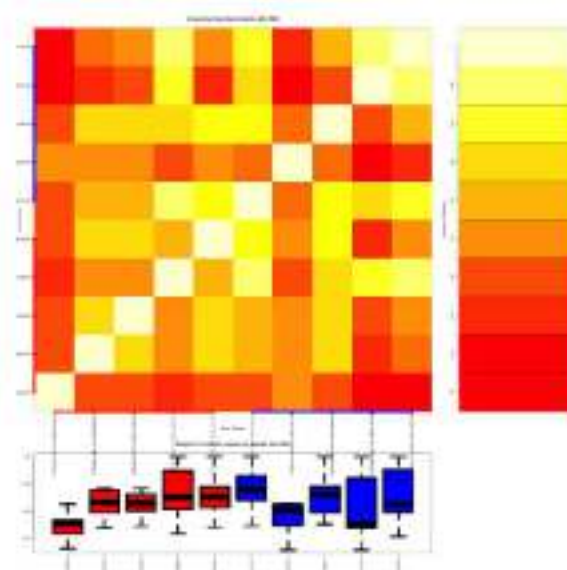
```

The result is shown below:

Table 3: Example of a  
7 samples in columns]

	TCGA-B9-A1FC-11A-52R-A18Q-07	TCGA-A2-A0CA-01A-53R-A11S-07	TCGA-A7-A13G-11A-51R-A15G-07
DPAGT1 1790	1255.21	6195.45	1489.0
PI4KB 5098	5093.00	9025.00	4044.0
KRTAP4-8 728274	0.00	0.00	0.0
DNAC24 120626	655.64	760.78	1073.6
HIRA 7290	1485.00	2014.00	2092.0
ZAP70 7535	145.00	205.00	25.0
SNHG10 293595	139.00	245.00	113.0
MEF2D 4209	2854.00	11776.00	3683.0
INTS5 80789	1329.00	5340.00	1064.0
EIF1B 10289	1722.00	4079.00	2389.0

The result from TCGAanalyze\_Preprocessing is shown below:



### TCGAanalyze\_DEA & TCGAanalyze\_LevelTab Differential expression analysis (DEA)

Perform DEA (Differential expression analysis) to identify differentially expressed genes (DEGs) using the `TCGAanalyze_DEA` function.

`TCGAanalyze_DEA` performs DEA using following functions from R `edgeR`:

1. `edgeR::DGEList` converts the count matrix into an `edgeR` object.
2. `edgeR::estimateCommonDisp` each gene gets assigned the same dispersion estimate.
3. `edgeR::exactTest` performs pair wise tests for differential expression between two groups.
4. `edgeR::topTags` takes the output from `exactTest()`, adjusts the raw p-values using the False Discovery Rate (FDR) correction, and returns the top differentially expressed genes.

This function receives as parameters:

- `mat1` The matrix of the first group (in the example group 1 is the normal samples).
- `mat2` The matrix of the second group (in the example group 2 is tumor samples)
- `Cond1type` Label for group 1
- `Cond2type` Label for group 2

After, we filter the output of `dataDEGs` by `abs(logFC) >= 1`, and uses the `TCGAanalyze_LevelTab` function to create a table with DEGs (differentially expressed genes), log Fold Change (FC), false discovery rate (FDR), the gene expression level for samples in `Cond1type`, and `Cond2type`, and Delta value (the difference of gene expression between the two conditions multiplied logFC).

```
# Eneustrum analysis using gene expression data
# TCGA samples from IlluminaHiSeqRNASeqV2 with type read.gene.results
# save(dataDEGs, geneInfo, file = "dataGeneExpression.rda")
library(TCGAbiolinks)

# Diff. expression analysis (DEA)
```

```
dataDEGs <- TCGAanalyze_DEA(dataFile1[,samplesNT], dataFile1[,samplesTT],
                             "Normal", "Tumor")

# DEGs filter by abs(logFC) >= 1
dataDEGsFilt <- dataDEGs[abs(dataDEGs$logFC) >= 1,]

# DEGs table with expression values in normal and tumor samples
dataDEGsFiltLevel <- TCGAanalyze_LevelTab(dataDEGsFilt, "Tumor", "Normal",
                                           dataFile1[,samplesNT], dataFile1[,samplesTT])
```

The result is shown below:

Table 4: Table DEGs after DEA

mRNA	logFC	FDR	Tumor	Normal	Delta
FN1	2.88	1.295151e-19	347787.48	41234.32	1001017.3
COL1A1	1.77	1.680844e-08	358010.32	89293.72	633086.3
C4orf7	3.20	2.820474e-50	87821.35	2132.70	456425.4
COL1A2	1.40	9.490478e-05	273395.44	91241.32	293242.9
GAPDH	1.32	3.290670e-05	179057.44	69663.00	238255.5
CLEC3A	6.79	7.971033e-74	27257.15	259.60	185150.6
IGFBP5	1.94	1.065717e-04	126186.88	53323.32	158674.6
CPR1	4.27	3.044021e-37	37001.76	2637.72	157068.8
CARTPT	6.72	1.021371e-72	21700.95	215.16	135872.8
DCD	7.26	1.017983e-80	15941.20	81.80	149306.3

### TCGAanalyze\_EAcomplete & TCGAvisualize\_EAbarplot: Enrichment Analysis

Researchers, in order to better understand the underlying biological processes, often want to retrieve a functional profile of a set of genes that might have an important role. This can be done by performing an enrichment analysis.

We will perform an enrichment analysis in gene sets using the `TCGAanalyze_EAcomplete` function. Given a set of genes that are up-regulated under certain conditions, an enrichment analysis will find identify clusters of genes or proteins that are over-represented using annotations for that gene set.

To view the results you can use the `TCGAvisualize_EAbarplot` function as shown below

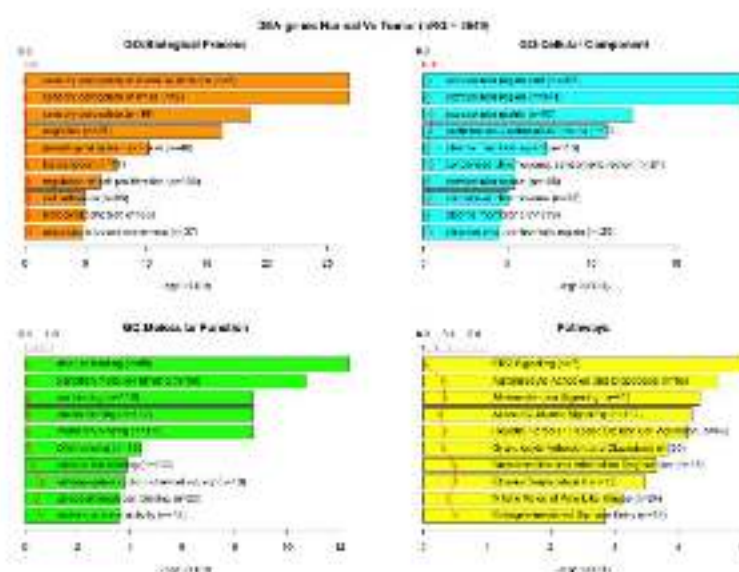
```
library(TCGAbiolinks)
# Enrichment Analysis EA
# Gene Ontology (GO) and Pathway enrichment by gene list
GeneList <- rownames(dataDEGsFiltLevel)

system.time(fEA <- TCGAanalyze_EAcomplete(TFname="DEs genes Normal Vs Tumor", GeneList))

# Enrichment Analysis EA (TCGAvisualize)
# Gene Ontology (GO) and Pathway enrichment barplot

TCGAvisualize_EAbarplot(tc = rownames(dataDEGsFilt),
                        GOETab = ansEA$GOBP,
                        GOCTab = ansEA$GOCC,
                        GEPTab = ansEA$GOEP,
                        PathTab = ansEA$BP,
                        ECCTab = GeneList,
                        xBar = 10)
```

The result is shown below:



### TCGAnalyze\_survival Survival Analysis: Cox Regression and dnet package

When analyzing survival times, different problems come up than the ones discussed so far. One question is how do we deal with subjects dropping out of a study. For example, assume that we test a new cancer drug. While some subjects die, others may believe that the new drug is not effective, and decide to drop out of the study before the study is finished. A similar problem would be faced when we investigate how long a machine lasts before it breaks down.

Using the clinical data, it is possible to create a survival plot with the function `TCGAnalyze_survival` as follows:

```
clin.gbm <- TCGAquery_clinic('gbm', 'clinical_patient')
clin.lgg <- TCGAquery_clinic('lgg', 'clinical_patient')

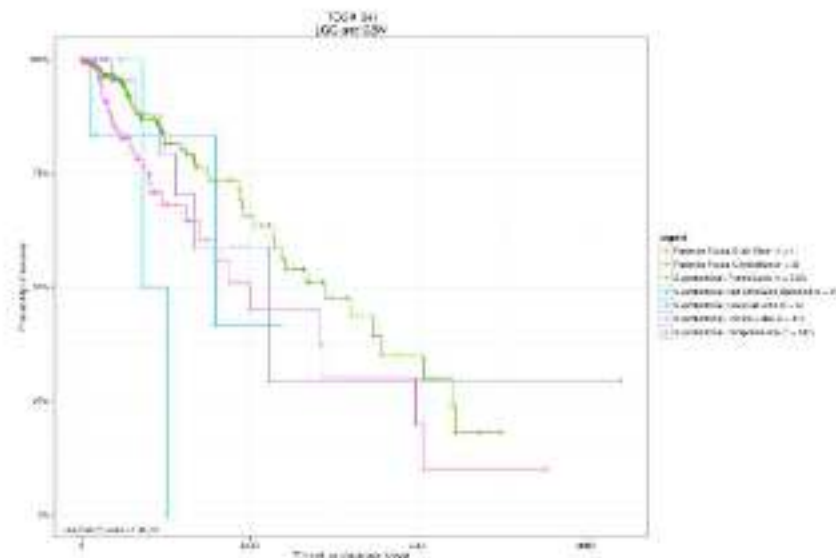
TCGAnalyze_survival(plyr::rbind.fill(clin.lgg, clin.gbm),
  "radiation_therapy",
  main = "TCGA lgg/ovca and gbm", height = 10, width=10)
```

The arguments of `TCGAnalyze_survival` are:

- **clinical\_patient** TCGA Clinical patient with the information `days_to_death`
- **clusterCol** Column with groups to plot. This is a mandatory field, the caption will be based in this column
- **legend** Legend title of the figure
- **cutoff.xlim** This parameter will be a limit in the x axis. That means, that patients with `days_to_death > cutoff` will be set to `Alive`
- **main** main title of the plot
- **ylab** y-axis text of the plot
- **xlab** x-axis text of the plot
- **filename** The name of the pdf file
- **color** Define the colors of the lines

The result is shown below:





```
library(TCGAbiolinks)
# Survival Analysis BR

clinical_patient_Cancer <- TCGAquery_clinic("brca","clinical_patient")
dataBRCAcomplete <- log2(BRCA_rnaseqv2)

tokenStop<- 1

tabSurvSKcomplete <- NULL

for( i in 1: round(nrow(dataBRCAcomplete)/100)){
  message( paste( i, "of ", round(nrow(dataBRCAcomplete)/100)) )
  tokenStart <- tokenStop
  tokenStop <-100+i
  tabSurvSK<-TCGAanalyze_SurvivalKM(clinical_patient_Cancer,dataBRCAcomplete,
                                     SampleList = rownames(dataBRCAcomplete)(tokenStart:tokenStop),
                                     Survresult = F,ThreshUp=0.57,ThreshDown=0.33)

  tabSurvSKcomplete <- rbind(tabSurvSKcomplete,tabSurvKN)
}

tabSurvKcomplete <- tabSurvKcomplete[tabSurvKcomplete$positive < 0.01,]
tabSurvKcomplete <- tabSurvKcomplete[!is.na(tabSurvKcomplete$RNA),]
rownames(tabSurvKcomplete) <- tabSurvKcomplete$RNA
tabSurvSKcomplete <- tabSurvSKcomplete[,-1]
tabSurvSKcomplete <- tabSurvSKcomplete[order(tabSurvKcomplete$positive, decreasing = F),]

tabSurvKcompleteBESs <- tabSurvKcomplete[rownames(tabSurvKcomplete) %in% dataBRCAcomplete$level1$BES,]
The result is shown below:
```

Table 5: Table KM-survival genes after SA

	pvalue	Cancer Deaths	Cancer Deaths with Top	Cancer Deaths with Down	Mean Time Top	Mean Ti
DCTPP1	6.294170e-06	66	45	20	19.31	
APCO	9.390193e-06	65	49	16	11.49	
LOC387646	1.035097e-06	69	48	21	7.97	
PCK1	1.196577e-06	71	49	22	15.66	
CCNE2	2.100342e-06	65	49	17	11.97	
CCDC75	2.920614e-06	74	45	28	9.57	
PCD3	3.035949e-06	69	23	46	12.30	
FAM160B	3.575856e-06	68	25	43	6.82	
MMP28	3.752351e-06	70	17	53	8.55	
ADHFE1	3.907103e-06	67	22	45	9.04	

## TCGAvisualize: Visualize results from analysis functions with TCGA's data.

You can easily visualize results from some following functions:

### TCGAvisualize\_PCA: Principal Component Analysis plot for differentially expressed genes

In order to understand better our genes, we can perform a PCA to reduce the number of dimensions of our gene set. The function `TCGAvisualize_PCA` will plot the PCA for different groups.

The parameters of this function are:

- **dataFilt**: The expression matrix after normalization and quantile filter
- **dataDEGsFiltLevel**: The `TCGAnalyze_LeveFilt` output
- **ntopgenes**: number of DEGs genes to plot in PCA

```
library(TCGAbiolinks)
```

```
# normalization of genes
```

```
dataNorm <- TCGAbiolinks::TCGAnalyze_Normalization(dataBREA, geninfo)
```

```
# quantile filter of genes
```

```
dataFilt <- TCGAnalyze_Filtering(dataNorm, 0.25)
```

```
# Principal Component Analysis plot for ntop selected genes
```

```
TCGAvisualize_PCA(dataFilt, dataDEGsFiltLevel, ntopgenes = 200)
```

```
# boxplot of normalized data
```

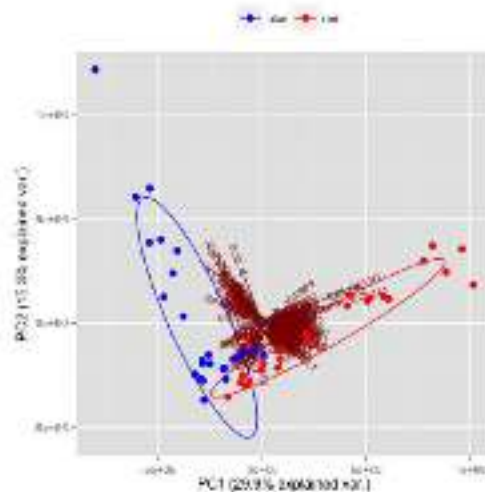
```
#sampleGenes <- rownames(dataFilt)[dataFilt$logFC >= 1,][1:20]
```

```
#boxplot(log(dataBREA[sampleGenes,]), las = 2)
```

```
#boxplot(log(dataFilt[sampleGenes,]), las = 2)
```

The result is shown below:

PCA top 200 Up and down differ genes between Normal vs Tumor



### TCGAvizualize\_SurvivalCoxNET Survival Analysis: Cox Regression and dnet package

TCGAvizualize\_SurvivalCoxNET can help an user to identify a group of survival genes that are significant from univariate Kaplan Meier Analysis and also for Cox Regression. It shows in the end a network build with community of genes with similar range of pvalues from Cox regression (same color) and that interaction among these genes is already validated in literatures using the STRING database (version 9.1).

```
library(TCGAbiolinks)
# Survival Analysis SE

clinical_patient_Cancer <- TCGAquery_clinic("brca","clinical_patient")
dataBRCAcomplete <- log2(BRCA_rawsig2)

tokenStop<- 1

tabSurvComplete <- NULL

for( i in 1: round(nrow(dataBRCAcomplete)/100)){
  message( paste( i, "of ", round(nrow(dataBRCAcomplete)/100)))
  tokenStart <- tokenStop
  tokenStop <- 100+i
  tabSurvSE<-TCGAanalyze_SurvivalEN(clinical_patient_Cancer,
                                   dataBRCAcomplete,
                                   GeneList = rownames(dataBRCAcomplete)(tokenStart:tokenStop),
                                   Survresult = F,ThreshTop=0.67,ThreshDown=0.33)

  tabSurvComplete <- rbind(tabSurvComplete,tabSurvSE)
}
```

```

tabSurvKcomplete <- tabSurvKcomplete[tabSurvKcomplete$psvalue < 0.05,]
tabSurvKcomplete <- tabSurvKcomplete[!duplicated(tabSurvKcomplete$uRNA),]
rownames(tabSurvKcomplete) <- tabSurvKcomplete$uRNA
tabSurvKcomplete <- tabSurvKcomplete[, -1]
tabSurvKcomplete <- tabSurvKcomplete[order(tabSurvKcomplete$psvalue, decreasing=T),]

tabSurvKcompleteDEGs <- tabSurvKcomplete[rownames(tabSurvKcomplete) %in% dataDEGs$H.L.level$uRNA,]

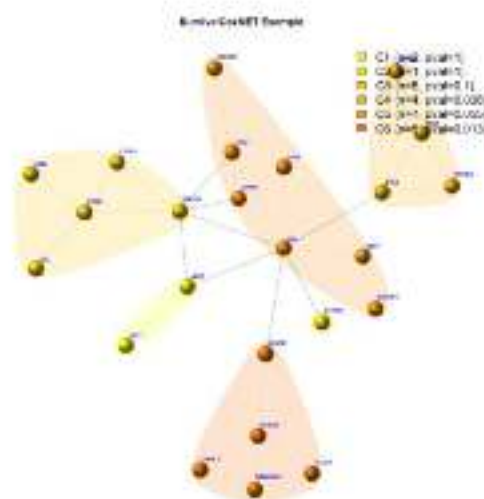
tfList <- KEGenes[KEGenes$Family == "transcription regulator", "Gene"]
tabSurvKcomplete_onlyTF <- tabSurvKcomplete[rownames(tabSurvKcomplete) %in% tfList,]

TabCoaSet <- TCGAvisualize_SurvivalCoaNET(clinical_patient, Cancer, DataSRCcomplete,
                                           CoaList = rownames(tabSurvKcomplete_onlyTF),
                                           scoreConfidence = 700, titlePlot = "TCGAvisualize_SurvivalCoaNET Example")

```

In particular the survival analysis with kaplan mair and Cox regression allow user to reduce the list of genes significant for survival. And using 'best' pipeline with 'TCGAvisualize\_SurvivalCoaNET' function the user can further filter those genes according some already validated interaction according STRING database. This is important because the user can have an idea about the biology inside the survival discrimination and further investigate in a sub group of genes that are working in a synergistic effect influencing the risk of survival. In the following picture the user can see some community of genes with same color and survival profiles.

The result is shown below:



## TCGA Downstream Analysis some workflows and pipelines



## Downstream Analysis n.1

After preparing the gene expression from TCGA data using the `TCGAprepare` function, you can do a normalization of genes using the function `TCGAanalyze_Normalization`, do a quantile filter of genes with the `TCGAanalyze_Filtering` function.

`TCGAanalyze_Normalization` allows user to normalize mRNA transcripts and miRNA, using R `EDASeq` package. Normalization for RNA-Seq: Numerical and graphical summaries of RNA-Seq read data. Within-lane normalization procedures to adjust for GC-content effect (or other gene-level effects) on read counts: loess robust linear regression, global-scaling, and full-quantile normalization (Risso, David and Schwan, Kafaja and Sherlock, Gysin and Dudoit, Sandrine 2011). Between-lane normalization procedures to adjust for distributional differences between lanes (e.g., sequencing depth): global scaling and full-quantile normalization (Bullard, James H and Purdom, Elizabeth and Hansen, Kasper D and O'dot, Sandrine 2016).

For instance returns all mRNA or miRNA with mean across all samples, higher than the threshold defined quantile mean across all samples.

Also, in order to classify your samples (tissue) you can use the `TCGAquery_SampleTypes` function, the typeSample "NT" will return the "Solid Tissue Normal" samples, while the typeSample "TP" will return "Primary Solid Tumor" samples.

```
# Downstream analysis using gene expression data
# TCGA samples from IlluminaHiSeq_RNASeq with type rnae.genes.matrix

library(TCGAbiolinks)

# dataBRCA is TCGAbiolinks package is a table from TCGA BRCA [10 samples] and genes from
# BRCAMatrix <- TCGAprepare(query,"dataBRCA") from above example
# dataBRCA <- BRCAMatrix

# normalization of genes
dataNorm <- TCGAbiolinks::TCGAanalyze_Normalization(dataBRCA, genome)

# quantile filter of genes
dataFilt <- TCGAbiolinks::TCGAanalyze_Filtering(dataNorm, 0.25)

# selection of normal samples "NT"
samplesNT <- TCGAquery_SampleTypes(colnames(dataFilt), typesample = c("NT"))

# selection of tumor samples "TP"
samplesTP <- TCGAquery_SampleTypes(colnames(dataFilt), typesample = c("TP"))
```

## Downstream Analysis n.2 IlluminaHiSeq\_RNASeq data

You can easily search TCGA samples, download and prepare a matrix of gene expression.

```
# Query platform IlluminaHiSeq_RNASeq within a list of barcode
query <- TCGAquery(tumor = "breast", platform = "IlluminaHiSeq_RNASeq", level = "R")

# You can define a list of samples to query and download providing relative TCGA barcodes.
listSamples <- TCGAquery_samplexfilter(query)

# Download only first 6 samples for test.

TCGAdownload(query, path = "dataBRCA", type = "gene.quantification",
              samples = listSamples$IlluminaHiSeq_RNASeq[1:6])
```

```
# Prepare expression matrix with gene id in rows and samples (barcodes) in columns
# rownames results as values
BSCAMatrix <- TCGAprepare(query,'dataBscA',type = "gene quantification")
```

### Downstream Analysis n.3 LGG and GBM Integration (Heatmap and Cluster)

```
library(TCGAbiolinks)
library(genstat)
library(clue)

BSCAmmseqV2 <- dataBSCA
BSCAmmseqV2HeatVar <- varFilter(BSCAmmseqV2, var.func = IQR, var.cutoff = 0.75,
                                rstat.by.quantile = TRUE)

aData <- t(BSCAmmseqV2HeatVar)
ddist <- dist(aData, method = "euclidean")
clhc <- hclust(ddist, method = "ward.D")

plot(knn, labels = FALSE, main = "HCA Hancer cluster dendrogram all samples",
     xlab = "Samples with relative group color", sub = "")

rect.hclust(slc, k=3, border="red")
tabCluster <- as.matrix(cutree(slc, k = 3))
colnames(tabCluster) <- "Cluster"
tabCluster<-cbind(Sample = rownames(tabCluster),Color = rownames(tabCluster), tabCluster)
tabCluster<-as.data.frame(tabCluster)
tabCluster<-tabCluster[order(tabCluster$Cluster,decreasing = FALSE),]
tabCluster<-as.data.frame(tabCluster)
tabCluster$Color<-as.character(tabCluster$Color)

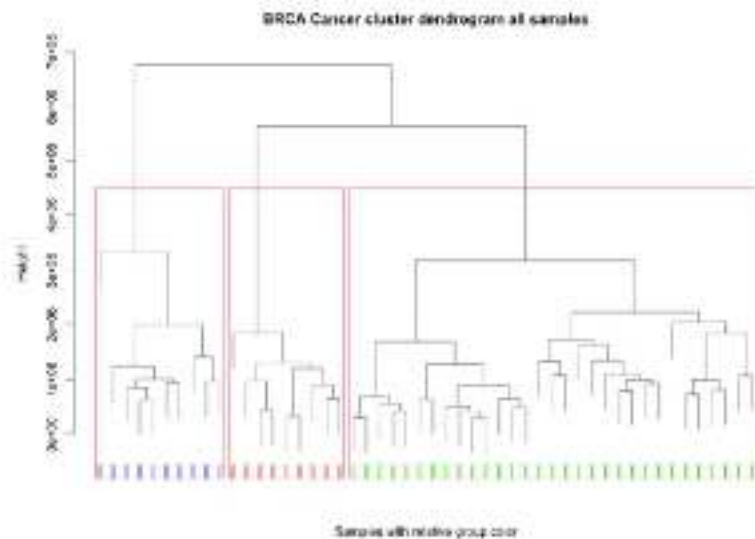
ccol <- palette()[1 : 10]

for(i in 1:3){
  tabCluster[tabCluster[, "Cluster"] == cc, "Color"] <- ccol[i]
}

tabCluster <- tabCluster[sBcLabels, ]

rug(which(tabCluster[sBcLabels, "Color"] == "blue"), col = "blue", lwd = 3)
rug(which(tabCluster[sBcLabels, "Color"] == "green"), col = "green", lwd = 3)
rug(which(tabCluster[sBcLabels, "Color"] == "red"), col = "red", lwd = 3)

The result is shown below:
```



```
library(TCGAbiolinks)

## Differential analysis
GroupBlueData <- BRCAAnaseqV2[, as.character(tabCluster[tabCluster$Color ==
  'blue', 'Sample'])]
GroupGreenData <- BRCAAnaseqV2[, as.character(tabCluster[tabCluster$Color ==
  'green', 'Sample'])]
GroupRedData <- BRCAAnaseqV2[, as.character(tabCluster[tabCluster$Color ==
  'red', 'Sample'])]

DEallData <- TCGAanalyze_DEA(cbind(GroupGreenData, GroupBlueData), GroupBlueData,
  'GroupBlue', 'GroupBlue')
DEallGreen <- TCGAanalyze_DEA(cbind(GroupBlueData, GroupRedData), GroupGreenData,
  'GroupOther', 'GroupGreen')
DEallRed <- TCGAanalyze_DEA(cbind(GroupBlueData, GroupGreenData), GroupRedData,
  'GroupOther', 'GroupRed')

dataDEGs <- TCGAanalyze_DEA(dataFilt[, samplesWT], dataFilt[, samplesTP], "Normal",
  "Tumor")

# DEGs filter by abs(logFC) >=1
dataDEGsFilt <- dataDEGs[abs(dataDEGs$logFC) >= 1, ]

dataDEGsFiltLevel <- TCGAanalyze_LevelTab(dataDEGsFilt, "Tumor", "Normal", dataFilt[,
  samplesTP], dataFilt[, samplesWT])

DEallData.new <- TCGAanalyze_LevelTab(DEallData, "GroupBlue", "GroupOther", GroupBlueData,
```

```

cbind(GroupGreenData, GroupRedData), typeOrder = TRUE)
DEGsGreenLevel <- TCGAanalyze_LevelTab(DEGsGreenS, "GroupGreenS", "GroupOther",
  GroupGreenData, cbind(GroupBlueData, GroupRedData), typeOrder = TRUE)
DEGsRedLevel <- TCGAanalyze_LevelTab(DEGsRed, "GroupRed", "GroupOther", GroupRedData,
  cbind(GroupBlueData, GroupGreenData), typeOrder = TRUE)

blueDEGs <- DEGsBlueLevel[DEGsBlueLevel$FDR < 0.01 & DEGsBlueLevel$logFC >=
  1, ]
blueDEGs <- blueDEGs[order(blueDEGs$FDR), ]
greenDEGs <- DEGsGreenLevel[DEGsGreenLevel$FDR < 0.01 & DEGsGreenLevel$logFC >=
  1, ]
greenDEGs <- greenDEGs[order(greenDEGs$FDR), ]
redDEGs <- DEGsRedLevel[DEGsRedLevel$FDR < 0.01 & DEGsRedLevel$logFC >=
  1, ]
redDEGs <- redDEGs[order(redDEGs$FDR), ]

blueDEGsSpec <- blueDEGs[setdiff(rownames(blueDEGs), union(rownames(greenDEGs),
  rownames(redDEGs))), ]
greenDEGsSpec <- greenDEGs[setdiff(rownames(greenDEGs), union(rownames(blueDEGs),
  rownames(redDEGs))), ]
redDEGsSpec <- redDEGs[setdiff(rownames(redDEGs), union(rownames(blueDEGs),
  rownames(greenDEGs))), ]

blueDEGsSpec <- blueDEGsSpec[1:50, ]
greenDEGsSpec <- greenDEGsSpec[1:50, ]
redDEGsSpec <- redDEGsSpec[1:50, ]

tabCluster <- tabCluster[order(tabCluster$Color), ]

MfiltQuantileOrdered <- BICArnaseqV2(c(rownames(blueDEGsSpec), rownames(greenDEGsSpec),
  rownames(redDEGsSpec)), rownames(tabCluster))

HActivity <- t(MfiltQuantileOrdered)

HActivity <- HActivity
thresholdquantile <- 0.75
HActivity[HActivity >= quantile(HActivity, thresholdquantile)] <- quantile(HActivity,
  thresholdquantile)

summary(as.vector(HActivity))
quantile(HActivity, 0.15)
quantile(HActivity, 0.85)
HActivity[HActivity <= quantile(HActivity, 0.15)] <- quantile(HActivity,
  0.15)
HActivity[HActivity >= quantile(HActivity, 0.85)] <- quantile(HActivity,
  0.85)

column_annotation <- matrix("", nrow = nrow(HActivity), ncol = 1)
column_annotation[, 1] <- tabCluster$Color

row_annotation <- matrix("", nrow = 1, ncol = ncol(HActivity))
row_annotation[, 1] <- c(rep("blue", nrow(blueDEGsSpec)), rep("green",
  nrow(greenDEGsSpec)), rep("red", nrow(redDEGsSpec)))

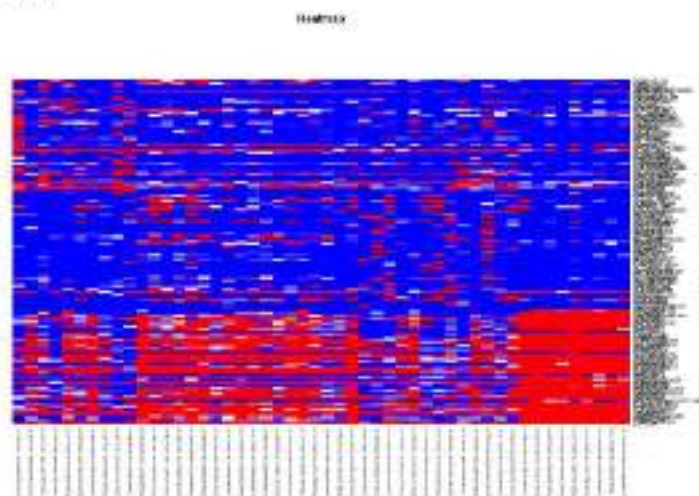
```



```
library("GEO")

png("BRCA_heatmap.png", width = 1000, height = 800)
heatmap.3(t(BPactivity), ColSideColors = column_annotation, RowSideColors = row_annotation,
  key = FALSE, Colv = NA, Rowv = NA,
  main = "none",
  dend = dendrogram(750),
  dendrogram = "none",
  colwidth = 58, labCol = NA,
  margins = c(1, 8), xlab.height.fraction = 0.25, keysize = 1.4, colrow = 1.6)
dev.off()
```

The result is shown below:



#### Downstream Analysis n.4 DNA methylation analysis

Some downstream analysis from DNA methylation data can be done with TCGAdata package. An example is shown below. Firstly, we search, download and prepare data from the HumanMethylation450 platform for the GBM tumor and also get the clinical information from the patients. In this case, we will have a SummarizedExperiment object, where the rows are the probes and the columns the samples. For more information about this object you can take a look in the documentation with the command ?SummarizedExperiment.

```
library(TCGAdataLinks)

# Getting the data
query <- TCGAquery(tumor = "gbm", platform = "HumanMethylation450", level = 2)
TCGAdownload(query, path=".")
data <- TCGAprepare(query, dir = ".", mode = T)
clinical <- TCGAquery_clinic("gbm", "clinical_patient")
```

```

4Preprocessing
# Remove probes with NA level
data <- subset(data, subset=(rowSums(is.na(assay(data)))==0))
As an example, we divided the data into groups in order to analyze the data.
# random split of patients into groups
clinical$group <- c(rep("group1",nrow(clinical)/4),
                    rep("group2",nrow(clinical)/4),
                    rep("group3",nrow(clinical)/4),
                    rep("group4",nrow(clinical)-3*(floor(nrow(clinical)/4))))

colData(data)$group <- c(rep("group1",ncol(data)/2), rep("group2",ncol(data)/2))

```

### TCGAvizualize\_meanMethylation: Sample Mean DNA Methylation Analysis

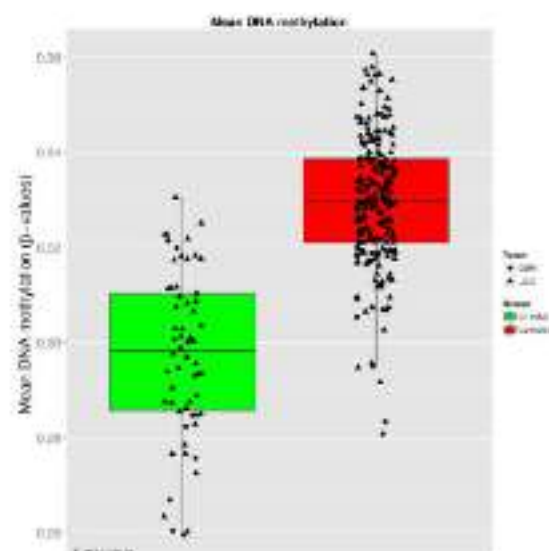
Using the data and calculating the mean DNA methylation per group, it is possible to create a mean DNA methylation heatmap with the function `TCGAvizualize_meanMethylation` as follows:

```
TCGAvizualize_meanMethylation(data,"group")
```

The arguments of `TCGAvizualize_meanMethylation` are:

- **data** SummarizedExperiment object obtained from `TCGAprepape`
- **groupCol** Column in `colData(data)` that defines the groups. If no column defined a column called "Patients" will be used
- **subgroupCol** Column in `colData(data)` that defines the subgroups.
- **shapes** Shape vector of the subgroups. It must have the size of the levels of the subgroups. Example: `shapes = c(21,23)` if for two levels
- **filename** The name of the pdf that will be saved
- **subgroup.legend** Name of the subgroup legend. **DEFAULT:** `subgroupCol`
- **group.legend** Name of the group legend. **DEFAULT:** `groupCol`
- **color** vector of colors to be used in graph
- **title** main title in the plot
- **ylab** y axis text in the plot
- **print.pvalue** Print p value for two groups in the plot
- **slab** x axis text in the plot
- **labels** Labels of the groups

The result is shown below:



### TCGAanalyze\_DMR: Differentially methylated regions Analysis

We will search for differentially methylated CpG sites using the `TCGAanalyze_DMR` function. In order to find these regions we use the beta values (methylation values ranging from 0.0 to 1.0) to compare two groups.

Firstly, it calculates the difference between the mean DNA methylation of each group for each probes.

Secondly, it calculates the p value using the wilcoxon test adjusting by the Benjamini-Hochberg method. The default parameters was set to require a minimum absolute beta values difference of 0.2 and a p value adjusted of  $< 0.01$ .

After these analysis, we save a volcano plot (x axis: diff mean methylation, y axis: significance) that will help the user identify the differentially methylated CpG sites and return the object with the calculus in the `rowRanges`.

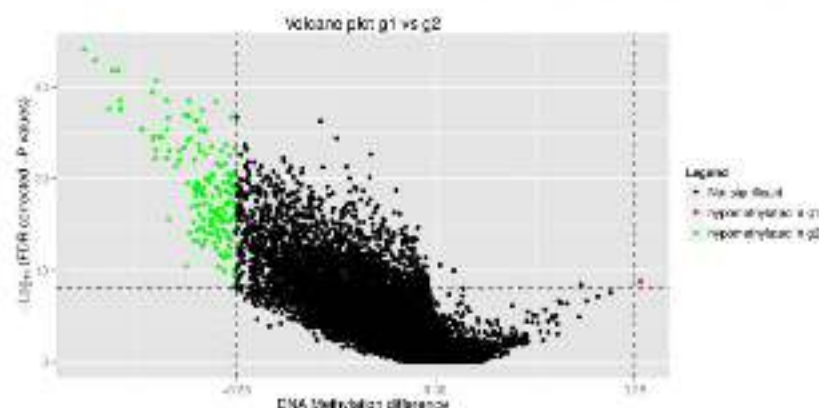
The arguments of `volcanoPlot` are:

- **data** SummarizedExperiment obtained from the `TCGAPrepare`
- **groupCol** Columns with the groups inside the SummarizedExperiment object. (This will be obtained by the function `colData(data)`)
- **group1** In case our object has more than 2 groups, you should set the name of the group
- **group2** In case our object has more than 2 groups, you should set the name of the group
- **filename** pdf filename. Default: volcano.pdf
- **legend** Legend title
- **color** vector of colors to be used in graph
- **title** main title. If not specified it will be "Volcano plot (group1 vs group2)"
- **ylab** y axis text
- **xlab** x axis text
- **xlim** x limits to cut image
- **ylim** y limits to cut image
- **label** vector of labels to be used in the figure. Example: c("1" = "Not Significant", "2" = "Hypermethylated in group1", "3" = "Hypermethylated in group2")

- `p.cut` p-value threshold. Default: 0.01
- `diffmean.cut` diffmean threshold. Default: 0.2
- `adj.method` Adjusted method for the p-value calculation
- `paired` Wilcoxon paired parameter. Default: FALSE
- `overwrite` Overwrite the p-values and diffmean values if already in the object for both groups? Default: FALSE

```
data <- TCGAanalyze_DMR(data, groupCol = "cluster_seth", subgroupCol = "disease",
  group.legend = "Groups", subgroup.legend = "Tumor",
  print.pvalue = TRUE)
```

The output will be a plot such as the figure below. The green dots are the probes that are hypomethylated in group 1 compared to group 2, while the red dots are the hypermethylated probes in group 1 compared to group 2.



Also, the `TCGAanalyze_DMR` function will save the plot as pdf and return the same `SummarizedExperiment` that was given as input with the values of p-value, p-value adjusted, diffmean and the group it belongs in the graph (non significant, hypomethylated, hypermethylated) in the rowRanges. The columns will be (where group1 and group2 are the names of the groups):

- `diffmean.group1.group2`
- `p.value.group1.group2`
- `p.value.adj.group1.group2`
- `status.group1.group2`

This values can be view/accessed using the `rowRanges` accessor (`rowRanges(data)`).

**Observations:** Calling the same function again, with the same arguments will only plot the results, as it was already calculated. With you want to have them recalculated, please set `overwrite` to TRUE or remove the calculated columns.

### TCGAvisualize\_starburst: Analyzing expression and methylation together

The starburst plot is proposed to combine information from two volcano plots, and is applied for a study of DNA methylation and gene expression. In order to reproduce this plot, we will use the `TCGAvisualize_starburst` function.

The function creates Starburst plot for comparison of DNA methylation and gene expression. The  $\log_{10}$  (FDR corrected P value) is plotted for beta value for DNA methylation (x axis) and gene expression (y axis) for each gene. The black dashed line shows the FDR-adjusted P value of 0.01.

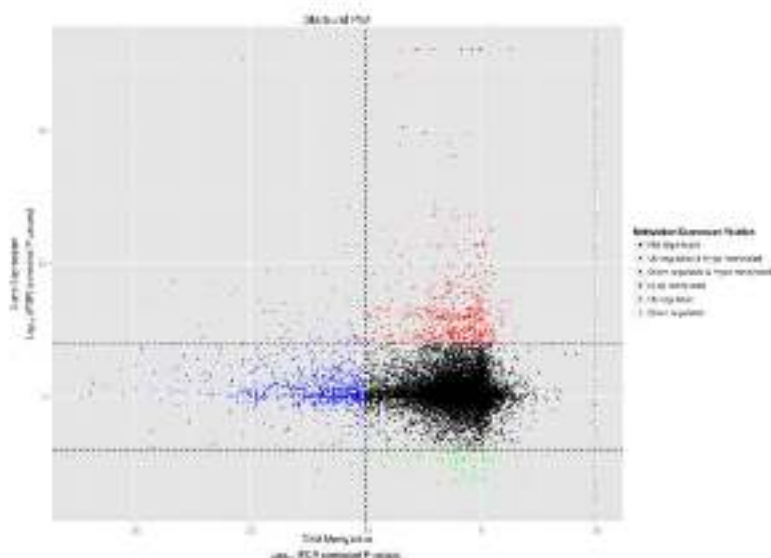
The parameters of this function are:



- **met**: SummarizedExperiment with methylation data obtained from the TCGAprepare and processed by TCGAanalyze\_BMR function. Expected colData columns: diffmean and pvalue.ac
- **exp**: Matrix with expression data obtained from the TCGAanalyze\_BER function. Expected colData columns: logFC, FDR
- **filename**: pdf filename
- **legend**: legend title
- **color**: vector of colors to be used in graph
- **label**: vector of labels to be used in graph
- **title**: main title
- **ylab**: y axis text
- **xlab**: x axis text
- **xlim**: x limits to cut image
- **ylim**: y limits to cut image
- **p.cut**: p value cut off
- **group1**: The name of the group 1. Obs: Column pvalue.ac|group1|group2 should exist
- **group2**: The name of the group 2. Obs: Column pvalue.ac|group1|group2 should exist

```
result <- TCGAanalyze_htlinks(met, exp, "g1", "g2", p.cut = 0.02)
```

As result the function will plot the figure below and return a matrix with The Gene\_symbol and its status in relation to expression (up regulated/down regulated) and methylation (Hyper/Hypo methylated).



## TCGAinvestigate: Searching questions, answers and literature

### TCGAinvestigate: Find most studied TFs in pubmed

Find most studied TFs in pubmed related to a specific cancer, disease, or tissue

```
# First perform DEGs with TCGAanalyze
# See previous section
library(TCGAdata)

# Select only transcription factors (TFs) from DEGs
TFs <- RAGenes[isGeneSymbol == 'transcription regulator',]
TFs_inDEGs <- intersect(TFsGenes, dataDEGsFiltLevel$GENE)
dataDEGsFiltLevelTFs <- dataDEGsFiltLevel[TFs_inDEGs,]

# Order table DEGs TFs according to Delta decrease
dataDEGsFiltLevelTFs <- dataDEGsFiltLevelTFs[order(dataDEGsFiltLevelTFs$Delta, decreasing = TRUE),]

# Find Pubmed of TF studied related to cancer
tblDEGsTFPubmed <- TCGAinvestigate("breast", dataDEGsFiltLevelTFs, topgenes = 10)

The result is shown below:
```

Table 6: Table with more studied TF in pubmed related to a specific cancer

mRNA	logFC	FDR	Tumor	Normal	Delta	Pubmed	PMID
MUC7	2.46	0	38436.56	6469.40	94523.36	897	26016502; 25589764; 25882681; 25873671;
FGS	-2.46	0	14080.32	66543.04	34627.41	313	26011749; 25556606; 25824586; 25788836;
MDM2	1.41	0	16132.29	6959.92	22024.14	441	26042562; 26001071; 25814109; 25893170;
GATA3	1.18	0	29394.60	8304.72	66410.03	180	26026330; 26008346; 25994056; 25905123;
FOXA1	1.45	0	16136.96	5378.58	23465.63	167	26008346; 25995731; 25994056; 25762479;
ECR1	-2.44	0	16073.08	74547.58	30275.29	77	25703326; 24080816; 24742492; 24675510;
TCF1	1.43	0	17765.96	6260.08	26470.39	13	25798844; 25889165; 25162636; 21937081;
MAGEB1	1.18	0	26850.16	8244.32	24633.09	5	24225485; 23884293; 22935435; 21618623;
PTTG	1.72	0	19220.12	44292.32	26194.62	5	25945513; 23216712; 21913217; 20427976;
ILF2	1.27	0	22250.32	7354.44	26245.23	0	0

## TCGAassoc: Searching questions, answers and literature

The TCGAassoc function has two types of searches, one that searches for most downloaded packages in CRAN or BioConductor and one that searches the most related question in biostar.

### TCGAassoc with BioConductor

Find most downloaded packages in CRAN or BioConductor

```
library(TCGAdata)

# Define a list of package to find number of downloads
listPackage <- c("limma", "edgeR", "survcomp")

tblPackage <- TCGAassoc(siteToFind = "bioconductor.org", listPackage)

# define a keyword to find in support.bioconductor.org returning a table with suggested packages
tblPackageKey <- TCGAassoc(siteToFind = "support.bioconductor.org", keyInfo = "key")

The result is shown below:
```

Table 7: Table with number of downloads about a list of packages

Package	NumberDownload
limma	70789
edgeR	33534
survcomp	3661

Table 8: Find most related question in support bioconductor.org with keyword = tugs

question	BiostarSite	PackageSuggested
A: Calculating Ibd Using R Package	/55481/	TIN
A: How To Identify Rotamer States From A Pdb ?	/96579/	SIM
A: Pathway Analysis In R	/14316/	sgPathway
A: Mys Question – Consensus	/17535/	sgPathway
A: How to read data file in Rsumtools R package?	/97978/	Rsumtools
A: Best Practices/Software To Calculate Ka/Ka Ratio	/5817/	ks
A: Trouble With Local Psblast	/79246/	ks
A: R Package For Annotations Of Genomic Regions	/43313/	ks
A: Question About Mafin Methylation Array	/89357/	LEA;MEDIPS
A: Find Out The Genes That Correspond To My Coordinates	/47526/	ChIPpeakAnno
Mina Sequents Using Biomart R Package	/96700/	biomaRt
A: Annotating Expression Profile Data	/60594/	AnnotationDbt;LEA
A: How to generate a Venn diagram	/102393	0
A: CNV calling for illumina 50k array	/109029	0
A: Error could not find function "headmap.2"	/105843	0
A: Extracting ProbeSet IDs from CEL files	/135042	0
A: Bam to nucleotide frequencies	/100798	0
A: Gene Regulatory Network using micro array data	/121070	0
A: R programming question: insert alternately	/139129	0
A: Ignoring NA on Each Side of the Chromosome	/146513	0
** MISSING **	NA	0
A: GCseq rat genome	/139732	0

## TCGAbiolinks with Biostar

Find most related question in biostar

```
library(TCGAbiolinks)
```

```
# Find most related question in biostar with TCGA
```

```
tabPackage1 <- TCGAbiolinks::getInfoFromBio("biostars.org", keyInfo = "TCGA")
```

```
# Find most related question in biostar with package
```

```
tabPackage2 <- TCGAbiolinks::getInfoFromBio("biostars.org", keyInfo = "package")
```

The result is shown below:

Table 9: Find most related question in biostar with TCGA

question	BiostarSite	PackageSuggested
A: Question About Toga Snp Array Data	/89141/	LEA;PROcess;ROC

question	BioStarsSite	PackageSuggested
A: Cnv Data	/95763/	DNACopy,HELP
A: Cnv Data	/95763/	DNACopy,HELP
A: Where To Find Test Datasets For Data Classification Problems	/60664/	convert,GEOquery,LEA,rMAT,roar,SIM
A: How to get public cancer RNA seq data?	/134370/	0
A: Microarray And Epigenomic Data For Same Cancer Cell Line?	/95224/	0

Table 10: Find more related question in bioStar web package

question	BioStarsSite	PackageSuggested
A: Calculating Irid Using R Package	/55481/	TIN
A: How To Identify Refactor Status From A Path?	/66579/	SIM
A: Pathway Analysis In R	/14326/	sigPathway
A: Ngs Question -- Consensus	/17535/	sigPathway
A: How to read .bam file in Rsamtools R package?	/54948/	Rsamtools
A: Fast Practices/Softwares To Calculate Ka/Ka Ratio	/58174/	ks
A: Trouble With Local Filecast	/79246/	ks
A: R Package For Annotations Of Circos Plots	/43313/	ks
A: Question About Medip Methylation Array	/89357/	LEA,MEDIPS
A: Find Out The Genes That Correspond To My Coordinates	/47826/	ChIPpeakAnno
Mini Sequences Using BiomaRt R Package	/96706/	biomaRt
A: Annotating Expression Profile Data	/60564/	AnnotationDbi,LEA
A: How to generate a Venn diagram	/102393/	0
A: CNV calling for Illumina bulk array	/109029/	0
A: Error: could not find function "heatmap.2"	/105843/	0
A: Extracting ProbeSet IDs from .CEL files	/135842/	0
A: Sum to nucleotide frequencies	/163798/	0
A: Gene Regulatory Network using micro array data	/121070/	0
A: R programming question: insert alternately	/139129/	0
A: Ignoring NAs on Each Side of the Chromosome	/146513/	0
*** MISSING ***	NA	0
A: RCase rat genomic	/135732/	0

## Session Information

```

sessioninfo()
## R version 3.2.1 (2015-05-16)
## Platform: x86_64-apple-darwin13.4.0 (64-bit)
## Running under: OS X 10.10.4 (Yosemite)
##
## locale:
## [1] pt_BR.UTF-8/pt_BR.UTF-8/pt_BR.UTF-8/C/pt_BR.UTF-8/pt_BR.UTF-8
##
## attached base packages:
## [1] grid      stats4     parallel  stats      graphics  grDevices  utils
## [8] datasets  methods    base
##

```



```

## other attached packages:
## [1] png_0.1-7           SummarizedExperiment_0.8.2
## [3] Biobase_2.29.1       GenomicRanges_1.21.17
## [5] GenomicInfoDb_1.5.9   IRanges_2.3.17
## [7] RAVenture_0.7.12      RmcExperiment_0.15.5
## [9] UCSCChromLink_0.99.1  RinnStyle_1.2.6
##
## loaded via a namespace (and not attached):
## [1] atom_1.0.0
## [2] edgeR_3.11.2
## [3] splines_3.2.1
## [4] R.utils_2.1.0
## [5] highr_0.6
## [6] azoom.light_2.0.2
## [7] latticeExtra_0.6-25
## [8] xlsxjars_0.8.1
## [9] coin_1.0-24
## [10] Rsumtools_1.20.14
## [11] yaml_2.1.13
## [12] RSQLite_1.0.0
## [13] lattice_0.20-33
## [14] limma_5.25.14
## [15] downloader_0.4
## [16] chrom_2.3-47
## [17] digest_0.6.9
## [18] RColorBrewer_1.1-2
## [19] Rwasm_0.9.1
## [20] rstat_0.2.0
## [21] colorspace_1.2-6
## [22] Matrix_1.2-2
## [23] rtracktools_0.2.6
## [24] R.oo_1.19.0
## [25] plyr_1.8.3
## [26] XML_3.96-1.3
## [27] devtools_1.9.0
## [28] ShortRead_1.27.5
## [29] BiocSnp_2.25.1
## [30] genefilter_1.51.0
## [31] xliblsc_1.15.0
## [32] xtable_1.7-4
## [33] rvmnorm_1.0-3
## [34] xlsxio_0.2.5
## [35] supraRex_1.7.2
## [36] BioParallel_1.3.47
## [37] gis2r_0.10.1
## [38] annotate_1.47.4
## [39] ggplot2_1.0.1
## [40] GenomicFeatures_1.21.13
## [41] hexbin_1.27.0
## [42] proto_0.2-10
## [43] survival_2.38-3
## [44] magrittr_1.5
## [45] reshape_0.2.1
## [46] evaluate_0.7

```

```

## [47] GGalaxy_0.5.0
## [48] R.methodsS3_1.7.0
## [49] nlme_3.1-121
## [50] MASS_7.3-43
## [51] xlsx_0.1.1
## [52] rpart_4.3.2
## [53] graph_1.47.2
## [54] test_3.2.1
## [55] data.table_1.5.6
## [56] formatR_1.2
## [57] matrixStats_0.14.2
## [58] stringi_1.0.0
## [59] xlsx_0.6.7
## [60] runsoll_0.4.2
## [61] AnnotationDbi_1.31.17
## [62] limma_3.11.7
## [63] rversions_1.0.0
## [64] Biostrings_2.37.2
## [65] DESeq_1.21.0
## [66] futile.logger_1.4.1
## [67] RUnit_1.95-4.7
## [68] rJava_0.2.15
## [69] igraph_1.0.1
## [70] bitops_1.0-5
## [71] rtracktools_0.7.1
## [72] dplyr_1.0.7
## [73] gtable_0.1.2
## [74] DEI_0.5.1
## [75] reshape_0.3.5
## [76] rsvgs2_0.1.1
## [77] curl_0.8.1
## [78] R6_2.1.0
## [79] reshape2_1.4.1
## [80] DESeq2_2.3.2
## [81] GenomicAlignments_1.6.12
## [82] knitr_1.10.5
## [83] rtracklayer_1.28.13
## [84] futile.options_1.0.0
## [85] Rgraphviz_2.15.0
## [86] ape_3.3
## [87] rJava_0.9-7
## [88] Tcdh.Hsapiens.GRCh38.knownGene_3.1.3
## [89] rstattools_0.2-21
## [90] stringr_0.5-5
## [91] Rcpp_0.12.0
## [92] ggplot2_1.47.0

```

## References

- Bullard, James H and Purdom, Elizabeth and Hansen, Kasper D and Dudoit, Sandrine. 2010. "Evaluation of Statistical Methods for Normalization and Differential Expression in mRNA Seq Experiments."
- Huay, Wolfgang and Carey, Vincent J and Gattaman, Robert and Anders, Simon and Carbon, Marc and Carvalho,

Berlton, S and Brown, Hector, Corrada and Davis, Sean and Gatto, Laurent and Grise, Thomas and others. 2015. "Orchestrating High-Throughput Genomic Analysis with Bioconductor."

Fries, Davide and Schwartz, Katja and Sherlock, Gavin and Dudar, Sanderine. 2011. "GC-Content Normalization for RNA-Seq Data."