

# curatedCRCData

*Princy Parsana, Markus Riester, Curtis Huttenhower, Levi Waldron*

2013

## Contents

1	curatedCRCData: Clinically Annotated Data for the colorectal Cancer Transcriptome . . . . .	1
2	Load data sets . . . . .	1
3	Load datasets based on rules. . . . .	2
4	Non-unique gene symbols. . . . .	6
A	Available Clinical Characteristics . . . . .	7
B	Summarizing the List of ExpressionSets . . . . .	7
C	For non-R users . . . . .	10
D	Session Info . . . . .	10

## 1 curatedCRCData: Clinically Annotated Data for the colorectal Cancer Transcriptome

---

This package represents a manually curated data collection for gene expression meta-analysis of patients with colorectal cancer. This resource provides uniformly prepared microarray data with curated and documented clinical metadata. It allows a computational user to efficiently identify studies and patient subgroups of interest for analysis and to run such analyses immediately without the challenges posed by harmonizing heterogeneous microarray technologies, study designs, expression data processing methods, and clinical data formats.

In this vignette, we give a short tour of the package and will show how to use it efficiently.

## 2 Load data sets

---

Loading a single dataset is very easy. First we load the package:

## curatedCRCDData

```
> library(curatedCRCDData)
```

To get a listing of all the datasets, use the `data` function:

```
> data(package="curatedCRCDData")
```

Now to load a single dataset, we use the `data` function again:

```
> data(TCGA.COAD_eset)
> TCGA.COAD_eset

ExpressionSet (storageMode: lockedEnvironment)
assayData: 17814 features, 130 samples
  element names: exprs
protocolData: none
phenoData
  sampleNames: TCGA.AA.3520 TCGA.AA.3532 ... TCGA.A6.2685 (130 total)
  varLabels: unique_patient_ID alt_sample_name ...
  uncurated_author_metadata (59 total)
  varMetadata: labelDescription
featureData
  featureNames: 15E1.2 2'-PDE ... ZZZ3 (17814 total)
  fvarLabels: probeset gene
  fvarMetadata: labelDescription
experimentData: use 'experimentData(object)'
pubMedIds: 22810696
Annotation: agilent-014850 whole human genome microarray 4x44k g4112f
```

The datasets are provided as Bioconductor `ExpressionSet` objects and we refer to the Bioconductor documentation for users unfamiliar with this data structure.

## 3 Load datasets based on rules

For a meta-analysis, we typically want to filter datasets and patients to get a population of patients we are interested in. We provide a short but powerful R script that does the filtering and provides the data as a list of `ExpressionSet` objects. One can use this script within R by first sourcing a config file which specifies the filters, like the minimum numbers of patients in each dataset. It is also possible to filter samples by annotation, for example to remove early stage and normal samples.

```
> source(system.file("extdata",
+ "patientselection_all.config", package="curatedCRCDData"))
> ls()

[1] "TCGA.COAD_eset"      "keep.common.only"   "meta.required"
[4] "min.number.of.events" "min.sample.size"    "quantile.cutoff"
[7] "rescale"             "strict.checking"
```

## curatedCRCDData

See what the values of these variables we have loaded are. The variable names are fairly descriptive, but note that “rule.1” is a character vector of length 2, where the first entry is the name of a clinical data variable, and the second entry is a Regular Expression providing a requirement for that variable. Any number of rules can be added, with increasing identifiers, e.g. “rule.2”, “rule.3”, etc.

Here `strict.checking` is `FALSE`, meaning that samples not annotated for the variables in these rules are allowed to pass the filter. If `strict.checking == TRUE`, samples missing this annotation will be removed.

```
> sapply(ls(), get)

$TCGA.COAD_eset
ExpressionSet (storageMode: lockedEnvironment)
assayData: 17814 features, 130 samples
  element names: exprs
protocolData: none
phenoData
  sampleNames: TCGA.AA.3520 TCGA.AA.3532 ... TCGA.A6.2685 (130 total)
  varLabels: unique_patient_ID alt_sample_name ...
    uncurated_author_metadata (59 total)
  varMetadata: labelDescription
featureData
  featureNames: 15E1.2 2'-PDE ... ZZZ3 (17814 total)
  fvarLabels: probeset gene
  fvarMetadata: labelDescription
experimentData: use 'experimentData(object)'
pubMedIds: 22810696
Annotation: agilent-014850 whole human genome microarray 4x44k g4112f

$keep.common.only
[1] FALSE

$meta.required
NULL

$min.number.of.events
[1] 0

$min.sample.size
[1] 1

$quantile.cutoff
[1] 0

$rescale
[1] FALSE

$strict.checking
[1] FALSE
```

Now that we have defined the sample filter, we create a list of `ExpressionSets` by sourcing the `createEsetList.R` file:

## curatedCRCDData

```
> source(system.file("extdata", "createEsetList.R", package =
+ "curatedCRCDData"))

2020-10-29 10:33:28 INFO::Inside script createEsetList.R - inputArgs =
2020-10-29 10:33:28 INFO::None provided
2020-10-29 10:33:28 INFO::Loading curatedCRCDData 2.22.0
2020-10-29 10:34:20 INFO::Clean up the esets.
2020-10-29 10:34:20 INFO::including GSE11237_eset
2020-10-29 10:34:20 INFO::including GSE12225.GPL3676_eset
2020-10-29 10:34:20 INFO::including GSE12945_eset
2020-10-29 10:34:21 INFO::including GSE13067_eset
2020-10-29 10:34:21 INFO::including GSE13294_eset
2020-10-29 10:34:21 INFO::including GSE14095_eset
2020-10-29 10:34:21 INFO::including GSE14333_eset
2020-10-29 10:34:21 INFO::including GSE16125.GPL5175_eset
2020-10-29 10:34:21 INFO::including GSE17536_eset
2020-10-29 10:34:21 INFO::including GSE17537_eset
2020-10-29 10:34:21 INFO::including GSE17538.GPL570_eset
2020-10-29 10:34:22 INFO::including GSE18105_eset
2020-10-29 10:34:22 INFO::including GSE2109_eset
2020-10-29 10:34:22 INFO::including GSE21510_eset
2020-10-29 10:34:22 INFO::including GSE21815_eset
2020-10-29 10:34:22 INFO::including GSE24549.GPL5175_eset
2020-10-29 10:34:22 INFO::including GSE24550.GPL5175_eset
2020-10-29 10:34:22 INFO::including GSE2630_eset
2020-10-29 10:34:22 INFO::including GSE26682.GPL570_eset
2020-10-29 10:34:23 INFO::including GSE26682.GPL96_eset
2020-10-29 10:34:23 INFO::including GSE26906_eset
2020-10-29 10:34:23 INFO::including GSE27544_eset
2020-10-29 10:34:23 INFO::including GSE28702_eset
2020-10-29 10:34:23 INFO::including GSE3294_eset
2020-10-29 10:34:23 INFO::including GSE33113_eset
2020-10-29 10:34:23 INFO::including GSE39582_eset
2020-10-29 10:34:23 INFO::including GSE3964_eset
2020-10-29 10:34:23 INFO::including GSE4045_eset
2020-10-29 10:34:24 INFO::including GSE4526_eset
2020-10-29 10:34:24 INFO::including GSE45270_eset
2020-10-29 10:34:24 INFO::including TCGA.COAD_eset
2020-10-29 10:34:24 INFO::including TCGA.READ_eset
2020-10-29 10:34:24 INFO::including TCGA.RNASeqV2.READ_eset
2020-10-29 10:34:24 INFO::including TCGA.RNASeqV2_eset
2020-10-29 10:34:24 INFO::Ids with missing data: GSE2630_eset, GSE3294_eset, TCGA.COAD_eset, TCGA.READ_eset
```

It is also possible to run the script from the command line and then load the R data file within R:

```
R --vanilla "--args patientselection.config crc.eset.rda tmp.log" < createEsetList.R
```

Now we have 34 datasets with samples that passed our filter in a list of `ExpressionSets` called `esets`:

## curatedCRCData

```
> names(esets)

 [1] "GSE11237_eset"           "GSE12225.GPL3676_eset"
 [3] "GSE12945_eset"          "GSE13067_eset"
 [5] "GSE13294_eset"          "GSE14095_eset"
 [7] "GSE14333_eset"          "GSE16125.GPL5175_eset"
 [9] "GSE17536_eset"          "GSE17537_eset"
[11] "GSE17538.GPL570_eset"    "GSE18105_eset"
[13] "GSE2109_eset"           "GSE21510_eset"
[15] "GSE21815_eset"          "GSE24549.GPL5175_eset"
[17] "GSE24550.GPL5175_eset"  "GSE2630_eset"
[19] "GSE26682.GPL570_eset"   "GSE26682.GPL96_eset"
[21] "GSE26906_eset"          "GSE27544_eset"
[23] "GSE28702_eset"          "GSE3294_eset"
[25] "GSE33113_eset"          "GSE39582_eset"
[27] "GSE3964_eset"           "GSE4045_eset"
[29] "GSE4526_eset"           "GSE45270_eset"
[31] "TCGA.COAD_eset"         "TCGA.READ_eset"
[33] "TCGA.RNASeqV2.READ_eset" "TCGA.RNASeqV2_eset"
```

## 4 Non-unique gene symbols

In the standard version of `curatedCRCDData` (the version available on Bioconductor), we collapse manufacturer probesets to official HGNC symbols using the Biomart database. Some probesets are mapped to multiple HGNC symbols in this database. For these probesets, we provide all the symbols. For example `220159_at` maps to `ABCA11P` and `ZNF721` and we provide `ABCA11P///ZNF721` as probeset name. If you have an array of gene symbols for which you want to access the expression data, "ABCA11P" would not be found in `curatedCRCDData` in this example. The following function will create a new `ExpressionSet` in which both `ZNF721` and `ABCA11P` are features with identical expression data:

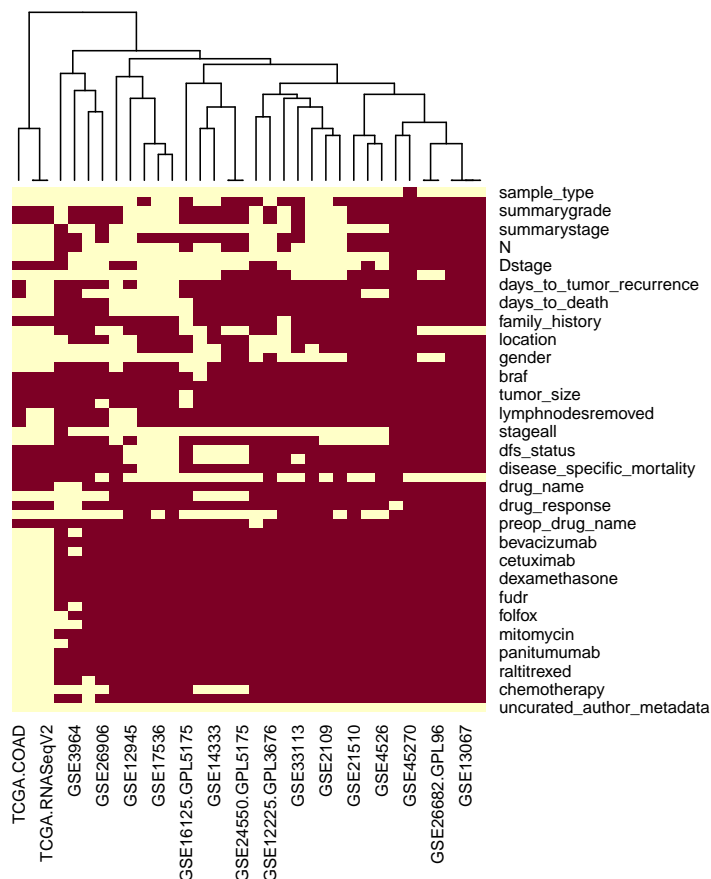
```
> expandProbesets <- function (eset, sep = "///")
+ {
+   x <- lapply(featureNames(eset), function(x) strsplit(x, sep)[[1]])
+   eset <- eset[order(sapply(x, length)), ]
+   x <- lapply(featureNames(eset), function(x) strsplit(x, sep)[[1]])
+   idx <- unlist(sapply(1:length(x), function(i) rep(i, length(x[[i]]))))
+   xx <- !duplicated(unlist(x))
+   idx <- idx[xx]
+   x <- unlist(x)[xx]
+   eset <- eset[idx, ]
+   featureNames(eset) <- x
+   eset
+ }
> X <- TCGA.COAD_eset[head(grep("AA", featureNames(TCGA.COAD_eset))), ]
> exprs(X)[,1:3]

      TCGA.AA.3520 TCGA.AA.3532 TCGA.AA.3553
AAAS      -0.72125    -1.51150    -1.01250
AACS       0.02225     0.82375    -0.08500
AADAC      0.02775    -0.42900     1.12525
AADACL1    1.06600     1.85550     0.92800
AADACL2    0.08750    -0.61825    -0.47525
AADACL3    0.38100     0.31100     0.33950

> exprs(expandProbesets(X))[,1:3]

      TCGA.AA.3520 TCGA.AA.3532 TCGA.AA.3553
AAAS      -0.72125    -1.51150    -1.01250
AACS       0.02225     0.82375    -0.08500
AADAC      0.02775    -0.42900     1.12525
AADACL1    1.06600     1.85550     0.92800
AADACL2    0.08750    -0.61825    -0.47525
AADACL3    0.38100     0.31100     0.33950
```

## A Available Clinical Characteristics



**Figure 1: Available clinical annotation**  
 This heatmap visualizes for each curated clinical characteristic (rows) the availability in each dataset (columns). Red indicates that the corresponding characteristic is available for at least one sample in the dataset.

## B Summarizing the List of ExpressionSets

This example provides a table summarizing the datasets being used, and is useful when publishing analyses based on curatedCRCDData. First, define some useful functions for this purpose:

```
> source(system.file("extdata", "summarizeEsets.R", package =
+ "curatedCRCDData"))
```

Optionally write this table to file, for example ( replace myfile <- tempfile() with something like myfile <- "nicetable.csv" )

## curatedCRCData

```
> (myfile <- tempfile())  
[1] "/tmp/Rtmpi9lh2H/file1807585fffa2"  
> write.table(summary.table, file=myfile, row.names=FALSE, quote=TRUE, sep=",")
```



Study	PMID	N samples
GSE11237	18653328	23
GSE12225.GPL3676	18959792	2
GSE12945	19399471	2
GSE13067	19088021	4
GSE13294	19088021	155
GSE14095	21680303	189
GSE14333	19996206	290
GSE16125.GPL5175	19672874	36
GSE17536	19914252	177
GSE17537	19914252	55
GSE17538.GPL570	19914252	232
GSE18105	20162577	111
GSE2109	PMID unknown	427
GSE21510	21270110	148
GSE21815	21862635	141
GSE24549.GPL5175	21619627	83
GSE24550.GPL5175	21619627	90
GSE2630	17390049	16
GSE26682.GPL570	21300766	156
GSE26682.GPL96	21300766	155
GSE26906	22496922	90
GSE27544	PMID unknown	22
GSE28702	22095227	83
GSE3294	16773188	24
GSE33113	22496204	96
GSE39582	23700391	566
GSE3964	16542501	29
GSE4045	16819509	37
GSE4526	19016304	36
GSE45270	PMID unknown	13
TCGA.COAD	22810696	130
TCGA.READ	22810696	51
TCGA.RNASeqV2.READ	22810696	6
TCGA.RNASeqV2	22810696	195

Table 1: Datasets provided by curatedCRCData

## C For non-R users

---

If you are not doing your analysis in R, and just want to get some data you have identified from the curatedCRCDData manual, here is a simple way to do it. For one dataset:

```
> library(curatedCRCDData)
> library(affy)
> data(TCGA.COAD_eset)
> write.csv(exprs(TCGA.COAD_eset), file="TCGA.COAD_eset_exprs.csv")
> write.csv(pData(TCGA.COAD_eset), file="TCGA.COAD_eset_clindata.csv")
```

Or for several datasets:

```
> data.to.fetch <- c("TCGA.COAD_eset", "GSE37317_eset")
> for (onedata in data.to.fetch){
+   print(paste("Fetching", onedata))
+   data(list=onedata)
+   write.csv(exprs(get(onedata)), file=paste(onedata, "_exprs.csv", sep=""))
+   write.csv(pData(get(onedata)), file=paste(onedata, "_clindata.csv", sep=""))
+ }
```

## D Session Info

---

- R version 4.0.3 (2020-10-10), x86\_64-pc-linux-gnu
- Locale: LC\_CTYPE=en\_US.UTF-8, LC\_NUMERIC=C, LC\_TIME=en\_US.UTF-8, LC\_COLLATE=C, LC\_MONETARY=en\_US.UTF-8, LC\_MESSAGES=en\_US.UTF-8, LC\_PAPER=en\_US.UTF-8, LC\_NAME=C, LC\_ADDRESS=C, LC\_TELEPHONE=C, LC\_MEASUREMENT=en\_US.UTF-8, LC\_IDENTIFICATION=C
- Running under: Ubuntu 18.04.5 LTS
- Matrix products: default
- BLAS: /home/biocbuild/bbs-3.12-bioc/R/lib/libRblas.so
- LAPACK: /home/biocbuild/bbs-3.12-bioc/R/lib/libRlapack.so
- Base packages: base, datasets, grDevices, graphics, methods, parallel, stats, utils
- Other packages: Biobase 2.50.0, BiocGenerics 0.36.0, BiocParallel 1.24.0, affy 1.68.0, curatedCRCDData 2.22.0, genefilter 1.72.0, logging 0.10-108, mgcv 1.8-33, nlme 3.1-150, survival 3.2-7, sva 3.38.0, xtable 1.8-4
- Loaded via a namespace (and not attached): AnnotationDbi 1.52.0, BiocManager 1.30.10, BiocStyle 2.18.0, DBI 1.1.0, IRanges 2.24.0, Matrix 1.2-18, R6 2.5.0, RSQLite 2.2.1, Rcpp 1.0.5, S4Vectors 0.28.0, XML 3.99-0.5, affyio 1.60.0, annotate 1.68.0, bit 4.0.4, bit64 4.0.5, blob 1.2.1, compiler 4.0.3, crayon 1.3.4, digest 0.6.27, edgeR 3.32.0, evaluate 0.14, grid 4.0.3, htmltools 0.5.0, httr 1.4.2, knitr 1.30, lattice 0.20-41, limma 3.46.0, locfit 1.5-9.4, matrixStats 0.57.0, memoise 1.1.0, preprocessCore 1.52.0, rlang 0.4.8, rmarkdown 2.5, splines 4.0.3, stats4 4.0.3, tools 4.0.3, vctrs 0.3.4, xfun 0.18, yaml 2.2.1, zlibbioc 1.36.0