

MetaGxBreast: a package for breast cancer gene expression analysis

Michael Zon¹, Deena M.A. Gendoo^{1,2}, Natchar Ratanasirigulchai¹,
Gregory Chen², Levi Waldron^{3,4}, and Benjamin Haibe-Kains^{*1,2}

¹Bioinformatics and Computational Genomics Laboratory, Princess
Margaret Cancer Center, University Health Network, Toronto,
Ontario, Canada

²Department of Medical Biophysics, University of Toronto, Toronto,
Canada

³Department of Biostatistics and Computational Biology,
Dana-Farber Cancer Institute, Boston, MA, USA

⁴Department of Biostatistics, Harvard School of Public Health,
Boston, MA, USA

November 4, 2025

Contents

| | | |
|----------|--|----------|
| 1 | Installing the Package | 2 |
| 2 | Loading Datasets | 2 |
| 3 | Obtaining Sample Counts in Datasets | 3 |
| 4 | Assess Phenotype Data | 3 |
| 5 | Session Info | 5 |

*benjamin.haibe.kains@utoronto.ca

1 Installing the Package

The MetaGxBreast package is a compendium of Breast Cancer datasets. The package is publicly available and can be installed from Bioconductor into R version 3.5.0 or higher.

To install the MetaGxBreast package from Bioconductor:

```
> if (!requireNamespace("BiocManager", quietly = TRUE))
+   install.packages("BiocManager")
> BiocManager::install("MetaGxBreast")
```

2 Loading Datasets

First we load the MetaGxBreast package into the workspace.

To load the packages into R and obtain some datasets, please use the following commands:

```
> library(MetaGxBreast)
> esets <- MetaGxBreast::loadBreastEsets(loadString = c("CAL", "DFHCC", "DFHCC2",
+   "DFHCC3", "DUKE", "DUKE2", "EMC2"))[[1]]
```

This will load 7 of the 37 expression datasets. Users can modify the parameters of the function to restrict datasets that do not meet certain criteria for loading. Also note that `loadString = "majority"` will load 37 of the 39 datasets. The larger metabric and tcga studies need to be loaded separately by altering the `loadString` variable to include the string metabric or tcga. Some example parameters are shown below:

Datasets: Retain only genes that are common across all platforms loaded
(default = FALSE)

Datasets: Retain studies with a minimum sample size (default = 0)

Datasets: Retain studies with a minimum number of genes (default = 0)

Datasets: Retain studies with a minimum number of survival events (default = 0)

Datasets: Remove duplicate samples (default = TRUE)

3 Obtaining Sample Counts in Datasets

To obtain the number of samples per dataset, run the following:

```
> library(Biobase)
> numSamples <- vapply(seq_along(esets), FUN=function(i, esets){
+   length(sampleNames(esets[[i]]))
+ }, numeric(1), esets=esets)
> SampleNumberSummaryAll <- data.frame(NumberOfSamples = numSamples,
+                                     row.names = names(esets))
> total <- sum(SampleNumberSummaryAll[, "NumberOfSamples"])
> SampleNumberSummaryAll <- rbind(SampleNumberSummaryAll, total)
> rownames(SampleNumberSummaryAll)[nrow(SampleNumberSummaryAll)] <- "Total"
> require(xtable)
> print(xtable(SampleNumberSummaryAll, digits = 2), floating = FALSE)
```

| | NumberOfSamples |
|--------|-----------------|
| CAL | 118.00 |
| DFHCC | 115.00 |
| DFHCC2 | 83.00 |
| DFHCC3 | 40.00 |
| DUKE | 169.00 |
| DUKE2 | 154.00 |
| EMC2 | 204.00 |
| Total | 883.00 |

4 Assess Phenotype Data

We can also obtain a summary of the phenotype data (pData) for each expression dataset. Here, we assess the proportion of samples in every datasets that contain a specific pData variable.

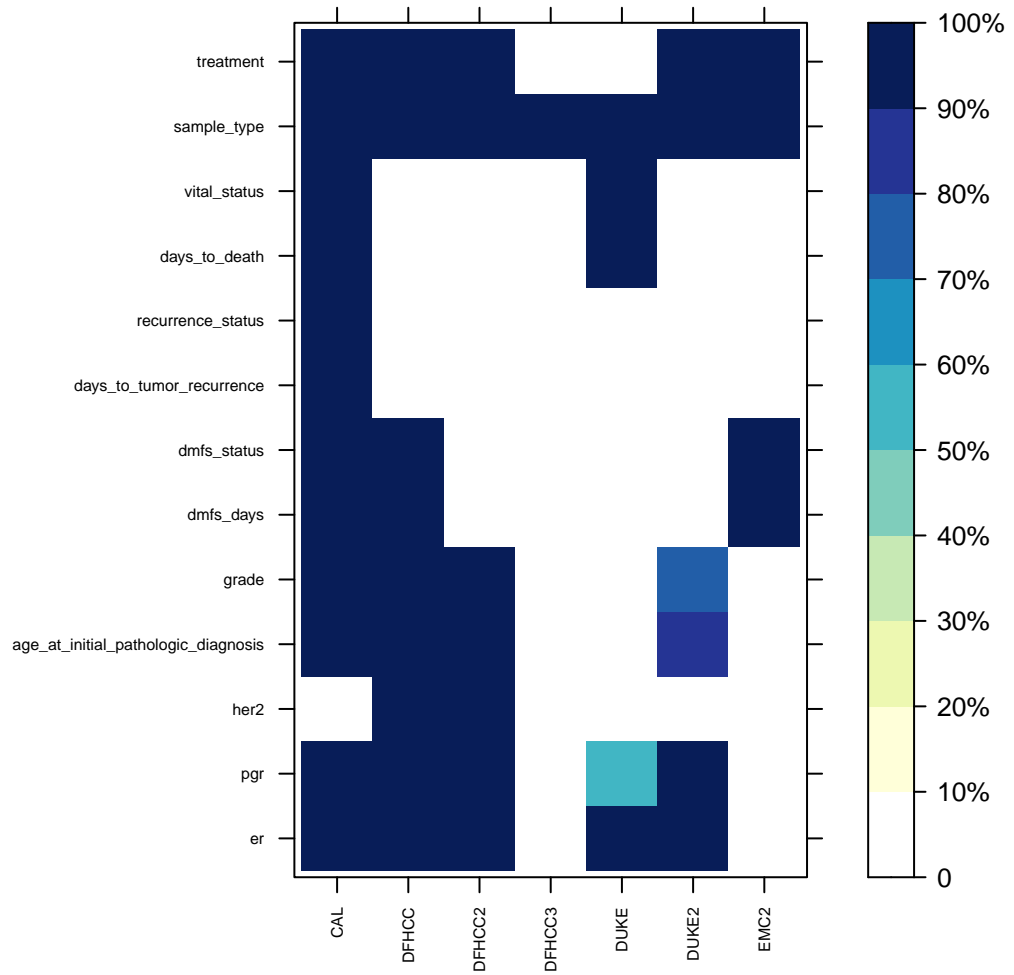
```
> #pData Variables
> pDataID <- c("er", "pgr", "her2", "age_at_initial_pathologic_diagnosis",
+             "grade", "dmfs_days", "dmfs_status", "days_to_tumor_recurrence",
+             "recurrence_status", "days_to_death", "vital_status",
+             "sample_type", "treatment")
> pDataPercentSummaryTable <- NULL
> pDataSummaryNumbersTable <- NULL
> pDataSummaryNumbersList <- lapply(esets, function(x)
+   vapply(pDataID, function(y) sum(!is.na(pData(x)[,y])), numeric(1)))
```

```

> pDataPercentSummaryList <- lapply(esets, function(x)
+   vapply(pDataID, function(y)
+     sum(!is.na(pData(x)[,y]))/nrow(pData(x)), numeric(1))*100)
> pDataSummaryNumbersTable <- sapply(pDataSummaryNumbersList, function(x) x)
> pDataPercentSummaryTable <- sapply(pDataPercentSummaryList, function(x) x)
> rownames(pDataSummaryNumbersTable) <- pDataID
> rownames(pDataPercentSummaryTable) <- pDataID
> colnames(pDataSummaryNumbersTable) <- names(esets)
> colnames(pDataPercentSummaryTable) <- names(esets)
> pDataSummaryNumbersTable <- rbind(pDataSummaryNumbersTable, total)
> rownames(pDataSummaryNumbersTable)[nrow(pDataSummaryNumbersTable)] <- "Total"
> # Generate a heatmap representation of the pData
> pDataPercentSummaryTable <- t(pDataPercentSummaryTable)
> pDataPercentSummaryTable <- cbind(Name=(rownames(pDataPercentSummaryTable))
+   ,pDataPercentSummaryTable)
> nba<-pDataPercentSummaryTable
> gradient_colors <- c("#ffffff", "#ffffd9", "#edf8b1", "#c7e9b4", "#7fcdbb",
+   "#41b6c4", "#1d91c0", "#225ea8", "#253494", "#081d58")
> library(lattice)
> nbamat<-as.matrix(nba)
> rownames(nbamat) <- nbamat[,1]
> nbamat <- nbamat[,-1]
> Interval <- as.numeric(c(10,20,30,40,50,60,70,80,90,100))
> levelplot(nbamat,col.regions=gradient_colors,
+   main="Available Clinical Annotation",
+   scales=list(x=list(rot=90, cex=0.5),
+     y= list(cex=0.5),key=list(cex=0.2)),
+   at=seq(from=0,to=100,length=10),
+   cex=0.2, ylab="", xlab="", lattice.options=list(),
+   colorkey=list(at=as.numeric(factor(c(seq(from=0, to=100, by=10))))),
+     labels=as.character(c( "0%", "10%", "20%", "30%", "40%", "50%",
+       "60%", "70%", "80%", "90%", "100%"),
+     cex=0.2,font=1,col="brown",height=1,
+     width=1.4), col=(gradient_colors)))
>

```

Available Clinical Annotation



5 Session Info

- R Under development (unstable) (2025-10-20 r88955), x86_64-pc-linux-gnu
- Locale: LC_CTYPE=en_US.UTF-8, LC_NUMERIC=C, LC_TIME=en_GB, 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

- Time zone: America/New_York
- TZcode source: system (glibc)
- Running under: Ubuntu 24.04.3 LTS
- Matrix products: default
- BLAS: /home/biocbuild/bbs-3.23-bioc/R/lib/libRblas.so
- LAPACK:
/usr/lib/x86_64-linux-gnu/lapack/liblapack.so.3.12.0
- Base packages: base, datasets, grDevices, graphics, methods, stats,
utils
- Other packages: AnnotationHub 4.1.0, Biobase 2.71.0,
BiocFileCache 3.1.0, BiocGenerics 0.57.0, ExperimentHub 3.1.0,
MetaGxBreast 1.31.0, dbplyr 2.5.1, generics 0.1.4, lattice 0.22-7,
xtable 1.8-4
- Loaded via a namespace (and not attached): AnnotationDbi 1.73.0,
BiocManager 1.30.26, BiocVersion 3.23.1, Biostrings 2.79.1,
DBI 1.2.3, DelayedArray 0.37.0, GenomicRanges 1.63.0,
IRanges 2.45.0, KEGGREST 1.51.0, Matrix 1.7-4,
MatrixGenerics 1.23.0, R6 2.6.1, RSQLite 2.4.3, S4Arrays 1.11.0,
S4Vectors 0.49.0, Seqinfo 1.1.0, SparseArray 1.11.1,
SummarizedExperiment 1.41.0, XVector 0.51.0, abind 1.4-8, bit 4.6.0,
bit64 4.6.0-1, blob 1.2.4, cachem 1.1.0, cli 3.6.5, compiler 4.6.0,
crayon 1.5.3, curl 7.0.0, dplyr 1.1.4, fastmap 1.2.0, filelock 1.0.3,
glue 1.8.0, grid 4.6.0, httr 1.4.7, httr2 1.2.1, impute 1.85.0,
lifecycle 1.0.4, magrittr 2.0.4, matrixStats 1.5.0, memoise 2.0.1,
pillar 1.11.1, pkgconfig 2.0.3, png 0.1-8, purrr 1.1.0, rappdirs 0.3.3,
rlang 1.1.6, stats4 4.6.0, tibble 3.3.0, tidyselect 1.2.1, tools 4.6.0,
vctrs 0.6.5, withr 3.0.2, yaml 2.3.10