

# MIGSA: Getting pbcmc datasets

**Juan C Rodriguez**

CONICET

Universidad Católica de Córdoba

Universidad Nacional de Córdoba

**Cristóbal Fresno**

Instituto Nacional de Medicina Genómica

**Andrea S Llera**

CONICET

Fundación Instituto Leloir

**Elmer A Fernández**

CONICET

Universidad Católica de Córdoba

Universidad Nacional de Córdoba

---

## Abstract

In this vignette we are going to show how we got the RData *pbcmcData.RData* which can be loaded via the **MIGSAdata** package using `data(pbcmcData)`.

*Keywords:* singular enrichment analysis, over representation analysis, gene set enrichment analysis, functional class scoring, big omics data.

---

## 1. Getting the data

Following we give the used code to download this data and their PAM50 subtypes.

```
> library(limma);
> library(pbcmc);
> # datasets included in BioConductor repository
> libNames <- c("mainz", "nki", "transbig", "unt", "upp", "vdx");
> # let's load them!
> pbcmcData <- lapply(libNames, function(actLibName) {
+   print(actLibName);
+
+   # the pbcmc package provides an easy way to download and classify them
+   actLib <- loadBCDataset(Class=PAM50, libname=actLibName, verbose=FALSE);
+   actLibFilt <- filtrate(actLib, verbose=FALSE);
+   actLibFilt <- classify(actLibFilt, std="none", verbose=FALSE);
+   actSubtypes <- classification(actLibFilt)$subtype;
+
+   # get the expression matrix and the annotation
+   actExprs <- exprs(actLib);
+   actAnnot <- annotation(actLib);
+ })
```

```

+   # we recommend working allways with Entrez IDs, let's match them with
+   # expression matrix rownames (and modify them)
+   if (all(actAnnot$probe == rownames(actExprs))) {
+       actExprs <- actExprs[!is.na(actAnnot$EntrezGene.ID),];
+       actAnnot <- actAnnot[!is.na(actAnnot$EntrezGene.ID),];
+       rownames(actExprs) <- as.character(actAnnot$EntrezGene.ID);
+   } else {
+       matchedEntrez <- match(rownames(actExprs), actAnnot$probe);
+       # all(rownames(actExprs) %in% actAnnot$probe == !is.na(matchedEntrez));
+
+       stopifnot(all(
+           actAnnot$probe[!is.na(matchedEntrez)] ==
+           rownames(actExprs)[!is.na(matchedEntrez)]));
+
+       actExprs <- actExprs[!is.na(matchedEntrez),];
+       actAnnot <- actAnnot[!is.na(matchedEntrez),];
+       stopifnot(all(actAnnot$probe == rownames(actExprs)));
+       actExprs <- actExprs[!is.na(actAnnot$EntrezGene.ID),];
+       actAnnot <- actAnnot[!is.na(actAnnot$EntrezGene.ID),];
+       rownames(actExprs) <- as.character(actAnnot$EntrezGene.ID);
+   }
+
+   # average repeated genes expression
+   actExprs <- avereps(actExprs);
+
+   stopifnot(all(colnames(actExprs) == names(actSubtypes)));
+   # filtrate only these two conditions
+   actExprs <- actExprs[, actSubtypes %in% c("Basal", "LumA")];
+   actSubtypes <- as.character(
+       actSubtypes[actSubtypes %in% c("Basal", "LumA")]);
+
+   return(list(geneExpr=actExprs, subtypes=actSubtypes));
+ })

```

```

[1] "mainz"
[1] "nki"
[1] "transbig"
[1] "unt"
[1] "upp"
[1] "vdx"

```

```
> names(pbcmcData) <- libNames;
```

And let's check it is the same data.

```

> # save the just created pbcmcData to newPbcmcData
> newPbcmcData <- pbcmcData;

```

```
> library(MIGSadata);
> # and load the MIGSadata one.
> data(pbcmcData);
> all.equal(newPbcmcData, pbcmcData);
```

```
[1] TRUE
```

## Session Info

```
> sessionInfo()
```

```
R version 3.5.1 Patched (2018-07-24 r75008)
Platform: x86_64-w64-mingw32/x64 (64-bit)
Running under: Windows Server 2012 R2 x64 (build 9600)
```

```
Matrix products: default
```

```
locale:
```

```
[1] LC_COLLATE=C
[2] LC_CTYPE=English_United States.1252
[3] LC_MONETARY=English_United States.1252
[4] LC_NUMERIC=C
[5] LC_TIME=English_United States.1252
```

```
attached base packages:
```

```
[1] stats4      parallel  stats      graphics  grDevices  utils      datasets
[8] methods     base
```

```
other attached packages:
```

```
[1] pbcmc_1.9.0          genefu_2.14.0        AIMS_1.14.0
[4] e1071_1.7-0          iC10_1.4.2           iC10TrainingData_1.3.1
[7] pamr_1.55            cluster_2.0.7-1      biomaRt_2.38.0
[10] mclust_5.4.1         survcomp_1.32.0      prodlim_2018.04.18
[13] survival_2.43-1      edgeR_3.24.0         MIGSadata_1.5.0
[16] MIGSA_1.6.0          mGSZ_1.0             ismev_1.42
[19] mgcv_1.8-25          nlme_3.1-137         MASS_7.3-51
[22] limma_3.38.0         GSA_1.03             BiocParallel_1.16.0
[25] GSEABase_1.44.0      graph_1.60.0         annotate_1.60.0
[28] XML_3.98-1.16        AnnotationDbi_1.44.0 IRanges_2.16.0
[31] S4Vectors_0.20.0     Biobase_2.42.0       BiocGenerics_0.28.0
```

```
loaded via a namespace (and not attached):
```

```
[1] amap_0.8-16          colorspace_1.3-2
[3] class_7.3-14         futile.logger_1.4.3
[5] breastCancerTRANSBIG_1.19.0 gg dendro_0.1-20
```

[7] bit64_0.9-7	splines_3.5.1
[9] SuppDists_1.1-9.4	breastCancerVDX_1.19.0
[11] G0.db_3.7.0	compiler_3.5.1
[13] httr_1.3.1	G0stats_2.48.0
[15] assertthat_0.2.0	Matrix_1.2-14
[17] lazyeval_0.2.1	formatR_1.5
[19] prettyunits_1.0.2	tools_3.5.1
[21] bindrcpp_0.2.2	gtable_0.2.0
[23] glue_1.3.0	Category_2.48.0
[25] reshape2_1.4.3	dplyr_0.7.7
[27] Rcpp_0.12.19	RJSONIO_1.3-0
[29] breastCancerNKI_1.19.0	stringr_1.3.1
[31] org.Hs.eg.db_3.7.0	scales_1.0.0
[33] hms_0.4.2	RBGL_1.58.0
[35] lambda.r_1.2.3	breastCancerUPP_1.19.0
[37] memoise_1.1.0	gridExtra_2.3
[39] breastCancerMAINZ_1.19.0	ggplot2_3.1.0
[41] stringi_1.2.4	RSQLite_2.1.1
[43] rmeta_3.0	genefilter_1.64.0
[45] permute_0.9-4	lava_1.6.3
[47] rlang_0.3.0.1	pkgconfig_2.0.2
[49] matrixStats_0.54.0	bitops_1.0-6
[51] lattice_0.20-35	purrr_0.2.5
[53] bindr_0.1.1	labeling_0.3
[55] survivalROC_1.0.3	cowplot_0.9.3
[57] bit_1.1-14	tidyselect_0.2.5
[59] AnnotationForge_1.24.0	plyr_1.8.4
[61] magrittr_1.5	R6_2.3.0
[63] snow_0.4-3	bootstrap_2017.2
[65] DBI_1.0.0	pillar_1.3.0
[67] RCurl_1.95-4.11	tibble_1.4.2
[69] crayon_1.3.4	futile.options_1.0.1
[71] KernSmooth_2.23-15	progress_1.2.0
[73] locfit_1.5-9.1	grid_3.5.1
[75] data.table_1.11.8	blob_1.1.1
[77] vegan_2.5-3	Rgraphviz_2.26.0
[79] digest_0.6.18	breastCancerUNT_1.19.0
[81] xtable_1.8-3	munsell_0.5.0

**Affiliation:**

Juan C Rodriguez & Elmer A Fernández  
 Bioscience Data Mining Group  
 Facultad de Ingeniería  
 Universidad Católica de Córdoba - CONICET

X5016DHK Córdoba, Argentina

E-mail: [jcrodriguez@bdmg.com.ar](mailto:jcrodriguez@bdmg.com.ar), [efernandez@bdmg.com.ar](mailto:efernandez@bdmg.com.ar)

URL: <http://www.bdmg.com.ar/>