

CBW pathways Workshops - example R notebooks

Ruth Isserlin

2023-04-14

Contents

1	Index	5
2	CBW Workshop example R Notebooks	7
3	Setup	9
3.1	Install R and RStudio	9
3.2	Docker [Optional]	10
3.3	Install Docker	10
4	Run g:profiler from R	13
4.1	Run g:profiler with supplied genesets	14
4.2	Create an output file of the results - Generic enrichment Map file from g:profiler gmt	17
4.3	Run g:profiler with your own genesets	18
4.4	Upload the gmt file to gprofiler	18
4.5	Create enrichment Results files	20
4.6	Create an output file of the results - Generic enrichment Map file from Baderlab gmt	20
5	Run GSEA from within R	23
5.1	Load in required libraries	23
5.2	Configurable Parameters	23
5.3	Download the latest pathway definition file	24
5.4	Run GSEA	25

6	Create GMT file from Ensembl	27
6.1	Load Libraries	28
6.2	Set up Biomart connection	28
6.3	Get species GO annotations	29
6.4	Format results into GMT file	30

Chapter 1

Index

Chapter 2

CBW Workshop example R Notebooks

Do you want to run the pathways and network analysis from R instead of doing everything manually as demonstrated in the workshop?

Everything (almost!) that was discussed in the lectures and practicals can be done computationally through R.

We are using the **bookdown** package (Xie 2023) in this Workshop R Notebooks book, which was built on top of R Markdown and **knitr** (Xie 2015).

Chapter 3

Setup

3.1 Install R and RStudio

As with many open source projects, **R** is a constantly evolving language with regular updates. There is a major release once a year with patch releases throughout the year. Often scripts and packages will work from one release to the next (ignoring pesky warnings that a package was compiled on a previous version of **R** is common) but there are exceptions. Some newer packages will only work on the latest version of **R** so sometimes the choice of upgrading or not using a new package might present themselves. Often, the amount of packages and work that is needed to upgrade is not realized until the process has begun. This is where docker demonstrates its most valuable features. You can create a new instance based on the latest release of **R** and all your needed packages without having to change any of your current settings.

In order to use these notebooks supplied here you need to have:

- **R** installed on your computer and
- a list of packages. (including BiocManager, BiomaRt, gprofiler2, GSA)

Each notebook in this set will check for the required packages and install them if they are missing so at the base level you need to just have **R** installed.

There are many different ways you can use and setup **R**.

1. By simply installing **R** you can use it directly but
2. it is highly recommended that you also install and use RStudio which is an Integrated development environment (IDE) for **R**. You cannot just download RStudio and use it. It requires an installation of **R**.

You don't need to install **R** and RStudio though. You can also use **R** and RStudio through docker. **I highly recommend using docker instead**

3.2 Docker [Optional]

Changing versions and environments are a continuing struggle with bioinformatics pipelines and computational pipelines in general. An analysis written and performed a year ago might not run or produce the same results when it is run today. Recording package and system versions or not updating certain packages rarely work in the long run.

One the best solutions to reproducibility issues is containing your workflow or pipeline in its own coding environment where everything from the operating system, programs and packages are defined and can be built from a set of given instructions. There are many systems that offer this type of control including:

- Docker.
- Singularity

“A container is a standard unit of software that packages up code and all its dependencies so the application runs quickly and reliably from one computing environment to another.” (“What Is a Container?” n.d.)

Why are containers great for Bioinformatics?

- allows you to create environments to run bioinformatics pipelines.
- create a consistent environment to use for your pipelines.
- test modifications to the pipeline without disrupting your current set up.
- Coming back to an analysis years later and there is no need to install older versions of packages or programming languages. Simply create a container and re-run.

3.3 Install Docker

1. Download and install docker desktop.
2. Follow slightly different instructions for Windows or MacOS/Linux

3.3.1 Windows

- it might prompt you to install additional updates (for example - <https://docs.microsoft.com/en-us/windows/wsl/install-win10#step-4---download-the-linux-kernel-update-package>) and require multiple restarts of your system or docker.
- launch docker desktop app.
- Open windows Power shell

- navigate to directory on your system where you plan on keeping all your code. For example: C:\USERS\risserlin\cbw_workshop_code
- Run the following command: (the only difference with the windows command is the way the current directory is written. \${PWD} instead of "\$(pwd)")

```
docker run -e PASSWORD=changeit --rm \
-v ${PWD}:/home/rstudio/projects -p 8787:8787 \
risserlin/workshop_base_image
```

- Windows defender firewall might pop up with warning. Click on *Allow access*.
- In docker desktop you see all containers you are running and easily manage them.

3.3.2 MacOS / Linux

- Open Terminal
- navigate to directory on your system where you plan on keeping all your code. For example: /Users/risserlin/bcb420_code
- Run the following command: (the only difference with the windows command is the way the current directory is written. \${PWD} instead of "\$(pwd)")

```
docker run -e PASSWORD=changeit --rm \
-v "$(pwd)"/:/home/rstudio/projects -p 8787:8787 \
--add-host "localhost:My.IP.address"
risserlin/workshop_base_image
```


Chapter 4

Run g:profiler from R

Detailed instructions on how to run g:Profiler programmatically from R

The parameters are set in the params option on this notebook but you can also manually set them here.

```
# for example - working_dir <- "./generated_data"
working_dir <- params$working_dir

data_dir <- params$data_dir

# for example - species <- "horse"
genelist_file <- params$genelist_file

# max size of the genesets for example - 350
max_gs_size <- params$max_gs_size

# max size of the genesets for example - 3
min_gs_size <- params$min_gs_size

#min intersection between your genelist and the geneset - for example 3
min_intersection <- params$min_intersection

# organism parameter used for g:profiler. First letter of first word in species name followed by
# the second word
# for example - hsapiens
organism <- params$organism

#use library
tryCatch(expr = { library("gprofiler2")},
         error = function(e) {
```

```

      install.packages("gprofiler2")),
      finally = library("gprofiler2"))

tryCatch(expr = { library("GSA")},
      error = function(e) {
        install.packages("GSA")},
      finally = library("GSA"))

```

Create or set a directory to store all the generated results

```

if(!dir.exists(params$working_dir)){
  dir.create(params$working_dir)
}

```

Load in the set of genes that we will be running g:profiler with

```

#load in the file
current_genelist <- read.table(file =
                                file.path(data_dir, genelist_file),
                                header = FALSE,
                                sep = "\t", quote = "",
                                stringsAsFactors = FALSE)

query_set <- current_genelist$V1

```

With regards to pathway sets there are two options when using g:Profiler -

- Use the genesets that are supplied by g:Profiler
- Upload your own genesets.

The most common reasons for supplying your own genesets is the ability to use up to date annotations or in-house annotations that might not be available in the public sphere yet. One of the greatest features of g:Profiler is that it is updated on a regular basis and most of the previous versions are available online on the gprofiler archive.

The gprofiler2 -g:Profiler R implementation is a wrapper for the web version. You require an internet connection to get enrichment results.

4.1 Run g:profiler with supplied genesets

For detailed descriptions of all the parameters that can be specified for the g:profiler function see [here](#)

For this query we are specifying -

- query - the set of genes of interest, as loaded in from the Supplementary_Table1_Cancer_drivers.txt file.
- significant - set to FALSE because we want g:Profiler to return all the results not just the ones that it deems significant by its predetermined threshold.
- ordered_query - set to TRUE because for this set of genes they are ordered in order of their significance
- correction_method - set to fdr. by default g:Profiler uses g:Scs
- organism - set to "hsapiens" for homo sapiens. Organism names are constructed by concatenating the first letter of the name and the family name (according to gprofiler2 documentation)
- source - the geneset source databases to use for the analysis. We recommend using GO biological process (GO:BP), WikiPathways (WP) and Reactome (Reac) but there are additional sources you can add (GO molecular function or cellular component(GO:MF, GO:CC), KEGG, transcription factors (TF), microRNA targets (MIRNA), corum complexes (CORUM), Human protein atlas (HPA), Human phenotype ontology (HP))

```
gprofiler_results <- gost(query = query_set ,
                          significant=FALSE,
                          ordered_query = TRUE,
                          exclude_iea=FALSE,
                          correction_method = "fdr",
                          organism = organism,
                          source = c("REAC", "WP", "GO:BP"))
```

```
#get the gprofiler results table
enrichment_results <- gprofiler_results$result

enrichment_results[1:5,]
```

```
##      query significant      p_value term_size query_size intersection_size
## 1 query_1          TRUE 1.426353e-37     5653         121             103
## 2 query_1          TRUE 3.391992e-36     5882         121             103
## 3 query_1          TRUE 7.172333e-36     3097         121              81
## 4 query_1          TRUE 2.511953e-35     5724         121             101
## 5 query_1          TRUE 7.298888e-35     3540         121              84
## precision    recall    term_id source
## 1 0.8512397 0.01822041 GO:0031323  GO:BP
## 2 0.8512397 0.01751105 GO:0080090  GO:BP
## 3 0.6694215 0.02615434 GO:0031325  GO:BP
## 4 0.8347107 0.01764500 GO:0051171  GO:BP
## 5 0.6942149 0.02372881 GO:0010604  GO:BP
##
##                                     term_name effective_domain_size
## 1                regulation of cellular metabolic process          21128
```

## 2	regulation of primary metabolic process	21128
## 3	positive regulation of cellular metabolic process	21128
## 4	regulation of nitrogen compound metabolic process	21128
## 5	positive regulation of macromolecule metabolic process	21128
##	source_order	parents
## 1	7549	G0:0019222, G0:0044237, G0:0050794
## 2	18924	G0:0019222, G0:0044238
## 3	7551	G0:0009893, G0:0031323, G0:0044237, G0:0048522
## 4	14399	G0:0006807, G0:0019222
## 5	4369	G0:0009893, G0:0043170, G0:0060255

Filter the table to include just the columns that are required for the generic enrichment map file results GEM. Restrict the results to just the ones that have at least min_gs_size and less than max_gs_size terms and min_intersection size include only the term_id, term_name, p_value (and p_value again because the p_value is actually the corrected p-value. The output file does not contain the nominal p_value. For down stream analysis though it is expected to have both a p-value and a q-value so just duplicate the q-value as both p-value and q-value)

```
# filter by params defined above
enrichment_results <- subset(enrichment_results, term_size >= min_gs_size &
                             term_size <= max_gs_size &
                             intersection_size >= min_intersection ,
                             select = c(term_id, term_name, p_value, p_value ))
```

In order to create a proper Generic enrichment results file we will need a copy of the gmt file used by g:Profiler. (also to create an Enrichment map).

Download the gmt file used for this analysis from g:profiler

```
#the link to the gmt file is static no matter what version
gprofiler_gmt_url <- "https://biit.cs.ut.ee/gprofiler//static/gprofiler_full_hsapiens.1"

#get version info gprofiler as the gmt file is always associated with a specific version
gprofiler_version <- get_version_info(organism=organism)

gprofiler_gmt_filename <- file.path(working_dir,
                                    paste("gprofiler_full", organism,
                                           gprofiler_version$gprofiler_version,
                                           ".name.gmt"))

download.file(url = gprofiler_gmt_url, destfile = gprofiler_gmt_filename)

#load in the g:profiler geneset file
```


4.2. CREATE AN OUTPUT FILE OF THE RESULTS - GENERIC ENRICHMENT MAP FILE FROM G:PROFILER

```
capt_output <- capture.output(genesets_gprofiler <- GSA.read.gmt(filename = gprofiler_gmt_filename))

names(genesets_gprofiler$genesets) <- genesets_gprofiler$geneset.names

# Given:
# query_genes - genes used for enrichment analysis (or as query)
#
# returns - the genes that overlap with the query set and part of the given
#           genesets
getGenesetGenes <- function(query_genes, subset_genesets){
  genes <- lapply(subset_genesets, FUN=function(x){intersect(x, query_genes)})

  # For each of the genes collapse to the comma separate text
  genes_collapsed <- unlist(lapply(genes, FUN=function(x){paste(x, collapse = ",")}))
  genes_collapsed_df <- data.frame(term_id = names(genes), genes = genes_collapsed, stringsAsFactors = FALSE)
  return(genes_collapsed_df)
}
```

4.2 Create an output file of the results - Generic enrichment Map file from g:profiler gmt

The file requires -

- name
- description
- p-value
- q-value
- phenotyp
- list of genes (overlap of query set and original geneset)

The list of genes needs to be calculated using the gmt file and original query set. For each geneset found in the result find the overlap between the set of genes that are a part of the geneset and the query set.

```
if(nrow(enrichment_results) > 0){

  #add phenotype to the results
  enrichment_results <- cbind(enrichment_results, 1)

  # Add the genes to the genesets
  subset_genesets <- genesets_gprofiler$genesets[which(genesets_gprofiler$geneset.names %
```

```

    genes <- getGenesetGenes(query_set, subset_genesets)

    enrichment_results <- merge(enrichment_results, genes, by.x=1, by.y=1)

    colnames(enrichment_results) <- c("name", "description", "p-value", "q-value", "padj")
  }

#output the enrichment map file
write.table(enrichment_results, file = file.path(working_dir, "gprofiler_GEM_using_gprofiler"),
            row.names = FALSE, col.names = TRUE, quote=FALSE)

```

4.3 Run g:profiler with your own genesets

Download the latest Bader lab genesets

4.4 Upload the gmt file to gprofiler

In order to use your own genesets with g:Profiler you need to upload the file to their server first. The function will return an ID that you need to specify in the organism parameter of the g:Profiler gost function call.

```

custom_gmt <- upload_GMT_file(gmtfile=dest_gmt_file)

## Your custom annotations ID is gp__4jKc_cgF5_T7M
## You can use this ID as an 'organism' name in all the related enrichment tests again

## Just use: gost(my_genes, organism = 'gp__4jKc_cgF5_T7M')

```

For this query we are specifying -

- query - the set of genes of interest, as loaded in from the Supplementary_Table1_Cancer_drivers.txt file.
- significant - set to FALSE because we want g:Profiler to return all the results not just the ones that it deems significant by its predetermined threshold.
- ordered_query - set to TRUE because for this set of genes they are ordered in order of their significance
- correction_method - set to fdr. by default g:Profiler uses g:Scs
- organism - set to the custom_gmt ID (for this run it is - gp__4jKc_cgF5_T7M) that we received when we uploaded our geneset file.

```
gprofiler_results_custom <- gost(query = query_set ,
                                significant=FALSE,
                                ordered_query = TRUE,
                                exclude_iea=FALSE,
                                correction_method = "fdr",
                                organism = custom_gmt
                                )
```

```
## Detected custom GMT source request
```

```
#get the gprofiler results table
enrichment_results_customgmt <- gprofiler_results_custom$result

enrichment_results_customgmt[1:5,]
```

```
##      query significant      p_value term_size query_size intersection_size
## 1 query_1          TRUE 6.730012e-37      4645         110             94
## 2 query_1          TRUE 2.579228e-35      4747         110             93
## 3 query_1          TRUE 2.579228e-35      4881         110             94
## 4 query_1          TRUE 4.280880e-34      5230         110             95
## 5 query_1          TRUE 2.069151e-33      3435         110             81
##      precision      recall
## 1 0.8545455 0.02023681
## 2 0.8454545 0.01959132
## 3 0.8545455 0.01925835
## 4 0.8636364 0.01816444
## 5 0.7363636 0.02358079
##
##                                     term_id
## 1                      REGULATION OF CELLULAR METABOLIC PROCESS%GOBP%GO:0031323
## 2          REGULATION OF NITROGEN COMPOUND METABOLIC PROCESS%GOBP%GO:0051171
## 3                      REGULATION OF PRIMARY METABOLIC PROCESS%GOBP%GO:0080090
## 4          REGULATION OF MACROMOLECULE METABOLIC PROCESS%GOBP%GO:0060255
## 5 REGULATION OF NUCLEOBASE-CONTAINING COMPOUND METABOLIC PROCESS%GOBP%GO:0019219
##
##                                     source
## 1 Human_GOBP_AllPathways_no_GO_iea_April_02_2023_symbol
## 2 Human_GOBP_AllPathways_no_GO_iea_April_02_2023_symbol
## 3 Human_GOBP_AllPathways_no_GO_iea_April_02_2023_symbol
## 4 Human_GOBP_AllPathways_no_GO_iea_April_02_2023_symbol
## 5 Human_GOBP_AllPathways_no_GO_iea_April_02_2023_symbol
##
##                                     term_name
## 1                      regulation of cellular metabolic process
## 2          regulation of nitrogen compound metabolic process
## 3                      regulation of primary metabolic process
## 4          regulation of macromolecule metabolic process
```

```
## 5 regulation of nucleobase-containing compound metabolic process
## effective_domain_size source_order parents
## 1          18525          18677      NULL
## 2          18525          11634      NULL
## 3          18525          11073      NULL
## 4          18525          10657      NULL
## 5          18525          5967       NULL
```

Filter the table to include just the columns that are required for the generic enrichment map file results GEM. Restrict the results to just the ones that have at least `min_gs_size` and less than `max_gs_size` terms and `min_intersection_size` include only the `term_id`, `term_name`, `p_value` (and `p_value` again because the `p_value` is actually the corrected p-value. The output file does not contain the nominal `p_value`. For down stream analysis though it is expected to have both a p-value and a q-value so just duplicate the q-value as both p-value and q-value)

```
# filter by params defined above
enrichment_results_customgmt <- subset(enrichment_results_customgmt,
                                       term_size >= min_gs_size &
                                       term_size <= max_gs_size &
                                       intersection_size >= min_intersection ,
                                       select = c(term_id,term_name,p_value,p_value ))
```

4.5 Create enrichment Results files

In order to use our results down stream in the Enrichment map we need to generate results files that we can pass to Enrichment Map.

Load in the GMT file

4.6 Create an output file of the results - Generic enrichment Map file from Baderlab gmt

The file requires -

- name
- description
- p-value
- q-value
- phenotyp
- list of genes (overlap of query set and original geneset)

4.6. CREATE AN OUTPUT FILE OF THE RESULTS - GENERIC ENRICHMENT MAP FILE FROM BADERLAB

The list of genes needs to be calculated using the gmt file and original query set.
For each geneset found in the result find the overlap between the set of genes that are a part of the geneset and the query set.

```
if(nrow(enrichment_results_customgmt) >0){  
  
    #add phenotype to the results  
    enrichment_results_customgmt <- cbind(enrichment_results_customgmt,1)  
  
    # Add the genes to the genesets  
    subset_genesets <- genesets_baderlab_genesets$genesets[which(genesets_baderlab_genesets  
  
    genes <- getGenesetGenes(query_set, subset_genesets)  
  
    enrichment_results_customgmt <- merge(enrichment_results_customgmt,genes,by.x=1, by.y=1  
  
    colnames(enrichment_results_customgmt) <- c("name","description","p-value","q-value","p  
  
}  
  
#output the enrichment map file  
write.table(enrichment_results_customgmt,  
            file = file.path(working_dir, "gprofiler_GEM_using_gprof_gmt.txt"),  
            row.names = FALSE, col.names = TRUE,quote = FALS
```


Chapter 5

Run GSEA from within R

This notebook is based largely on the original notebook published with EnrichmentMap Protocol(Reimand et al. 2019)

There is no package to run the original algorithm of GSEA(Subramanian et al. 2005) in R. There are many packages that have been published to imitate the process but none are recognized by The GSEA team.

Revised GSEA
R script

The original GSEA R script from 2005 was revised in 2019 to run on current versions of R. This [updated version is available on GitHub](#). The original script is available from our [Archived Downloads](#) page. Note that neither of these GSEA R scripts are actively supported by the GSEA-MSigDB team; we recommend use of the GSEA software provided above.

5.1 Load in required libraries

```
#install required R and bioconductor packages
tryCatch(expr = { library("RCurl")},
         error = function(e) { install.packages("RCurl")},
         finally = library("RCurl"))
```

5.2 Configurable Parameters

In order to run GSEA automatically through the notebook you will need to download the gsea jar from [here](#). Specify the exact path to the gsea jar in the parameters in order to automatically compute enrichments using GSEA.

If you are running this notebook using the baderlab workshop docker image then the image comes pre-installed with the gsea jar that you can use to run gsea directly in the docker. The path to the GSEA jar in the docker is - /home/rstudio/GSEA_4.3.2/gsea-cli.sh

In order to run GSEA automatically you need to specify the path to the gsea jar file. The gsea_jar needs to be the full path to the GSEA 4.3.2 directory that you downloaded from GSEA. for example /Users/johnsmith/GSEA_4.3.2/gsea-cli.sh

For each variable defined there is an example of how you can set the variable using parameters so you can run re-use your notebook with varying parameters and run completely from the command line. For more info see [here](#)

```
#path to GSEA jar
gsea_jar <- params$gsea_jar --> if you want to define parameters for the notebook and
gsea_jar <- "/home/rstudio/GSEA_4.3.2/gsea-cli.sh"
```

Set the working directory as the directory to the directory where you downloaded all protocol files. For example /User/JohnSmith/EMProtocolFiles/data

```
#directory where all the data files are found. For example - ./data/
working_dir <- params$working_dir --> if you want to define parameters for the notebook
working_dir <- "./data/"

#The name to give the analysis in GSEA - for example Mesen_vs_Immuno
analysis_name <- params$analysis_name --> if you want to define parameters for the notebook
analysis_name <- "Mesen_vs_Immuno"

#rank file to use in GSEA analysis.
#For example - MesenchymalvsImmunoreactive_edger_ranks.rnk
rnk_file <- params$rnk_file --> if you want to define parameters for the notebook and
rnk_file <- "MesenchymalvsImmunoreactive_edger_ranks.rnk"

run_gsea <- true/false
# This parameter is for the compilation of the notebook.
run_gsea <- params$run_gsea --> if you want to define parameters for the notebook and
run_gsea <- FALSE
```

5.3 Download the latest pathway definition file

Only Human, Mouse, Rat, and Woodchuck gene set files are currently available on the baderlab downloads site. If you are working with a species other than human (and it is either rat,mouse or woodchuck) change the gmt_url below to the correct species. Check [here](#) to see all available species.

To create your own GMT file using Ensembl see [Create GMT file from Ensembl](#)


```
gmt_url = "http://download.baderlab.org/EM_Genesets/current_release/Human/symbol/"

#list all the files on the server
filenames = getURL(gmt_url)
tc = textConnection(filenames)
contents = readLines(tc)
close(tc)

#get the gmt that has all the pathways and does not include terms inferred from electronic annotations
#start with gmt file that has pathways only
rx = gregexpr("(?<=<a href=\")(.*GOBP_AllPathways_no_GO_iea.*)(.gmt)(?=\">)",
  contents, perl = TRUE)
gmt_file = unlist(regmatches(contents, rx))

dest_gmt_file <- file.path(working_dir,gmt_file )

#check if this gmt file already exists
if(!file.exists(dest_gmt_file)){
  download.file(
    paste(gmt_url,gmt_file,sep=""),
    destfile=dest_gmt_file
  )
}
```

5.4 Run GSEA

(GSEA)[<http://software.broadinstitute.org/gsea/index.jsp>] is a stand alone java program with many customizable options. It can be easily run through its integrated user interface. To make this a seamless pipeline we can run GSEA from the command line with a set of options. Any of the supplied options can be customized and there are many additional options that can be specified. For more details see (here)[http://software.broadinstitute.org/gsea/doc/GSEAUsguideTEXT.htm#_Running_GSEA_from]

In the below command the following options have been specified:

- `rnk` - path to the rank file
- `gmx` - path to the gene set definition (gmt) file
- `collapse` - true/false indicates whether the expression/rnk file needs to be collapsed from probes to gene symbols
- `nperm` - number of permutations
- `scoring_scheme` -

- rpt_label - name of the directory with output
- rnd_seed - random seed to use
- set_max - maximum size for individual gene sets. In GSEA interface this is set to 500 but we prefer to use a more stringent setting of 200.
- set_min - minimum size for individual gene sets
- zip_report - true/false to zip output directory
- out - directory where to place the result directory.

```
if(run_gsea){  
  command <- paste("",gsea_jar,  
                    "GSEAPreRanked -gmx", dest_gmt_file,  
                    "-rnk" ,file.path(working_dir,rnk_file),  
                    "-collapse false -nperm 1000 -scoring_scheme weighted",  
                    "-rpt_label ",analysis_name,  
                    " -plot_top_x 20 -rnd_seed 12345 -set_max 200",  
                    " -set_min 15 -zip_report false ",  
                    "-out" ,working_dir,  
                    " > gsea_output.txt",sep=" ")  
  system(command)  
}
```

Chapter 6

Create GMT file from Ensembl

The Baderlab geneset download site is an updated resource for geneset files from GO, Reactome, WikiPathways, Pathbank, NetPath, HumanCyc, IOB, ... many others that can be used in g:Profiler or GSEA and many other enrichment tools that support the gmt format.

Unfortunately genesets are only supplied for:

- Human
- Mouse
- Rat
- Woodchuck

If you are working in a different species you will need to generate your own gmt file. The best way to do this is through ensembl. Ensembl doesn't have annotations for all the pathway databases listed above but it has annotations for most species from GO.

The parameters are set in the params option on this notebook but you can also manually set them here.

```
# for example - working_dir <- "./generated_data"
working_dir <- params$working_dir

# for example - species <- "horse"
species <- params$species

# for example - ensembl_dataset <- "ecaballus_gene_ensembl"
ensembl_dataset <- params$ensembl_dataset
```

```

#use library
#make sure biocManager is installed
tryCatch(expr = { library("BiocManager")},
  error = function(e) {
    install.packages("BiocManager")},
  finally = library("BiocManager"))

tryCatch(expr = { library("biomaRt")},
  error = function(e) {
    BiocManager::install("biomaRt")},
  finally = library("biomaRt"))

```

6.1 Load Libraries

Create or set a directory to store all the generated results

```

if(!dir.exists(params$working_dir)){
  dir.create(params$working_dir)
}

```

6.2 Set up Biomart connection

Connect to Biomart

```
ensembl <- useMart("ensembl")
```

Figure out which dataset you want to use - for some species there might be a few datasets to choose from. Not all of the datasets have common names associated with them. For example, if you search for 'yeast' nothing will be returned but if you look for *Saccharomyces* or *cerevisiae* you will be able to find it.

```

all_datasets <- listDatasets(ensembl)

#get all the datasets that match our species definition
all_datasets[grepl(all_datasets$description,
  pattern=species,
  ignore.case = TRUE),]

```

##	dataset	description
----	---------	-------------

```
## 60      ecaballus_gene_ensembl      Horse genes (EquCab3.0)
## 76      hcomes_gene_ensembl      Tiger tail seahorse genes (H_comes_QL1_v1)
## 164 rferrumequinum_gene_ensembl Greater horseshoe bat genes (mRhiFer1_v1.p)
##          version
## 60      EquCab3.0
## 76      H_comes_QL1_v1
## 164     mRhiFer1_v1.p
```

If you know the ensembl dataset that you want to use you can specify it in the parameters above or grab from the above table the dataset of the species that you are interested in.

```
ensembl = useDataset(ensembl_dataset, mart=ensembl)
```

6.3 Get species GO annotations

Get the GO annotations for our species

```
go_annotation <- getBM(attributes = c("external_gene_name",
                                     "ensembl_gene_id",
                                     "ensembl_transcript_id",
                                     "go_id",
                                     "name_1006",
                                     "namespace_1003",
                                     "go_linkage_type"),
                      filters=list(biotype='protein_coding'), mart=ensembl);

#get just the go biological process subset
#####
# Get rid of this line if you want to include all of go and not just biological process
#####
go_annotation_bp <- go_annotation[which(
  go_annotation$namespace_1003 == "biological_process"),]

#compute the unique pathway sets
go_pathway_sets <- aggregate(go_annotation_bp[,1:5],
                             by = list(go_annotation_bp$go_id),
                             FUN = function(x){list(unique(x))})

#unlist the go descriptions
go_pathway_sets$name_1006 <- apply(go_pathway_sets,1,FUN=function(x){
  paste(gsub(unlist(x$name_1006),pattern= "\\\"",
             replacement = ""),collapse = "")})
```

There are two identifiers that you can choose from in the above table * external_symbols * ensembl_ids

Each of these is stored as a list in the dataframe. In order to convert it to the right format for the gmt file we need to convert the list to string of tab delimited strings. (unfortunately there is no straightforward way to write out a dataframe's column of lists.)

```
go_pathway_sets[1:3,"external_gene_name"]
```

```
## [[1]]
## [1] "MEF2A"      "SLC25A36" "OPA1"      "MGME1"      "SLC25A33" "TYMP"      "AKT3"
## [8] "PIF1"
##
## [[2]]
## [1] "GNRH1" "GNRH2" "LIN9"
##
## [[3]]
## [1] "ERCC6" "ERCC8" "LIG4"  "APLF"  "APTX"  "XRCC1" "SIRT1" "XNDC1"
```

```
go_pathway_sets[1:3,"ensembl_gene_id"]
```

```
## [[1]]
## [1] "ENSECAG000000011593" "ENSECAG000000010094" "ENSECAG000000024248"
## [4] "ENSECAG000000012675" "ENSECAG000000016862" "ENSECAG00000001072"
## [7] "ENSECAG000000019722" "ENSECAG000000005316"
##
## [[2]]
## [1] "ENSECAG000000010664" "ENSECAG000000039220" "ENSECAG000000014325"
##
## [[3]]
## [1] "ENSECAG000000014160" "ENSECAG000000018335" "ENSECAG00000003257"
## [4] "ENSECAG000000013246" "ENSECAG000000012674" "ENSECAG000000014127"
## [7] "ENSECAG000000013909" "ENSECAG000000042118"
```

6.4 Format results into GMT file

Convert column of lists to a tab delimited string of gene names

```
go_pathway_sets$collapsed_genenames <- apply(go_pathway_sets,1,
                                              FUN=function(x){
        paste(gsub(unlist(x$external_gene_name),pattern= "\\\"",
                    replacement = ""),collapse = "\\t")
      })
```

Convert column of lists to a tab delimited string of gene names

```
go_pathway_sets$collapsed_ensemblids <- apply(go_pathway_sets,1,
                                              FUN=function(x){
        paste(gsub(unlist(x$ensembl_gene_id),pattern= "\\\"",
                      replacement = "\""),collapse = "\\t")
      })
```

The format of the GMT file is described https://software.broadinstitute.org/cancer/software/gsea/wiki/index.php/Data_formats and consists of rows with the following

- Name
- Description
- tab delimited list of genes a part of this geneset

Write out the gmt file with genenames

```
gmt_file_genenames <- go_pathway_sets[,c("Group.1","name_1006",
                                          "collapsed_genenames")]
colnames(gmt_file_genenames)[1:2] <- c("name","description")

gmt_genenames_filename <- file.path(params$working_dir, paste(species,ensembl_dataset,"GO_geneset",
                                                                "name_1006",
                                                                "collapsed_genenames",
                                                                "gmt"))

write.table(x = gmt_file_genenames,file = gmt_genenames_filename,
           quote = FALSE,sep = "\\t",row.names = FALSE,
           col.names=TRUE)
```

Write out the gmt file with ensembl ids

```
gmt_file_ensemblids <- go_pathway_sets[,c("Group.1","name_1006",
                                          "collapsed_ensemblids")]
colnames(gmt_file_ensemblids)[1:2] <- c("name","description")

gmt_ensemblids_filename <- file.path(params$working_dir, paste(species,ensembl_dataset,"GO_geneset",
                                                                "name_1006",
                                                                "collapsed_ensemblids",
                                                                "gmt"))

write.table(x = gmt_file_ensemblids,file = gmt_ensemblids_filename,
           quote = FALSE,sep = "\\t",row.names = FALSE,
           col.names=TRUE)
```

Reimand, Jüri, Ruth Isserlin, Veronique Voisin, Mike Kucera, Christian Tannus-Lopes, Asha Rostamianfar, Lina Wadi, et al. 2019. "Pathway Enrichment Analysis and Visualization of Omics Data Using g: Profiler, GSEA, Cytoscape and EnrichmentMap." *Nature Protocols* 14 (2): 482–517.

Subramanian, Aravind, Pablo Tamayo, Vamsi K Mootha, Sayan Mukherjee, Benjamin L Ebert, Michael A Gillette, Amanda Paulovich, et al. 2005.

- “Gene Set Enrichment Analysis: A Knowledge-Based Approach for Interpreting Genome-Wide Expression Profiles.” *Proceedings of the National Academy of Sciences* 102 (43): 15545–50.
- “What Is a Container?” n.d. *Docker*. <https://www.docker.com/resources/what-container>.
- Xie, Yihui. 2015. *Dynamic Documents with R and Knitr*. 2nd ed. Boca Raton, Florida: Chapman; Hall/CRC. <http://yihui.name/knitr/>.
- . 2023. *Bookdown: Authoring Books and Technical Documents with r Markdown*. <https://CRAN.R-project.org/package=bookdown>.