

Gene set enrichment analysis

Kim Dill-McFarland, kadm@uw.edu

version July 12, 2021

Contents

Overview	1
Prior to the workshop	1
Load data	2
Gene sets	2
Broad MSigDB	2
Hypergeometric enrichment	3
Enrichment exercises	8
Gene set enrichment analysis (GSEA)	8
Additional resources	8
Groups	8
Online	9
R session	9

Overview

In this workshop, we introduce gene set analysis relevant to RNA-sequencing data. In it, we cover:

- Broad Molecular Signatures Database (MSigDB) gene sets
- hypergeometric enrichment with `clusterProfiler`
- gene set enrichment analysis (GSEA) with `fgsea`

During the workshop, we will build an R script together, which will be posted as ‘live_notes’ after the workshop here.

Prior to the workshop

Please install R, RStudio, and the following packages. See the setup instructions for more details.

```
#Data manipulation
library(tidyverse)
#Broad gene set database
library(msigdb)
#Hypergeo enrichment
library(clusterProfiler)
```

```
#GSEA
library(fgsea)
```

Load data

```
#RNAseq expression data
load("data/RSTR_data_clean_subset.RData")
#Linear model results
model.results <- read_csv("data/RSTR.Mtb.model.results.anno.csv") %>%
  #Filter results for Mtb vs media condition
  filter(variable == "conditionTB")

##
## -- Column specification -----
## cols(
##   group = col_character(),
##   model = col_character(),
##   gene = col_character(),
##   variable = col_character(),
##   pval = col_double(),
##   sigma = col_double(),
##   FDR = col_double(),
##   ensembl_gene_id = col_character(),
##   CHR = col_character(),
##   start = col_double(),
##   end = col_double(),
##   entrezgene_id = col_character()
## )
```

Gene sets

Definition: two or more genes with shared function, structure, localization, or any other defined grouping variable

Use: glean interpretable, biologically-relevant meaning from a list of genes of interest

Gene sets can be as general as “cytosolic” to as specific as “up-regulated in blood vessel cells in response to wound in *Roy et al 2007*”. They can be defined by an individual research study (as in the latter) up to generally agreed upon knowledge formed from an entire body of literature (like the former). They also vary in size (2 to thousands), are not exclusive (one gene is in many gene sets), and many are redundant (transmembrane, membrane, lipid bilayer...).

It is up to you, the researcher, to select what group of gene sets you want to test. More on this below. Keep in mind, this is definitely not a case for “more is better”. For example, getting 200 significant gene sets from your 250 significant genes does not help you much.

Broad MSigDB

The Broad Institute has collated many useful gene sets in their Molecular Signatures Database (MSigDB). This includes the following used in this workshop (descriptions modified from Broad).

- Hallmark: summarize and represent specific well-defined biological processes. Generated by a computational methodology based on identifying overlaps between gene sets in other MSigDB collections and retaining genes that display coordinate expression

- Total: 50
- Examples: glycolysis, inflammatory response, apoptosis
- Curated (C2)
 - Canonical pathways: from pathway databases including KEGG, REACTOME, etc. Canonical representations of a biological process compiled by domain experts
 - Total: 2922
 - Examples: Caspase pathway, signaling by NOTCH1, DNA repair
- Gene ontology GO (C5)
 - Biological process: molecular-level activities performed by gene products
 - Total: 7481
 - Examples: viral life cycle, vitamin D biosynthetic process, mitochondrial calcium-ion transmembrane transport

Hallmark is the most general group and is usually where we start our analyses. It helps reduce the complexity of a large gene list down to easily understood terms. It also provides some context when looking at more specific gene sets as you can often group the other sets under a more general Hallmark term. However, because Hallmark is the most general, it is also the most incomplete with 50 gene sets representing just under 4400 genes.

Curated and gene ontology sets are more specific than Hallmark, though they contain a range of levels. They are particularly useful in adding details to a Hallmark term or defining testable hypotheses for follow-up experiments in the lab. They contain a lot more genes than Hallmark but many sets are redundant or overlapping. If there are many significant Hallmark terms, we often do not run these more specific sets as it gives too many significant results for interpretation.

There are several other sets in MSigDB that are not described here. Some of these are too specific to be of use in our global RNA-sequencing analyses (C2 chemical and genetic perturbations or C7 immunologic signatures) while other covers areas not relevant to our experimental design (C1 chromosome position, C3 gene regulation, C4/C6 cancer-oriented).

Next we will use gene sets to unravel how monocytes respond to *M. tuberculosis* infection.

Hypergeometric enrichment

This analysis is a simple enrichment. We ask if our list of significant genes is enriched in gene sets more than random chance would allow.

Get gene set data

First, we extract the lists of genes in Broad terms from the package `msigdb`. We format them as a data frames so we can use the `tidyverse`.

```
#Hallmark
H <- as.data.frame(msigdb(species = "Homo sapiens",
                          category = "H"))

#Canonical pathways
#Note that we have to call each database and combine them together
C2.CP <- as.data.frame(msigdb(species = "Homo sapiens",
                              category = "C2",
                              subcategory = "CP:BIOCARTA")) %>%
  bind_rows(as.data.frame(msigdb(species = "Homo sapiens",
                                category = "C2",
                                subcategory = "CP:KEGG")) %>%
    bind_rows(as.data.frame(msigdb(species = "Homo sapiens",
                                  category = "C2",
                                  subcategory = "CP:PID")) %>%
```

```

bind_rows(as.data.frame(msigdbr(species = "Homo sapiens",
                                category = "C2",
                                subcategory = "CP:REACTOME"))))
#Gene ontology biological process
C5.BP <- as.data.frame(msigdbr(species = "Homo sapiens",
                                category = "C5", subcategory = "GO:BP"))

```

Here, we see the start of the Hallmark data.

```

head(H)

##   gs_cat gs_subcat      gs_name gene_symbol entrez_gene
## 1      H          HALLMARK_ADIPOGENESIS      ABCA1         19
## 2      H          HALLMARK_ADIPOGENESIS      ABCB8        11194
## 3      H          HALLMARK_ADIPOGENESIS      ACAA2        10449
## 4      H          HALLMARK_ADIPOGENESIS      ACADL          33
## 5      H          HALLMARK_ADIPOGENESIS      ACADM          34
## 6      H          HALLMARK_ADIPOGENESIS      ACADS          35
##   ensembl_gene human_gene_symbol human_entrez_gene human_ensembl_gene gs_id
## 1 ENSG00000165029      ABCA1         19      ENSG00000165029 M5905
## 2 ENSG00000197150      ABCB8        11194      ENSG00000197150 M5905
## 3 ENSG00000167315      ACAA2        10449      ENSG00000167315 M5905
## 4 ENSG00000115361      ACADL          33      ENSG00000115361 M5905
## 5 ENSG00000117054      ACADM          34      ENSG00000117054 M5905
## 6 ENSG00000122971      ACADS          35      ENSG00000122971 M5905
##   gs_pmid gs_geoid gs_exact_source gs_url
## 1
## 2
## 3
## 4
## 5
## 6
##
##   gs_description
## 1 Genes up-regulated during adipocyte differentiation (adipogenesis).
## 2 Genes up-regulated during adipocyte differentiation (adipogenesis).
## 3 Genes up-regulated during adipocyte differentiation (adipogenesis).
## 4 Genes up-regulated during adipocyte differentiation (adipogenesis).
## 5 Genes up-regulated during adipocyte differentiation (adipogenesis).
## 6 Genes up-regulated during adipocyte differentiation (adipogenesis).

```

Define significant genes

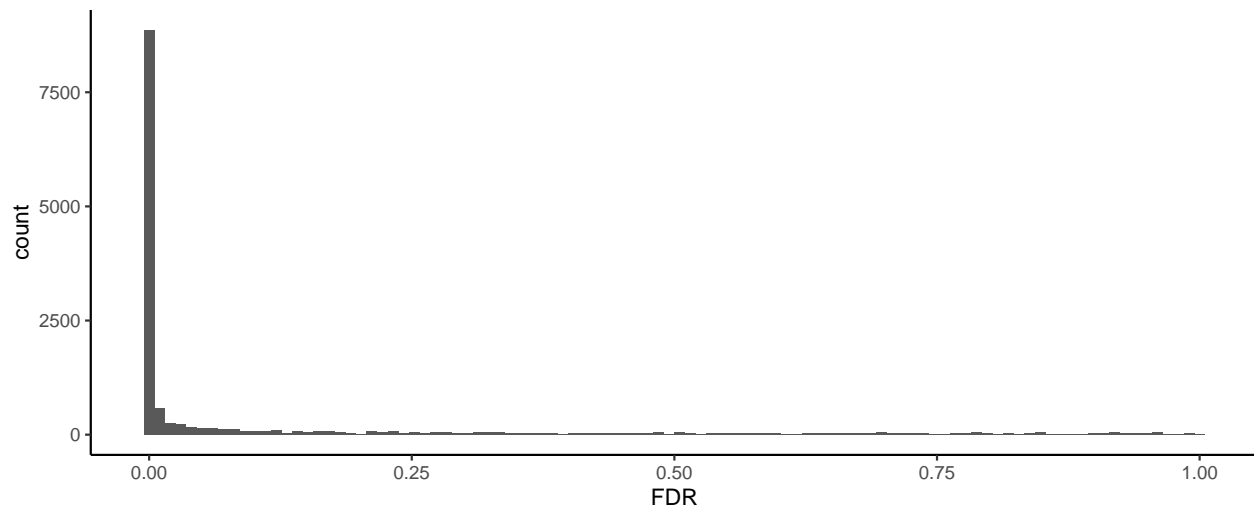
Next, we must define our “significant genes”. This type of enrichment is binary in that it only considers if a gene is significant or not. There is no granularity beyond yes/no. Thus, it is very important to consider our FDR cutoff as it is an arbitrary value that can dramatically impact results.

Let’s plot our FDR values as a histogram to help us determine a cutoff. We see that the data are heavily skewed. This is not uncommon when you have a large biological perturbation like bacterial infection.

```

ggplot(model.results, aes(x=FDR)) +
  geom_histogram(bins=100) +
  theme_classic()

```



In fact, there is a large proportion of FDR values that are zero.

```
table(model.results$FDR == 0)
```

```
##
## FALSE TRUE
## 10476 3502
```

These zeroes are actually a result of limitations in our linear modeling function. After a certain point (in this case, $1E-16$), the function stops computing and saves the P-value as zero to save time.

Here, we will use $FDR < 1E-16$ to see the genes that change the most with Mtb infection. *This is not necessarily what I'd recommend for these data.* It just for workshop purposes.

```
signif <- model.results %>%
  filter(FDR <= 1E-16)
```

Run enrichment

Now we can run the enrichment in `clusterProfiler`. We can use any gene ID in our dataset; it just has to match the gene IDs in the database.

Here, we run the enrichment using ENTREZ ID and the Hallmark database. Note that you must select ONLY the columns needed from the Hallmark database. `enricher()` isn't smart enough to know which columns to use on its own.

```
#Select unique ENSEMBL IDs for significant genes
signif.entrez <- unique(signif$entrezgene_id)
#Select matching ID column in database
H.entrez <- select(H, gs_name, entrez_gene)
#enrichment
enrich.H <- enricher(gene = signif.entrez, TERM2GENE = H.entrez)
```

Similarly, we could use ENSEMBL IDs.

```
#Select unique ENSEMBL IDs for significant genes
signif.ensembl <- unique(signif$ensembl_gene_id)
#Select matching ID column in database
H.ensembl <- select(H, gs_name, ensembl_gene)
#enrichment
enrich.H <- enricher(gene = signif.ensembl, TERM2GENE = H.ensembl)
```

Importantly, these will yield slightly different results, because the ID databases are not perfect. We must just accept that this is life and choose whichever ID we have / prefer. Since the expression data use ENSEMBL, we'll move forward with that.

```
length(signif.ensembl)
```

```
## [1] 3492
```

```
length(signif.entrez)
```

```
## [1] 3496
```

Extract enrichment results

clusterProfiler outputs results in an S4 object.

```
class(enrich.H)
```

```
## [1] "enrichResult"
```

```
## attr(,"package")
```

```
## [1] "DOSE"
```

We briefly mentioned this data type in intro R. It works very similarly to S3 (like the expression data) except you use @ to extract data. For example, the enrichment results are in a data frame.

```
head(enrich.H@result)
```

```
##
##                                     ID
## HALLMARK_TNFA_SIGNALING_VIA_NFKB   HALLMARK_TNFA_SIGNALING_VIA_NFKB
## HALLMARK_INFLAMMATORY_RESPONSE     HALLMARK_INFLAMMATORY_RESPONSE
## HALLMARK_INTERFERON_GAMMA_RESPONSE HALLMARK_INTERFERON_GAMMA_RESPONSE
## HALLMARK_IL2_STAT5_SIGNALING        HALLMARK_IL2_STAT5_SIGNALING
## HALLMARK_P53_PATHWAY                HALLMARK_P53_PATHWAY
## HALLMARK_APOPTOSIS                 HALLMARK_APOPTOSIS
##                                     Description GeneRatio
## HALLMARK_TNFA_SIGNALING_VIA_NFKB   HALLMARK_TNFA_SIGNALING_VIA_NFKB 164/1283
## HALLMARK_INFLAMMATORY_RESPONSE     HALLMARK_INFLAMMATORY_RESPONSE 112/1283
## HALLMARK_INTERFERON_GAMMA_RESPONSE HALLMARK_INTERFERON_GAMMA_RESPONSE 114/1283
## HALLMARK_IL2_STAT5_SIGNALING        HALLMARK_IL2_STAT5_SIGNALING 89/1283
## HALLMARK_P53_PATHWAY                HALLMARK_P53_PATHWAY 86/1283
## HALLMARK_APOPTOSIS                 HALLMARK_APOPTOSIS 74/1283
##                                     BgRatio    pvalue    p.adjust
## HALLMARK_TNFA_SIGNALING_VIA_NFKB   227/4892 5.621591e-50 2.810796e-48
## HALLMARK_INFLAMMATORY_RESPONSE     222/4892 2.541007e-15 6.352517e-14
## HALLMARK_INTERFERON_GAMMA_RESPONSE 284/4892 9.402307e-08 1.567051e-06
## HALLMARK_IL2_STAT5_SIGNALING        216/4892 6.864771e-07 8.580963e-06
## HALLMARK_P53_PATHWAY                214/4892 3.648903e-06 3.648903e-05
## HALLMARK_APOPTOSIS                 182/4892 1.083569e-05 9.029741e-05
##                                     qvalue
## HALLMARK_TNFA_SIGNALING_VIA_NFKB   2.011938e-48
## HALLMARK_INFLAMMATORY_RESPONSE     4.547065e-14
## HALLMARK_INTERFERON_GAMMA_RESPONSE 1.121679e-06
## HALLMARK_IL2_STAT5_SIGNALING        6.142163e-06
## HALLMARK_P53_PATHWAY                2.611846e-05
## HALLMARK_APOPTOSIS                 6.463394e-05
##
## HALLMARK_TNFA_SIGNALING_VIA_NFKB   ENSG00000109321/ENSG00000162772/ENSG00000070961/ENSG00000086062/ENSG00000109321
## HALLMARK_INFLAMMATORY_RESPONSE     ENSG00000162772/ENSG00000070961/ENSG00000086062/ENSG00000109321
```

```
## HALLMARK_INTERFERON_GAMMA_RESPONSE
## HALLMARK_IL2_STAT5_SIGNALING
## HALLMARK_P53_PATHWAY
## HALLMARK_APOPTOSIS
##                                     Count
## HALLMARK_TNFA_SIGNALING_VIA_NFKB    164
## HALLMARK_INFLAMMATORY_RESPONSE      112
## HALLMARK_INTERFERON_GAMMA_RESPONSE  114
## HALLMARK_IL2_STAT5_SIGNALING         89
## HALLMARK_P53_PATHWAY                  86
## HALLMARK_APOPTOSIS                   74
```

This includes

```
colnames(enrich.H@result)
```

```
## [1] "ID"          "Description" "GeneRatio"   "BgRatio"     "pvalue"
## [6] "p.adjust"    "qvalue"      "geneID"      "Count"
```

- ID: Unique name for gene set
- Description: Descriptive text for gene set, often the same as ID
- GeneRatio: Significant genes in gene set / Significant genes in all gene sets tested
- BgRatio: Total genes in gene set / Total genes in all gene sets tested
- pvalue: Significance
- p.adjust: FDR adjust significance
- qvalue: Q-value used in P-value calculation
- geneID: vector of significant gene IDs in gene set. Separated by /
- Count: Significant genes in gene set

These column names aren't the most informative and R doesn't understand the ratios as numbers.

```
class(enrich.H@result$GeneRatio)
```

```
## [1] "character"
```

So we'll do some additional formatting. We refer to a gene set as a "term" and the Hallmark database as a "category" for brevity.

```
enrich.H.df <- enrich.H@result %>%
  #separate ratios into 2 columns of data
  separate(BgRatio, into=c("size.term", "size.category"), sep="/") %>%
  separate(GeneRatio, into=c("size.overlap.term", "size.overlap.category"),
           sep="/") %>%
  #convert to numeric
  mutate_at(vars("size.term", "size.category",
                 "size.overlap.term", "size.overlap.category"),
            as.numeric) %>%
  #Calculate k/K
  mutate("k.K"=size.overlap.term/size.term)
```

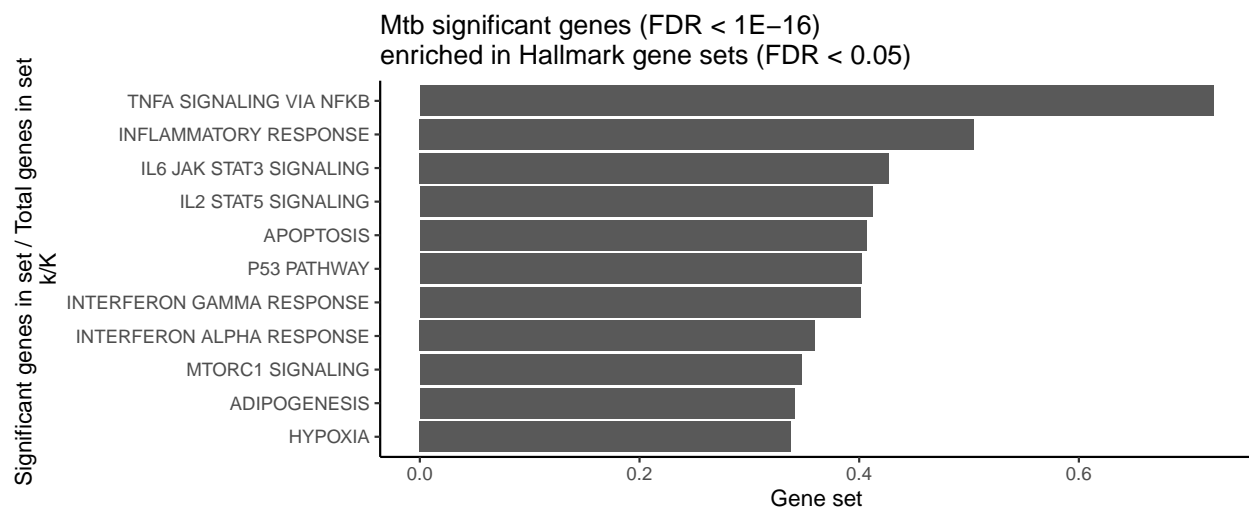
Visualize significant enrichment

The most common visualization for enrichment is k/K as it is the ratio of significant genes relative to total genes in the gene set. We also subset to significant enrichment at an defined FDR. Here, let's use $FDR < 0.05$. This is actually pretty strict of enrichment and people often go up to $FDR < 0.2$.

```
enrich.H.df %>%
  filter(p.adjust <= 0.05) %>%
  #Beautify descriptions by removing _ and HALLMARK
```

```
mutate(Description = gsub("HALLMARK_", "", Description),
       Description = gsub("_", " ", Description)) %>%

ggplot(aes(x=reorder(Description, k.K), #Reorder gene sets by k/K values
            y=k.K)) +
  geom_col() +
  theme_classic() +
  #Some more customization to pretty it up
  #Flip x and y so long labels can be read
  coord_flip() +
  #fix labels
  labs(x="Significant genes in set / Total genes in set \nk/K",
       y="Gene set",
       title = "Mtb significant genes (FDR < 1E-16)\nenriched in Hallmark gene sets (FDR < 0.05)")
```



Enrichment exercises

1. Run enrichment of the current significant gene list against the canonical pathways C2.CP database. Do you get more or fewer significant gene sets? Is this what you'd expect given what you know about C2?
2. Run Hallmark enrichment on genes significant at FDR < 0.01. Do you get more or fewer significant gene sets? Why might this be the case?
3. Further modify the final ggplot. Add color, change the FDR cutoff, do whatever your heart desires!

Gene set enrichment analysis (GSEA)

Additional resources

Groups

- Rladies Seattle Not just for ladies! A pro-actively inclusive R community with both in-person and online workshops, hangouts, etc.
- TidyTuesday A weekly plotting challenge
- R code club Dr. Pat Schloss opened his lab's coding club to remote participation.
- Seattle useR Group

Online

- R cheatsheets also available in RStudio under Help > Cheatsheets
- The Carpentries

R session

```
sessionInfo()
```

```
## R version 4.0.2 (2020-06-22)
## Platform: x86_64-apple-darwin17.0 (64-bit)
## Running under: macOS 10.16
##
## Matrix products: default
## BLAS: /Library/Frameworks/R.framework/Versions/4.0/Resources/lib/libRblas.dylib
## LAPACK: /Library/Frameworks/R.framework/Versions/4.0/Resources/lib/libRlapack.dylib
##
## locale:
## [1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
##
## attached base packages:
## [1] stats      graphics  grDevices  utils      datasets  methods   base
##
## other attached packages:
## [1] fgsea_1.14.0      clusterProfiler_3.16.1 msigdb_7.4.1
## [4] forcats_0.5.1     stringr_1.4.0          dplyr_1.0.7
## [7] purrr_0.3.4       readr_1.4.0            tidyr_1.1.3
## [10] tibble_3.1.2      ggplot2_3.3.5          tidyverse_1.3.1
##
## loaded via a namespace (and not attached):
## [1] colorspace_2.0-2  ggribes_0.5.3          ellipsis_0.3.2
## [4] qvalue_2.20.0     fs_1.5.0               rstudioapi_0.13
## [7] farver_2.1.0      urltools_1.7.3         graphlayouts_0.7.1
## [10] ggrepel_0.9.1     bit64_4.0.5            scatterpie_0.1.6
## [13] AnnotationDbi_1.50.3 fansi_0.5.0            lubridate_1.7.10
## [16] xml2_1.3.2        splines_4.0.2          cachem_1.0.5
## [19] GOSemSim_2.14.2   knitr_1.33             polyclip_1.10-0
## [22] jsonlite_1.7.2    broom_0.7.8            GO.db_3.11.4
## [25] dbplyr_2.1.1      ggforce_0.3.3          BiocManager_1.30.16
## [28] compiler_4.0.2    httr_1.4.2             rvcheck_0.1.8
## [31] backports_1.2.1   assertthat_0.2.1       Matrix_1.3-4
## [34] fastmap_1.1.0     cli_3.0.0              tweenr_1.0.2
## [37] htmltools_0.5.1.1 prettyunits_1.1.1      tools_4.0.2
## [40] igraph_1.2.6      gtable_0.3.0           glue_1.4.2
## [43] reshape2_1.4.4    DO.db_2.9              fastmatch_1.1-0
## [46] Rcpp_1.0.7        enrichplot_1.8.1       Biobase_2.48.0
## [49] cellranger_1.1.0  vctrs_0.3.8            babelgene_21.4
## [52] ggraph_2.0.5      xfun_0.24              rvest_1.0.0
## [55] lifecycle_1.0.0   DOSE_3.14.0            europepmc_0.4
## [58] MASS_7.3-54       scales_1.1.1           tidygraph_1.2.0
## [61] hms_1.1.0         parallel_4.0.2         RColorBrewer_1.1-2
## [64] yaml_2.2.1        memoise_2.0.0          gridExtra_2.3
## [67] downloader_0.4    triebeard_0.3.0        stringi_1.6.2
## [70] RSQLite_2.2.7     highr_0.9              S4Vectors_0.26.1
```

## [73]	BiocGenerics_0.34.0	BiocParallel_1.24.1	rlang_0.4.11
## [76]	pkgconfig_2.0.3	evaluate_0.14	lattice_0.20-44
## [79]	labeling_0.4.2	cowplot_1.1.1	bit_4.0.4
## [82]	tidyselect_1.1.1	plyr_1.8.6.9000	magrittr_2.0.1
## [85]	R6_2.5.0	IRanges_2.22.2	generics_0.1.0
## [88]	DBI_1.1.1	pillar_1.6.1	haven_2.4.1
## [91]	withr_2.4.2	modelr_0.1.8	crayon_1.4.1
## [94]	utf8_1.2.1	rmarkdown_2.9	viridis_0.6.1
## [97]	progress_1.2.2	grid_4.0.2	readxl_1.3.1
## [100]	data.table_1.14.0	blob_1.2.1	reprex_2.0.0
## [103]	digest_0.6.27	gridGraphics_0.5-1	stats4_4.0.2
## [106]	munsell_0.5.0	viridisLite_0.4.0	ggplotify_0.0.7
