

Dorothea tutorial

Xiaosai Yao

12 November 2022

Package

epiregulon 1.0.16

Contents

1	Load regulon	2
2	Load scRNA-seq data	2
3	Calculate activity	3
4	Visualize activity	4
5	Pathway enrichment	7
6	Session Info	9

Epiregulon also supports transcription factor activity inference when users only have scRNA-seq. After all, multiome or scATAC-seq data is still relatively rare. To enable TF activity inference on scRNA-seq, users can supply a pre-constructed gene regulatory network. [Dorothea](#) provides both human and mouse pre-constructed gene regulatory networks based on curated experimental and computational data. In this vignette, we bypass the regulon construction step and go straight to calculate TF activity from a Dorothea GRN.

1 Load regulon

Dorothea assigns confidence level to its regulons with A being the most confident (i.e. supported by multiple lines of evidence) and E being the least confident. For this demo, we further trim the regulons to only 4 TFs.

```
library(dorothea)
data(dorothea_mm, package = "dorothea")
regulon <- dorothea_mm

#trim regulon
genes_to_plot <- c("Foxa1", "Neurod1", "Pdx1", "Arx")
regulon <- regulon[regulon$tf %in% genes_to_plot, ]
```

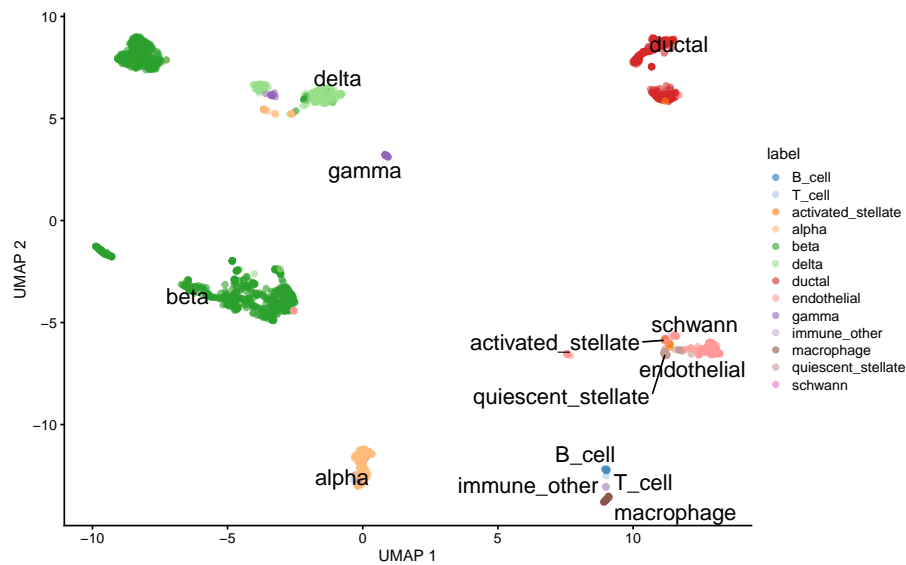
2 Load scRNA-seq data

We download the raw counts of a mouse pancreas data set from [scRNAseq](#). We add normalized logcounts, perform dimension reduction and visualize the embeddings using [scater](#).

```
library(scRNAseq)
library(scater)

sce <- BaronPancreasData('mouse')
sce <- logNormCounts(sce)
sce <- runPCA(sce)
sce <- runUMAP(sce)

plotUMAP(sce, colour_by = "label", text_by = "label")
```



3 Calculate activity

Even though Dorothea provides weights under the `mor` column, we achieved superior performance if we recompute the weights based on the correlation between `tf` and target gene expression based on our own data. We performed 2 steps, the first step is to add weights to the Dorothea regulons and the second step is to estimate the TF activity by taking the weighted average of the target gene expression.

```
library(epiregulon)

#Add weights to regulon
regulon.ms <- addWeights(regulon = regulon,
                        expMatrix = sce,
                        clusters = sce$label,
                        BPPARAM = BiocParallel::MulticoreParam())

## computing correlation of the regulon...
##
```

	0%
=====	25%
=====	50%
=====	75%
=====	100%

```
#Calculate activity
score.combine <- calculateActivity(sce,
```

```

regulon = regulon.ms,
mode = "weight",
method = "weightedMean")

```

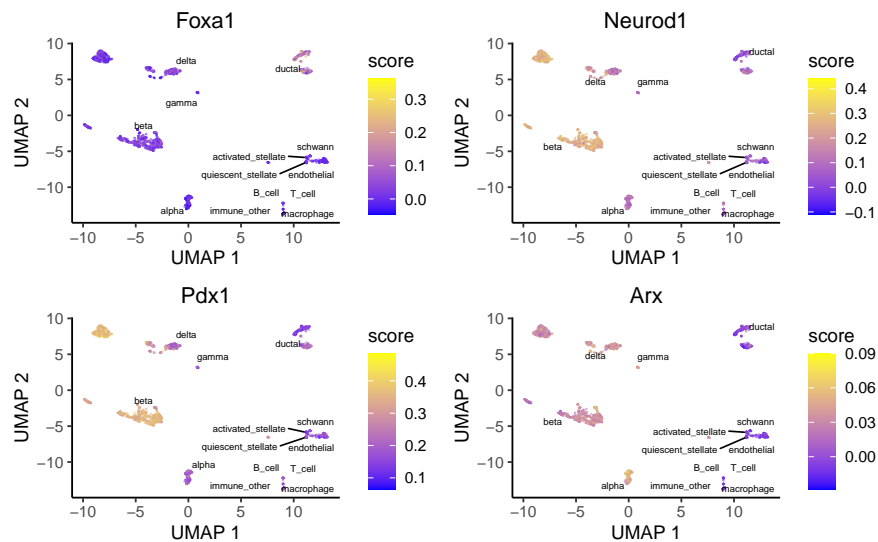
4 Visualize activity

Finally we visualize the TF activity by either UMAP, violin plots or bubble plots. We confirm the activity of known lineage factors Pdx1 and Neurod1 in beta cells, Arx in alpha cells and Foxa1 in ductal cells.

```

# plot umap
plotActivityDim(sce = sce,
  activity_matrix = score.combine,
  tf = genes_to_plot,
  legend.label = "score",
  point_size = 0.1,
  dimtype = "UMAP",
  label = "label",
  combine = TRUE,
  text_size = 2)

```

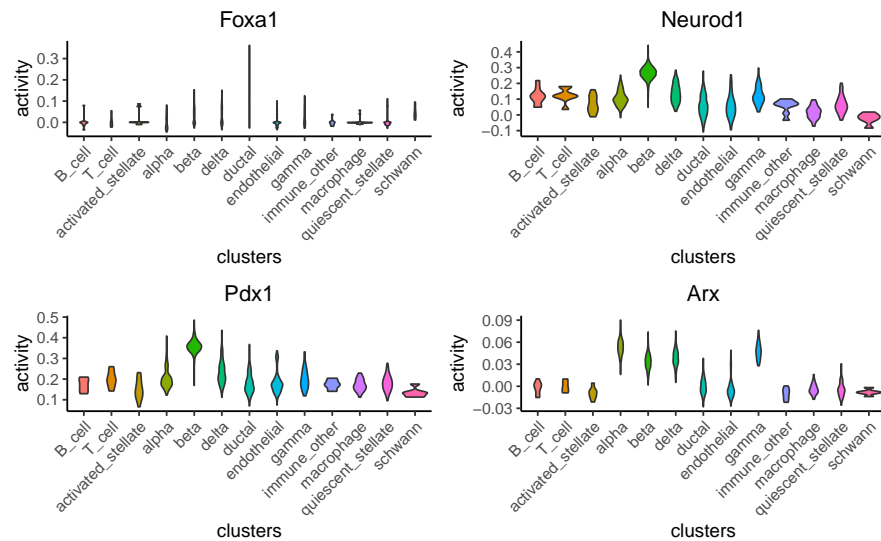


```

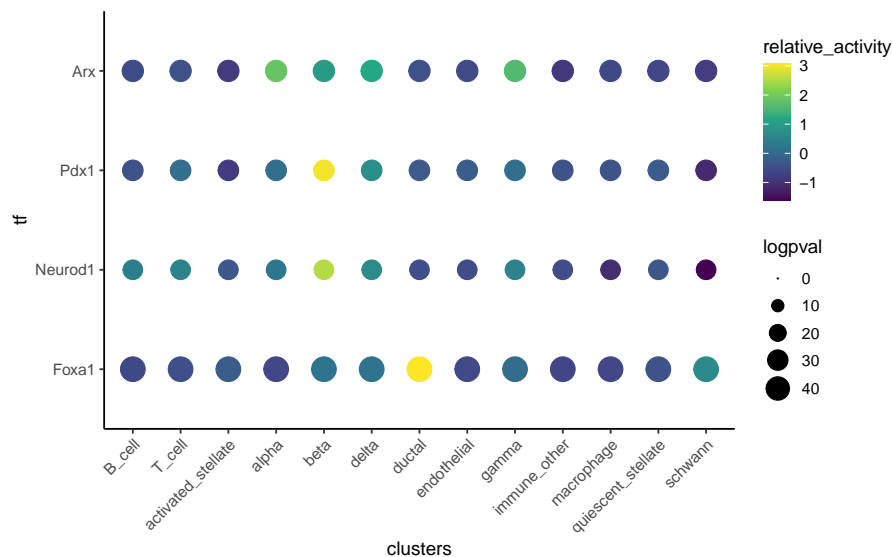
# plot violin plot
plotActivityViolin(score.combine,
  tf = genes_to_plot,
  clusters = sce$label)

```

Dorothea tutorial



```
# plot Bubble plot
plotBubble(score.combine,
            tf = genes_to_plot,
            clusters = sce$label)
```

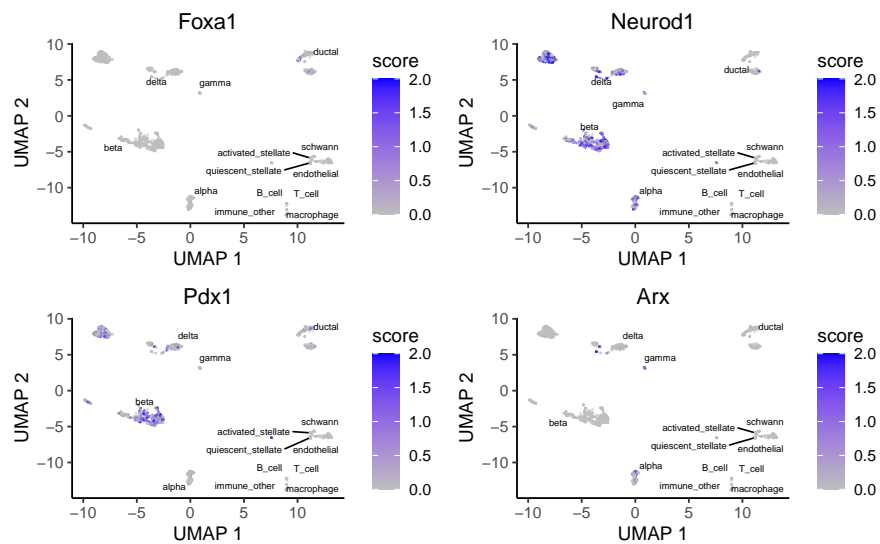


We can adapt the epiregulon package to plot gene expression. When compared against TF activity, gene expression of Foxa1 and Arx has noisy signals and high dropout rates. Epiregulon enhances the signal to noise ratio of TF activity and better resolves lineage differences.

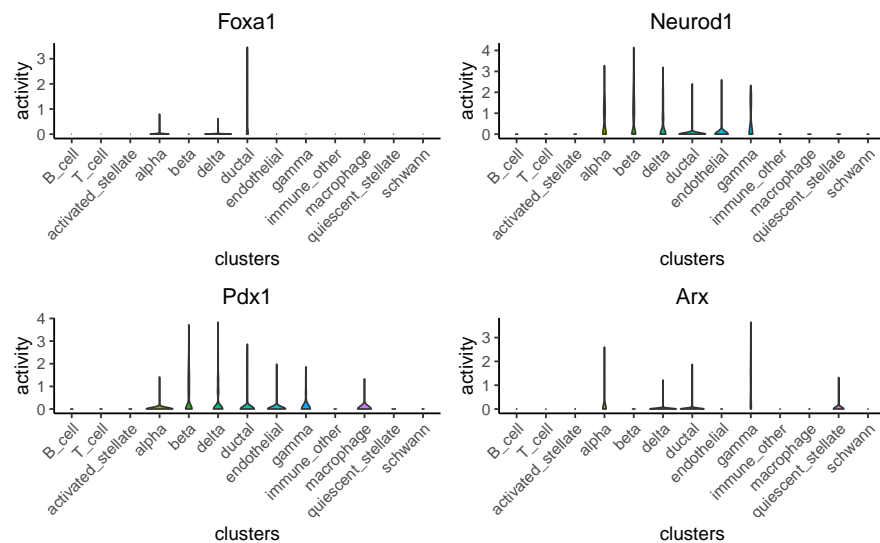
```
# plot umap
plotActivityDim(sce = sce,
                activity_matrix = logcounts(sce)[genes_to_plot,],
                tf = genes_to_plot,
                legend.label = "score",
```

Dorothea tutorial

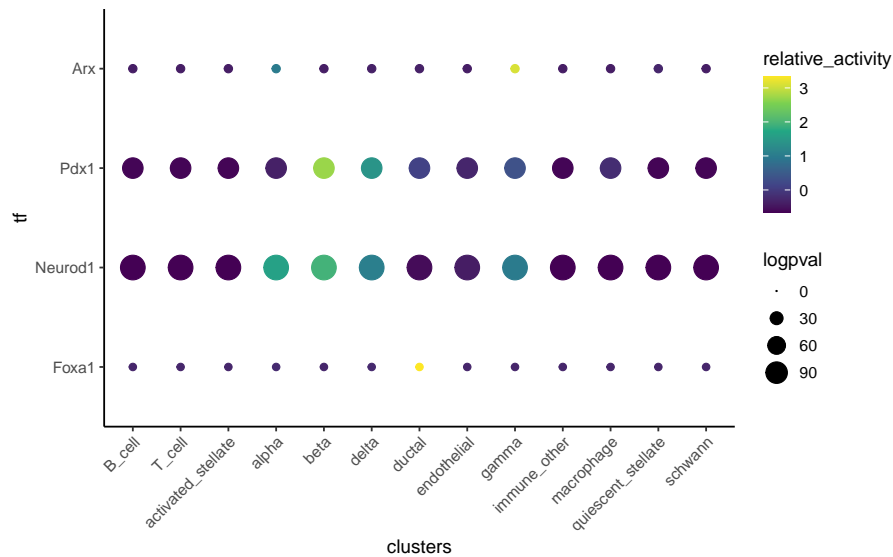
```
point_size = 0.1,
dimtype = "UMAP",
label = "label",
combine = TRUE,
text_size = 2,
colors = c("gray", "blue"),
limit = c(0,2))
```



```
# plot violin plot
plotActivityViolin(logcounts(sce)[genes_to_plot,],
  tf = genes_to_plot,
  clusters = sce$label)
```



```
# plot Bubble plot
plotBubble(logcounts(sce)[genes_to_plot,],
           tf = genes_to_plot,
           clusters = sce$label)
```



5 Pathway enrichment

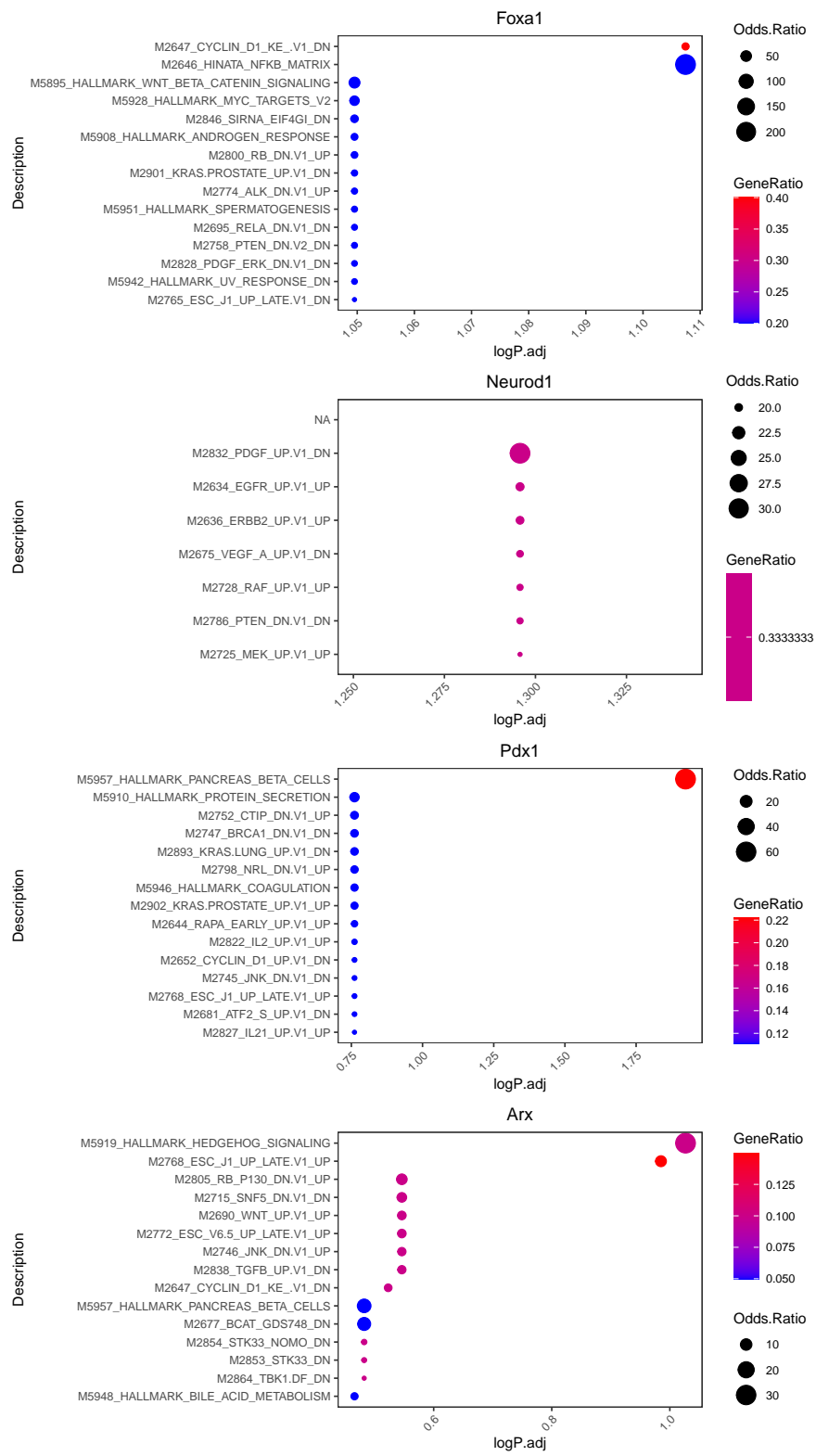
Sometimes it is useful to understand what pathways are enriched in the regulons. We take the highly correlated target genes of a regulon and perform geneset enrichment using the `enricher` function from [clusterProfiler](#).

```
#retrieve genesets
H <- EnrichmentBrowser::getGenesets(org = "mmu", db = "msigdb", cat = "H", gene.id.type = "SYMBOL" )
## Using cached version from 2022-11-12 19:27:57
C6 <- EnrichmentBrowser::getGenesets(org = "mmu", db = "msigdb", cat = "C6", gene.id.type = "SYMBOL" )
## Using cached version from 2022-11-12 19:28:02

#combine genesets and convert genesets to be compatible with enricher
gs <- c(H,C6)
gs.list <- do.call(rbind,lapply(names(gs), function(x)
  {data.frame(gs = x, genes = gs[[x])}))

enrichresults <- regulonEnrich(genes_to_plot,
                              regulon = regulon.ms,
                              corr = "weight",
                              corr_cutoff = 0.5,
                              genesets = gs.list)

#plot results
enrichPlot(results = enrichresults, ncol = 1)
```



6 Session Info

```

sessionInfo()
## R version 4.2.0 (2022-04-22)
## Platform: x86_64-pc-linux-gnu (64-bit)
## Running under: Ubuntu 18.04.6 LTS
##
## Matrix products: default
## BLAS: /usr/local/lib/R/lib/libRblas.so
## LAPACK: /usr/local/lib/R/lib/libRlapack.so
##
## Random number generation:
## RNG: L'Ecuyer-CMRG
## Normal: Inversion
## Sample: Rejection
##
## locale:
## [1] LC_CTYPE=en_US.UTF-8 LC_NUMERIC=C LC_TIME=C LC_COLLATE=C
## [5] LC_MONETARY=C LC_MESSAGES=C LC_PAPER=C LC_NAME=C
## [9] LC_ADDRESS=C LC_TELEPHONE=C LC_MEASUREMENT=C LC_IDENTIFICATION=C
##
## attached base packages:
## [1] parallel grid stats4 stats graphics grDevices utils datasets
## [9] methods base
##
## other attached packages:
## [1] scater_1.26.0 scuttle_1.8.0 scRNAseq_2.12.0
## [4] dorothea_1.10.0 nabor_0.5.0 rhdf5_2.42.0
## [7] Rcpp_1.0.9 Matrix_1.5-3 data.table_1.14.4
## [10] stringr_1.4.0 plyr_1.8.8 magrittr_2.0.3
## [13] gtable_0.3.1 gtools_3.9.3 gridExtra_2.3
## [16] ArchR_1.0.2 ggplot2_3.4.0 msigdb_7.5.1
## [19] epiregulon_1.0.16 SingleCellExperiment_1.20.0 SummarizedExperiment_1.29.1
## [22] Biobase_2.58.0 GenomicRanges_1.50.1 GenomeInfoDb_1.34.3
## [25] IRanges_2.32.0 S4Vectors_0.36.0 BiocGenerics_0.44.0
## [28] MatrixGenerics_1.10.0 matrixStats_0.62.0 BiocStyle_2.26.0
##
## loaded via a namespace (and not attached):
## [1] rappdirs_0.3.3 rtracklayer_1.58.0
## [3] tidyr_1.2.1 bit64_4.0.5
## [5] knitr_1.40 irlba_2.3.5.1
## [7] DelayedArray_0.24.0 AnnotationFilter_1.22.0
## [9] KEGGREST_1.38.0 RCurl_1.98-1.9
## [11] generics_0.1.3 GenomicFeatures_1.50.2
## [13] ScaledMatrix_1.6.0 cowplot_1.1.1
## [15] RSQLite_2.2.18 shadowtext_0.1.2
## [17] artificer.mae_1.3.4 bit_4.0.4
## [19] enrichplot_1.18.0 base64url_1.4
## [21] gp.cache_1.7.1 xml2_1.3.3
## [23] httpuv_1.6.6 assertthat_0.2.1

```

```
## [25] genomitory_2.1.6      viridis_0.6.2
## [27] xfun_0.31              hms_1.1.2
## [29] promises_1.2.0.1       babelgene_22.9
## [31] evaluate_0.18          restfulr_0.0.15
## [33] progress_1.2.2         fansi_1.0.3
## [35] dbplyr_2.2.1           Rgraphviz_2.42.0
## [37] igraph_1.3.5           DBI_1.1.3
## [39] ellipsis_0.3.2        purrr_0.3.5
## [41] dplyr_1.0.10          backports_1.4.1
## [43] bookdown_0.30         annotate_1.76.0
## [45] biomaRt_2.54.0        sparseMatrixStats_1.10.0
## [47] artificer.matrix_1.3.7 vctrs_0.5.0
## [49] ensemblDb_2.22.0      Cairo_1.6-0
## [51] dsdb.plus_1.3.2       cachem_1.0.6
## [53] withr_2.5.0           ggforce_0.4.1
## [55] HD0.db_0.99.0         checkmate_2.1.0
## [57] GenomicAlignments_1.34.0 metacommons_1.9.0
## [59] treeio_1.22.0         prettyunits_1.1.1
## [61] MultiAssayExperiment_1.24.0 scran_1.26.0
## [63] cluster_2.1.3         DOSE_3.23.3
## [65] ExperimentHub_2.6.0   BiocBaseUtils_1.1.0
## [67] ape_5.6-2             lazyeval_0.2.2
## [69] crayon_1.5.1          edgeR_3.40.0
## [71] pkgconfig_2.0.3       labeling_0.4.2
## [73] tweenr_2.0.2          ProtGenerics_1.30.0
## [75] nlme_3.1-160          vipor_0.4.5
## [77] rlang_1.0.6           lifecycle_1.0.3
## [79] artificer.schemas_0.99.2 downloader_0.4
## [81] filelock_1.0.2        artificer.base_1.3.19
## [83] BiocFileCache_2.6.0   rsvd_1.0.5
## [85] AnnotationHub_3.6.0   polyclip_1.10-4
## [87] GSVA_1.46.0           graph_1.76.0
## [89] aplot_0.1.8           Rhdf5lib_1.20.0
## [91] beeswarm_0.4.0        rjson_0.2.21
## [93] png_0.1-7             viridisLite_0.4.1
## [95] artificer.ranges_1.3.4 bitops_1.0-7
## [97] artificer.se_1.3.4    getPass_0.2-2
## [99] gson_0.0.9            rhdf5filters_1.10.0
## [101] EnrichmentBrowser_2.28.0 Biostrings_2.66.0
## [103] blob_1.2.3            DelayedMatrixStats_1.20.0
## [105] qvalue_2.30.0         gridGraphics_0.5-1
## [107] beachmat_2.14.0       scales_1.2.1
## [109] memoise_2.0.1         GSEABase_1.60.0
## [111] zlibbioc_1.44.0       compiler_4.2.0
## [113] scatterpie_0.1.8      BiocIO_1.8.0
## [115] dqrng_0.3.0           tinytex_0.42
## [117] RColorBrewer_1.1-3    KEGGgraph_1.58.0
## [119] Rsamtools_2.14.0      cli_3.4.1
## [121] XVector_0.38.0        patchwork_1.1.2
## [123] ArtifactDB_1.9.5      MASS_7.3-58.1
## [125] tidyselect_1.2.0      stringi_1.7.6
```

```
## [127] yaml_2.3.5
## [129] BiocSingular_1.14.0
## [131] ggrepel_0.9.2
## [133] fastmatch_1.1-3
## [135] rstudioapi_0.13
## [137] metapod_1.6.0
## [139] ggraph_2.1.0
## [141] BiocManager_1.30.19
## [143] shiny_1.7.3
## [145] BiocVersion_3.16.0
## [147] httr_1.4.3
## [149] rsconnect_0.8.28
## [151] XML_3.99-0.12
## [153] uwot_0.1.14
## [155] statmod_1.4.37
## [157] graphlayouts_0.8.3
## [159] xtable_1.8-4
## [161] ggtree_3.6.2
## [163] ggfun_0.0.8
## [165] R6_2.5.1
## [167] pillar_1.8.1
## [169] glue_1.6.2
## [171] clusterProfiler_4.6.0
## [173] BiocNeighbors_1.16.0
## [175] codetools_0.2-18
## [177] gp.version_1.5.0
## [179] lattice_0.20-45
## [181] curl_4.3.2
## [183] ggbeeswarm_0.6.0
## [185] GO.db_3.16.0
## [187] dsassembly_1.7.6
## [189] rmarkdown_2.18
## [191] GenomeInfoDbData_1.2.9
## [193] reshape2_1.4.4
GOSemSim_2.24.0
locfit_1.5-9.6
bcellViper_1.34.0
tools_4.2.0
bluster_1.8.0
farver_2.1.1
digest_0.6.29
FNN_1.1.3.1
later_1.3.0
gp.auth_1.7.0
AnnotationDbi_1.60.0
colorspace_2.0-3
splines_4.2.0
yulab.utils_0.0.5
tidytree_0.4.1
ggplotify_0.1.0
jsonlite_1.8.3
tidygraph_1.2.2
ShadowArray_1.7.1
mime_0.12
htmltools_0.5.3
fastmap_1.1.0
BiocParallel_1.32.1
interactiveDisplayBase_1.36.0
fgsea_1.24.0
utf8_1.2.2
tibble_3.1.8
genomitory.schemas_0.99.0
artificer.sce_1.3.4
limma_3.54.0
dsdb.schemas_0.99.1
munSELL_0.5.0
HDF5Array_1.26.0
```