

Epiregulon

Shang-Yang Chen* Xiaosai Yao†

March 29th, 2022

Introduction

Gene regulatory networks model the underlying gene regulation hierarchies that drive gene expression and observed phenotypes. The main function of the epiregulon R package is to infer TF activity in single cells by constructing a gene regulatory network (regulons). This is achieved through integration of scATAC-seq and scRNA-seq data and incorporation of public bulk TF ChIP-seq data. Links between regulatory elements and their target genes are established by computing correlations between chromatin accessibility and gene expressions.

Current prerequisite for running epiregulon is a ArchR project with pre-computed peak and gene matrices. It is also expected that LSI dimensionality reduction and integration with an scRNA-seq dataset has been performed. The scATAC-seq experiment can be either paired or unpaired with the scRNA-seq dataset as long as they were already integrated by ArchR. The final output of epiregulon is a matrix of TF activities where rows are individual TFs and columns are single cell indexes.

Alternatively, users can now supply peak, gene, and dimensional reduction matrices derived from a MultiAssayExperiment object. This is to be compatible with future GPSA multiome workflow. Epiregulon implements a custom algorithm to derive a more stringent set of P2G correlations compared to ArchR.

In this vignette we demonstrate the workflow of epiregulon along with some visualization functionalities using the tutorial datasets from ArchR development team. In this dataset, scRNaseq and scATACseq were unpaired and integrated by the `addGeneIntegrationMatrix` function.

Installation

Epiregulon is currently available on R/dev

```
library(epiregulon)
```

If you would like to install from gitlab,

```
devtools::install_gitlab(repo='scwg/gene-transcriptional-network/activity-inference/Epiregulon',
                         auth_token = "<gitlab token>",
                         host = "https://code.roche.com/" )
library(epiregulon)
```

*chens179@gene.com

†yaox19@gene.com

Data preparation

Please refer to the full ArchR manual for instructions

Before running Epiregulon, the following analyses need to be completed: 1. Obtain a peak matrix on scATAC-seq by using addGroupCovariates > addReproduciblePeakSet > addPeakMatrix. See chapter 10 from ArchR manual 2. RNA-seq integration. a. For unpaired scATAC-seq, use addGeneIntegrationMatrix. See chapter 8 from ArchR manual b. For multiome data, use addGeneExpressionMatrix. See multiome tutorial 3. Perform dimensionality reduction from with either single modalities or joint scRNA-seq and scATAC-seq using addCombinedDims

To verify that all the necessary matrices are present,

```
library(ArchR)
archR_project_path="/gstore/project/lineage/sam/heme_GRN/OUTPUT"
proj<- loadArchRProject(path = archR_project_path, showLogo = FALSE)
#> Successfully loaded ArchRProject!
getAvailableMatrices(proj)
#> [1] "GeneIntegrationMatrix" "GeneScoreMatrix"           "PeakMatrix"
#> [4] "TileMatrix"
```

Quick start

Retrieve bulk TF ChIP-seq binding sites

First, we retrieve the information of TF binding sites collected from Cistrome and ENCODE ChIP-seq, which are hosted on Genomitory. Currently, human genomes hg19 and hg38 and mouse genome mm10 are available

```
grl <- getTFMotifInfo(genome = "hg19")
head(grl)
#> GRangesList object of length 6:
#> $ADNP
#> GRanges object with 20474 ranges and 0 metadata columns:
#>   seqnames      ranges strand
#>   <Rle>      <IRanges>  <Rle>
#>   [1] chr1    565199-565456   *
#>   [2] chr1    569272-569544   *
#>   [3] chr1    895823-896145   *
#>   [4] chr1    946692-947524   *
#>   [5] chr1    993405-993727   *
#>   ...     ...
#>   [20470] chrX  153402026-153403267  *
#>   [20471] chrX  153729859-153730943  *
#>   [20472] chrX  153979887-153980719  *
#>   [20473] chrX  154297459-154298291  *
#>   [20474] chrX  154446717-154447848  *
#> -----
#>   seqinfo: 52 sequences from an unspecified genome; no seqlengths
#>
#> ...
#> <5 more elements>
```

Link ATACseq peaks to target genes

Next, we compute peak to gene correlations using the calculateP2G function from ArchR package. The user would need to supply a path to an ArchR project that already contains the peak matrix, gene expression matrix and Latent semantic indexing (LSI) dimensionality reduction. The example project shown here utilizes the tutorial datasets provided by the ArchR development team.

```
# path to ArchR project
p2g <- calculateP2G(ArchR_path = archR_project_path)
#> Setting ArchRLogging = FALSE
#> 2022-04-21 12:17:33 : Getting Available Matrices, 0 mins elapsed.
#> 2022-04-21 12:17:41 : Filtered Low Prediction Score Cells (684 of 10250, 0.067), 0.002 mins elapsed.
#> 2022-04-21 12:17:41 : Computing KNN, 0.007 mins elapsed.
#> 2022-04-21 12:17:53 : Identifying Non-Overlapping KNN pairs, 0.215 mins elapsed.
#> 2022-04-21 12:17:56 : Identified 494 Groupings!, 0.261 mins elapsed.
#> 2022-04-21 12:17:56 : Getting Group RNA Matrix, 0.265 mins elapsed.
#> 2022-04-21 12:19:13 : Getting Group ATAC Matrix, 1.537 mins elapsed.
#> 2022-04-21 12:20:15 : Normalizing Group Matrices, 2.569 mins elapsed.
#> 2022-04-21 12:20:19 : Finding Peak Gene Pairings, 2.647 mins elapsed.
#> 2022-04-21 12:20:20 : Computing Correlations, 2.655 mins elapsed.
#> 2022-04-21 12:20:26 : Completed Peak2Gene Correlations!, 2.759 mins elapsed.

head(p2g)
#>   idxATAC Chrom idxRNA      Gene Correlation
#> 1     6    chr1     9    ISG15  0.5308352
#> 3     24   chr1     6    KLHL17  0.5698560
#> 2     24   chr1    14 TNFRSF18  0.5651328
#> 5     25   chr1     9    ISG15  0.5895757
#> 4     25   chr1    14 TNFRSF18  0.5680771
#> 6     37   chr1     8    HES4   0.6679620
```

Alternatively, users can now also supply peak, gene, and dimensional reduction matrices derived from a MultiAssayExperiment object. This is compatible with future GPSA multiome workflow. Epiregulon implements a custom algorithm to derive a more stringent set of P2G correlations compared to ArchR.

```
# load the MAE object
heme <- readRDS("/gstore/project/archr_importer/heme_grn_unpaired_data/mae.rds")
# peak matrix
peakmatrix <- heme[["PeakMatrix"]]
# expression matrix
expmatrix <- heme[["GeneIntegrationMatrix"]]
rownames(expmatrix) <- rowData(expmatrix)$name
# dimensional reduction matrix
reducedDim <- SingleCellExperiment::reducedDims(heme[['TileMatrix500']])[[["IterativeLSI"]]]

p2g <- calculateP2G(peakmatrix, expmatrix, reducedDim)

head(p2g)
```

Add TF motif binding to peaks

The next step is to add the TF motif binding information by overlapping the regions of the peak matrix with the bulk chip-seq database loaded in 2. The user can supply either an archR project path and this function will retrieve the peak matrix, or a peakMatrix in the form of a Granges object or RangedSummarizedExperiment.

```

overlap <- addTFMotifInfo(p2g, grl, archR_project_path = archR_project_path)
#> Successfully loaded ArchRProject!
#> Computing overlap...
#> Success!
head(overlap)
#>      idxATAC idxTF      tf
#> 522       6   427    MAX
#> 523       6   601  POLR2G
#> 524       6   785    TAL1
#> 525       6   796   TCF12
#> 526       6   809    TERC
#> 2200     24       3    AFF1

```

Generate regulons

A long format dataframe, representing the inferred regulons, is then generated. The dataframe consists of three columns:

- tf (transcription factor)
- target gene
- peak to gene correlation between tf and target gene

```

regulon <- getRegulon(p2g, overlap, aggregate=TRUE)
head(regulon)
#>      tf target      corr
#> 1  AFF1 A2M-AS1 0.6021363
#> 2 ARID4A A2M-AS1 0.6021363
#> 3 ARID4B A2M-AS1 0.6021363
#> 4 ASH2L A2M-AS1 0.6021363
#> 5 ATF1 A2M-AS1 0.6021363
#> 6 ATF2 A2M-AS1 0.6021363

```

Epiregulon outputs two different correlations. The first, termed “corr”, is the correlation between chromatin accessibility of regulatory elements vs expression of target genes calculated by ArchR. The second, termed “weight”, can be generated by the addWeights function, which compute the correlation between gene expressions of TF vs expressions of target genes, shown below. The user is required to supply the clustering or batch labels of the scRNA-seq dataset when running addWeights. “Weight” is the preferred metric for calculating activity.

load scRNA-seq data

```
sce=readRDS("/gstore/project/lineage/sam/heme_GRN/scRNA-Granja-2019.rds")
```

Trim regulon for demonstration purposes

```

TFs <- c("FOXA1", "GATA3", "SOX9", "SPI1")
regulon <- regulon %>% dplyr::filter(tf %in% TFs)
nrow(regulon)
#> [1] 15336

regulon.w <- addWeights(regulon,
                         sce=sce,
                         cluster_factor="BioClassification",
                         block_factor=NULL,
                         corr=TRUE,

```

```

    MI=FALSE,
    BPPARAM=BiocParallel::MulticoreParam())
#> calculating average expression across clusters...
#> computing correlation of the regulon...
#> if the standard deviation for TF expression is zero, the derived weight will be NA...
#> /
head(regulon.w)
#>      tf target      corr      weight
#> 3 FOXA1 A4GALT 0.7514745 -0.04072622
#> 7 FOXA1 AAGAB 0.5301808  0.01432174
#> 8 FOXA1 AAK1 0.6251292 -0.11678985
#> 11 FOXA1 AAMDC 0.5667211  0.06322510
#> 13 FOXA1 AAMP 0.5984585 -0.14155255
#> 16 FOXA1 AAR2 0.5086169  0.37346716

```

Calculate TF activity

Finally, the activities for a specific TF in each cell are computed by averaging expressions of target genes linked to the TF weighted by the correlation variable of user's choice.

$$y = \frac{1}{n} \sum_{i=1}^n x_i * corr_i$$

where y is the activity of a TF for a cell n is the total number of targets for a TF x_i is the log count expression of target i where i in $\{1,2,\dots,n\}$ $corr_i$ is the weight of TF and target i

```

score.combine <- calculateActivity(sce, regulon.w, "weight", method="weightedMean",
                                    assay = "logcounts")
#> calculating TF activity from regulon using weightedmean
#> /
head(score.combine[,1:10])
#> 4 x 10 Matrix of class "dgeMatrix"
#>      CD34_32_R5:AAACCTGAGTATCGAA-1 CD34_32_R5:AAACCTGAGTCGTTG-1
#> FOXA1          0.006467344          0.014336410
#> GATA3         -0.021012397         -0.027471145
#> SOX9          0.002726760         -0.001256817
#> SPI1          -0.003549379          0.010048047
#>      CD34_32_R5:AAACCTGGTTCACAA-1 CD34_32_R5:AAACGGGAGCTTCGCG-1
#> FOXA1          0.002326760          0.0040956059
#> GATA3         -0.016458312         -0.0165527279
#> SOX9          0.003799146         -0.0007582545
#> SPI1          -0.004201749         -0.0051176003
#>      CD34_32_R5:AAACGGGAGGGAGTAA-1 CD34_32_R5:AAACGGGAGTTACGGG-1
#> FOXA1          0.002576713          0.001192130
#> GATA3         -0.013892506         -0.011573274
#> SOX9          0.007332733          0.005315836
#> SPI1          -0.001195408         -0.005240490
#>      CD34_32_R5:AAACGGGCAAGTAATG-1 CD34_32_R5:AAACGGGCAAGTTCTG-1
#> FOXA1         -0.001796598          0.002849936
#> GATA3         -0.005744173         -0.012281146
#> SOX9          0.006684796          0.003803802
#> SPI1          -0.002644292         -0.003329163
#>      CD34_32_R5:AAACGGGCACAGACAG-1 CD34_32_R5:AAACGGGGTAACGTTC-1

```

#> FOXA1	0.003780386	0.0032500098
#> GATA3	-0.024982442	-0.0191374106
#> SOX9	0.001013378	0.0004819318
#> SPI1	-0.008007367	-0.0019797585

Differential TF activity test

We can next determine which TFs exhibit differential activities across cell clusters/groups via the findDifferentialActivity function. This function depends on findMarkers function from scran package and allow the same parameters.

```
da_list <- findDifferentialActivity(score.combine, sce$BioClassification, pval.type="some",
                                      direction="up", test.type= "t")
```

getSigGenes compiles the different test results into a single dataframe and enables user to supply their desired cutoffs for significance and variable to order by.

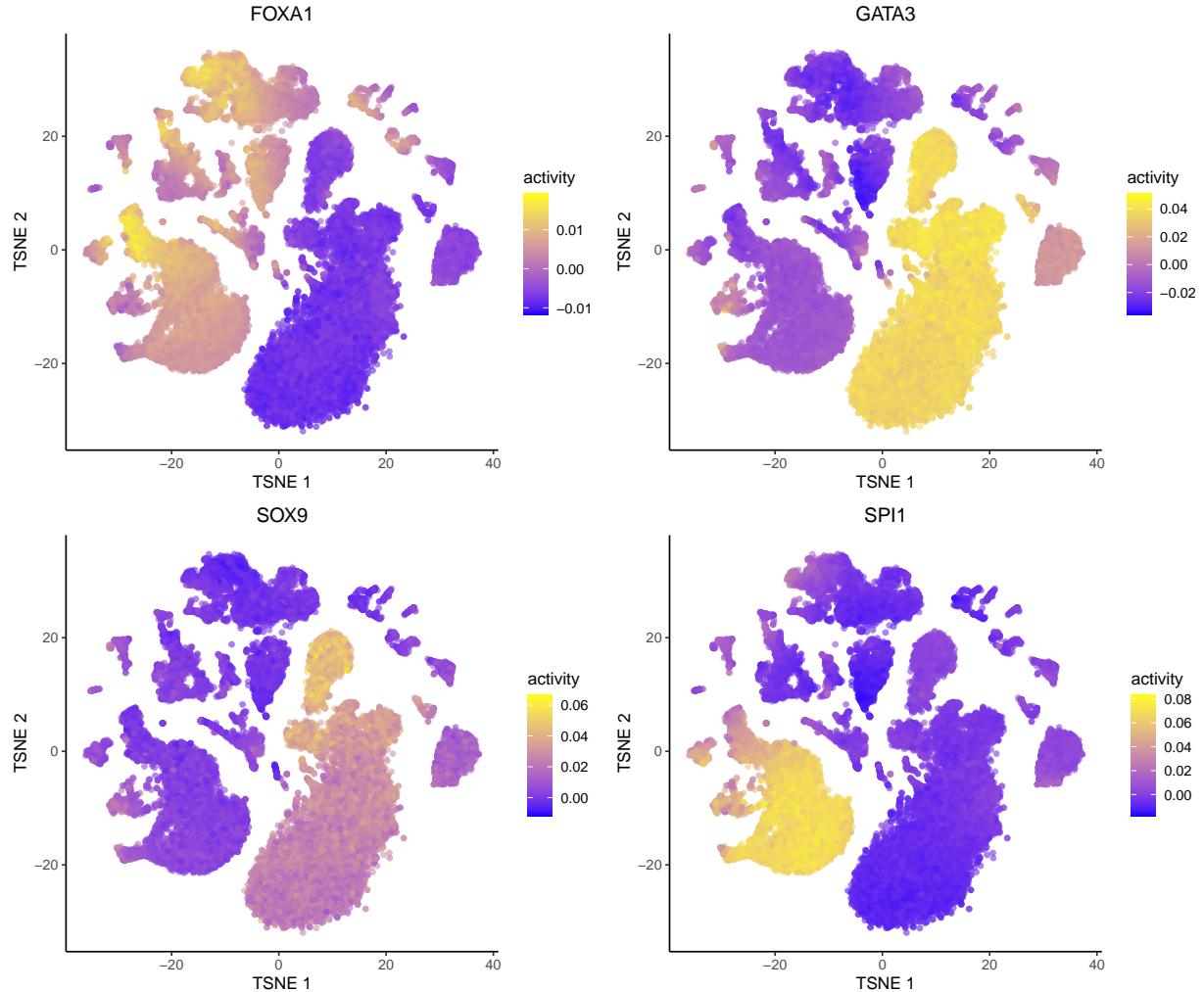
```
markers <- getSigGenes(da_list, fdr_cutoff = 0.05)
#> Using a logFC cutoff of 0 for class 1
#> Using a logFC cutoff of 0 for class 2
#> Using a logFC cutoff of 0 for class 3
#> Using a logFC cutoff of 0 for class 4
#> Using a logFC cutoff of 0 for class 5
#> Using a logFC cutoff of 0 for class 6
#> Using a logFC cutoff of 0 for class 7
#> Using a logFC cutoff of 0 for class 8
#> Using a logFC cutoff of 0 for class 9
#> Using a logFC cutoff of 0 for class 10
#> Using a logFC cutoff of 0.1 for class 11
#> Using a logFC cutoff of 0.1 for class 12
#> Using a logFC cutoff of 0.1 for class 13
#> Using a logFC cutoff of 0 for class 14
#> Using a logFC cutoff of 0 for class 15
#> Using a logFC cutoff of 0 for class 16
#> Using a logFC cutoff of 0 for class 17
#> Using a logFC cutoff of 0 for class 18
#> Using a logFC cutoff of 0 for class 19
#> Using a logFC cutoff of 0 for class 20
#> Using a logFC cutoff of 0 for class 21
#> Using a logFC cutoff of 0 for class 22
#> Using a logFC cutoff of 0 for class 23
#> Using a logFC cutoff of 0 for class 24
#> Using a logFC cutoff of 0 for class 25
#> Using a logFC cutoff of 0 for class 26
head(markers)
#>
#>      p.value          FDR summary.logFC      class      tf
#> 1  2.196874e-123 8.787498e-123  0.002982791 02_Early.Eryth FOXA1
#> 2  1.089715e-84  4.358859e-84   0.004170094 03_Late.Eryth FOXA1
#> 3  5.227473e-37  2.090989e-36   0.004031504 04_Early.Baso FOXA1
#> 4  6.075948e-213 2.430379e-212  0.007563786 05_CMP.LMPP FOXA1
#> 5  6.748446e-62  2.699379e-61   0.002395315 06_CLP.1 FOXA1
#> 21 0.000000e+00  0.000000e+00   0.022055115 07_GMP  SPI1
```

Visualizing TF activities

Epiregulon also provides multiple options for visualizing the inferred TF activities.

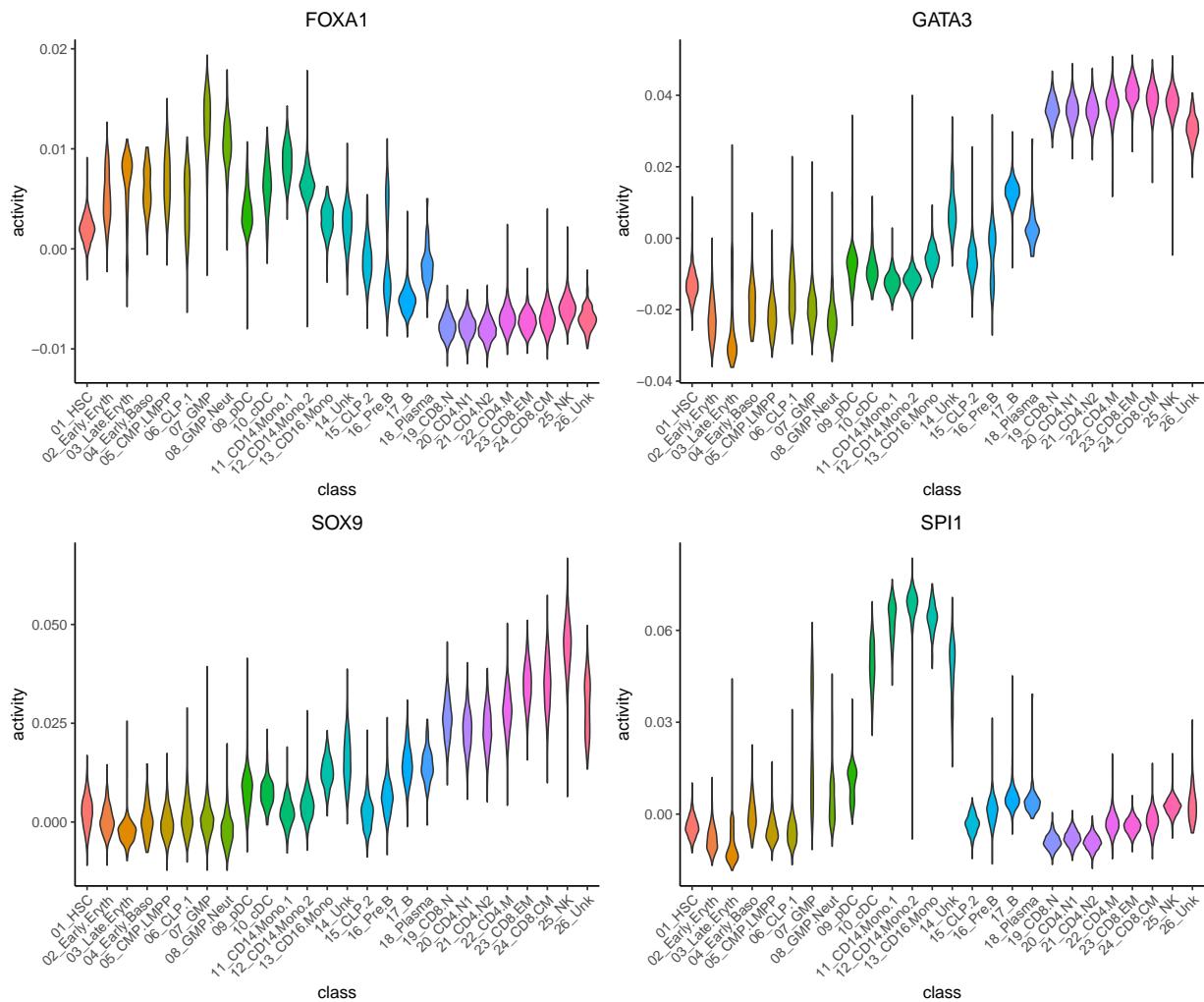
tSNE or UMAP plots:

```
plotActivityDim(sce, score.combine, c("FOXA1", "GATA3", "SOX9", "SPI1"), "TSNE", combine = T)
```



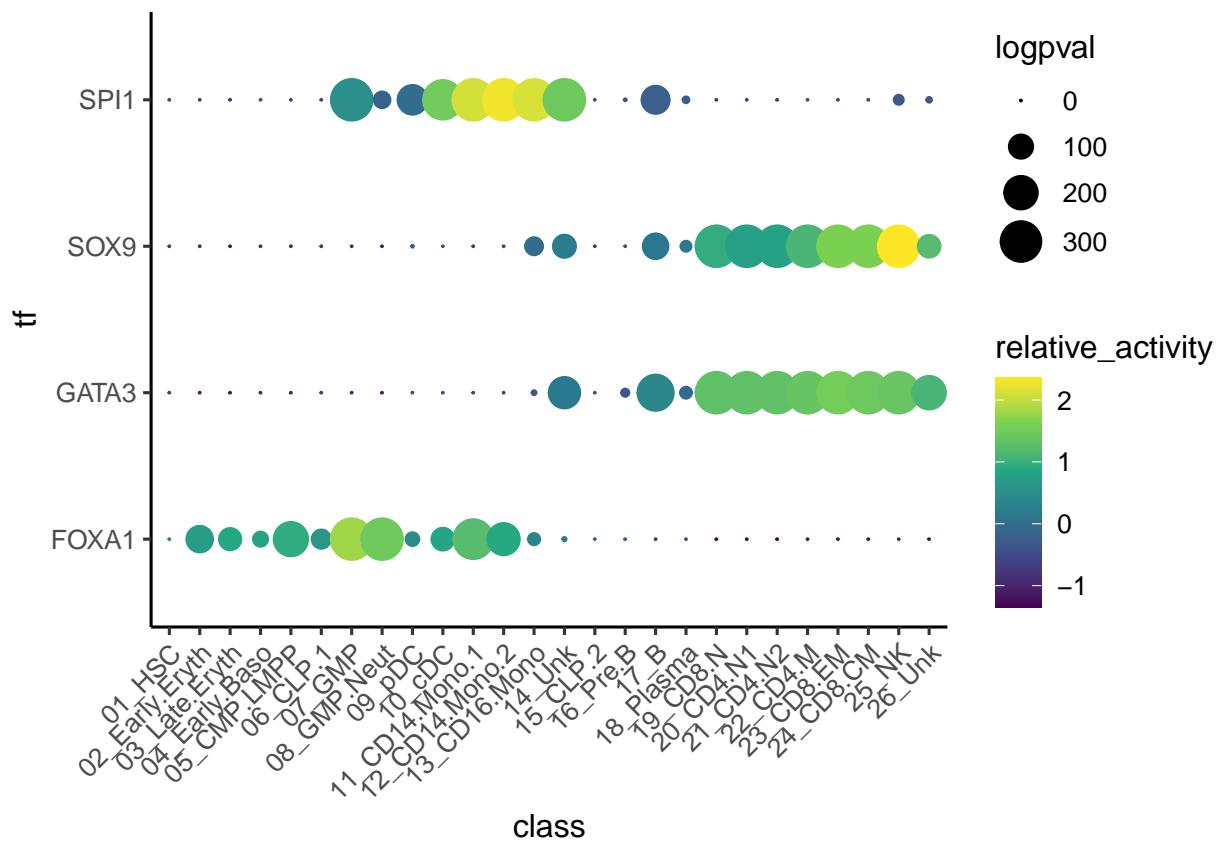
Violin plots:

```
plotActivityViolin(score.combine, c("FOXA1", "GATA3", "SOX9", "SPI1"), sce$BioClassification)
```



Bubble plot:

```
plotBubble(as.matrix(score.combine), c("FOXA1", "GATA3", "SOX9", "SPI1"),
sce$BioClassification, bubblesize = "FDR")
```



Session Info

```
sessionInfo()
#> R Under development (unstable) (2022-02-23 r81801)
#> Platform: x86_64-pc-linux-gnu (64-bit)
#> Running under: Ubuntu 18.04.5 LTS
#>
#> Matrix products: default
#> BLAS: /usr/local/lib/R/lib/libRblas.so
#> LAPACK: /usr/local/lib/R/lib/libRlapack.so
#>
#> locale:
#> [1] LC_CTYPE=en_US.UTF-8          LC_NUMERIC=C           LC_TIME=C
#> [4] LC_COLLATE=C                 LC_MONETARY=C        LC_MESSAGES=C
#> [7] LC_PAPER=C                  LC_NAME=C            LC_ADDRESS=C
#> [10] LC_TELEPHONE=C             LC_MEASUREMENT=C   LC_IDENTIFICATION=C
#>
#> attached base packages:
#> [1] parallel    stats       graphics   grDevices  utils      datasets
#> [8] methods     base
#>
#> other attached packages:
#> [1] nabor_0.5.0               epiregulon_1.0.2
```

```

#> [3] SingleCellExperiment_1.17.2 ArchR_1.0.1
#> [5] magrittr_2.0.3           rhdf5_2.39.6
#> [7] Matrix_1.4-0            data.table_1.14.2
#> [9] SummarizedExperiment_1.25.3 Biobase_2.55.2
#> [11] GenomicRanges_1.47.6   GenomeInfoDb_1.31.7
#> [13] IRanges_2.29.1         S4Vectors_0.33.17
#> [15] BiocGenerics_0.41.2   MatrixGenerics_1.7.0
#> [17] matrixStats_0.62.0    ggplot2_3.3.5
#>
#> loaded via a namespace (and not attached):
#> [1] ggbeeswarm_0.6.0        colorspace_2.0-3
#> [3] ellipsis_0.3.2          scuttle_1.5.2
#> [5] bluster_1.5.1          XVector_0.35.0
#> [7] BiocNeighbors_1.13.0   rstudioapi_0.13
#> [9] farver_2.1.0            ggrepel_0.9.1
#> [11] bit64_4.0.5            fansi_1.0.3
#> [13] codetools_0.2-18      sparseMatrixStats_1.7.0
#> [15] cachem_1.0.6          knitr_1.38
#> [17] scater_1.23.6          jsonlite_1.8.0
#> [19] cluster_2.1.2          compiler_4.2.0
#> [21] httr_1.4.2             dqrng_0.3.0
#> [23] backports_1.4.1       assertthat_0.2.1
#> [25] fastmap_1.1.0          limma_3.51.8
#> [27] cli_3.2.0              BiocSingular_1.11.0
#> [29] htmltools_0.5.2        tools_4.2.0
#> [31] rsud_1.0.5              igraph_1.3.1
#> [33] gp.auth_1.5.2          gtable_0.3.0
#> [35] glue_1.6.2              GenomeInfoDbData_1.2.8
#> [37] gp.cache_1.5.4          reshape2_1.4.4
#> [39] dplyr_1.0.8             Rcpp_1.0.8.3
#> [41] vctrs_0.3.8             rhdf5filters_1.7.0
#> [43] DelayedMatrixStats_1.17.0 xfun_0.30
#> [45] stringr_1.4.0           beachmat_2.11.0
#> [47] lifecycle_1.0.1          irlba_2.3.5
#> [49] gtools_3.9.2             statmod_1.4.36
#> [51] edgeR_3.37.1            getPass_0.2-2
#> [53] zlibbioc_1.41.0          scales_1.2.0
#> [55] yaml_2.3.5              memoise_2.0.1
#> [57] gridExtra_2.3            metacommons_1.7.2
#> [59] stringi_1.7.6            RSQLite_2.2.12
#> [61] highr_0.9                artificer.ranges_1.1.0
#> [63] ScaledMatrix_1.3.0       scran_1.23.1
#> [65] filelock_1.0.2           BiocParallel_1.29.21
#> [67] rlang_1.0.2              pkgconfig_2.0.3
#> [69] bitops_1.0-7              evaluate_0.15
#> [71] lattice_0.20-45          purrr_0.3.4
#> [73] Rhdf5lib_1.17.3          labeling_0.4.2
#> [75] patchwork_1.1.1          cowplot_1.1.1
#> [77] bit_4.0.4                tidyselect_1.1.2
#> [79] plyr_1.8.7               R6_2.5.1
#> [81] generics_0.1.2            base64url_1.4
#> [83] metapod_1.3.0            artificer.base_1.1.5
#> [85] DelayedArray_0.21.2      DBI_1.1.2

```

```
#> [87] pillar_1.7.0           withr_2.5.0
#> [89] RCurl_1.98-1.6        tibble_3.1.6
#> [91] crayon_1.5.0          utf8_1.2.2
#> [93] rmarkdown_2.13         viridis_0.6.2
#> [95] locfit_1.5-9.5        grid_4.2.0
#> [97] genomitory_1.99.3     blob_1.2.3
#> [99] gp.version_1.5.0      digest_0.6.29
#> [101] munsell_0.5.0        beeswarm_0.4.0
#> [103] viridisLite_0.4.0    ArtifactDB_1.7.4
#> [105] vipor_0.4.5
```