

SCRIP simulation for scRNA-seq data

Fei Qin

Last updated: 08/17/2021

1. Introduction to SCRIP method

SCRIP proposed two frameworks based on Gamma-Poisson and Beta-Gamma-Poisson distribution for simulating scRNA-seq data. Both Gamma-Poisson and Beta-Gamma-Poisson distribution model the over dispersion of scRNA-seq data. Specifically, Beta-Gamma-Poisson model was used to model bursting effect. The dispersion was accurately simulated by fitting the mean-BCV dependency using generalized additive model (GAM). Other key characteristics of scRNA-seq data including library size, zero inflation and outliers were also modeled by SCRIP. With its flexible modeling, SCRIP enables various application for different experimental designs and goals including DE analysis, clustering analysis, trajectory-based analysis and bursting analysis

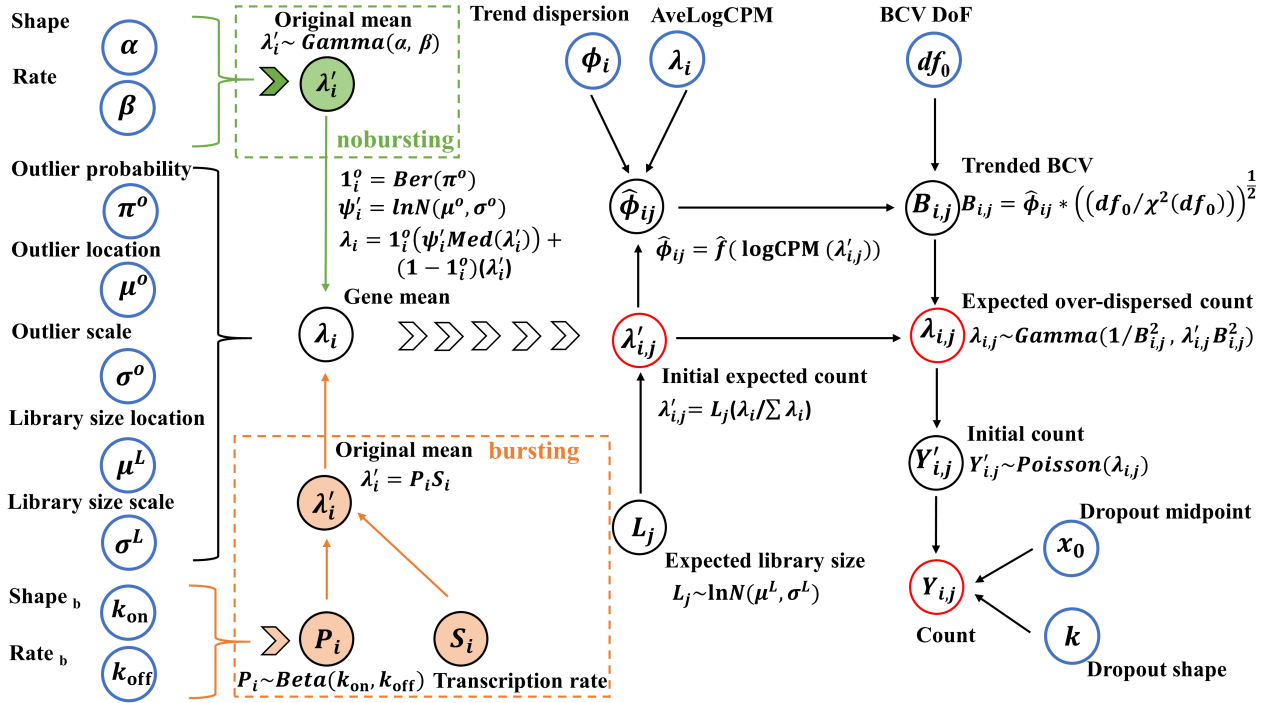


Figure 1: SCRIP framework

2. Installation

```
BiocManager::install("splatter")

library(devtools)
install_github("thecailab/SCRIP")
```

3. Quick start

Assuming you already have a count matrix for scRNA-seq data, and you want to simulation data based on it. Only a few steps are needed to creat a simulation data using SCRIP.

A dataset from Xin data is used for example.

```
library(splatter)
library(Biobase)
library(SCRIP)

EMTAB.eset <- EMTABesethealthy
expre_data = exprs(EMTAB.eset)
pheno_data = pData(EMTAB.eset)

dim(expre_data)
```

```
## [1] 25453 1097
```

```
expre_data[1:6,1:6]
```

```
##           AZ_A10 AZ_A11 AZ_A12 AZ_A2 AZ_A5 AZ_A6
## SGIP1          0      0      0    32      0      0
## AZIN2          0      0      0      0      0      0
## CLIC4          3      0      0      1      0      0
## AGBL4          0      0      0      0      0      0
## NECAP2         0      0      0      0      0      0
## SLC45A1        0      0      0      0      0      0
```

```
head(pheno_data)
```

```
##           sampleID SubjectName cellTypeID cellType
## AZ_A10           1   Non T2D 1           5   delta
## AZ_A11           1   Non T2D 1           2   alpha
## AZ_A12           1   Non T2D 1           5   delta
## AZ_A2            1   Non T2D 1           9   gamma
## AZ_A5            1   Non T2D 1           6   ductal
## AZ_A6            1   Non T2D 1           2   alpha
```

```
# To simplify computation here, only acinar celltype in sample 5 was utilized here.
sub_expre_data=expre_data[1:10000,which((pheno_data$sampleID %in% c(5)) &
                                         (pheno_data[,4] == "acinar"))]
params <- splatEstimate(sub_expre_data)
```

```

## $start.arg
## $start.arg$shape
## [1] 0.6954772
##
## $start.arg$rate
## [1] 0.08910949
##
##
## $fix.arg
## NULL
##
## $start.arg
## $start.arg$meanlog
## [1] 12.03125
##
## $start.arg$sdlog
## [1] 1.056239
##
##
## $fix.arg
## NULL
##
## $start.arg
## $start.arg$meanlog
## [1] 5.256435
##
## $start.arg$sdlog
## [1] 1.276158
##
##
## $fix.arg
## NULL

```

```

sim_trend <- SCRIPsimu(data=sub_expre_data, params=params, mode="GP-trendedBCV")
sim_trend

```

```

## class: SingleCellExperiment
## dim: 10000 80
## metadata(13): Params method ... batch.facScale bcv.shrink
## assays(5): BatchCellMeans BaseCellMeans CellMeans TrueCounts counts
## rownames(10000): Gene1 Gene2 ... Gene9999 Gene10000
## rowData names(4): Gene BaseGeneMean OutlierFactor GeneMean
## colnames(80): Cell1 Cell2 ... Cell179 Cell180
## colData names(3): Cell Batch ExpLibSize
## reducedDimNames(0):
## mainExpName: NULL
## altExpNames(0):

```

4 Single cell type simulation

4.1 Parameter estimation

SCRIP utilized the estimation strategy from splatter, but also provided more parameters (Fold change, dropout rates, library size, BCV degree of freedom) to serve different experimental designs (i.e. Simulation for differential expression analysis, clustering analysis and trajectory analysis). Detailed description about other parameters will be shown in other sections of this document.

4.2 Simulation

The default mode in SCRIP for simulation is “GP-trendedBCV”. You can also choose other modes (“GP-commonBCV”, “BGP-commonBCV”, “BP”, “BGP-trendedBCV”) in the SCRIPsimu() function. For single cell type simulation, you have to set the “method” as “single”, which was default in SCRIPsimu() function.

4.2.1 GP-commonBCV

GP-commonBCV is the model used by splatter. GP-commonBCV applied the Gamma-Poisson mixture model with mean-BCV dependency fitted by a common BCV across genes.

```
##### GP-commonBCV model/Splatter #####  
#####  
sim_GPcommon <- SCRIPsimu(data=sub_expre_data, params=params, mode="GP-commonBCV")  
sim_GPcommon
```

```
## class: SingleCellExperiment  
## dim: 10000 80  
## metadata(13): Params method ... batch.facScale bcv.shrink  
## assays(5): BatchCellMeans BaseCellMeans CellMeans TrueCounts counts  
## rownames(10000): Gene1 Gene2 ... Gene9999 Gene10000  
## rowData names(4): Gene BaseGeneMean OutlierFactor GeneMean  
## colnames(80): Cell1 Cell2 ... Cell179 Cell180  
## colData names(3): Cell Batch ExpLibSize  
## reducedDimNames(0):  
## mainExpName: NULL  
## altExpNames(0):
```

4.2.2 GP-trendedBCV

GP-trendedBCV is the major model of SCIRP. which used the Gamma-Poisson mixture model with mean-BCV dependency fitted by GAM.

```
##### GP-trendedBCV model #####  
#####  
sim_GPtrend <- SCRIPsimu(data=sub_expre_data, params=params, mode="GP-trendedBCV")
```

4.2.3 BP

BP is the model used for simulating bursting effect using Beta-Poisson mixture distribution without considering BCV effect.

```
##### BP-commonBCV model #####
#####
sim_BP <- SCRIPsimu(data=sub_expre_data, params=params, mode="BP")
```

```
## $start.arg
## $start.arg$shape1
## [1] 0.16309
##
## $start.arg$shape2
## [1] 1.954278
##
##
## $fix.arg
## NULL
```

4.2.4 BGP-commonBCV

BP-commonBCV is the model used for simulating bursting effect with Beta-Gamma-Poisson mixture distribution. The mean-BCV dependency was fitted by a common BCV across genes.

```
##### BP-commonBCV model #####
#####
sim_BGPcommon <- SCRIPsimu(data=sub_expre_data, params=params, mode="BGP-commonBCV")
```

```
## $start.arg
## $start.arg$shape1
## [1] 0.16309
##
## $start.arg$shape2
## [1] 1.954278
##
##
## $fix.arg
## NULL
```

4.2.5 BGP-trendedBCV

BP-trendedBCV is the model used for simulating bursting effect with Beta-Gamma-Poisson mixture distribution. The mean-BCV dependency was fitted by a GAM.

```
##### BP-trendedBCV model #####
#####
sim_BGPtrend <- SCRIPsimu(data=sub_expre_data, params=params, mode="BGP-trendedBCV")
```

```
## $start.arg
## $start.arg$shape1
## [1] 0.16309
##
## $start.arg$shape2
## [1] 1.954278
##
##
```

```
##  
## $fix.arg  
## NULL
```

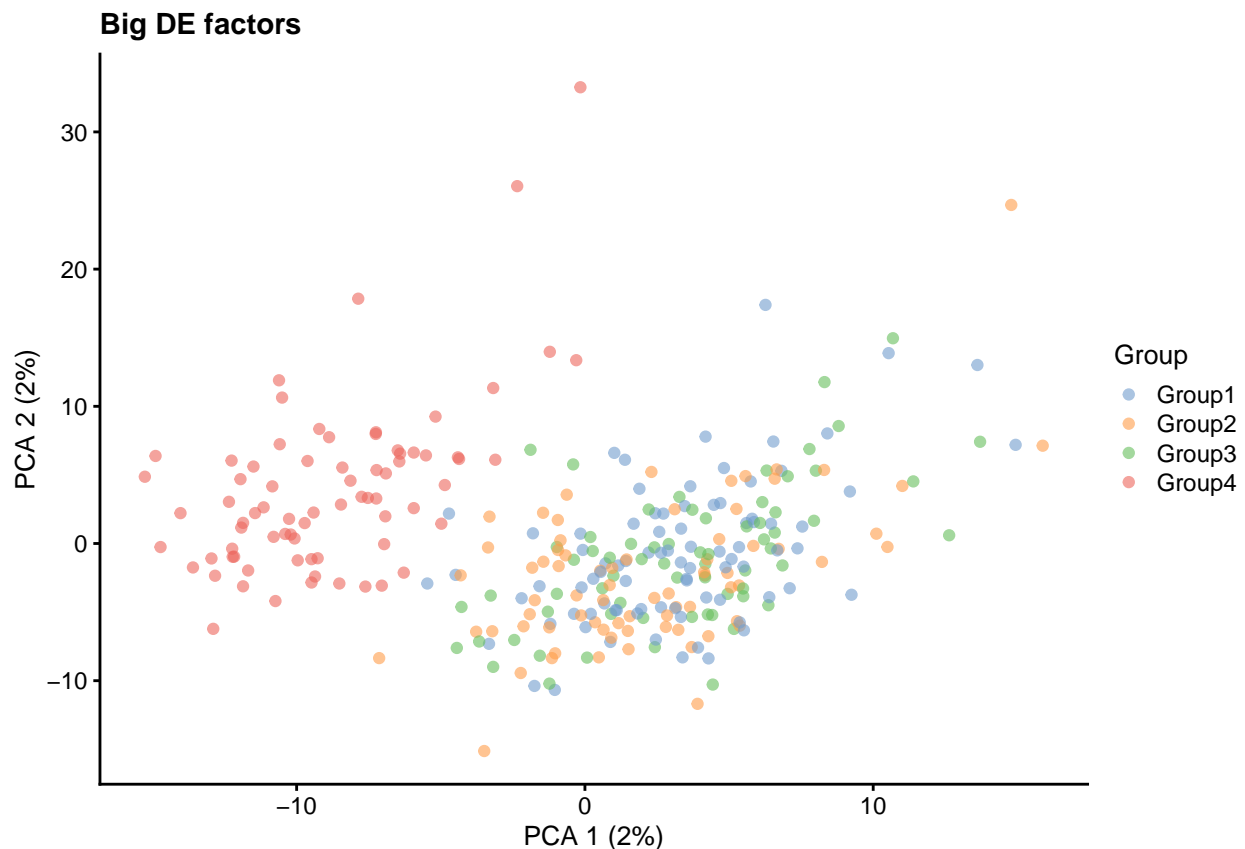
5 Group simulation

Group simulation is useful for studying different experimental conditions, especially for differential expression (DE) analysis. To serve different applications in scRNA-seq analysis, SCRIP provides flexible simulation. It can simulate scRNA-seq data with different parameters from multiple cell groups (i.e. cell types), which is useful for evaluating the detection of global characteristics such as clustering. It also allows simulation of group difference in a single cell group, which is useful for evaluating typical DE analysis methods.

5.1 Basic group simulation

DEGs were simulated using multiplicative differential expression factors from a log-normal distribution with parameters including number of genes (nGenes), the path-specific proportion of DE genes (de.prob), the proportion of down-regulated DE genes (de.downProb), DE location factor (de.facLoc) and DE scale factor (de.facScale).

```
sim.SCRIP2 <- SCRIPsimu(data=sub_expre_data, params=params, method="groups",  
  batchCells=300, group.prob = c(0.25, 0.25, 0.25, 0.25),  
  de.prob = c(0.2, 0.2, 0.2, 0.2),  
  de.downProb = c(0.5, 0.5, 0.5, 0.5),  
  de.facLoc = c(0.2, 0.3, 0.4, 0.5),  
  de.facScale=c(0.2, 0.2, 0.2, 0.2))  
  
library(scater)  
sim2 <- logNormCounts(sim.SCRIP2)  
sim2 <- runPCA(sim2)  
plotPCA(sim2, colour_by = "Group") + ggtitle("Big DE factors")
```



5.2 Group simulation with batch effect

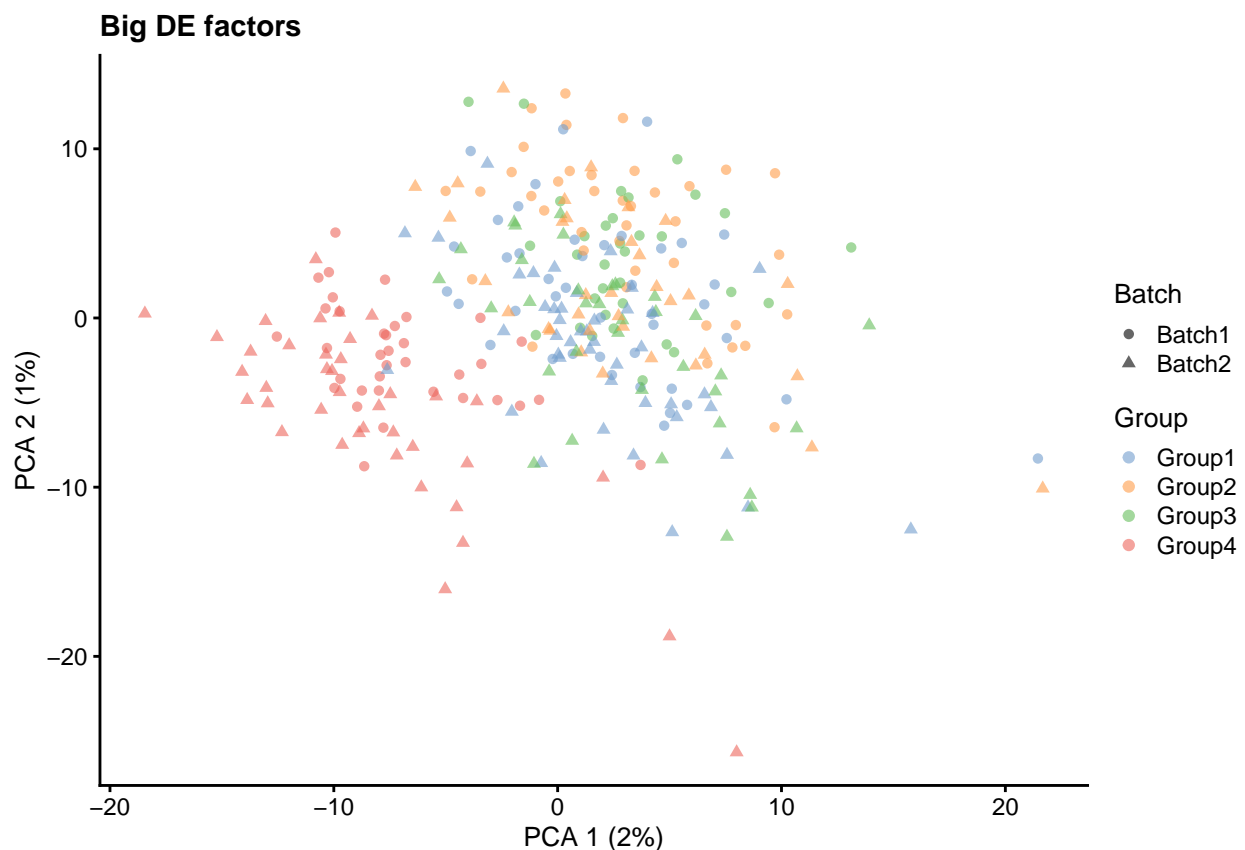
Batch effect factors are also generated from a log-normal distribution with parameters including batchCells, batch.facLoc and batch.facScale.

batchCells: number of cells for each batch

batch.facLoc: Batch location factor in log-normal distribution for batch factor

batch.facScale: Batch scale factor in log-normal distribution for batch factor

```
sim.SCRIP3 <- SCRIPsimu(data=sub_expre_data, params=params, method="groups",
  batchCells=c(150, 150),
  batch.facLoc = c(0.1, 0.1),
  batch.facScale = c(0.1, 0.1),
  group.prob = c(0.25, 0.25, 0.25, 0.25),
  de.prob = c(0.2, 0.2, 0.2, 0.2),
  de.downProb = c(0.5, 0.5, 0.5, 0.5),
  de.facLoc = c(0.2, 0.3, 0.4, 0.5),
  de.facScale=c(0.2, 0.2, 0.2, 0.2))
sim3 <- logNormCounts(sim.SCRIP3)
sim3 <- runPCA(sim3)
plotPCA(sim3, colour_by = "Group", shape_by="Batch") + ggtitle("Big DE factors")
```



5.3 Group simulation for clustering

To simulate multi-cell type scRNA-seq data, we estimated hyperparameters from real data and simulate data for each cell type separately. SCRIP preserved cluster discriminative genes, which is not capable in

available methods. We identified variably expressed genes (VEGs) from each cell type in real data. For these VEGs, true cell mean was set to their expression means in real data to simulate cell-type discriminative genes. Whereafter, aforementioned simulation steps (simulating true cell mean; simulating true counts) will be followed to complete the simulation.

```
### Xin ###
EMTAB.eset <- EMTABesethealthy
library(Biobase)
library(edgeR)
library(stats)
expre_data = exprs(EMTAB.eset)
pheno_data=pData(EMTAB.eset)
expre_data <- as.matrix(expre_data[1:2000,])

expre_data[1:6,1:6]
```

```
##          AZ_A10 AZ_A11 AZ_A12 AZ_A2 AZ_A5 AZ_A6
## SGIP1         0      0      0     32      0      0
## AZIN2         0      0      0      0      0      0
## CLIC4         3      0      0      1      0      0
## AGBL4         0      0      0      0      0      0
## NECAP2        0      0      0      0      0      0
## SLC45A1       0      0      0      0      0      0
```

```
head(pheno_data)
```

```
##      sampleID SubjectName cellTypeID cellType
## AZ_A10         1   Non T2D 1           5   delta
## AZ_A11         1   Non T2D 1           2   alpha
## AZ_A12         1   Non T2D 1           5   delta
## AZ_A2          1   Non T2D 1           9   gamma
## AZ_A5          1   Non T2D 1           6   ductal
## AZ_A6          1   Non T2D 1           2   alpha
```

```
table(pheno_data$cellType)
```

```
##
##          acinar          alpha          beta
##          112          443          171
##      co-expression          delta          ductal
##          26          59          135
##      endothelial          epsilon          gamma
##          13          5          75
##          mast          MHC class II          PSC
##          4          1          23
##      unclassified unclassified endocrine
##          1          29
```

```
CTlist <- c("acinar","alpha","beta")
```

```
params <- splatEstimate(expre_data)
```

```
## $start.arg
## $start.arg$shape
## [1] 0.673965
##
## $start.arg$rate
## [1] 0.0813325
##
##
## $fix.arg
## NULL
##
## $start.arg
## $start.arg$meanlog
## [1] 9.733591
##
## $start.arg$sdlog
## [1] 1.170375
##
##
## $fix.arg
## NULL
```

```
library(Seurat)
final.list <- simu_cluster(expre_data = expre_data, pheno_data = pheno_data, CT=CTlist,
                           mode="GP-trendedBCV", nfeatures=100)
```

```
## $start.arg
## $start.arg$shape
## [1] 0.673965
##
## $start.arg$rate
## [1] 0.0813325
##
##
## $fix.arg
## NULL
##
## $start.arg
## $start.arg$meanlog
## [1] 9.733591
##
## $start.arg$sdlog
## [1] 1.170375
##
##
## $fix.arg
## NULL
```

```
final <- final.list$final
CT.infor <- final.list$CT.infor

library(ggplot2)
Group.sim <- function(data){
```

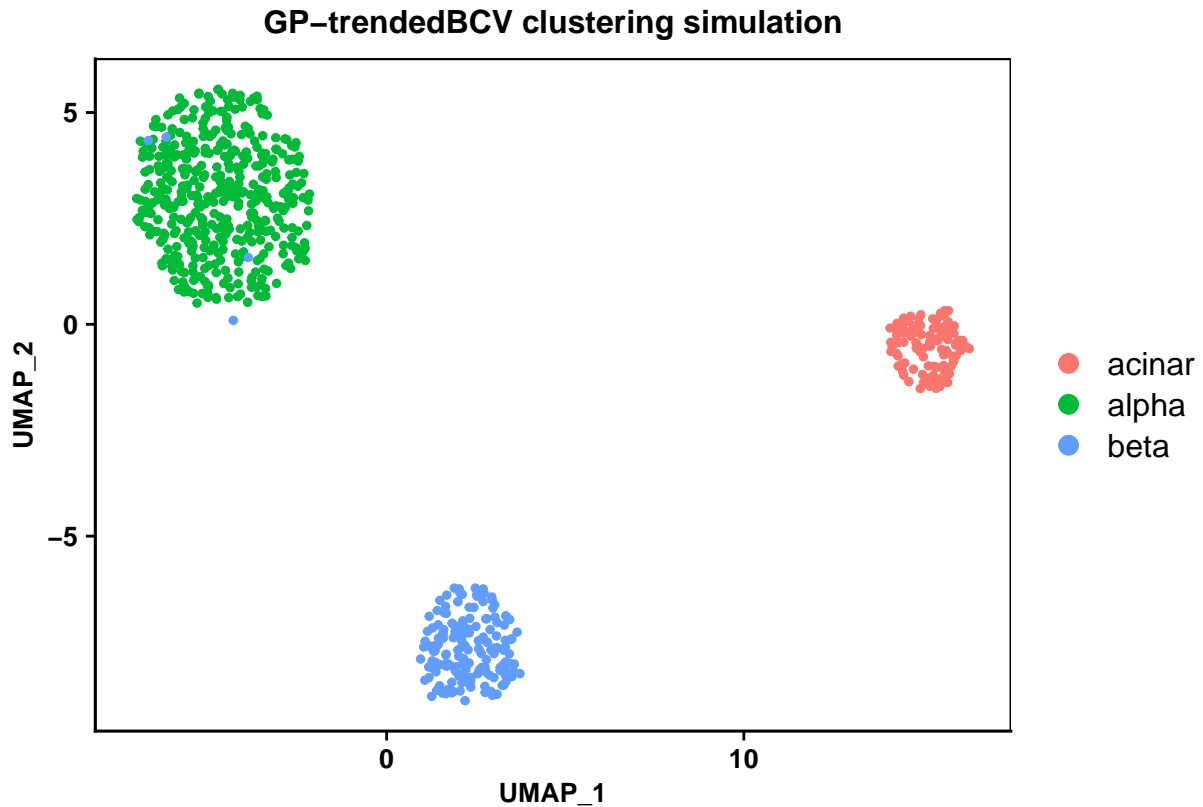
```

library(dplyr)
library(Seurat)
library(cowplot)
pbmc <- CreateSeuratObject(counts = data, project = "MuSiC", min.cells = 3,
                           min.features = 200)
pbmc <- NormalizeData(pbmc, normalization.method = "LogNormalize", scale.factor = 10000)
pbmc <- FindVariableFeatures(pbmc, selection.method = "vst", nfeatures = 100)
all.genes <- rownames(pbmc)
pbmc <- ScaleData(pbmc, features = all.genes)
pbmc <- RunPCA(pbmc, features = VariableFeatures(object = pbmc))
pbmc <- RunUMAP(pbmc, dims = 1:10)
pbmc <- RunTSNE(pbmc)
return(pbmc)
}

pbmc <- Group.sim(data=final)
pbmc@meta.data$Celltype <- CT.infor
umap <- DimPlot(pbmc, reduction = "umap", group.by = "Celltype", pt.size=1)+
  labs(title="GP-trendedBCV clustering simulation")+
  theme(plot.title = element_text(face="bold",size=12, hjust=0.5),
        legend.text=element_text(size=12),
        axis.text=element_text(size=10,face="bold"),
        axis.title.x = element_text(size=10,face="bold"),
        axis.title.y = element_text(size=10,face="bold"),
        panel.background = element_rect(fill='white',colour="black"))

umap

```



5.4 Group simulation for DE analysis

To simulate single-cell type scRNA-seq data, SCRIP enabled to simulate with different DE rates, fold change, dropout rate, library size, BCV.df and BCV.shrink. You change these parameters in `simu_DE()` function, which will finally export the simulation results for group1 and group2, as well as the true DEGs.

```
## Here I will use Tung dataset (given in package) to give examples about how to change
## parameter for evaluating DE analysis methods.
```

```
counts <- as.matrix(Tung)
```

```
library(Biobase)
```

```
## To simplify computation, only 1000 genes and 100 cells from Tung dataset will be use
## for estimation and simulation
```

```
expre_data_sub <- counts[1:1000,1:100]
```

```
params1 <- splatEstimate(expre_data_sub)
```

```
## $start.arg
```

```
## $start.arg$shape
```

```
## [1] 0.6022893
```

```
##
```

```
## $start.arg$rate
```

```
## [1] 0.2100649
```

```
##
```

```
##
```

```
## $fix.arg
```

```
## NULL
##
## $start.arg
## $start.arg$meanlog
## [1] 8.267739
##
## $start.arg$sdlog
## [1] 0.3449362
##
##
## $fix.arg
## NULL
##
## $start.arg
## $start.arg$meanlog
## [1] 4.342277
##
## $start.arg$sdlog
## [1] 0.3501836
##
##
## $fix.arg
## NULL
```

- 1) Fold change: FC enables us to control the difference between two groups. Larger fold change value means larger pre-defined difference between two groups; it is expected that we can find better performance for each DE analysis tool.

```
##### Change the parameter of fold change (FC) #####
#####
resFC <- simu_DE(expre_data = sub_expre_data, params = params1, nDE=50, FC=1.4)
## simu_DE() function will export the simulation results for group1 and group2,
## as well as the true DEGs.
resFC
```

```
## [[1]]
## class: SingleCellExperiment
## dim: 1000 100
## metadata(14): Params method ... batch.facScale bcv.shrink
## assays(5): BatchCellMeans BaseCellMeans CellMeans TrueCounts counts
## rownames(1000): Gene1 Gene2 ... Gene999 Gene1000
## rowData names(4): Gene BaseGeneMean OutlierFactor GeneMean
## colnames(100): Cell1 Cell2 ... Cell99 Cell100
## colData names(3): Cell Batch ExpLibSize
## reducedDimNames(0):
## mainExpName: NULL
## altExpNames(0):
##
## [[2]]
## class: SingleCellExperiment
## dim: 1000 100
## metadata(14): Params method ... batch.facScale bcv.shrink
## assays(5): BatchCellMeans BaseCellMeans CellMeans TrueCounts counts
```

```
## rownames(1000): Gene1 Gene2 ... Gene999 Gene1000
## rowData names(4): Gene BaseGeneMean OutlierFactor GeneMean
## colnames(100): Cell1 Cell2 ... Cell199 Cell100
## colData names(3): Cell Batch ExpLibSize
## reducedDimNames(0):
## mainExpName: NULL
## altExpNames(0):
##
## [[3]]
## [1] "gene341" "gene102" "gene316" "gene314" "gene206" "gene630" "gene474"
## [8] "gene908" "gene337" "gene931" "gene202" "gene813" "gene441" "gene446"
## [15] "gene653" "gene947" "gene470" "gene533" "gene900" "gene522" "gene338"
## [22] "gene780" "gene806" "gene485" "gene4" "gene526" "gene428" "gene348"
## [29] "gene64" "gene128" "gene223" "gene773" "gene907" "gene677" "gene687"
## [36] "gene963" "gene895" "gene216" "gene878" "gene471" "gene319" "gene21"
## [43] "gene449" "gene705" "gene569" "gene131" "gene798" "gene537" "gene703"
## [50] "gene204"
```

- 2) Dropout rate: dropout rate controls the zero proportion in scRNA-seq data, since zero inflation is one of the main characteristics in scRNA-seq data.

```
##### Change the parameter of dropout rates (Dropout_rate) #####
#####
resDrop <- simu_DE(expre_data = expre_data_sub, params = params1, nDE=50, FC=1.4,
                  Dropout_rate=0.2)
```

- 3) Library size: library size can be used to control the expression level for each gene, as well as the zero proportions in the scRNA-seq data.

```
##### Change the parameter of library size (libsize) #####
#####
resLib <- simu_DE(expre_data = expre_data_sub, params = params1, nDE=50, FC=1.4,
                  libsize=5000)
```

- 4) BCV.df: BCV.df enables us to change the variation of BCV values. Smaller variation is given to BCV values when BCV.df gets larger.

```
##### Change the parameter of BCV degree of freedom (pre.bcv.df) #####
#####
resDf <- simu_DE(expre_data = expre_data_sub, params = params1, nDE=50, FC=1.4,
                  pre.bcv.df=5)
```

- 5) BCV.shrink: This parameter is introduced to amplify (> 1) or shrink (< 1) BCV values, which mean it can control the BCV levels.

```
##### Change the parameter of BCV shrinkage factor (bcv.shrink) #####
#####
resShrink <- simu_DE(expre_data = expre_data_sub, params = params1, nDE=50, FC=1.4,
                     bcv.shrink = 1.5)
```

Here I take fold change as an example to introduce how to use the simulation data to evaluate DE analysis methods (i.e. edgeR, DESeq2, limma voom and MAST):

```

sim=resFC[[1]]
simDE=resFC[[2]]
genenameDE=resFC[[3]]

exps <- counts(sim)
expsDE <- counts(simDE)
counts <- cbind(exps,expsDE)

colnames(counts) <- paste0("cell",1:ncol(counts))
rownames(counts) <- paste0("gene",1:nrow(counts))

truth <- rownames(counts)
truth[which(rownames(counts) %in% genenameDE)] <- "DE"
truth[-which(rownames(counts) %in% genenameDE)] <- "common"
truth <- as.factor(truth)

```

As an application of SCRIP, we simulated scRNA-seq data using SCRIP and compared four DE analysis tools including edgeR, DESeq2, Limma-voom and MAST. To better account for zero inflation of scRNA-seq data in DE analysis, a weighting strategy of ZINB-WaVE was also integrated with edgeR, DESeq2 and Limma-voom.

ZINB WAVE: ZINB-WaVE is used to compute gene- and cell-specific weights, based on zero-inflated negative binomial model, to “unlock” bulk RNA-seq tools for single-cell applications.

```

#####
##### ZINB WAVE #####
library(zinbwave)
library(DESeq2)
group <- factor(c(rep(1,ncol(exps)),rep(2,ncol(expsDE))))
coldata <- data.frame(condition=group)
rownames(coldata) <- colnames(counts)
fluidigm <- DESeqDataSetFromMatrix(countData = counts,
                                   colData = coldata,
                                   design = ~ condition)
zinb <- zinbFit(fluidigm, K=2, epsilon=1000)
fluidigm_zinb <- zinbwave(fluidigm, fitted_model = zinb, K = 2, epsilon=1000,
                          observationalWeights = TRUE)
ZINB.Wave.weight <- assay(fluidigm_zinb, "weights")

```

5.4.1 edgeR

edgeR was originally designed for bulk RNA sequencing data and it used an overdispersed Poisson model and an Empirical Bayes model to account for both biological and technical variability.

```

#####
##### edgeR #####
library(edgeR)
dgList <- DGEList(assay(fluidigm))
countsPerMillion <- cpm(dgList)
countCheck <- countsPerMillion > 1
keep <- which(rowSums(countCheck) >= 2)
dgList <- dgList[keep,]
dgList <- calcNormFactors(dgList, method="TMM")

```

```

design <- model.matrix(~condition, data = colData(fluidigm))
dgList <- estimateDisp(dgList, design)
fit <- glmFit(dgList, design)
lrt <- glmLRT(fit, coef=2)
edgeR.res <- as.data.frame(topTags(lrt,n=nrow(lrt$coefficients)))
edgeR.res <- edgeR.res[rownames(counts),]

```

5.4.2 edgeR ZINB Wave

To apply edgeR for scRNA-seq data, we computed gene- and cell-specific weights using ZINB Wave approach. Then such weights were considered in the DE analysis.

```

#####
##### edgeR ZINB WAVE #####

library(edgeR)
dge <- DGEList(assay(fluidigm_zinb))
dge <- calcNormFactors(dge)

design <- model.matrix(~condition, data = colData(fluidigm))
dge$weights <- weights
dge <- estimateDisp(dge, design)
fit <- glmFit(dge, design)
lrt <- glmWeightedF(fit, coef = 2)
edgeR.weighted.res <- as.data.frame(topTags(lrt,n=nrow(lrt$coefficients)))
edgeR.weighted.res <- edgeR.weighted.res[rownames(counts),]

```

5.4.3 DESeq2

DESeq2, originally designed for bulk RNA-seq data, used a shrinkage estimation approach for dispersions and fold change to improve the stability and interpretability of estimates in DE analysis

```

#####
##### DESeq2 #####

library(DESeq2)
dds <- DESeqDataSet(fluidigm, design = ~ condition)
dds <- DESeq(dds, sfType="poscounts", useT=TRUE, minmu=1e-6)

DESeq2.res <- results(dds)
DESeq2.res <- as.data.frame(DESeq2.res)
DESeq2.res <- DESeq2.res[rownames(counts),]

```

5.4.4 DESeq2 ZINB Wave

```

#####
##### DESeq2 ZINB WAVE #####

library(DESeq2)
dds <- DESeqDataSet(fluidigm_zinb, design = ~ condition)
dds <- DESeq(dds, sfType="poscounts", useT=TRUE, minmu=1e-6)

```



```
DESeq2.weighted.res <- results(dds)
DESeq2.weighted.res <- as.data.frame(DESeq2.weighted.res)
DESeq2.weighted.res <- DESeq2.weighted.res[rownames(counts),]
```

5.4.5 limma voom

Voom model estimated the mean-variance relationship for the log-counts, which then generated a precision weight for each observation. Then these weights were considered into the limma empirical Bayes analysis pipeline.

```
#####
##### limma voom #####
library(edgeR)
d0 <- DGEList(counts)
d0 <- calcNormFactors(d0)

mm <- model.matrix(~0 + group)
y <- voom(d0, mm, plot = F)

fit <- lmFit(y, mm)

contr <- makeContrasts(group1 - group2, levels = colnames(coef(fit)))
tmp <- contrasts.fit(fit, contr)
tmp <- eBayes(tmp)
top.table <- topTable(tmp, sort.by = "none", n = Inf)
limma.res <- top.table
limma.res <- limma.res[rownames(counts),]
```

5.4.6 limma voom ZINB Wave

```
#####
##### limma voom ZINB WAVE #####
library(edgeR)
d0 <- DGEList(counts)
d0 <- calcNormFactors(d0)

mm <- model.matrix(~0 + group)
y <- voom(d0, mm, plot = F)
y$weights <- weights*y$weights

fit <- lmFit(y, mm)

contr <- makeContrasts(group1 - group2, levels = colnames(coef(fit)))
tmp <- contrasts.fit(fit, contr)
tmp <- eBayes(tmp)
top.table <- topTable(tmp, sort.by = "none", n = Inf)
limma.weighted.res <- top.table
limma.weighted.res <- limma.weighted.res[rownames(counts),]
```

5.4.7 MAST

MAST was designed specifically for scRNA-seq data. It developed a two-part, generalized linear model to parameterizes the characteristics in scRNA-seq data.

```
#####
##### MAST #####

group <- factor(c(rep("Group0",ncol(exps)),rep("Group1",ncol(expsDE))))
group <- relevel(group,"Group0")
coldata <- data.frame(condition=group)
rownames(coldata) <- colnames(counts)
fluidigm <- DESeqDataSetFromMatrix(countData = counts,
                                   colData = coldata,
                                   design = ~ condition)

library(MAST)
library(data.table)
sca <- FromMatrix(counts,data.frame(condition=group),
                  data.frame(Geneid=rownames(counts)),check_sanity = FALSE)
zlmCond <- zlm(~condition, sca)

summaryCond <- summary(zlmCond, doLRT='conditionGroup1')
print(summaryCond, n=4)
summaryDt <- summaryCond$datatable
fcHurdle <- merge(summaryDt[contrast=='conditionGroup1' & component=='H',.(primerid, `Pr(>Chisq)`)],
                  #hurdle P values
                  summaryDt[contrast=='conditionGroup1' & component=='logFC',
                             .(primerid, coef, ci.hi, ci.lo)],
                  by='primerid') #logFC coefficients
MAST.res=fcHurdle
MAST.res <- MAST.res[rownames(counts),]
```

For each method, true positive rates and false positive rates of DE detections can be calculated based on the knowledge of true and false DEGs. Receiver operating characteristic (ROC) curve and the areas under the ROC curve (AUC) can be used as measures to assess the performance of these DE analysis methods.

6 Cell lineage simulation

RNA transcriptomes data can be used to trace cell divisions and migrations over time, which is also referred to as cell lineage analysis. The developmental order of cells can be represented by a quantitative measure, pseudotime. To simulate pseudotime effect, we used the strategy of Splat based on differential expression process in cell developmental paths. Similarly to group simulation, in such differential expression process, we first defined the start and end expression level for each path. DEGs were simulated using multiplicative differential expression factors from a log-normal distribution with parameters including number of genes (nGenes), the path-specific proportion of DE genes (de.prob), the proportion of down-regulated DE genes (de.downProb), DE location factor (de.facLoc) and DE scale factor (de.facScale). Specifically, a few other parameters are provided to control the pattern of the lineages.

- 1) path.nSteps
 - number of steps between the start point and end point for each path
- 2) path.skew
 - Controls how likely cells are from the start or end point of the path. Cells are more likely to come from end point when path.skew gets close to 0, while cells are more likely to come from start point when path.skew gets close to 1.
- 3) path.from
 - Control the relationship between these paths.

```
counts <- as.matrix(Tung)

counts <- counts[1:10000,1:200]
params <- splatEstimate(counts)
```

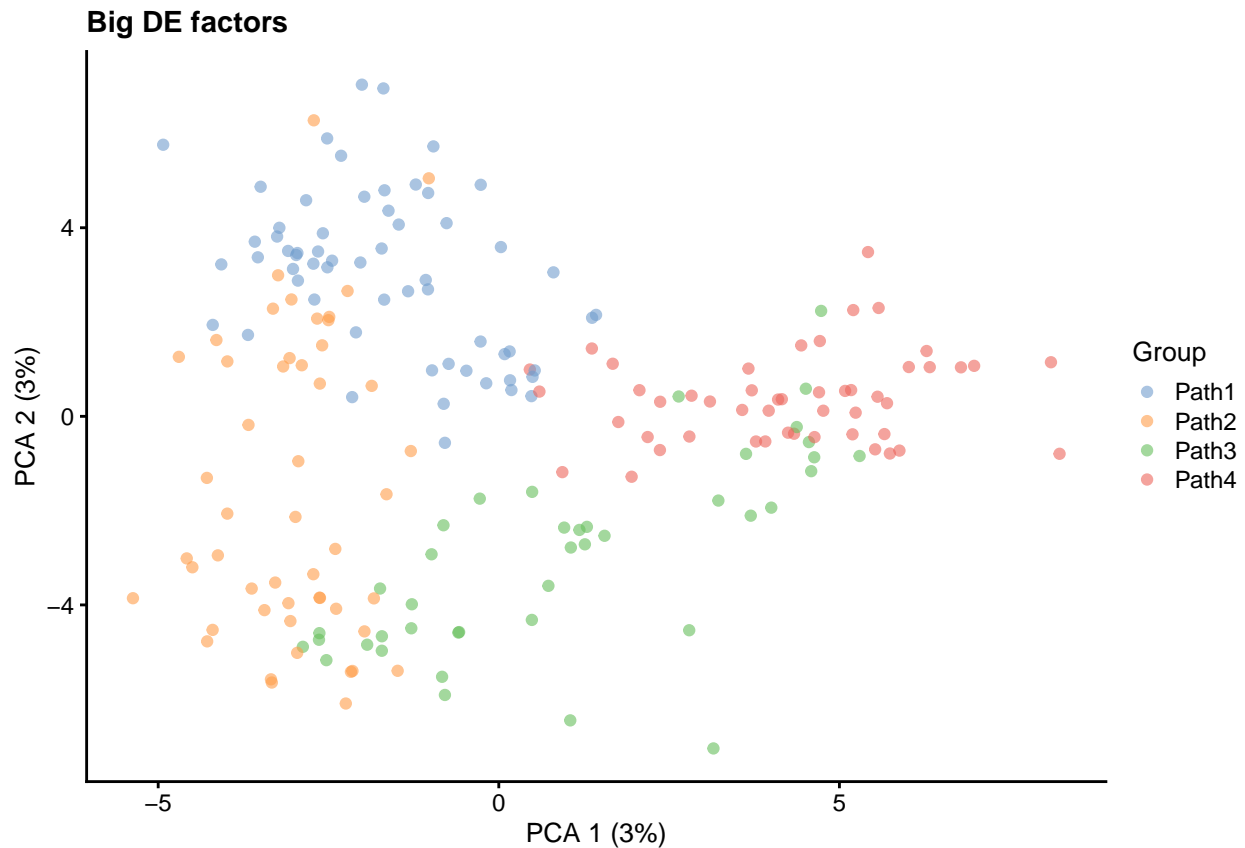
```
## $start.arg
## $start.arg$shape
## [1] 0.6324803
##
## $start.arg$rate
## [1] 0.2723197
##
##
## $fix.arg
## NULL
##
## $start.arg
## $start.arg$meanlog
## [1] 10.40535
##
## $start.arg$sdlog
## [1] 0.390699
##
##
## $fix.arg
## NULL
##
## $start.arg
## $start.arg$meanlog
## [1] 4.67681
```

```
##
## $start.arg$sdlog
## [1] 0.3149738
##
##
## $fix.arg
## NULL
```

You can either use the parameters estimated from real data to simulate data or pre-define those parameters by yourself.

```
### Use parameters estimated from real data to simulate data and you only need
### to define the group proportions (group.prob) and path directions (path.from).
sim.SCRIP1 <- SCRIPsimu(data=counts, params=params, method="paths",
                      group.prob = c(0.25, 0.25, 0.25, 0.25),
                      path.from = c(0, 1, 2, 3))

library(scater)
sim1 <- logNormCounts(sim.SCRIP1)
sim1 <- runPCA(sim1)
plotPCA(sim1, colour_by = "Group") + ggtitle("Big DE factors")
```



```
### Predefine parameters to simulate data
sim.SCRIP2 <- SCRIPsimu(data=counts, params=params, method="paths", batchCells=300,
                      group.prob = c(0.25, 0.25, 0.25, 0.25),
```

```

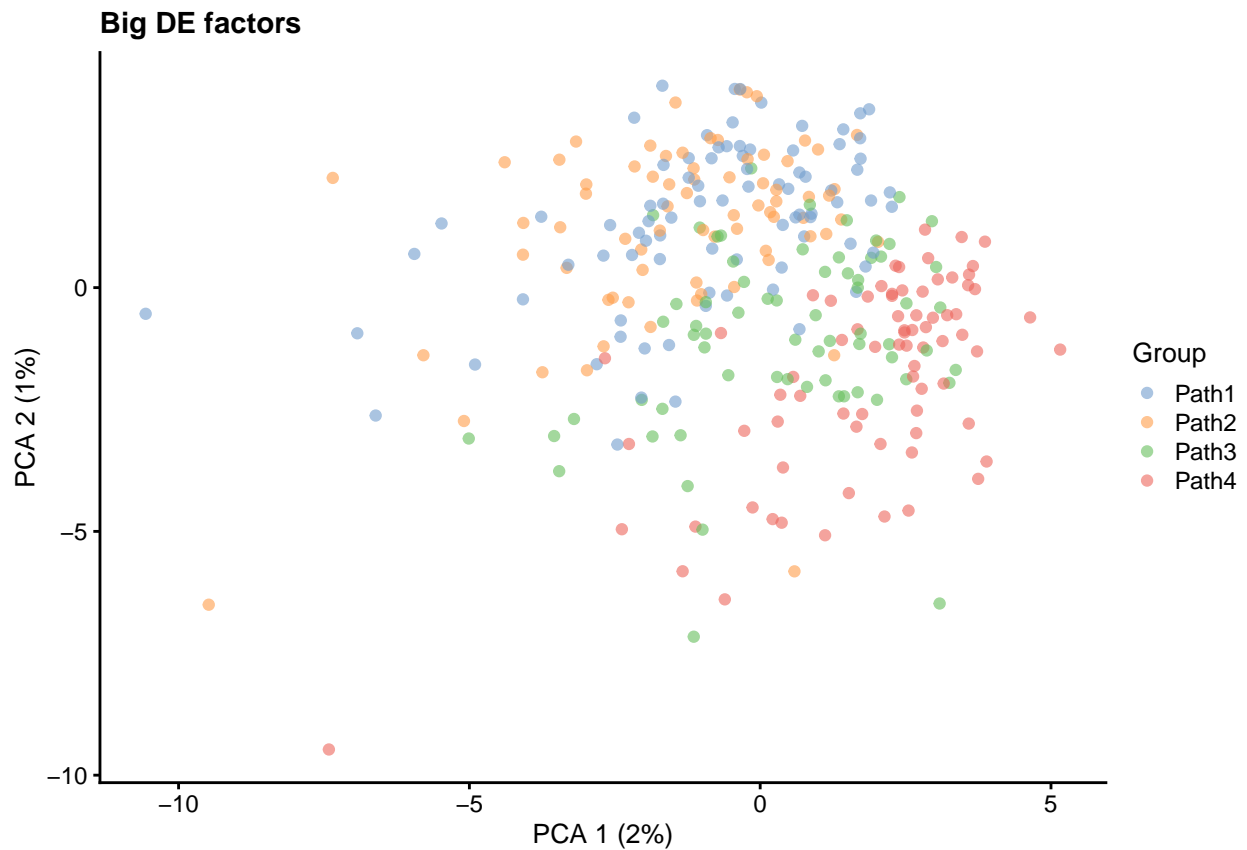
de.prob = c(0.2, 0.2, 0.2, 0.2),
de.downProb = c(0.5, 0.5, 0.5, 0.5),
de.facLoc = c(0.2, 0.2, 0.2, 0.2),
de.facScale=c(0.2, 0.2, 0.2, 0.2),
path.nSteps=c(100, 100, 100, 100),
path.from = c(0, 1, 2, 3),
path.skew=c(0.5,0.5,0.5,0.5))

```

```

library(scater)
sim2 <- logNormCounts(sim.SCRIP2)
sim2 <- runPCA(sim2)
plotPCA(sim2, colour_by = "Group") + ggtitle("Big DE factors")

```



7 Trajectory DE simulation

To mitigate the side effect of non-significant genes during analysis, it is essential to find genes that are differentially expressed within lineage or between lineages, referred as trajectory DE analysis. SCRIP enables evaluation of trajectory DE analysis methods with simulated DEGs for within lineage or between lineages. Between-lineage and within-lineage DEGs were simulated separately, with multiplicative differential expression factors generated from a log-normal distribution described above.

7.1 Within-lineage trajectory DE simulation

For within-lineage simulation, 20% genes were also randomly selected as within-lineage differentially expressed with other within-lineage DE parameters `de.downProb`, `de.facLoc` and `de.facScale`.

```
library(SCRIP)
library(splatter)
counts <- as.matrix(Tung)

counts <- counts[1:10000,1:200]
params <- splatEstimate((counts))
```

```
## $start.arg
## $start.arg$shape
## [1] 0.6324803
##
## $start.arg$rate
## [1] 0.2723197
##
##
## $fix.arg
## NULL
##
## $start.arg
## $start.arg$meanlog
## [1] 10.40535
##
## $start.arg$sdlog
## [1] 0.390699
##
##
## $fix.arg
## NULL
##
## $start.arg
## $start.arg$meanlog
## [1] 4.67681
##
## $start.arg$sdlog
## [1] 0.3149738
##
##
## $fix.arg
## NULL
```

```
sim.SCRIP.within <- SCRIPsimu(data=counts, params=params, method="paths",
                             de.prob = 0.2, de.downProb=0.5,
                             de.facLoc=1.2, de.facScale=0.1)
```

```
sim <- sim.SCRIP.within
All.genes <- rownames(counts(sim))
Truepaths <- rowData(sim)
DE.within.path1 <- Truepaths$DEFacPath1
True.within <- rep("common",10000)
True.within[which(DE.within.path1 != 1)] <- "DE"
table(True.within)
```

```
## True.within
## common      DE
##    7978    2022
```

7.1.1 Monocle3

Monocle3 test whether any of the genes change over time in their expression by fitting a generalized linear model.

```
library(monocle3)
expression_matrix <- counts(sim.SCRIP.within)
cell_metadata <- data.frame(libsize = apply(expression_matrix,2,sum))
rownames(cell_metadata) <- colnames(expression_matrix)
gene_annotation <- data.frame(gene_short_name=rownames(expression_matrix))
rownames(gene_annotation) <- rownames(expression_matrix)

cds.within <- new_cell_data_set(expression_matrix,
                               cell_metadata = cell_metadata,
                               gene_metadata = gene_annotation)

cds.within <- preprocess_cds(cds.within, num_dim = 100, method = "PCA")

cds.within <- reduce_dimension(cds.within, preprocess_method = "PCA",
                              reduction_method = "UMAP")
cds.within <- cluster_cells(cds.within, reduction_method="UMAP")
cds.within <- learn_graph(cds.within)

# We find all the cells that are close to the starting point
closest_vertex <- cds.within@principal_graph_aux[["UMAP"]]$pr_graph_cell_proj_closest_vertex
closest_vertex <- as.matrix(closest_vertex[colnames(cds.within), ])
closest_vertex <- as.numeric(names(which.max(table(closest_vertex))))
mst <- principal_graph(cds.within)$UMAP
root_pr_nodes <- igrph::V(mst)$name[closest_vertex]

# We compute the trajectory
cds.within <- order_cells(cds.within, root_pr_nodes = root_pr_nodes)

# plot_cells(cds.within)

gene_fits <- fit_models(cds.within, model_formula_str = "~pseudotime")
```

```
fit_coefs <- coefficient_table(gene_fits)

fit_coefs <- fit_coefs[which(fit_coefs$term != "(Intercept)"),]
```

7.1.2 tradeSeq

tradeSeq implements the `associationTest()` function to test whether the average gene expression is significantly changing along pseudotime.

```
library(magrittr)
# Get the closest vertice for every cell
y_to_cells <- principal_graph_aux(cds.within)$UMAP$pr_graph_cell_proj_closest_vertex %>%
  as.data.frame()
y_to_cells$cells <- rownames(y_to_cells)
y_to_cells$Y <- y_to_cells$V1

# Get the root vertices
# It is the same node as above
root <- cds.within@principal_graph_aux$UMAP$root_pr_nodes

# Get the other endpoints
endpoints <- names(which(igraph::degree(mst) == 1))
endpoints <- endpoints[!endpoints %in% root]

# For each endpoint
cellWeights <- lapply(endpoints, function(endpoint) {
  # We find the path between the endpoint and the root
  path <- igraph::shortest_paths(mst, root, endpoint)$vpath[[1]]
  path <- as.character(path)
  # We find the cells that map along that path
  df <- y_to_cells[y_to_cells$Y %in% path, ]
  df <- data.frame(weights = as.numeric(colnames(cds.within) %in% df$cells))
  colnames(df) <- endpoint
  return(df)
}) %>% do.call(what = 'cbind', args = .) %>%
  as.matrix()
rownames(cellWeights) <- colnames(cds.within)
pseudotime <- matrix(pseudotime(cds.within), ncol = ncol(cellWeights),
  nrow = ncol(cds.within), byrow = FALSE)

library(tradeSeq)
sce.within <- fitGAM(counts = expression_matrix ,
  pseudotime = pseudotime,
  cellWeights = cellWeights)
```

```
library(tradeSeq)
## Association of gene expression with pseudotime
assoRes <- associationTest(sce.within)
```


7.2 Between-lineage trajectory DE simulation

Between-lineage DE factors: Between-lineage DE factors are generated from a log-normal distribution with parameters including number of genes (nGenes), the path-specific proportion of DE genes (de.prob), the proportion of down-regulated DE genes (de.downProb), location factor (de.facLoc) and scale factor (de.facScale).

```
library(SCRIP)
library(splatter)
counts <- as.matrix(Tung)

counts <- counts[1:10000,1:400]
params <- splatEstimate(counts)
```

```
## $start.arg
## $start.arg$shape
## [1] 0.5919538
##
## $start.arg$rate
## [1] 0.2868007
##
##
## $fix.arg
## NULL
##
## $start.arg
## $start.arg$meanlog
## [1] 10.35125
##
## $start.arg$sdlog
## [1] 0.4409724
##
##
## $fix.arg
## NULL
##
## $start.arg
## $start.arg$meanlog
## [1] 4.996054
##
## $start.arg$sdlog
## [1] 0.257332
##
##
## $fix.arg
## NULL
```

```
sim.SCRIP.between <- SCRIPsimu(data=counts, params=params, method="paths",
                                group.prob = c(0.5,0.5),
                                de.prob = c(0.0,0.2), de.downProb=c(0.5,0.5),
                                de.facLoc=c(1.2,1.2), de.facScale=c(0.1,0.1),
                                path.from = c(0, 1),path.nSteps = c(100, 100))
```

```

Truepaths <- rowData(sim.SCRIP.between)
DE.between.path <- Truepaths$DEFacPath2
True.between <- rep("common",10000)
True.between[which(DE.between.path != 1)] <- "DE"
table(True.between)

```

```

## True.between
## common      DE
##    7984    2016

```

7.2.1 monocle3

The `graph_test()` in `monocle3` uses a statistic from spatial autocorrelation analysis called Moran's I to identify the DEGs.

```

library(monocle3)
expression_matrix <- counts(sim.SCRIP.between)
cell_metadata <- data.frame(libsize = apply(expression_matrix,2,sum))
rownames(cell_metadata) <- colnames(expression_matrix)
gene_annotation <- data.frame(gene_short_name=rownames(expression_matrix))
rownames(gene_annotation) <- rownames(expression_matrix)

cds <- new_cell_data_set(expression_matrix,
                        cell_metadata = cell_metadata,
                        gene_metadata = gene_annotation)

cds <- preprocess_cds(cds, num_dim = 100, method = "PCA")

cds <- reduce_dimension(cds,preprocess_method = "PCA",
                      reduction_method = "UMAP")
cds <- cluster_cells(cds,reduction_method="UMAP")
cds <- learn_graph(cds,use_partition = FALSE)

# We find all the cells that are close to the starting point
closest_vertex <- cds@principal_graph_aux[["UMAP"]]$pr_graph_cell_proj_closest_vertex
closest_vertex <- as.matrix(closest_vertex[colnames(cds), ])
closest_vertex <- as.numeric(names(which.max(table(closest_vertex))))
mst <- principal_graph(cds)$UMAP
root_pr_nodes <- igraph::V(mst)$name[closest_vertex]

# We compute the trajectory
cds <- order_cells(cds, root_pr_nodes = root_pr_nodes)
plot_cells(cds, color_cells_by = "pseudotime")

pr_graph_test_res <- graph_test(cds, neighbor_graph="knn", cores=8)

```

7.2.2 tradeSeq

The function `patternTest()` implements a statistical model to identify DEGs which have significantly different expression patterns.

```

library(magrittr)
# Get the closest vertice for every cell
y_to_cells <- principal_graph_aux(cds)$UMAP$pr_graph_cell_proj_closest_vertex %>%
  as.data.frame()
y_to_cells$cells <- rownames(y_to_cells)
y_to_cells$Y <- y_to_cells$V1

# Get the root vertices
root <- cds@principal_graph_aux$UMAP$root_pr_nodes

# Get the other endpoints
endpoints <- names(which(igraph::degree(mst) == 1))
endpoints <- endpoints[!endpoints %in% root]

# For each endpoint
cellWeights <- lapply(endpoints, function(endpoint) {
  # We find the path between the endpoint and the root
  path <- igraph::shortest_paths(mst, root, endpoint)$vpath[[1]]
  path <- as.character(path)
  # We find the cells that map along that path
  df <- y_to_cells[y_to_cells$Y %in% path, ]
  df <- data.frame(weights = as.numeric(colnames(cds) %in% df$cells))
  colnames(df) <- endpoint
  return(df)
}) %>% do.call(what = 'cbind', args = .) %>%
  as.matrix()

rownames(cellWeights) <- colnames(cds)
pseudotime <- matrix(pseudotime(cds), ncol = ncol(cellWeights),
                     nrow = ncol(cds), byrow = FALSE)

library(tradeSeq)
sce <- fitGAM(counts = expression_matrix,
              pseudotime = pseudotime,
              cellWeights = cellWeights)

library(tradeSeq)
patternRes <- patternTest(sce)

```