

# scGPS introduction

Quan Nguyen and Michael Thompson

2018-12-17

## Contents

<b>1. Installation instruction</b>	<b>1</b>
<b>2. A simple workflow of the scGPS:</b>	<b>2</b>
2.1 Create scGPS objects . . . . .	2
2.2 Run prediction . . . . .	2
2.3 Summarise results . . . . .	2
<b>3. A complete workflow of the scGPS:</b>	<b>3</b>
3.1 Identify clusters in a dataset using CORE . . . . .	3
3.1 Identify clusters in a dataset using SCORE (Stable Clustering at Optimal REsolution) . . . . .	3
3.2 Visualise all cluster results in all iterations . . . . .	5
3.4 Compare clustering results with other dimensional reduction methods (e.g., CIDR) . . . . .	5
3.5 Find gene markers and annotate clusters . . . . .	5
<b>4. Relationship between clusters within one sample or between two samples</b>	<b>7</b>
4.1 Start the scGPS prediction to find relationship between clusters . . . . .	7
4.2 Display summary results for the prediction . . . . .	12
4.3 Plot the relationship between clusters in one sample . . . . .	13
4.3 Plot the relationship between clusters in two samples . . . . .	15

## 1. Installation instruction

```
# Prior to installing scGPS you need to install the SummarizedExperiment
# bioconductor package as the following
# source('https://bioconductor.org/biocLite.R') biocLite('SummarizedExperiment')

# To install scGPS from github (Depending on the configuration of the local
# computer or HPC, possible custom C++ compilation may be required - see
# installation trouble-shootings below)
devtools::install_github("IMB-Computational-Genomics-Lab/scGPS")

# for C++ compilation trouble-shooting, manual download and installation can be
# done from github

git clone https://github.com/IMB-Computational-Genomics-Lab/scGPS

# then check in scGPS/src if any of the precompiled (e.g. those with *.so and
# *.o) files exist and delete them before recompiling

# create a Makevars file in the scGPS/src with one line: PKG_LIBS =
# $(LAPACK_LIBS) $(BLAS_LIBS) $(FLIBS)

# then with the scGPS as the R working directory, manually recompile scGPS in R
```

```
# using devtools to load and install functions
devtools::document()
#load the package to the workspace
devtools::load_all()
```

## 2. A simple workflow of the scGPS:

*The purpose of this workflow is to solve the following task: given a mixed population with known subpopulations, estimate transition scores between these subpopulation*

### 2.1 Create scGPS objects

```
# load mixed population 1 (loaded from sample1 dataset, named it as day2)
# setwd('/Users/quan.nguyen/Documents/Powell_group_MacQuan/AllCodes/scGPS/vignettes/')
devtools::load_all()

day2 <- sample1
mixedpop1 <- NewscGPS(ExpressionMatrix = day2$dat2_counts, GeneMetadata = day2$dat2geneInfo,
  CellMetadata = day2$dat2_clusters)

# load mixed population 2 (loaded from sample2 dataset, named it as day5)
day5 <- sample2
mixedpop2 <- NewscGPS(ExpressionMatrix = day5$dat5_counts, GeneMetadata = day5$dat5geneInfo,
  CellMetadata = day5$dat5_clusters)
```

### 2.2 Run prediction

```
# select a subpopulation
c_selectID <- 1
# load gene list (this can be any lists of user selected genes)
genes <- GeneList
genes <- genes$Merged_unique
# load cluster information
cluster_mixedpop1 <- colData(mixedpop1)[,1]
cluster_mixedpop2 <- colData(mixedpop2)[,1]
#run training
LSOLDA_dat <- bootstrap_scGPS(nboots = 2, mixedpop1 = mixedpop1,
  mixedpop2 = mixedpop2, genes = genes, c_selectID = c_selectID, listData = list(),
  cluster_mixedpop1 = cluster_mixedpop1,
  cluster_mixedpop2 = cluster_mixedpop2)
```

### 2.3 Summarise results

```
# display the list of result information in the LASOLDA_dat object
names(LSOLDA_dat)
LSOLDA_dat$ElasticNetPredict
```

```

LSOLDA_dat$LDAPredict

# summary results LDA
summary_prediction_lda(LSOLDA_dat = LSOLDA_dat, nPredSubpop = 4)

# summary results Lasso to show the percent of cells classified as cells belonging
summary_prediction_lasso(LSOLDA_dat = LSOLDA_dat, nPredSubpop = 4)

# summary accuracy to check the model accuracy in the leave-out test set
summary_accuracy(object = LSOLDA_dat)

# summary maximum deviance explained by the model
summary_deviance(object = LSOLDA_dat)

```

### 3. A complete workflow of the scGPS:

*The purpose of this workflow is to solve the following task: given an unknown mixed population, find clusters and estimate relationship between clusters*

#### 3.1 Identify clusters in a dataset using CORE

*(skip this step if clusters are known)*

```

# find clustering information in an expression data using CORE
day5 <- sample2
cellnames <- colnames(day5$dat5_counts)
cluster <- day5$dat5_clusters
cellnames <- data.frame("Cluster"=cluster, "cellBarcodes" = cellnames)
mixedpop2 <- NewscGPS(ExpressionMatrix = day5$dat5_counts, GeneMetadata = day5$dat5geneInfo, CellMetadata = cellnames)

CORE_cluster <- CORE_scGPS(mixedpop2, remove_outlier = c(0), PCA=FALSE)

```

#### 3.1 Identify clusters in a dataset using SCORE (Stable Clustering at Optimal Resolution)

*(skip this step if clusters are known) (SCORE aims to get stable subpopulation results, by introducing bagging aggregation and bootstrapping to the CORE algorithm)*

```

# find clustering information in an expression data using SCORE
day5 <- sample2
cellnames <- colnames(day5$dat5_counts)
cluster <- day5$dat5_clusters
cellnames <- data.frame("Cluster"=cluster, "cellBarcodes" = cellnames)
mixedpop2 <- NewscGPS(ExpressionMatrix = day5$dat5_counts, GeneMetadata = day5$dat5geneInfo, CellMetadata = cellnames)

SCORE_test <- CORE_scGPS_bagging(mixedpop2, remove_outlier = c(0), PCA=FALSE,
                                bagging_run = 20, subsample_proportion = .8)

#> [1] "Performing 1 round of filtering"
#> [1] "Identifying top variable genes"
#> [1] "Calculating distance matrix"

```



## 3.2 Visualise all cluster results in all iterations

```
##3.2.1 plot CORE clustering
plot_CORE(CORE_cluster$tree, CORE_cluster$Cluster) #plot all clustering bars
#extract optimal index identified by CORE_scGPS
key_height <- CORE_cluster$OptimalClust$KeyStats$Height
optimal_res <- CORE_cluster$OptimalClust$OptimalRes
optimal_index = which(key_height == optimal_res)
#plot one optimal clustering bar
plot_optimal_CORE(original_tree= CORE_cluster$tree,
                  optimal_cluster = unlist(CORE_cluster$Cluster[optimal_index]), shift = -2000)
# you can customise the cluster color bars (provide color_branch values)
plot_CORE(CORE_cluster$tree, CORE_cluster$Cluster, color_branch = c("#208eb7", "#6ce9d3", "#1c5e39", "#"))

##3.2.2 plot SCORE clustering
plot_CORE(SCORE_test$tree, list_clusters = SCORE_test$Cluster) #plot all clustering bars
#plot one stable optimal clustering bar
plot_optimal_CORE(original_tree= SCORE_test$tree,
                  optimal_cluster = unlist(SCORE_test$Cluster[SCORE_test$optimal_index]), shift = -100)
```

## 3.4 Compare clustering results with other dimensional reduction methods (e.g., CIDR)

```
library(cidr)
t <- CIDR_scGPS(expression.matrix=assay(mixedpop2))
p2 <- plotReduced_scGPS(t, color_fac = factor(colData(mixedpop2)[,1]), palletes = 1:length(unique(colData(mixedpop2)[,1])))
p2
```

## 3.5 Find gene markers and annotate clusters

```
#load gene list (this can be any lists of user-selected genes)
genes <- GeneList
genes <- genes$Merged_unique

#the gene list can also be objectively identified by differential expression analysis
#cluster information is required for findMarkers_scGPS. Here, we use CORE results.

colData(mixedpop2)[,1] <- unlist(SCORE_test$Cluster[SCORE_test$optimal_index])

suppressMessages(library(locfit))
suppressMessages(library(DESeq))

DEgenes <- findMarkers_scGPS(expression_matrix=assay(mixedpop2), cluster = colData(mixedpop2)[,1],
                           selected_cluster=unique(colData(mixedpop2)[,1]))

#> [1] "Start estimate dispersions for cluster 1..."
#> [1] "Done estimate dispersions. Start nbinom test for cluster 1..."
#> [1] "Done nbinom test for cluster 1 ..."
#> [1] "Adjust foldchange by subtracting basemean to 1..."
#> [1] "Start estimate dispersions for cluster 2..."
#> [1] "Done estimate dispersions. Start nbinom test for cluster 2..."
```

```

#> [1] "Done nbinom test for cluster 2 ..."
#> [1] "Adjust foldchange by subtracting basemean to 1..."
#> [1] "Start estimate dispersions for cluster 3..."
#> [1] "Done estimate dispersions. Start nbinom test for cluster 3..."
#> [1] "Done nbinom test for cluster 3 ..."
#> [1] "Adjust foldchange by subtracting basemean to 1..."
#> [1] "Start estimate dispersions for cluster 4..."
#> [1] "Done estimate dispersions. Start nbinom test for cluster 4..."
#> [1] "Done nbinom test for cluster 4 ..."
#> [1] "Adjust foldchange by subtracting basemean to 1..."

#the output contains dataframes for each cluster.
#the data frame contains all genes, sorted by p-values
names(DEgenes)
#> [1] "DE_Subpop1vsRemaining" "DE_Subpop2vsRemaining" "DE_Subpop3vsRemaining"
#> [4] "DE_Subpop4vsRemaining"

#you can annotate the identified clusters
DEgeneList_3vsOthers <- DEgenes$DE_Subpop3vsRemaining$id

#users need to check the format of the gene input to make sure they are consistent to
#the gene names in the expression matrix
DEgeneList_3vsOthers <-gsub("_.*", "", DEgeneList_3vsOthers )

#the following command saves the file "PathwayEnrichment.xlsx" to the working dir
#use 500 top DE genes
suppressMessages(library(DOSE))
suppressMessages(library(ReactomePA))
suppressMessages(library(clusterProfiler))
enrichment_test <- annotate_scGPS(DEgeneList_3vsOthers[1:500], pvalueCutoff=0.05, gene_symbol=TRUE)
#> [1] "Original gene number in geneList"
#> [1] 500
#> [1] "Number of genes successfully converted"
#> [1] 486

#the enrichment outputs can be displayed by running
dotplot(enrichment_test, showCategory=15)

```

Signaling by Receptor Tyrosi  
 Extracellular matrix o  
 Muscle  
 Striated Muscle (c  
 ECM pro  
 Degradation of the extracel  
 Smooth Muscle (c  
 Growth Factor (IGF) transport and uptake by Insulin-like Growth Factor Binding Protein  
 Collagen c  
 Cell junction o  
 Assembly of collagen fibrils and other multimeric  
 Collagen chain tr  
 Molecules associated with el  
 Cell-extracellular matrix i  
 Endosomal/Vacuol

## 4. Relationship between clusters within one sample or between two samples

*The purpose of this workflow is to solve the following task: given one or two unknown mixed population(s) and clusters in each mixed population, estimate and visualise relationship between clusters*

### 4.1 Start the scGPS prediction to find relationship between clusters

```

#select a subpopulation, and input gene list
c_selectID <- 1
genes = DEgenes$DE_Subpop1vsRemaining$id[1:500]
#format gene names
genes <- gsub("_.*", "", genes)

#run the test bootstrap with nboots = 2 runs

cluster_mixedpop1 <- colData(mixedpop1)[,1]
cluster_mixedpop2 <- colData(mixedpop2)[,1]

sink("temp")
LSOLDA_dat <- bootstrap_scGPS(nboots = 2, mixedpop1 = mixedpop1,
  mixedpop2 = mixedpop2, genes = genes, c_selectID = c_selectID, listData = list(),
  cluster_mixedpop1 = cluster_mixedpop1,
  cluster_mixedpop2 = cluster_mixedpop2)
#>
#> Call: glmnet(x = as.matrix(dataset[, -which(colnames(dataset) == "Cluster_class")]), y = y_cat
#>

```

```

#>      Df      %Dev  Lambda
#> [1,]  0 -2.563e-15  0.293400
#> [2,]  2  2.893e-02  0.280000
#> [3,]  4  6.172e-02  0.267300
#> [4,]  4  9.250e-02  0.255200
#> [5,]  5  1.217e-01  0.243600
#> [6,]  5  1.493e-01  0.232500
#> [7,]  6  1.750e-01  0.221900
#> [8,]  6  1.992e-01  0.211800
#> [9,]  6  2.217e-01  0.202200
#> [10,] 6  2.428e-01  0.193000
#> [11,] 6  2.624e-01  0.184200
#> [12,] 6  2.809e-01  0.175900
#> [13,] 7  2.983e-01  0.167900
#> [14,] 8  3.163e-01  0.160200
#> [15,] 8  3.339e-01  0.153000
#> [16,] 9  3.524e-01  0.146000
#> [17,] 10 3.716e-01  0.139400
#> [18,] 10 3.901e-01  0.133000
#> [19,] 11 4.079e-01  0.127000
#> [20,] 13 4.256e-01  0.121200
#> [21,] 13 4.425e-01  0.115700
#> [22,] 14 4.585e-01  0.110500
#> [23,] 16 4.742e-01  0.105400
#> [24,] 16 4.894e-01  0.100600
#> [25,] 16 5.039e-01  0.096060
#> [26,] 16 5.177e-01  0.091700
#> [27,] 16 5.309e-01  0.087530
#> [28,] 16 5.436e-01  0.083550
#> [29,] 16 5.556e-01  0.079750
#> [30,] 16 5.671e-01  0.076130
#> [31,] 16 5.782e-01  0.072670
#> [32,] 17 5.899e-01  0.069370
#> [33,] 20 6.017e-01  0.066210
#> [34,] 20 6.133e-01  0.063200
#> [35,] 20 6.243e-01  0.060330
#> [36,] 21 6.350e-01  0.057590
#> [37,] 22 6.454e-01  0.054970
#> [38,] 24 6.556e-01  0.052470
#> [39,] 26 6.657e-01  0.050090
#> [40,] 28 6.756e-01  0.047810
#> [41,] 29 6.853e-01  0.045640
#> [42,] 29 6.949e-01  0.043560
#> [43,] 31 7.042e-01  0.041580
#> [44,] 31 7.133e-01  0.039690
#> [45,] 33 7.221e-01  0.037890
#> [46,] 33 7.307e-01  0.036170
#> [47,] 35 7.394e-01  0.034520
#> [48,] 36 7.483e-01  0.032950
#> [49,] 38 7.569e-01  0.031460
#> [50,] 39 7.656e-01  0.030030
#> [51,] 38 7.740e-01  0.028660
#> [52,] 40 7.822e-01  0.027360

```



```

#> [53,] 40 7.903e-01 0.026120
#> [54,] 40 7.979e-01 0.024930
#> [55,] 40 8.053e-01 0.023800
#> [56,] 41 8.125e-01 0.022710
#> [57,] 40 8.194e-01 0.021680
#> [58,] 44 8.262e-01 0.020700
#> [59,] 45 8.329e-01 0.019760
#> [60,] 46 8.394e-01 0.018860
#> [61,] 44 8.458e-01 0.018000
#> [62,] 44 8.519e-01 0.017180
#> [63,] 45 8.577e-01 0.016400
#> [64,] 45 8.635e-01 0.015660
#> [65,] 44 8.690e-01 0.014940
#> [66,] 48 8.743e-01 0.014270
#> [67,] 49 8.796e-01 0.013620
#> [68,] 49 8.847e-01 0.013000
#> [69,] 49 8.895e-01 0.012410
#> [70,] 49 8.943e-01 0.011840
#> [71,] 50 8.988e-01 0.011300
#> [72,] 50 9.031e-01 0.010790
#> [73,] 50 9.073e-01 0.010300
#> [74,] 50 9.113e-01 0.009832
#> [75,] 51 9.151e-01 0.009386
#> [76,] 52 9.188e-01 0.008959
#> [77,] 52 9.223e-01 0.008552
#> [78,] 53 9.256e-01 0.008163
#> [79,] 53 9.289e-01 0.007792
#> [80,] 54 9.319e-01 0.007438
#> [81,] 55 9.349e-01 0.007100
#> [82,] 55 9.378e-01 0.006777
#> [83,] 56 9.405e-01 0.006469
#> [84,] 56 9.431e-01 0.006175
#> [85,] 58 9.456e-01 0.005894
#> [86,] 59 9.481e-01 0.005627
#> [87,] 60 9.504e-01 0.005371
#> [88,] 63 9.526e-01 0.005127
#> [89,] 63 9.547e-01 0.004894
#> [90,] 63 9.568e-01 0.004671
#> [91,] 62 9.587e-01 0.004459
#> [92,] 62 9.606e-01 0.004256
#> [93,] 62 9.624e-01 0.004063
#> [94,] 62 9.641e-01 0.003878
#> [95,] 62 9.657e-01 0.003702
#> [96,] 62 9.673e-01 0.003534
#> [97,] 62 9.687e-01 0.003373
#> [98,] 63 9.701e-01 0.003220
#> [99,] 64 9.715e-01 0.003073
#> [100,] 64 9.728e-01 0.002934
#> [1] "please check the lambda min output ..."
#> [1] "done bootstrap 1"
#>
#> Call: glmnet(x = as.matrix(dataset[, -which(colnames(dataset) == "Cluster_class")]),
#>
y = y_cat

```

```

#>      Df      %Dev   Lambda
#> [1,]  0 -2.563e-15  0.299100
#> [2,]  3  3.026e-02  0.285500
#> [3,]  3  5.969e-02  0.272500
#> [4,]  3  8.679e-02  0.260100
#> [5,]  3  1.119e-01  0.248300
#> [6,]  3  1.351e-01  0.237000
#> [7,]  3  1.568e-01  0.226300
#> [8,]  4  1.774e-01  0.216000
#> [9,]  4  1.981e-01  0.206200
#> [10,] 4  2.174e-01  0.196800
#> [11,] 5  2.365e-01  0.187800
#> [12,] 5  2.545e-01  0.179300
#> [13,] 5  2.714e-01  0.171200
#> [14,] 5  2.872e-01  0.163400
#> [15,] 6  3.026e-01  0.155900
#> [16,] 6  3.180e-01  0.148900
#> [17,] 7  3.333e-01  0.142100
#> [18,] 8  3.478e-01  0.135600
#> [19,] 8  3.616e-01  0.129500
#> [20,] 9  3.746e-01  0.123600
#> [21,] 9  3.874e-01  0.118000
#> [22,] 10 4.003e-01  0.112600
#> [23,] 10 4.128e-01  0.107500
#> [24,] 10 4.246e-01  0.102600
#> [25,] 10 4.359e-01  0.097940
#> [26,] 10 4.465e-01  0.093490
#> [27,] 10 4.567e-01  0.089240
#> [28,] 10 4.663e-01  0.085180
#> [29,] 11 4.756e-01  0.081310
#> [30,] 13 4.860e-01  0.077620
#> [31,] 14 4.967e-01  0.074090
#> [32,] 14 5.070e-01  0.070720
#> [33,] 16 5.171e-01  0.067510
#> [34,] 17 5.269e-01  0.064440
#> [35,] 17 5.363e-01  0.061510
#> [36,] 17 5.452e-01  0.058710
#> [37,] 20 5.543e-01  0.056050
#> [38,] 21 5.643e-01  0.053500
#> [39,] 21 5.742e-01  0.051070
#> [40,] 22 5.837e-01  0.048750
#> [41,] 22 5.928e-01  0.046530
#> [42,] 24 6.017e-01  0.044410
#> [43,] 25 6.112e-01  0.042400
#> [44,] 25 6.207e-01  0.040470
#> [45,] 26 6.299e-01  0.038630
#> [46,] 27 6.390e-01  0.036870
#> [47,] 27 6.481e-01  0.035200
#> [48,] 27 6.567e-01  0.033600
#> [49,] 31 6.654e-01  0.032070
#> [50,] 33 6.749e-01  0.030610
#> [51,] 35 6.847e-01  0.029220
#> [52,] 38 6.945e-01  0.027890

```

```

#> [53,] 38 7.045e-01 0.026630
#> [54,] 40 7.143e-01 0.025420
#> [55,] 43 7.241e-01 0.024260
#> [56,] 45 7.341e-01 0.023160
#> [57,] 47 7.441e-01 0.022110
#> [58,] 47 7.540e-01 0.021100
#> [59,] 48 7.635e-01 0.020140
#> [60,] 49 7.727e-01 0.019230
#> [61,] 54 7.820e-01 0.018350
#> [62,] 55 7.916e-01 0.017520
#> [63,] 54 8.006e-01 0.016720
#> [64,] 54 8.092e-01 0.015960
#> [65,] 56 8.176e-01 0.015240
#> [66,] 58 8.257e-01 0.014540
#> [67,] 59 8.336e-01 0.013880
#> [68,] 60 8.411e-01 0.013250
#> [69,] 60 8.483e-01 0.012650
#> [70,] 60 8.552e-01 0.012070
#> [71,] 62 8.618e-01 0.011530
#> [72,] 64 8.681e-01 0.011000
#> [73,] 64 8.742e-01 0.010500
#> [74,] 66 8.799e-01 0.010020
#> [75,] 67 8.854e-01 0.009569
#> [76,] 67 8.906e-01 0.009134
#> [77,] 68 8.956e-01 0.008719
#> [78,] 67 9.003e-01 0.008323
#> [79,] 67 9.048e-01 0.007944
#> [80,] 67 9.091e-01 0.007583
#> [81,] 69 9.132e-01 0.007238
#> [82,] 69 9.171e-01 0.006909
#> [83,] 70 9.209e-01 0.006595
#> [84,] 70 9.244e-01 0.006296
#> [85,] 69 9.278e-01 0.006010
#> [86,] 70 9.311e-01 0.005736
#> [87,] 70 9.342e-01 0.005476
#> [88,] 70 9.371e-01 0.005227
#> [89,] 72 9.399e-01 0.004989
#> [90,] 73 9.426e-01 0.004762
#> [91,] 73 9.452e-01 0.004546
#> [92,] 75 9.477e-01 0.004339
#> [93,] 77 9.500e-01 0.004142
#> [94,] 78 9.523e-01 0.003954
#> [95,] 78 9.544e-01 0.003774
#> [96,] 78 9.565e-01 0.003603
#> [97,] 78 9.585e-01 0.003439
#> [98,] 80 9.604e-01 0.003283
#> [99,] 80 9.621e-01 0.003133
#> [100,] 81 9.639e-01 0.002991
#> [1] "done bootstrap 2"

```

```
sink()
```

## 4.2 Display summary results for the prediction

```
#get the number of rows for the summary matrix
row_cluster <-length(unique(colData(mixedpop2)[,1]))

#summary results LDA to show the percent of cells classified as cells belonging by LDA classifier
summary_prediction_lda(LSOLDA_dat=LSOLDA_dat, nPredSubpop = row_cluster )
#>               V1               V2               names
#> 1             56.25             66.40625 LDA for subpop 1 in target mixedpop2
#> 2 44.16666666666667             35 LDA for subpop 2 in target mixedpop2
#> 3 44.2105263157895 27.3684210526316 LDA for subpop 3 in target mixedpop2
#> 4 48.2758620689655 34.4827586206897 LDA for subpop 4 in target mixedpop2

#summary results Lasso to show the percent of cells classified as cells belonging by Lasso classifier
summary_prediction_lasso(LSOLDA_dat=LSOLDA_dat, nPredSubpop = row_cluster)
#>               V1               V2
#> 1             53.515625             60.546875
#> 2 58.33333333333333 20.83333333333333
#> 3 62.1052631578947 14.7368421052632
#> 4 72.4137931034483 31.0344827586207
#>               names
#> 1 ElasticNet for subpop1 in target mixedpop2
#> 2 ElasticNet for subpop2 in target mixedpop2
#> 3 ElasticNet for subpop3 in target mixedpop2
#> 4 ElasticNet for subpop4 in target mixedpop2

# summary maximum deviance explained by the model during the model training
summary_deviance(object = LSOLDA_dat)
#> $allDeviance
#> [1] "0.9728" "0.6945"
#>
#> $DeviMax
#>      Dfd  Deviance  DEgenes
#> 1      0 -2.563e-15 genes_cluster1
#> 2      2  0.02893 genes_cluster1
#> 3      4  0.0925 genes_cluster1
#> 4      5  0.1493 genes_cluster1
#> 5      6  0.2809 genes_cluster1
#> 6      7  0.2983 genes_cluster1
#> 7      8  0.3339 genes_cluster1
#> 8      9  0.3524 genes_cluster1
#> 9     10  0.3901 genes_cluster1
#> 10     11  0.4079 genes_cluster1
#> 11     13  0.4425 genes_cluster1
#> 12     14  0.4585 genes_cluster1
#> 13     16  0.5782 genes_cluster1
#> 14     17  0.5899 genes_cluster1
#> 15     20  0.6243 genes_cluster1
#> 16     21  0.635 genes_cluster1
#> 17     22  0.6454 genes_cluster1
#> 18     24  0.6556 genes_cluster1
#> 19     26  0.6657 genes_cluster1
#> 20     28  0.6756 genes_cluster1
```

```

#> 21      29      0.6949 genes_cluster1
#> 22      31      0.7133 genes_cluster1
#> 23      33      0.7307 genes_cluster1
#> 24      35      0.7394 genes_cluster1
#> 25      36      0.7483 genes_cluster1
#> 26      38      0.774 genes_cluster1
#> 27      39      0.7656 genes_cluster1
#> 28      40      0.8194 genes_cluster1
#> 29      41      0.8125 genes_cluster1
#> 30      44      0.869 genes_cluster1
#> 31      45      0.8635 genes_cluster1
#> 32      46      0.8394 genes_cluster1
#> 33      48      0.8743 genes_cluster1
#> 34      49      0.8943 genes_cluster1
#> 35      50      0.9113 genes_cluster1
#> 36      51      0.9151 genes_cluster1
#> 37      52      0.9223 genes_cluster1
#> 38      53      0.9289 genes_cluster1
#> 39      54      0.9319 genes_cluster1
#> 40      55      0.9378 genes_cluster1
#> 41      56      0.9431 genes_cluster1
#> 42      58      0.9456 genes_cluster1
#> 43      59      0.9481 genes_cluster1
#> 44      60      0.9504 genes_cluster1
#> 45      62      0.9687 genes_cluster1
#> 46      63      0.9701 genes_cluster1
#> 47      64      0.9728 genes_cluster1
#> 48 remaining      1      DEgenes
#>
#> $LassoGenesMax
#> NULL

# summary accuracy to check the model accuracy in the leave-out test set
summary_accuracy(object = LSOLDA_dat)
#> [1] 86.16071 91.96429

```

### 4.3 Plot the relationship between clusters in one sample

Here we look at one example use case to find relationship between clusters within one sample or between two sample

```

#run prediction for 3 clusters
cluster_mixedpop1 <- colData(mixedpop1)[,1]
cluster_mixedpop2 <- as.numeric(as.vector(colData(mixedpop2)[,1]))

c_selectID <- 1
genes = DEgenes$DE_Subpop1vsRemaining$id[1:200] #top 200 gene markers distinguishing cluster 1
genes <- gsub("_.*", "", genes)

LSOLDA_dat1 <- bootstrap_scGPS(nboots = 2, mixedpop1 = mixedpop2, mixedpop2 = mixedpop2, genes=genes, c

c_selectID <- 2
genes = DEgenes$DE_Subpop2vsRemaining$id[1:200]

```

```

genes <- gsub("_.*", "", genes)
LSOLDA_dat2 <- bootstrap_scGPS(nboots = 2, mixedpop1 = mixedpop2, mixedpop2 = mixedpop2, genes=genes, c_
cluster_mixedpop2 = cluster_mixedpop2)

c_selectID <- 3
genes = DEgenes$DE_Subpop3vsRemaining$id[1:200]
genes <- gsub("_.*", "", genes)
LSOLDA_dat3 <- bootstrap_scGPS(nboots = 2, mixedpop1 = mixedpop2, mixedpop2 = mixedpop2, genes=genes, c_
cluster_mixedpop2 = cluster_mixedpop2)

c_selectID <- 4
genes = DEgenes$DE_Subpop4vsRemaining$id[1:200]
genes <- gsub("_.*", "", genes)
LSOLDA_dat4 <- bootstrap_scGPS(nboots = 2, mixedpop1 = mixedpop2, mixedpop2 = mixedpop2, genes=genes, c_
cluster_mixedpop2 = cluster_mixedpop2)

#prepare table input for sankey plot

LASSO_C1S2 <- reformat_LASSO(c_selectID=1, mp_selectID = 2, LSOLDA_dat=LSOLDA_dat1,
nPredSubpop = length(unique(colData(mixedpop2)[,1])),
Nodes_group = "#7570b3")

LASSO_C2S2 <- reformat_LASSO(c_selectID=2, mp_selectID =2, LSOLDA_dat=LSOLDA_dat2,
nPredSubpop = length(unique(colData(mixedpop2)[,1])),
Nodes_group = "#1b9e77")

LASSO_C3S2 <- reformat_LASSO(c_selectID=3, mp_selectID =2, LSOLDA_dat=LSOLDA_dat3,
nPredSubpop = length(unique(colData(mixedpop2)[,1])),
Nodes_group = "#e7298a")

LASSO_C4S2 <- reformat_LASSO(c_selectID=4, mp_selectID =2, LSOLDA_dat=LSOLDA_dat4,
nPredSubpop = length(unique(colData(mixedpop2)[,1])),
Nodes_group = "#00FFFF")

combined <- rbind(LASSO_C1S2, LASSO_C2S2, LASSO_C3S2, LASSO_C4S2 )
combined <- combined[is.na(combined$Value) != TRUE,]

nboots = 2
#links: source, target, value
#source: node, nodegroup
combined_D3obj <- list(Nodes=combined[, (nboots+3):(nboots+4)], Links=combined[, c((nboots+2):(nboots+1), n

library(networkD3)

Node_source <- as.vector(sort(unique(combined_D3obj$Links$Source)))
Node_target <- as.vector(sort(unique(combined_D3obj$Links$Target)))
Node_all <- unique(c(Node_source, Node_target))

#assign IDs for Source (start from 0)
Source <- combined_D3obj$Links$Source
Target <- combined_D3obj$Links$Target

```

```

for(i in 1:length(Node_all)){
  Source[Source==Node_all[i]] <-i-1
  Target[Target==Node_all[i]] <-i-1
}

combined_D3obj$Links$Source <- as.numeric(Source)
combined_D3obj$Links$Target <- as.numeric(Target)
combined_D3obj$Links$LinkColor <- combined$NodeGroup

#prepare node info
node_df <-data.frame(Node=Node_all)
node_df$id <-as.numeric(c(0, 1:(length(Node_all)-1)))

suppressMessages(library(dplyr))
Color <- combined %>% count(Node, color=NodeGroup) %>% select(2)
node_df$color <- Color$color

suppressMessages(library(networkD3))
p1<-sankeyNetwork(Links=combined_D3obj$Links, Nodes = node_df, Value = "Value", NodeGroup ="color", L
fontSize = 22 )
p1

#saveNetwork(p1, file = paste0(path, 'Subpopulation_Net.html'))
##R Setting Information
#sessionInfo()
#rmarkdown::render("/Users/quan.nguyen/Documents/Powell_group_MacQuan/AllCodes/scGPS/vignettes/vignette
#rmarkdown::render("/Users/quan.nguyen/Documents/Powell_group_MacQuan/AllCodes/scGPS/vignettes/vignette

```

### 4.3 Plot the relationship between clusters in two samples

Here we look at one example use case to find relationship between clusters within one sample or between two sample

```

#run prediction for 3 clusters
cluster_mixedpop1 <- colData(mixedpop1)[,1]
cluster_mixedpop2 <- as.numeric(as.vector(colData(mixedpop2)[,1]))
row_cluster <-length(unique(colData(mixedpop2)[,1]))

c_selectID <- 1
genes = DEgenes$DE_Subpop1vsRemaining$id[1:200] #top 200 gene markers distinguishing cluster 1
genes <- gsub("_.*", "", genes)

LSOLDA_dat1 <- bootstrap_scGPS(nboots = 1, mixedpop1 = mixedpop1, mixedpop2 = mixedpop2, genes=genes, c

c_selectID <- 2
genes = DEgenes$DE_Subpop2vsRemaining$id[1:200]
genes <- gsub("_.*", "", genes)
LSOLDA_dat2 <- bootstrap_scGPS(nboots = 1,mixedpop1 = mixedpop1, mixedpop2 = mixedpop2, genes=genes, c_
cluster_mixedpop2 = cluster_mixedpop2)

c_selectID <- 3
genes = DEgenes$DE_Subpop3vsRemaining$id[1:200]

```

```

genes <- gsub("_.*", "", genes)
LSOLDA_dat3 <- bootstrap_scGPS(nboots = 1, mixedpop1 = mixedpop1, mixedpop2 = mixedpop2, genes=genes, c_
  cluster_mixedpop2 = cluster_mixedpop2)

#prepare table input for sankey plot

LASSO_C1S1 <- reformat_LASSO(c_selectID=1, mp_selectID = 1, LSOLDA_dat=LSOLDA_dat1,
  nPredSubpop = row_cluster, Nodes_group = "#7570b3")

LASSO_C2S1 <- reformat_LASSO(c_selectID=2, mp_selectID = 1, LSOLDA_dat=LSOLDA_dat2,
  nPredSubpop = row_cluster, Nodes_group = "#1b9e77")

LASSO_C3S1 <- reformat_LASSO(c_selectID=3, mp_selectID = 1, LSOLDA_dat=LSOLDA_dat3,
  nPredSubpop = row_cluster, Nodes_group = "#e7298a")

combined <- rbind(LASSO_C1S1, LASSO_C2S1, LASSO_C3S1)
combined <- combined[is.na(combined$Value) != TRUE,]
combined_D3obj <- list(Nodes=combined[,4:5], Links=combined[,c(3,2,1)])

library(networkD3)

Node_source <- as.vector(sort(unique(combined_D3obj$Links$Source)))
Node_target <- as.vector(sort(unique(combined_D3obj$Links$Target)))
Node_all <- unique(c(Node_source, Node_target))

#assign IDs for Source (start from 0)
Source <- combined_D3obj$Links$Source
Target <- combined_D3obj$Links$Target

for(i in 1:length(Node_all)){
  Source[Source==Node_all[i]] <- i-1
  Target[Target==Node_all[i]] <- i-1
}

combined_D3obj$Links$Source <- as.numeric(Source)
combined_D3obj$Links$Target <- as.numeric(Target)
combined_D3obj$Links$LinkColor <- combined$NodeGroup

#prepare node info
node_df <- data.frame(Node=Node_all)
node_df$id <- as.numeric(c(0, 1:(length(Node_all)-1)))

suppressMessages(library(dplyr))
Color <- combined %>% count(Node, color=NodeGroup) %>% select(2)

n <- length(unique(node_df$Node))
Color = RColorBrewer::brewer.pal(n, "Set2")

node_df$color <- Color

suppressMessages(library(networkD3))
p1<-sankeyNetwork(Links =combined_D3obj$Links, Nodes = node_df, Value = "Value", NodeGroup ="color", L

```



```
fontSize = 22 )  
p1  
  
#saveNetwork(p1, file = paste0(path, 'Subpopulation_Net.html'))  
##R Setting Information  
#sessionInfo()  
#rmarkdown::render("/Users/quan.nguyen/Documents/Powell_group_MacQuan/AllCodes/scGPS/vignettes/vignette  
#rmarkdown::render("/Users/quan.nguyen/Documents/Powell_group_MacQuan/AllCodes/scGPS/vignettes/vignette
```