NIFA: Non-negative Independent Factor Analysis disentangles discrete and continuous sources of variation in scRNA-seq data

January 30, 2022

Contents

Data Normalization 1

Run NIFA 3

This walkthrough provides a quick start guide for using **NIFA**. In this walkthrough we start from a publicly available scRNA-seq dataset **SimKumar4easy**. This simulated scRNA-seq dataset is available via the bioconductor package *DuoClustering2018*. We first perform data normalization as described in the manuscript, then we apply **NIFA** and compare the the latent factors with the cell labels.

We first load the data SimKumar4easy from the bioconductor package DuoClustering2018.

```
if (!require("DuoClustering2018", character.only = T, warn.conflicts = F,
    quietly = T)) {
    BiocManager::install("DuoClustering2018")
}
library(DuoClustering2018)

library(rsvd)
library(pheatmap)
library(ggpubr)
```

rsvd, pheatmap and ggpubr are loaded only for the purpose of data normalization and visualization. We also load the **NIFA** package into the current working environment.

```
if (!require("NIFA", character.only = T, warn.conflicts = F, quietly = T)) {
    devtools::install_github("wgmao/NIFA")
}
library(NIFA)
```

Data Normalization

We incorporate one additional function normalize_barcode_sums_to_median() credited to https://github.com/hb-gitified/cellrangerRkit/blob/master/R/normalize.r.

```
normalize_barcode_sums_to_median <- function(gbm) {
    bc_sums <- colSums(gbm)
    median_sum <- median(bc_sums)
    return(sweep(gbm, 2, median_sum/bc_sums, "*"))
} #normalize_barcode_sums_to_median</pre>
```

We load the dataset SimKumar4easy and extract the scRNA-seq matrix (gene-by-cell) and cell labels.

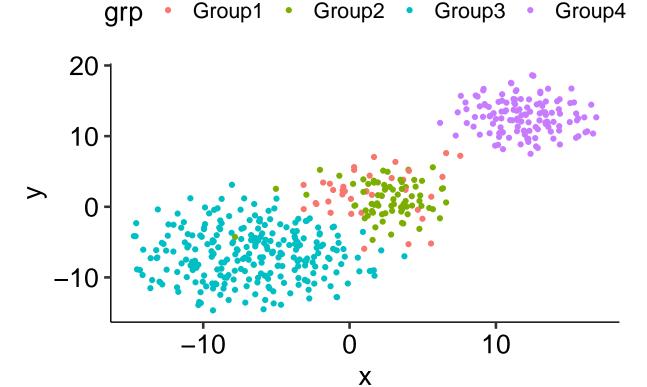
• rnaseq: this is the log transformed scRNA-seq matrix

• celltype: this encodes the cell labels. There are four types of cells in total, which are Group1, Group2, Group3 and Group4.

```
sce <- sce_full_SimKumar4easy(metadata = F)
rnaseq <- sce@assays$data$logcounts
celltype <- sce$Group</pre>
```

We visualize rnaseq in the tSNE space and color dots based on celltype.

```
dat.plot <- data.frame(x = sce@reducedDims$TSNE[, 1], y = sce@reducedDims$TSNE[,
    2], grp = celltype)
ggplot(dat.plot, aes(x = x, y = y, color = grp)) + geom_point() +
    theme_pubr(base_size = 20)</pre>
```



We filter out genes with low-expression levels (missing > 5% of cells)

```
use_genes <- which(apply(rnaseq == 0, 1, sum)/ncol(rnaseq) < 0.95)
rnaseq <- rnaseq[use_genes, ]</pre>
```

We then normalize the total gene expression of each cell to be the median value across all cells using the function normalize_barcode_sums_to_median().

```
rnaseq <- normalize_barcode_sums_to_median(rnaseq)</pre>
```

Finally we perform SVD smoothing by replacing rnaseq with the SVD approximation of the first 50 decomposed principal components.

```
svdres <- rsvd(tscale(rnaseq), k = 50)
rnaseq.norm <- svdres$u %*% diag(svdres$d[1:50]) %*% t(svdres$v)
rownames(rnaseq.norm) <- rownames(rnaseq)
colnames(rnaseq.norm) <- colnames(rnaseq)</pre>
```

Run NIFA

one_hot_encode() function is used to convert celltype with four cell types into one hot encoding (cell-by-celltype). We can easily compute the correlation between decomposed latent factors and one hot encoding to check the correspondence between latent factors and each cell group.

```
one_hot_encode <- function(x) {
    nrow <- length(x)
    cat <- names(table(x))
    ncol <- length(cat)
    res <- matrix(0, nrow = nrow, ncol = ncol)
    for (i in 1:ncol) {
        res[which(x == cat[i]), i] <- 1
    } #for i
    return(res)
} #one_hot_encode</pre>
```

We run NIFA by calling NIFA() function. Here the number of latent factors is set to be 4 and the number of mixture components associated with each latent factor is also set to be 4.

```
NIFA.res <- NIFA(tscale(rnaseq.norm), K = 4, M = 4, max.iter = 500,
    S_threshold = 6e-05, init = "sd", A.init = NULL, S.init = NULL,
    verbose = F, ref = one_hot_encode(celltype), beta_expect_flag = NULL,
    L1.sd = NULL, L2.sd = NULL, b_noise_prior = prod(dim(rnaseq.norm)) *
    5)</pre>
```

We calculate the Pearson correlations between latent factors and one hot encoding of all cell types and we visualize this correspondence using heatmap. There is a one-to-one correspondence between cell types and NIFA latent factors with Pearson correlation > 0.9.

