

# Citrus: a toolkit for single cell sequencing data analysis

Mengjie Chen, Xiang Zhou

April 15, 2016

**Citrus version:** 4.0

Please cite:  
Mengjie Chen, Xiang Zhou.  
Normalization of single cell RNA sequencing data using both control and target genes.

Contents

---

1	Introduction of <i>Citrus</i>	3
2	scPLS	3
2.1	Model and input . . . . .	3
2.2	Removing technical effect using spike-ins . . . . .	4
2.3	Removing cell cycle effect . . . . .	6
2.4	Other algorithms . . . . .	9

## 1 Introduction of *Citrus*

---

*Citrus* is a R/c++ package that is specialized in the analysis of RNA sequencing data from single cell studies. The current release features `scPLS`, a data normalization method that can control for confounding effects by borrowing information from both control and target genes. `scPLS` can be used to remove unwanted variation, such as technical effect or cell cycle effect.

## 2 `scPLS`

---

Like any other genomic sequencing experiment, scRNAseq studies are influenced by many factors that can introduce unwanted variation in the sequencing data and confound the down-stream analysis. Due to low capture efficiency and low amount of input material, such unwanted variation are exacerbated in scRNAseq experiments. As a result, adjusting for confounding factors and normalizing scRNAseq data is crucial for accurate estimation of gene expression levels and successful down-stream analysis. Yet some confounding effects are due to observable batches and can be adjusted for by including batch labels and technician ids as covariates. However, many confounding factors are hidden and are difficult or even impossible to measure. Common hidden confounding factors include various technical artifacts during library preparation and sequencing, and unwanted biological confounders such as cell cycle status. These hidden confounding factors can cause systematic bias, are notoriously difficult control for, and are the focus of the present study.

To mitigate the influence of hidden confounding factors, we develop `scPLS`, a latent factor based model. The genes are divided into two sets: a control set of genes that are used to infer the confounding factors and a target set of genes that are of primary interest. The confounding factors inferred from the control set are used to remove unwanted variation in the target genes for subsequent downstream analysis. For example, most scRNAseq studies add ERCC spike-in controls during the PCR amplification and sequencing steps. The spike-in controls can be used to capture the hidden confounding technical noise associated with the experimental procedures for normalizing genes of primary interest. Similarly, most scRNAseq studies include a set of control genes that are known to have varying expression levels across cell cycles. These cell cycle genes can be used to capture the unmeasured cell cycle status of each cell, which is further used to normalize target genes. Compared to existing methods, our approach taking into account the information from both control and target genes. In addition, our method can model other systematic biological variation and heterogeneity, which are often observed in the target genes. By incorporating such systematic heterogeneity, we can further improve the estimation of the confounding factors and the removal of unwanted variation. We have shown the effectiveness of our approach in our manuscript. This method is implemented in function `scPLS()`.

### 2.1 Model and input

We model the control and target gene sets jointly by the following `scPLS` model

$$\mathbf{x}_i \sim MVN(\mathbf{\Lambda}_x^T \mathbf{z}_i, \psi_x), \quad (1)$$

$$\mathbf{y}_i \sim MVN(\mathbf{\Lambda}_y^T \mathbf{z}_i + \mathbf{\Lambda}_u^T \mathbf{u}_i, \psi_y), \quad (2)$$

where for  $i$ -th individual cell,  $\mathbf{x}_i$  is a  $q$ -vector of expression level for  $q$  control genes;  $\mathbf{y}_i$  is a  $p$ -vector of expression level for  $p$  target genes;  $\mathbf{z}_i$  is a  $k_z$ -vector of confounding factors that influence both control

and target genes ( $k_z < q$ ); the confounding effects are represented by the  $q$  by  $k_z$  loading matrix  $\Lambda_x$  for the controls and the  $p$  by  $k_1$  loading matrix  $\Lambda_y$  for the target genes;  $\mathbf{u}_i$  is a  $k_u$ -vector of biological factors unique to target genes ( $k_u < p$ ); the biological factors could represent some underlying pathways or transcription factors that affect multiple genes;  $\Lambda_u$  is a  $p$  by  $k_u$  loading matrix;  $\psi_x$  is the diagonal covariance matrix for  $q$ -dimensional idiosyncratic error for  $x_i$ ;  $\psi_y$  is the diagonal covariance matrix for  $p$ -dimensional idiosyncratic error for  $y_i$ ; MVN denotes the multivariate normal distribution.

In the above model,  $x_i$  and  $y_i$  represent input data,  $k_1$  and  $k_2$  are prespecified number of latent factors,  $\Lambda_x$ ,  $\Lambda_y$ ,  $\Lambda_u$ ,  $\mathbf{z}_i$  and  $\mathbf{u}_i$  are parameters to be estimated and will be the output of our algorithm.

## 2.2 Removing technical effect using spike-ins

We use the cell cycle data published in Buettner et al as an example to illustrate the usage of `scPLS()` function. We selected top 500 non cell cycle genes, 100 cell cycle genes and 20 spike-ins with largest variance across samples.

```
library(Citrus)
data(CellCycleData)
table(CellCycleData$Label)

##
##      CellCycle NonCellCycle      Spikein
##           100           500           20
```

To remove technical effect, `scPLS()` takes inputs of a  $n \times q$  matrix of target gene expression, a  $n \times p$  matrix of spike-in, a pre-specified number of technical factors  $k_1$  and a pre-specified number of structured biological factors  $k_2$ . The numbers of factors  $k_1$  and  $k_2$  can be determined by model selection criteria later.

```
Control <- CellCycleData$Expression[, CellCycleData$Label ==
  "Spikein"]
Target <- CellCycleData$Expression[, CellCycleData$Label != "Spikein"]
system.time(res <- scPLS(Target, Control, k1 = 2, k2 = 5, iter = 100,
  method = "EM", Chunk = FALSE, center = TRUE))

##      user  system elapsed
##    10.73    4.11    14.94
```

The default algorithm is the “EM-in-Chunks” algorithm, which will randomly divide genes into different chunks, average the estimates from chunks and thus accelerate the computation. The default chunk size is 1000. We recommend to use the “EM-in-Chunks” algorithm only for large gene sets. When `Chunk = FALSE`, the “Naive EM” algorithm will be used. In addition, `scPLS()` will center the expression of each gene to mean of 0 unless setting `center = FALSE`.

```
summary(res)

##              Length Class      Mode
## Factor           3000 -none-   numeric
## Loading           910 -none-   numeric
## Likelihood         1 -none-   numeric
## Z                 364 -none-   numeric
```

```
## lambdaY          1200 -none-      numeric
## lambdaX           40 -none-      numeric
## psi_y            600 -none-      numeric
## Method            1 -none-      character
## Adjusted        109200 -none-     numeric
## VarianceSummary    6 data.frame list
```

`scPLS()` outputs the estimates for confounding factor matrix  $\Lambda_y$  (`lambdaY`), confounding loading matrix  $\mathbf{Z}$  (`Z`), structured factor matrix  $\Lambda_u$  (`Factor`), structured factor loading matrix  $\mathbf{U}$  in `2` (`Loading`). The likelihood for `2` is saved as `Likelihood`, which can be used for model comparison later. The expression matrix after adjusting for confounding effect is stored in `Adjusted`, which is calculated by  $Y - \Lambda_y \mathbf{Z}$ .

```
head(res$VarianceSummary)

##              Sample SampleConfounding SampleStructure
## ENSMUSG00000000127    1.04             0.002         0.088
## ENSMUSG00000000149    1.13             0.014         0.106
## ENSMUSG00000000276    1.19             0.008         0.061
## ENSMUSG00000000308    1.12             0.073         0.074
## ENSMUSG00000000325    1.19             0.003         0.183
## ENSMUSG00000000409    1.14             0.006         0.057
##              Model ModelConfounding ModelStructure
## ENSMUSG00000000127    1.03             0.002         0.080
## ENSMUSG00000000149    1.11             0.014         0.099
## ENSMUSG00000000276    1.18             0.008         0.067
## ENSMUSG00000000308    1.10             0.076         0.070
## ENSMUSG00000000325    1.16             0.003         0.163
## ENSMUSG00000000409    1.13             0.006         0.055
```

`VarianceSummary` stores total sample variance and sample technical variance as well as total variance and technical variance estimated by the fitted model for each gene. We can further quantify the proportion of variance explained by different sources as the following. For each gene, the expression variance is partitioned into three components: a component that is explained by technical factors, a component that is explained by structured biological factors, and the residual error variance. Genes are evenly divided into ten quantiles based on the sample variance.

```
Var <- res$VarianceSummary
ConProp <- round(Var["ModelConfounding"]/Var["Model"], 3)
StructuredProp <- round(Var["ModelStructure"]/Var["Model"], 3)
PropTable <- data.frame(ConProp, StructuredProp, 1 - ConProp -
  StructuredProp)
colnames(PropTable) <- c("Tech", "Structured", "Residue")
QuantileTable <- QuantileSummary(PropTable, quantiles = seq(0.1,
  1, by = 0.1), rankingby = unlist(Var["Sample"]))
head(QuantileTable)

##      Tech Structured Residue Quantile
## 1 0.0287      0.128   0.843        10
## 2 0.0336      0.113   0.853        20
## 3 0.0275      0.100   0.872        30
```

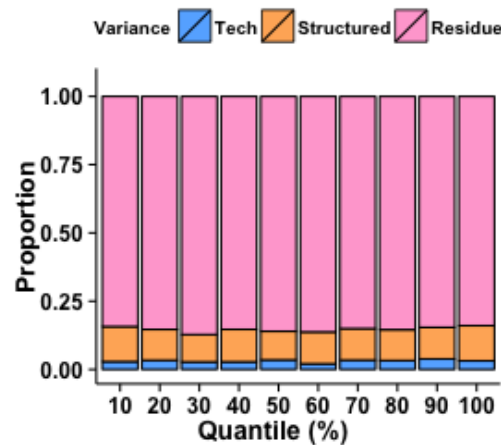


Figure 1: Decomposition of variance into three components. Genes are evenly divided into ten quantiles based on the sample variance.

```
## 4 0.0278      0.119      0.853      40
## 5 0.0350      0.105      0.860      50
## 6 0.0207      0.117      0.863      60
```

We can use barplot to visualize the result as in Fig 1.

```
library(reshape2)
library(ggplot2)
library(easyGgplot2)

df <- melt(QuantileTable, id.vars = "Quantile")
colnames(df)[3] <- "Proportion"
colnames(df)[2] <- "Variance"

ff <- ggplot2.barplot(data = df, xName = "Quantile", yName = "Proportion",
  groupName = "Variance", groupColors = c("#66B2FF", "#FFB266",
    "#FFAAD4"), position = position_stack(), backgroundColor = "white",
  color = "black", xtitle = "Quantile (%)", ytitle = "Proportion",
  mainTitle = "", removePanelGrid = TRUE, removePanelBorder = TRUE,
  axisLine = c(0.5, "solid", "black"), ylim = c(0, 1.05), legendPosition = "top",
  legendTextFont = c(10, "bold", "black"))
```

```
ff
```

## 2.3 Removing cell cycle effect

On the same dataset, we can further remove the cell cycle effect by treating cell cycle genes as control and non cell cycle genes as target.

```

Expression <- res$Adjusted
Label <- CellCycleData$Label[CellCycleData$Label != "Spikein"]
Control <- Expression[, Label == "CellCycle"]
Target <- Expression[, Label == "NonCellCycle"]
system.time(res2 <- scPLS(Target, Control, k1 = 2, k2 = 5, iter = 100,
  method = "EM", Chunk = FALSE, center = TRUE))

##      user      system elapsed
##      6.29       2.33      8.65

```

Similarly we can quantify the proportion of variance explained by cell cycle process.

```

Var <- res2$VarianceSummary
ConProp <- round(Var["ModelConfounding"]/Var["Model"], 3)
StructuredProp <- round(Var["ModelStructure"]/Var["Model"], 3)
PropTable <- data.frame(ConProp, StructuredProp, 1 - ConProp -
  StructuredProp)
colnames(PropTable) <- c("CellCycle", "Structured", "Residue")
median(PropTable$CellCycle)

## [1] 0.0585

```

```

QuantileTable <- QuantileSummary(PropTable, quantiles = seq(0.1,
  1, by = 0.1), rankingby = unlist(Var["Sample"]))
df <- melt(QuantileTable, id.vars = "Quantile")
colnames(df)[3] <- "Proportion"
colnames(df)[2] <- "Variance"

ff <- ggplot2.barplot(data = df, xName = "Quantile", yName = "Proportion",
  groupName = "Variance", groupColors = c("#66B2FF", "#FFB266",
    "#FFAAD4"), position = position_stack(), backgroundColor = "white",
  color = "black", xtitle = "Quantile (%)", ytitle = "Proportion",
  mainTitle = "", removePanelGrid = TRUE, removePanelBorder = TRUE,
  axisLine = c(0.5, "solid", "black"), ylim = c(0, 1.05), legendPosition = "top",
  legendTextFont = c(10, "bold", "black"))

```

```
ff
```

We can compare the expression of non cell cycle genes in original dataset and after removing both technical effect and cell cycle effect.

```

celltypes <- gsub("_[^<>]*", "", rownames(CellCycleData$Expression))
pcs <- prcomp(CellCycleData$Expression[, CellCycleData$Label ==
  "NonCellCycle"], center = TRUE)
PC1 <- pcs$x[, 1]
PC2 <- pcs$x[, 2]
df <- data.frame(PC1, PC2, "raw", celltypes)
colnames(df) <- c("PC1", "PC2", "Method", "Type")

ff1 <- ggplot2.scatterplot(data = df, xName = "PC1", yName = "PC2",
  groupName = "Type", size = 5, backgroundColor = "white",

```

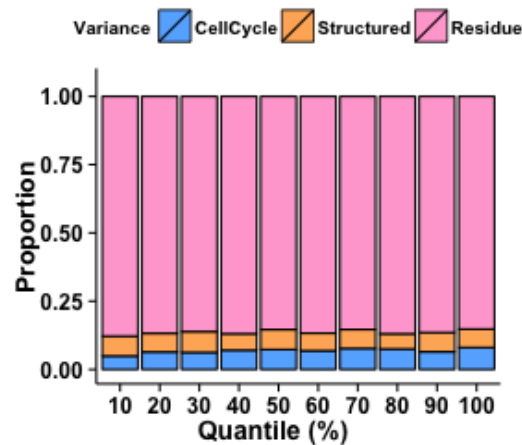


Figure 2: Decomposition of variance into three components. Genes are evenly divided into ten quantiles based on the sample variance.

```
groupColors = c("#66B2FF", "#FFB266", "#FFAAD4"), xtitle = "PC1",
ytitle = "PC2", mainTitle = "Raw", removePanelGrid = TRUE,
removePanelBorder = TRUE, setShapeByGroupName = TRUE, showLegend = FALSE,
mainTitleFont = c(25, "bold", "black"), xtitleFont = c(25,
  "bold", "black"), ytitleFont = c(25, "bold", "black"),
xTickLabelFont = c(25, "bold", "black"), yTickLabelFont = c(25,
  "bold", "black"))
```

```
pcs <- prcomp(res2$Adjusted, center = TRUE)
```

```
PC1 <- pcs$x[, 1]
```

```
PC2 <- pcs$x[, 2]
```

```
df <- data.frame(PC1, PC2, "rmcellcycle", celltypes)
```

```
colnames(df) <- c("PC1", "PC2", "Method", "Type")
```

```
ff2 <- ggplot2.scatterplot(data = df, xName = "PC1", yName = "PC2",
  groupName = "Type", size = 5, backgroundColor = "white",
  groupColors = c("#66B2FF", "#FFB266", "#FFAAD4"), xtitle = "PC1",
  ytitle = "PC2", mainTitle = "Cell cyle effect removed", removePanelGrid = TRUE,
  removePanelBorder = TRUE, setShapeByGroupName = TRUE, showLegend = TRUE,
  mainTitleFont = c(25, "bold", "black"), xtitleFont = c(25,
    "bold", "black"), ytitleFont = c(25, "bold", "black"),
  xTickLabelFont = c(25, "bold", "black"), yTickLabelFont = c(25,
    "bold", "black"), legendTitle = "Cell type", legendTitleFont = c(15,
    "bold", "black"), legendTextFont = c(15, "bold", "black"))
```

```
ff1
```

```
ff2
```



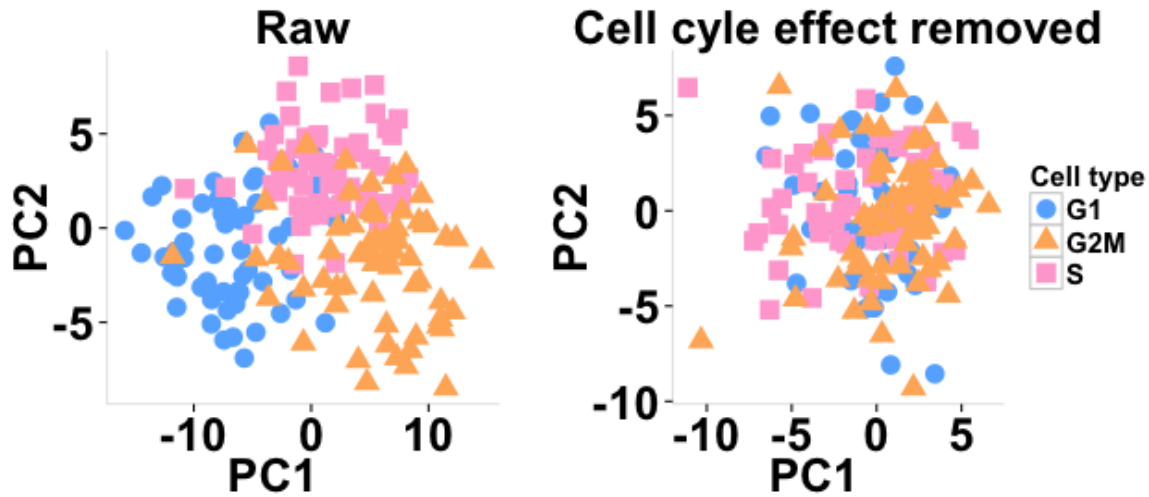


Figure 3: PCA analysis for the uncorrected data and scPLS corrected data. In the uncorrected data, there is a clear separation of cells by cell-cycle stage. Such separation of cells is no longer observed in the corrected data.

## 2.4 Other algorithms

`scPLS()` also provides other alternative algorithms to estimate the confounding factors. Users can select different algorithms by specifying the option `method`. `PCA` algorithm implements the initializer of our EM algorithm. The latent factors are estimated from a Singular Value Decomposition. Details of the algorithm can be found in []. Other choices all impose sparsity assumptions on the structured biological factor matrix to improve interpretability. These include “`EMSparse`” algorithm, with penalty on sparsity of the factor matrix pre-specified, “`EMSparseTraining`” algorithm with penalty on sparsity learned from training samples, “`EMparseNfold`” algorithm with penalty on sparsity learned from N fold cross-validation, and “`IBP`” algorithm with the sparse factor matrix modeled by an Indian Buffet Process prior.