

BUSseq: Batch Effects Correction with Unknown Subtypes for scRNA-seq data

User's Guide

Fangda Song*, Ga Ming Chan and Yingying Wei

The Chinese University of Hong Kong

*sfd1994895@gmail.com

January 20, 2019

Contents

1	Introduction	2
2	Data Preparation	2
3	Model Fitting	4
4	Estimated Cell Types, Batch and Cell-Specific Effects	7
5	Intrinsic Gene Identification	9
6	Corrected Read Count Data and Visualization	10
7	Model Selection using BIC	14

BUSseq: Batch Effects Correction with Unknown Subtypes for scRNA-seq data

User's Guide

1 Introduction

Single-cell RNA-sequencing (scRNA-seq) technologies enable the measurement of the transcriptome of individual cells, which provides unprecedented opportunities to discover cell types and understand cellular heterogeneity [1]. Despite their widespread applications, single-cell RNA-sequencing (scRNA-seq) experiments are still plagued by batch effects and dropout events.

One of the major tasks of scRNA-seq experiments is to identify cell types for a population of cells [1]. Therefore, the cell type of each individual cell is always unknown and is the target of inference. However, most existing methods for batch effects correction, such as Combat space [2] and the surrogate variable analysis (SVA)([3], [4]), are designed for bulk experiments and require knowledge of the subtype information, which corresponds to cell type information for scRNA-seq data, of each sample a priori.

Here, the R package *BUSseq* fits an interpretable Bayesian hierarchical model—the Batch Effects Correction with Unknown Subtypes for scRNA seq Data(BUSseq)—to correct batch effects in the presence of unknown cell types [5]. BUSseq is able to simultaneously correct batch effects, clusters cell types, and takes care of the count data nature, the overdispersion, the dropout events, and the cell-specific sequencing depth of scRNA-seq data. After correcting the batch effects with BUSseq, the corrected value can be used for downstream analysis as if all cells were sequenced in a single batch. BUSseq can integrate the read count matrices measured from different platforms and allow cell types to be measured in some but not all of the batches as long as the experimental design fulfills the conditions listed in [5].

This guide provides step-by-step instructions for applying the BUSseq model to correct batch effects and identify the unknown cell type indicators for each cell for scRNA-seq data.

2 Data Preparation

The input data should be an R list with the length equal to the batch number. Each element of the list is a read count matrix, where each row represents a gene and each column corresponds to a cell. Specifically, assuming there are three batches, the R list consists of three read count matrices with genes in rows and cells in columns. In the following, we provide examples to

BUSseq: Batch Effects Correction with Unknown Subtypes for scRNA-seq data User's Guide

illustrate the input formats, where `BUSseqfits_example` is a sample `BUSseqfits` object. We use the data list stored in this object as an example.

```
library(BUSseq)

#Input data is a list
CountData <- BUSseqfits_example$CountData_raw
class(CountData)

## [1] "list"

#The length of the list is three, so we have three batches
length(CountData)

## [1] 3

#Each element of the list is a matrix
class(CountData[[1]])

## [1] "matrix"

#In the matrix, each row is a gene, and each column corresponds to a cell
dim(CountData[[1]])

## [1] 1000 150

dim(CountData[[2]])

## [1] 1000 150

dim(CountData[[3]])

## [1] 1000 150

#Peek at the read counts
head(CountData[[1]][,1:4])

##      [,1] [,2] [,3] [,4]
## [1,]   68    0   54    5
## [2,]   17   27   83   56
## [3,]   53   74   72   94
## [4,]   40  110   22   10
## [5,]   54   16   74   69
## [6,]    0   97   50   31
```

BUSseq: Batch Effects Correction with Unknown Subtypes for scRNA-seq data

User's Guide

The example data `CountData` consist of three batches. In total, 1,000 genes are measured. The number of cells in each batch is 150, respectively. Because it is a simulated dataset, we actually know that all of the cells come from 4 cell types.

In a nutshell, the user can use the R list with length equal to the batch number as input. Note that the gene numbers of all batches need to be the same.

3 Model Fitting

Once we have prepared the input data and specified the cell type number, we are able to fit the BUSseq model, which requires the function `BUSseq_MCMC`.

The first argument, `Data`, of `BUSseq_MCMC` should be an R list where each element is a data matrix for a specific batch. In the matrix, each row corresponds to a gene or a genomic feature and each column corresponds to a cell.

The second argument, `n.celltypes`, is the number of cell types among all cells, which needs to be specified by the user in advance. As discussed later, if `n.celltypes` is unknown, we can vary the cell type number and use the Bayesian Information Criterion (BIC) to select the optimal number.

The third argument, `n.iterations`, is the total number of iterations of the MCMC algorithm for the posterior inference of the BUSseq model. The user can also set the number of burnin iterations by the argument, `n.burnin`. Given `n.iterations`, the default number of burnins is `n.iterations/2` iterations. The parameters are inferred by samples after the burnin iterations.

The forth argument, `working_dir`, is the user's directory to store the posterior samples. The default directory is the current directory.

The fifth argument, `showIteration`, lets the user decide whether `BUSseq_MCMC` displays the number of iterations that have been run. To obtain reproducible results, we highly recommend to set the argument `seed`.

BUSseq: Batch Effects Correction with Unknown Subtypes for scRNA-seq data User's Guide

```
#Conduct MCMC sampling and posterior inference for BUSseq model
BUSseqfits_res <- BUSseq_MCMC(Data = CountData, n.celltypes = 4,
                             n.iterations = 500, working_dir = ".",
                             showIteration = TRUE, seed = 123)

## Finish the 100-th iterations
## Finish the 200-th iterations
## Finish the 300-th iterations
## Finish the 400-th iterations
## Finish the 500-th iterations

## The MCMC sampling takes: 3.404 mins
## loading the posterior samples...
## calculating posterior means and posterior modes...
## calculating posterior means and posterior takes: 0.057 mins
## calculating BIC...
## calculating BIC takes: 0.045 mins

class(BUSseqfits_res)

## [1] "BUSseqfits"
```

The `summary` command provides an overview of the output object `BUSseqfits_res` from `BUSseq_MCMC`. `BUSseqfits_res` collects two lists of the read count matrices and all posterior estimates of parameters as well as the BIC value. These two lists consist of the raw observed read counts, the inferred underlying true read counts after imputing the dropout events. The posterior estimates contain the cell-type indicators for each cell, the cell-type proportions for each batch, the cell-type-specific mean expression levels, the location batch effects, the overdispersion parameters and the odds ratios of the logistic regression for dropout events.

```
summary(BUSseqfits_res)

## B = 3 batches
## G = 1000 genes
```

BUSseq: Batch Effects Correction with Unknown Subtypes for scRNA-seq data

User's Guide

```
## K = 4 cell types
## N = 450 cells in total
## n.record = 250 iterations are recorded.
## BUSseqfits is an R list that contains the following main elements:
##   BUSseqfits$w.est : the estimated cell type indicators, a list with length
##                      equal to B.
##   BUSseqfits$pi.est : the estimated cell type proportions across batches,
##                      a K by B matrix.
##   BUSseqfits$gamma.est : the estimated the coefficients of the logistic regression
##                      for the dropout events, a B by 2 matrix
##   BUSseqfits$alpha.est : the estimated log-scale baseline expression levels,
##                      a vector with length G.
##   BUSseqfits$beta.est : the estimated cell type effects, a G by K matrix.
##   BUSseqfits$delta.est : the estimated cell-specific effects, a list with
##                      length equal to B.
##   BUSseqfits$nu.est : the estimated location batch effects, a B by G matrix.
##   BUSseqfits$phi.est : the estimated overdispersion parameters, a B by G matrix.
##   BUSseqfits$BIC : the BIC, a scalar.
##   BUSseqfits$PPI.est : the posterior marginal probablity of being an intrinsic
##                      gene, a G by K matrix.
##   For more output values, please use "?BUSseq_MCMC"
##
```

4 Estimated Cell Types, Batch and Cell-Specific Effects

Our main interests are the estimation of the cell type for each cell and the estimation of the batch effects. We can call the `celltypes` function to extract the cell type information from `BUSseqfits_res`.

```
celltypes_est <- celltypes(BUSseqfits_res)

## Batch 1 cells' cell type indicators:  4,4,4...  ...
## Batch 2 cells' cell type indicators:  4,4,4...  ...
## Batch 3 cells' cell type indicators:  4,4,4...  ...
## The output format is a list with length equal to the batch number.
## Each element of the list is a cell type indicator vector in that batch.
```

There is a message from the function `celltypes` to remind the user of the format of `celltypes_est`. In this example, `celltypes_est` is a list of length three, corresponding to the three batches in the study. `celltypes_est[[1]]` shows the cell type for each of the 150 cells in batch one.

Similarly, you can call `location_batch_effects` and `overdispersions` functions to get the estimated location batch effects and batch-specific gene-specific overdispersion parameters. Notice that the first batch is taken as the reference batch, so its location batch effects are zeros for all genes.

```
location_batch_effects_est <- location_batch_effects(
                                BUSseqfits_res)

## The output format is a matrix.
## Each row represents a batch, and each column corresponds to a gene.

head(t(location_batch_effects_est))

##      [,1]      [,2]      [,3]
## [1,]    0 1.883665 3.064765
```

BUSseq: Batch Effects Correction with Unknown Subtypes for scRNA-seq data User's Guide

```
## [2,]    0 1.960276 2.931372
## [3,]    0 2.108560 3.018982
## [4,]    0 2.096778 3.228843
## [5,]    0 1.620678 3.080005
## [6,]    0 1.714777 2.881382

overdispersion_est <- overdispersions(BUSseqfits_res)

## The output format is a matrix.
##
## Each row represents a batch, and each column corresponds to a gene.

head(t(overdispersion_est))

##           [,1]      [,2]      [,3]
## [1,] 0.2594994 0.58109288 0.69905184
## [2,] 1.5180288 9.87009034 4.09068726
## [3,] 0.6085407 0.07304145 0.11348195
## [4,] 0.4305076 6.88936898 4.24084139
## [5,] 1.3748051 0.05600024 0.05938803
## [6,] 0.4493926 2.74961051 3.63539254
```

The cell-specific size effects is available using the `cell_effect_values`. Notice that the first element of each batch is 0 as the first cell in each batch is taken as the reference one.

```
cell_effects_est <- cell_effect_values(BUSseqfits_res)

## The output format is a list with length equal to the batch number.
## Each element of the list is a cell-specific size factor vector of that batch.

head(cell_effects_est[[1]])

## [1] 0.0000000 0.1966175 0.2026120 0.2079615 0.1891741 0.1811446
```

`celltype_mean_expression` function provides the estimated cell-type-specific mean expression levels. The estimates remove the technical artifacts, including the location batch effects and the cell-specific global effects, but retain the biological features for each cell type. Moreover, the cell type effects can be obtained by the `celltype_effects` function. Notice that the first cell type is taken as the baseline cell type implying all zeros in the first column of `celltype_effects_est`.

BUSseq: Batch Effects Correction with Unknown Subtypes for scRNA-seq data User's Guide

```
celltype_mean_expression_est <- celltype_mean_expression(BUSseqfits_example)

## The output format is a matrix.

## Each row represents a gene, and each column corresponds to a cell type.

head(celltype_mean_expression_est)

##           [,1]      [,2]      [,3]      [,4]
## [1,] 6.159967 6.144799 6.179064 45.66401
## [2,] 5.861053 5.931605 5.762307 47.78957
## [3,] 5.235860 5.219132 5.209213 43.81115
## [4,] 5.155947 5.169952 5.089094 40.40547
## [5,] 8.585969 5.607348 6.548128 38.49434
## [6,] 7.512649 7.303040 7.649924 51.10746

celltype_effects_est <- celltype_effects(BUSseqfits_res)

## The output format is a matrix.
##
## Each row represents a gene, and each column corresponds to a cell type.

head(celltype_effects_est)

##           [,1]      [,2]      [,3]      [,4]
## [1,]      0 -0.002465419  0.003095364 2.003239
## [2,]      0  0.011965474 -0.016991307 2.098478
## [3,]      0 -0.003200017 -0.005102163 2.124357
## [4,]      0  0.002712594 -0.013050974 2.058814
## [5,]      0 -0.426051516 -0.270950171 1.500382
## [6,]      0 -0.028297359  0.018107678 1.917343
```

5 Intrinsic Gene Identification

The intrinsic genes are the genes that differentiate cell types [6]. More specifically, a gene is an intrinsic gene if the gene is differentially expressed in at least one cell type compared with the first cell type. We use `intrinsic_genes_BUSseq` to identify the intrinsic genes by controlling the false discovery rate (FDR). In the simulated dataset, we set the first 250 genes as intrinsic genes, whereas the other genes have the same mean expression levels in all cells.

BUSseq: Batch Effects Correction with Unknown Subtypes for scRNA-seq data User's Guide

```
#obtain the intrinsic gene indicators
intrinsic_gene_indices <- intrinsic_genes_BUSseq(
    BUSseqfits_res,
    fdr_threshold = 0.05)

## The output format is a matrix.

## Each row represents a gene, and each column corresponds to a cell type from
2 to K

## 250 intrinsic genes are found.

## The output format is a vector implying the intrinsic gene indices.

#The estimated FDR, the first 250 genes are known intrinsic
#genes in the simulation setting.
false_discovery_ind <- !(intrinsic_gene_indices %in% 1:250)
fdr_est <- sum(false_discovery_ind)/length(intrinsic_gene_indices)
fdr_est

## [1] 0
```

Therefore, the true FDR is 0 less than the estimated FDR, 0.05.

6 Corrected Read Count Data and Visualization

The function `BUSseq_MCMC` not only conducts MCMC sampling and posterior inference, but also imputes the missing data caused by dropout events and corrects batch effects. The function `corrected_read_counts` calculates the corrected read count data of a `BUSseqfits` object. The message is a reminder of the output format, and the output is a `CountData` object. The `summary` command shows the number of batches, genes and cells in each batch.

```
corrected_countdata <- corrected_read_counts(BUSseqfits_res)

## correcting read counts...

## Correcting read counts takes: 0.085 mins

## The output format is a "CountData" object with length equal to the batch number.
```

BUSseq: Batch Effects Correction with Unknown Subtypes for scRNA-seq data

User's Guide

```
## Each element of the object is the corrected read count matrix.  
## In each matrix, each row represents a gene and each column corresponds to a cell.  
class(corrected_countdata)  
## [1] "CountData"  
summary(corrected_countdata)  
## There are 3 batches and 1000 genes.  
## Batch 1 contains 150 cells.  
## Batch 2 contains 150 cells.  
## Batch 3 contains 150 cells.  
##
```

Subsequently, we can compare the raw count data that suffer from batch effects and dropout events, the inferred true expression levels after imputing dropout events, and the corrected count data which remove the batch effects and impute the dropout events. The function `raw_read_counts` and `imputed_read_counts` gives the raw and imputed read count data of a `BUSseqfits` object, respectively. The outputs of them are also a `CountData` object.

```
raw_countdata <- raw_read_counts(BUSseqfits_res)  
## The output format is a "CountData" object with length equal to the batch number.  
## Each element of the object is the raw read count matrix.  
## In each matrix, each row represents a gene and each column corresponds to a cell.  
imputed_countdata <- imputed_read_counts(BUSseqfits_res)  
## The output format is a "CountData" object with length equal to the batch number.  
## Each element of the object is the imputed read count matrix.  
## In each matrix, each row represents a gene and each column corresponds to a cell.
```

The function `heatmap_data_BUSseq` plots the heatmap for the count data across batches for a `CountData` object. The images are saved in the local folder according to the argument `image_dir`. The image name can be modified by the argument `project_name`. The user can

BUSseq: Batch Effects Correction with Unknown Subtypes for scRNA-seq data User's Guide

specify the argument `gene_ind_set`, which is a vector of gene indices, selecting the genes to be displayed in the heatmap.

```
#generate the heatmap of raw read count data
heatmap_data_BUSseq(raw_countdata,
                     project_name = "BUSseq_raw_allgenes",
                     image_dir = "./heatmap")

## pdf
## 2

#display only the first 100 genes in the heatmap
heatmap_data_BUSseq(raw_countdata,
                     gene_set = 1:100,
                     project_name = "BUSseq_raw_100genes",
                     image_dir = "./heatmap")

## pdf
## 2
```

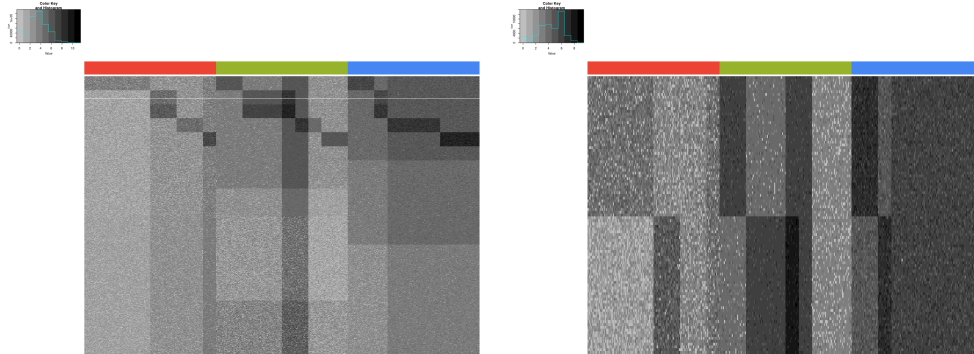


Figure 1: The heatmap of the raw count data of all genes and the first 100 genes

```
#generate the heatmap of imputed read count data
heatmap_data_BUSseq(imputed_countdata,
                     project_name = "BUSseq_imputed_allgenes",
                     image_dir = "./heatmap")

## pdf
## 2
```

BUSseq: Batch Effects Correction with Unknown Subtypes for scRNA-seq data User's Guide

```
#generate the heatmap of corrected read count data
heatmap_data_BUSseq(corrected_countdata,
                     project_name = "BUSseq_corrected_allgenes",
                     image_dir = "./heatmap")

## pdf
## 2
```

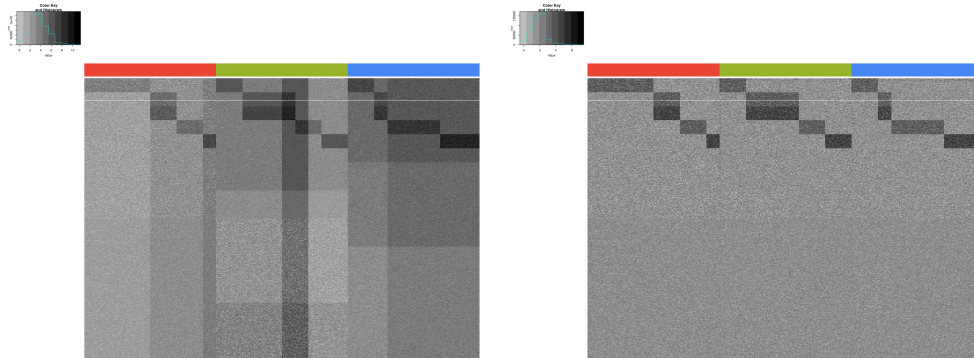


Figure 2: The heatmap for the imputed and corrected count data of all genes

In all these heatmaps, the top bar indicates the batch origin for each cell. Cells under the same color are from the same batch. The batch effects present in the raw data are correctly removed in the corrected count data, and only the biological variabilities are kept. We can also only display the identified intrinsic genes in the corrected count data by setting the argument `gene_set` as `intrinsic_gene_indices`.

```
#Only show the identified intrinsic genes
heatmap_data_BUSseq(corrected_countdata,
                     gene_set = intrinsic_gene_indices,
                     project_name = "BUSseq_corrected_intrinsic_genes",
                     image_dir = "./heatmap")

## pdf
## 2
```

BUSseq: Batch Effects Correction with Unknown Subtypes for scRNA-seq data User's Guide

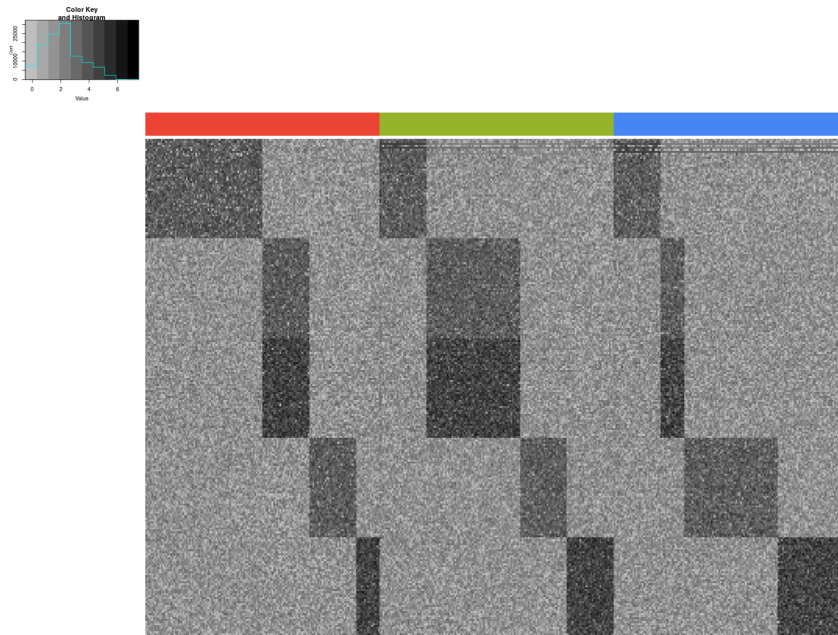


Figure 3: The heatmap for the corrected count data of identified intrinsic genes

7 Model Selection using BIC

If we have no prior knowledge about the cell type number, we can vary the argument `n.celltypes` in the function `BUSseq_MCMC`, e.g., from 2 to 10 and identify the underlying true cell type number K as the one that achieves the minimal BIC.

The user can obtain the BIC value from the `BUSseqfits_res` by the `BIC_BUSseq`.

```
BIC_val <- BIC_BUSseq(BUSseqfits_res)
## BIC is 4118092.49755456
## The output is a scalar.
```

In this example, the underlying true number of cell types is four. For an illustration, we vary the `n.celltypes` from 2 to 6.

BUSseq: Batch Effects Correction with Unknown Subtypes for scRNA-seq data

User's Guide

```
BIC_values <- NULL
for(k in 2:6){
  BUSseqfits_temp <- BUSseq_MCMC(Data = CountData, n.celltypes = k,
                                n.iterations = 500, working_dir = ".",
                                showIteration = FALSE, seed = 123)
  BIC_values <- c(BIC_values, BIC_BUSseq(BUSseqfits_temp))
}

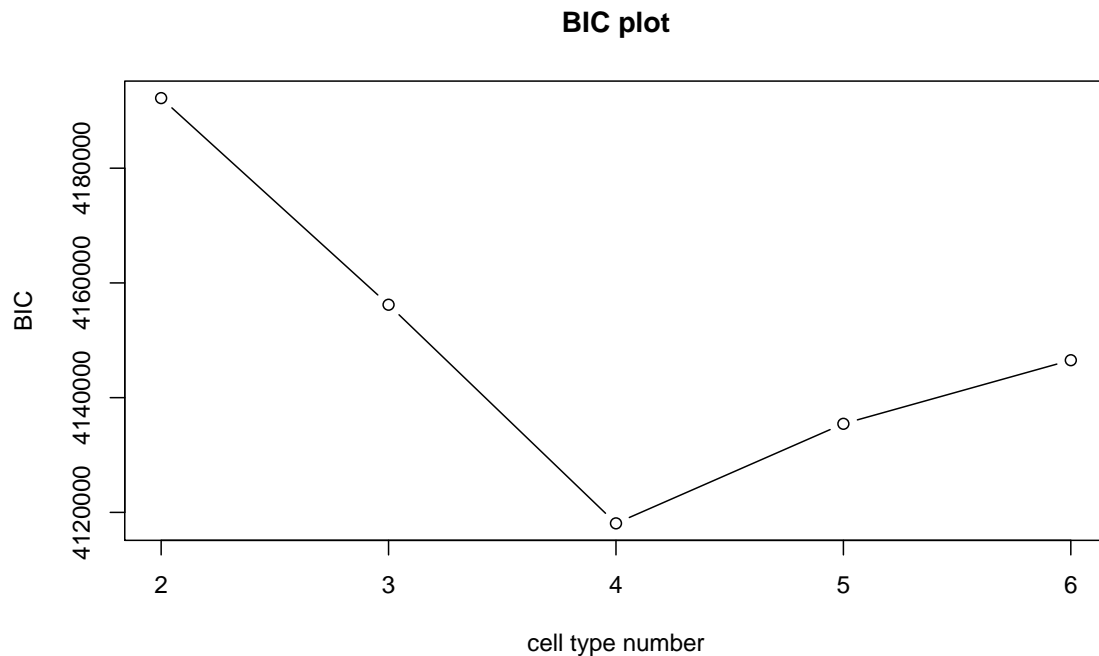
## The MCMC sampling takes: 3.286 mins
## loading the posterior samples...
## calculating posterior means and posterior modes...
## calculating posterior means and posterior takes: 0.045 mins
## calculating BIC...
## calculating BIC takes: 0.024 mins
## BIC is 4192216.51458923
## The output is a scalar.
## The MCMC sampling takes: 3.315 mins
## loading the posterior samples...
## calculating posterior means and posterior modes...
## calculating posterior means and posterior takes: 0.05 mins
## calculating BIC...
## calculating BIC takes: 0.034 mins
## BIC is 4156202.81753306
## The output is a scalar.
## The MCMC sampling takes: 3.309 mins
## loading the posterior samples...
## calculating posterior means and posterior modes...
## calculating posterior means and posterior takes: 0.058 mins
```

BUSseq: Batch Effects Correction with Unknown Subtypes for scRNA-seq data

User's Guide

```
## calculating BIC...
## calculating BIC takes: 0.045 mins
## BIC is 4118092.49755456
## The output is a scalar.
## The MCMC sampling takes: 3.268 mins
## loading the posterior samples...
## calculating posterior means and posterior modes...
## calculating posterior means and posterior takes: 0.062 mins
## calculating BIC...
## calculating BIC takes: 0.054 mins
## BIC is 4135445.65539535
## The output is a scalar.
## The MCMC sampling takes: 3.322 mins
## loading the posterior samples...
## calculating posterior means and posterior modes...
## calculating posterior means and posterior takes: 0.07 mins
## calculating BIC...
## calculating BIC takes: 0.066 mins
## BIC is 4146529.21180375
## The output is a scalar.
plot(2:6, BIC_values, xlab="cell type number", ylab="BIC", main="BIC plot", type="b")
```


BUSseq: Batch Effects Correction with Unknown Subtypes for scRNA-seq data User's Guide



The BIC attains the minimum at `n.celltypes=4`, thus correctly recovering the true cell type number.

References

- [1] Rhonda Bacher and Christina Kendzierski. Design and computational analysis of single-cell rna-sequencing experiments. *Genome Biology*, 17(1):63, 2016.
- [2] W Evan Johnson, Cheng Li, and Ariel Rabinovic. Adjusting batch effects in microarray expression data using empirical bayes methods. *Biostatistics*, 8(1):118–127, 2007.
- [3] Jeffrey T Leek and John D Storey. Capturing heterogeneity in gene expression studies by surrogate variable analysis. *PLoS Genetics*, 3(9):e161, 2007.
- [4] Jeffrey T Leek. svaseq: removing batch effects and other unwanted noise from sequencing data. *Nucleic Acids Research*, page gku864, 2014.
- [5] Fangda Song, Ga Ming Chan, and Yingying Wei. Flexible experimental designs for valid single-cell rna-sequencing experiments allowing batch effects correction. *Manuscript*.

BUSseq: Batch Effects Correction with Unknown Subtypes for scRNA-seq data User's Guide

- [6] Zhiguang Huo, Ying Ding, Silvia Liu, Steffi Oesterreich, and George Tseng. Meta-analytic framework for sparse k-means to identify disease subtypes in multiple transcriptomic studies. *Journal of the American Statistical Association*, 111(513):27–42, 2016.