# broadSeq Report

Rishi Das Roy, Rintu Kutum                    Gautam Ahuja, Santosh Adhikari

**Note**: All code was compiled using `R-4.3.3`, `Bioconductor version 3.18` and `Rtools43`. The operating system used were Windows 11 Version 23H2 and Ubuntu-22.04

We successfully installed the library on `R-4.3.3`. I worked on a Windows-based system, while Santosh worked on an Ubuntu system. Below is a detailed report of each error encountered during the installation and execution of the `broadSeq` package.

# 1 Installation - Windows 11

## 1.1 BiocManager

In windows, we started by installing the `BiocManager` package as follows:

```
if (!require("BiocManager")) {
  install.packages("BiocManager")
}
```

This was installed successfully without any errors.

## 1.2 devtools

Next, we installed the `devtools` package:

```
install.packages('devtools')
```

This was installed successfully without any errors.

## 1.3 broadSeq

Before installing `broadSeq`, we first had to install `rtools43` on Windows. We followed the steps mentioned here for Windows. We then attempted to install the development version of `broadSeq` with `vignettes` as follows:

```
devtools::install_github("dasroy/broadSeq", build_vignettes = TRUE)
```

However, this resulted in an error.
To resolve this, we first installed the `BiocStyle` package from `BiocManager`:

```
BiocManager::install('BiocStyle')
```

After this, we were able to install and browse the `vignettes`.

We then installed `broadSeq` from `BiocManager` as follows:

```
1  BiocManager::install("broadSeq")
```

This was installed successfully without any errors.

## 1.4  Other Libraries

We installed the `ggplot2` and `ggpubr` packages as mentioned in the `vignettes`. We are now ready to run the `broadSeq` library.

```
1  install.packages('ggplot2')
2  install.packages('ggpubr')
```

## 1.5  Vignettes

We also `knitted` the vignette seperately by downloading the `.Rmd` files and the images. It was compiled without error.

# 2  Installation - Ubuntu-22.04

## 2.1  BiocManager

In Linux system with Ubuntu- 22.04, the installation for the `BiocManager` package is as follows:

```
1  if (!require("BiocManager")) {
2    install.packages("BiocManager")
3  }
```

This was installed successfully without any errors.

## 2.2  devtools

Next, we proceeded to install the `devtools` package:

```
1  install.packages('devtools')
```

This was installed successfully without any errors.

## 2.3  broadSeq

While we attempted to install `broadSeq`, We attempted to install the development version of `broadSeq` with `vignettes` as follows:

```
1  devtools::install_github("dasroy/broadSeq", build_vignettes = TRUE)
```

However, this resulted in an error stating that BiocStyle library is missing.
To resolve this, we first installed the `BiocStyle` package from `BiocManager`:

```
1  BiocManager::install('BiocStyle')
```

Again while we attempted to install `broadSeq`, we got the error message stating that
`ERROR: dependency 'clusterProfiler' is not available for package 'broadSeq'`
along with the following `warning message`.

```
1   Warning messages:
2   1: In i.p(...) : installation of package 'ggtree' had non-zero exit
    ↪  status
3   2: In i.p(...) :
4     installation of package 'enrichplot' had non-zero exit status
5   3: In i.p(...) :
6     installation of package \\
7     'clusterProfiler' had non-zero exit status
8   4: In i.p(...) :
9     installation of package \\
10    '/tmp/RtmpvfE7g9/file164647e8fbdac/broadSeq_0.99.3.tar.gz' had
    ↪  non-zero exit status
```

We used the following command to resolve the encountered errors:

```
1   BiocManager::install("clusterProfiler")
2   BiocManager::install("ggtree")
```

After this, we were able to install and browse the `vignettes`.

We then installed `broadSeq` from `BiocManager` as follows:

```
1   BiocManager::install("broadSeq")
```

Before the final installation was over we also got the following warning message:
`Bioconductor version 3.18 (BiocManager 1.30.23), R 4.3.3 (2024-02-29)`
`Warning message:`
`package 'broadSeq' is not available for Bioconductor version '3.18'`
This was then installed successfully without any errors.

## 2.4  Other Libraries

We installed the `ggplot2` and `ggpubr` packages as mentioned in the `vignettes`. We are now ready to run the `broadSeq` library.

```
1   install.packages('ggplot2')
2   install.packages('ggpubr')
```

# 3 Testing

The broadSeq package simplifies the process of including many Bioconductor packages for RNA-seq data and evaluating their performance.

We followed the instructions in the vignettes to verify that we could run all the commands correctly.

## 3.1 Imports

We import both the broadSeq and ggplot2 libraries.

```
1  library(broadSeq)
2  library(ggplot2)
```

On loading the broadSeq package we get the following mesage:

```
1   Loading required package: dplyr
2
3   Attaching package: 'dplyr'
4
5   The following objects are masked from 'package:stats':
6
7       filter, lag
8
9   ...
10
11  The following objects are masked from 'package:matrixStats':
12
13      anyMissing, rowMedians
```

There were no errors encountered.

## 3.2 Reading the data

The broadSeq takes a SummarizedEnperiment as an input format and gives the output as a data.frame. We load the gene expression data in .rds format as shown:

```
1  se <- readRDS(system.file("extdata","rat_vole_mouseSE_salmon.rds",
   ↪  package = "broadSeq"))
2  SummarizedExperiment::assayNames(se)
```

The SummarizedExperiment output was as follows:

```
1  [1] "counts"       "abundance"    "avgTxLength"    "vst"
```

### 3.2.1 Sample Metadata

We then explored the sample metadata:

```
as.data.frame(colData(se)) %>%
  dplyr::count(stage, species) %>%
  tidyr::spread(stage, n)
se$stage <- factor(se$stage, levels = c("Bud", "Cap", "Late Cap",
→   "Bell"))
```

There were no errors encountered and we got expected results.

### 3.2.2 Filtering Low Expression Genes

We filtered out low expression genes and visualized the distribution to help in improving the signal-to-noise ratio. The first step involves visualizing the distribution of gene expression counts followed by applying a filter to remove genes with low expression levels.

```
assays(se)[["counts"]][,5] %>%
  ggpubr::ggdensity(y = "count") +
  ggplot2::geom_vline(xintercept = 10) +
  ggplot2::scale_x_log10()

keep <- (assays(se)[["counts"]] >= 3) %>% rowSums() >= 5
table(keep)
```

There were no errors encountered and we got expected results.

## 3.3 Normalization

Normalization was performed using CPM and TMM methods:

```
# CPM
se <- broadSeq::normalizeEdgerCPM(se, method = "none", cpm.log = TRUE)
SummarizedExperiment::assayNames(se)
# TMM
se <- broadSeq::normalizeEdgerCPM(se, method = "TMM", cpm.log = FALSE)
SummarizedExperiment::assayNames(se)
```

There were no errors encountered and we got expected results. Accessing the normalized data:

```
assays(se)[["counts"]][1:5, 1:5]
assays(se)[["TMM"]][1:5, 1:5]
assays(se)[["logCPM"]][1:5, 1:5]
```

There were no errors encountered and we got expected results.

## 3.4   Transformation

We applied various transformations to the data:

```
# VST
se <- broadSeq::transformDESeq2(se, method = "vst")

# Normalized counts transformation
se <- broadSeq::transformDESeq2(se, method = "normTransform")

# rlog
se <- broadSeq::transformDESeq2(se, method = "rlog")
SummarizedExperiment::assayNames(se)
```

There were no errors encountered and we got expected results.

### 3.4.1   Comparison of Transformations

We compared the transformations visually:

```
p <- broadSeq::sampleAssay_plot(se[, se$species=="Mouse" ],
                                assayName = "counts", fill = "stage",
                                yscale = "log2")+ rremove("x.text")
p1 <- broadSeq::sampleAssay_plot(se[, se$species=="Mouse"],
                                  assayName = "vst", fill = "stage")+
                                  ↪ rremove("x.text")
p2 <- broadSeq::sampleAssay_plot(se[, se$species=="Mouse"],
                                  assayName = "TMM", fill = "stage",
                                  yscale = "log10")+ rremove("x.text")
p3 <- broadSeq::sampleAssay_plot(se[, se$species=="Mouse"],
                                  assayName = "logCPM", fill = "stage")+
                                  ↪ rremove("x.text")
ggarrange(p,p1,p2,p3, common.legend = TRUE, labels = c("A","B","C"))
```

There were no errors encountered and we got expected results.
If the vsn package is available, we perform mean-variance plots:

```
if (requireNamespace("vsn", quietly = TRUE)) {
  library("vsn")
  x <- meanSdPlot(log2(assays(se[, se$species == "Rat"])[["counts"]] +
    ↪ 1), plot = FALSE)
  print(x$gg + ggtitle(label = "log2(n+1)"))
  x <- meanSdPlot(assays(se[, se$species == "Rat"])[["vst"]], plot =
    ↪ FALSE)
  print(x$gg + ggtitle(label = "Vst"))
  x <- meanSdPlot(assays(se[, se$species == "Rat"])[["logCPM"]], plot =
    ↪ FALSE)
  print(x$gg + ggtitle(label = "logCPM"))
}
```

There were no errors encountered and we got expected results.

## 3.5   Visualization of Gene Expression

Visualizations of gene expression were generated using various assays:

```r
broadSeq::assay_plot(se, feature = c("Shh"),
                     assayNames = c("counts","logCPM","vst","TMM"),
                     x = "stage", fill="species", add="dotplot", palette
                     ↪  = "npg")
broadSeq::genes_plot(se,
                     features = c("Shh","Edar"),
                     facet.by = "symbol",
                     x = "stage", assayName = "vst", fill="species",
                     ↪  palette = "jco")
```

Using predefined or custom color palettes based on journals:

```r
jco <- broadSeq::genes_plot(se[,se$species == "Mouse"],
                            features = c("Shh"), facet.by = "symbol",
                            ↪  assayName =  "logCPM",
                            x = "stage",  fill="stage", add="dotplot",
                            ↪  xlab = "",
                            title = "Journal of Clinical Oncology",
                            ↪  palette = "jco")

npg <- broadSeq::genes_plot(se[,se$species == "Mouse"],
                            features = c("Shh"), facet.by =
                            ↪  "symbol",assayName =  "logCPM",
                            x = "stage", fill="stage", add="dotplot",
                            ↪  xlab = "",
                            title = "Nature Publishing Group", palette =
                            ↪  "npg")

aaas <- broadSeq::genes_plot(se[,se$species == "Mouse"],
                             features = c("Shh"), facet.by = "symbol",
                             ↪  assayName = "logCPM",
                             x = "stage", fill="stage", add="dotplot",
                             ↪  xlab = "",
                             title = "Science", palette = "aaas")

nejm <- broadSeq::genes_plot(se[,se$species == "Mouse"],
                             features = c("Shh"), facet.by = "symbol",
                             ↪  assayName = "logCPM",
                             x = "stage", fill="stage", add="dotplot",
                             ↪  xlab = "",
                             title = "New England Journal of
                             ↪  Medicine",palette = "nejm")

ggarrange(jco+ggpubr::rotate_x_text(), npg+ggpubr::rotate_x_text(),
          aaas+ggpubr::rotate_x_text(),nejm+ggpubr::rotate_x_text(),
```

```
23          nrow = 1, common.legend = TRUE,legend = "none",
24          labels = c("A","B","C","D")) %>%
25     annotate_figure( top = text_grob("Color palette"))
```

There were no errors encountered and we got expected results.

## 3.6  Quality Control with Clustering

Quality control was performed using various clustering methods:

```
1   # MDS plot
2   broadSeq::plot_MDS(se, scaledAssay = "vst", ntop=500,
3                      color = "species", shape = "stage",
4                      ellipse=TRUE, legend = "bottom")
5   head(rowData(se))
6   # Hierarchical clustering
7   p_vst <- broadSeq::plotHeatmapCluster(
8       se,
9       scaledAssay = "vst",
10      annotation_col = c("species", "stage"),
11      annotation_row = c("Class","gene_biotype"),
12      ntop = 30, show_geneAs = "symbol",
13      cluster_cols = TRUE, cluster_rows = FALSE,
14      show_rownames = TRUE, show_colnames = FALSE,
15      main = "Top 30 variable gene vst"
16  )
```

### 3.6.1  PCA plot

```
1   # prcompTidy
2   computedPCA_logCPM <- broadSeq::prcompTidy(se, scaledAssay = "logCPM",
    ↪  ntop = 500)
3   computedPCA_vst <- broadSeq::prcompTidy(se, scaledAssay = "vst", ntop =
    ↪  500)
4   # Plot
5   # logCPM
6   plotAnyPC(computedPCA = computedPCA_logCPM,
7             x = 1, y = 2, color = "species", shape = "stage",
8             legend = "bottom")
9   # VST
10  pca_vst <- plotAnyPC(computedPCA = computedPCA_vst,
11                       x = 2, y = 3,  color = "species", shape = "stage",
12                       legend = "bottom")
13  # Other PCs
14  computedPCA_vst$eigen_values %>%
15    dplyr::filter(var_exp >= 2) %>%
16    ggbarplot(x="PC",y="var_exp", label = TRUE, label.pos = "out")
17  pca_vst_2_3 <-plotAnyPC(computedPCA = computedPCA_vst,
```

```
18                                  x = 2, y = 3,
19                                  color = "species", shape = "stage", legend =
       ↪    "bottom")
20  # Gene loading
21  computedPCA_vst %>% broadSeq::getFeatureLoadRanking(keep =
    ↪  c("symbol","Class")) %>% head()
22  computedPCA_vst$loadings %>% top_n(5,abs(PC2)  ) %>%
    ↪   dplyr::select(gene,PC2)
23  pca_vst_loading <- computedPCA_vst %>%
24    broadSeq::getFeatureLoadRanking(keep = c("symbol","Class"), topN = 50,
      ↪  pcs=1:10) %>%
25    dplyr::count(Class, PC) %>%
26    ggbarplot(
27      x = "PC", y = "n", fill = "Class",
28      legend = "bottom", palette =
        ↪   c("red","blue","orange","purple","white","grey")
29    )
30  # Biplot
31  pca_vst_bi <- broadSeq::biplotAnyPC(computedPCA = computedPCA_vst,
32                                      x = 1, y = 2, genesLabel = "symbol",
33                                      color = "species", shape = "stage",
34                                      legend = "bottom")
35  ggarrange(
36    ggarrange(pca_vst_bi+ggtitle(label =  ""),
37              pca_vst_2_3+ggtitle(label =  ""), common.legend = TRUE),
38    pca_vst_loading, nrow = 2)
39  # User defined genes
40  biplotAnyPC(computedPCA = computedPCA_vst,x = 2, y = 3,
41              color = "species", shape = "stage",
42              genes= computedPCA_vst$loadings %>%
43                top_n(5,abs(PC3)) %>% pull(gene),
44              genesLabel = "symbol")
45
46  ## Plot progression gene "Shh"
47  biplotAnyPC(computedPCA = computedPCA_vst,x = 2, y = 3,
48              color = "species", shape = "stage",
49              genes=c("Shh"),
50              genesLabel = "symbol")
```

There were no errors encountered and we got expected results.

## 3.7   Compare Differential Expression

We load the data again

```
1  se <- readRDS(system.file("extdata","rat_vole_mouseSE_salmon.rds",
   ↪  package = "broadSeq"))
2  se <- se[,colData(se)$species == "Mouse"]
```

```
3   # Gene information
4   head(rownames(se))
5   head(rowData(se))
6   # Sample information
7   head(colData(se))
8   table(colData(se)$stage)
```

### 3.7.1 Differential Expression

```
1   # Function pattern
2   result_Noiseq <-
3     use_NOIseq(se = se,
4               colData_id = "stage", control = "Bud", treatment = "Cap",
5               rank = TRUE,
6               r = 10)
7   head(result_Noiseq)
8   pg <- broadSeq::genes_plot(se, x = "stage", assayName =  "counts",
9                             features = result_Noiseq %>%
                              ↪  dplyr::filter(rank <5) %>% rownames(),
10                            fill="stage", facet.by = "symbol",
11                            palette="jco", add =
                              ↪  "dotplot")+rotate_x_text()
12  pg_sc <- ggscatter(result_Noiseq, x="Bud_mean", y="Cap_mean",color =
    ↪  "prob")+
13    scale_x_log10()+scale_y_log10()
14  pg+pg_sc
15  # Available methods
16  ?use_limma_trend(se, colData_id, control, treatment, rank = FALSE, ...)
    ↪  # limma
17  ?use_limma_voom(se, colData_id, control, treatment, rank = FALSE, ...)
18  ?use_edgeR_exact(se, colData_id, control, treatment, rank = FALSE, ...)
    ↪  # edgeR
19  ?use_edgeR_GLM(se, colData_id, control, treatment, rank = FALSE, ...)
20  ?use_deseq2(se, colData_id, control, treatment, rank = FALSE, ...) #
    ↪  deseq2
21  ?use_DELocal(se, colData_id, control, treatment, rank = FALSE, ...) #
    ↪  DELocal
22  ?use_NOIseq(se, colData_id, control, treatment, rank = FALSE, ...)  #
    ↪  noiseq
23  ?use_EBSeq(se, colData_id, control, treatment, rank = FALSE, ...) #
    ↪  EBSeq
24  ?use_SAMseq(se, colData_id, control, treatment, rank = FALSE, ...) #
    ↪  samseq
```

There were no errors encountered and we got expected results.

### 3.7.2 Compare DE results

```
1   funs <- list(limma_trend = use_limma_trend, limma_voom = use_limma_voom,
2                edgeR_exact = use_edgeR_exact, edgeR_glm = use_edgeR_GLM,
3                deseq2 = use_deseq2,
4                DELocal = use_DELocal, noiseq = use_NOIseq,
5                EBSeq = use_EBSeq)
6   multi_result <- broadSeq::use_multDE(
7     se = se,
8     deFun_list = funs, return.df = TRUE,
9     colData_id = "stage", control = "Bud", treatment = "Cap",
10    rank = TRUE)
11  head(multi_result)
12  colnames(multi_result)
13  # Similarity of methods
14  clusters <- multi_result %>% dplyr::select(ends_with("rank")) %>% t()
    ↪  %>% dist() %>% hclust()
15  plot(clusters,main =  "distance: Euclidean")
16  # Plots
17  # Volcano
18  multi_result %>% broadSeq::volcanoPlot(
19    pValName = "deseq2_padj",
20    lFCName = "deseq2_log2FoldChange",
21    labelName = "symbol",
22    palette = "lancet" ,
23    selectedLabel =
24      multi_result %>% dplyr::arrange(deseq2_padj) %>% pull(symbol) %>%
        ↪  head()
25  )
26  multi_result %>% broadSeq::volcanoPlot(
27    pValName = "deseq2_padj",
28    lFCName = "deseq2_log2FoldChange",
29    labelName = "symbol",
30    palette = c("purple","orange","grey"),
31    selectedLabel = list(criteria = "(`x` > 5 | `x` < -2) & (`y` > 10)")
32  ) +xlim(-7.5,7.5)
```

There were no errors encountered and we got expected results.

When ploting, there was a reoccuring warning at multiple steps:

```
1   Warning message:
2   In seq_len(computedPCA$pc_scores[, dottedArg$shape]) :
3     first element used of 'length.out' argument
```