

FRASER: Find RAre Splicing Events in RNA-seq

Christian Mertes¹, Ines Scheller¹, Julien Gagneur¹

¹ Technische Universität München, Department of Informatics, Garching, Germany

June 9, 2020

Abstract

Genetic variants affecting splicing are a major cause of rare diseases yet their identification remains challenging. Recently, detecting splicing defects by RNA sequencing (RNA-seq) has proven to be an effective complementary avenue to genomic variant interpretation. However, no specialized method exists for the detection of aberrant splicing events in RNA-seq data. Here, we addressed this issue by developing the statistical method *FRASER* (Find RAre Splicing Events in RNA-seq). *FRASER* detects splice sites de novo, assesses both alternative splicing and intron retention, automatically controls for latent confounders using a denoising autoencoder, and provides significance estimates using an over-dispersed count fraction distribution. *FRASER* outperforms state-of-the-art approaches on simulated data and on enrichments for rare near-splice site variants in 48 tissues of the GTEx dataset. Application to a previously analysed rare disease dataset led to a new diagnostic by reprioritizing an aberrant exon truncation in TAZ. Altogether, we foresee *FRASER* as an important tool for RNA-seq based diagnostics of rare diseases.

If you use *FRASER* in published research, please cite:

Mertes C, Scheller I, Yopez V, *et al.* **Detection of aberrant splicing events in RNA-seq data with FRASER**, biorXiv, 2019,
<https://doi.org/10.1101/2019.12.18.866830>

Package

FRASER 1.0.1

Contents

1	Introduction	3
2	Quick guide to <i>FRASER</i>	4
3	A detailed <i>FRASER</i> analysis	6
3.1	Data preparation	7
3.1.1	Creating a <i>FraserDataSet</i> and Counting reads	7
3.1.2	Creating a <i>FraserDataSet</i> from existing count matrices	10
3.2	Data preprocessing and QC	12
3.2.1	Filtering	12
3.2.2	Sample co-variation	15
3.3	Detection of aberrant splicing events.	16
3.3.1	Fitting the splicing model	16
3.3.2	Calling splicing outliers	18
3.3.3	Interpreting the result table	18
3.4	Finding splicing candidates in patients	19
3.5	Saving and loading a <i>FraserDataSet</i>	21
4	More details on <i>FRASER</i>	22
4.1	Correction for confounders	22
4.1.1	Finding the dimension of the latent space	23
4.2	P-value calculation	24
4.3	Z-score calculation	26
4.4	Result visualization	26
	References	29
5	Session Info	30

1 Introduction

FRASER (Find RAre Splicing Evens in RNA-seq) is a tool for finding aberrant splicing events in RNA-seq samples. It works on the splice metrics ψ_5 , ψ_3 and θ to be able to detect any type of aberrant splicing event from exon skipping over alternative donor usage to intron retention. To detect these aberrant events, **FRASER** uses a similar approach as the **OUTRIDER** package that aims to find aberrantly expressed genes and makes use of an autoencoder to automatically control for confounders within the data. **FRASER** also uses this autoencoder approach and models the read count ratios in the ψ values by fitting a beta binomial model to the ψ values obtained from RNA-seq read counts and correcting for apparent co-variations across samples. Similarly as in **OUTRIDER**, read counts that significantly deviate from the distribution are detected as outliers. A scheme of this approach is given in Figure 1.

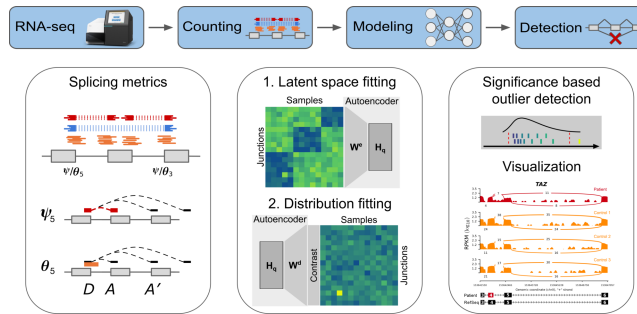


Figure 1: The **FRASER splicing outlier detection workflow**

The workflow starts with RNA-seq aligned reads and performs splicing outlier detection in three steps. First (left column), a splice site map is generated in an annotation-free fashion based on RNA-seq split reads. Split reads supporting exon-exon junctions as well as non-split reads overlapping splice sites are counted. Splicing metrics quantifying alternative acceptors (ψ_5), alternative donors (ψ_3) and splicing efficiencies at donors (θ_5) and acceptors (θ_3) are computed. Second (middle column), a statistical model is fitted for each splicing metric that controls for sample covariations (latent space fitting using a denoising autoencoder) and overdispersed count ratios (beta-binomial distribution). Third (right column), outliers are detected as data points significantly deviating from the fitted models. Candidates are then visualized with a genome browser.

FRASER uses the following splicing metrics as described by Pervouchine et al[1]: we compute for each sample, for donor D (5' splice site) and acceptor A (3' splice site) the ψ_5 and ψ_3 values, respectively, as:

$$\psi_5(D, A) = \frac{n(D, A)}{\sum_{A'} n(D, A')} \quad 1$$

and

$$\psi_3(D, A) = \frac{n(D, A)}{\sum_{D'} n(D', A)} \quad 2$$

FRASER: Find RAre Splicing Events in RNA-seq

where $n(D, A)$ denotes the number of split reads spanning the intron between donor D and acceptor A and the summands in the denominators are computed over all acceptors found to splice with the donor of interest (Equation 1), and all donors found to splice with the acceptor of interest (Equation 2). To not only detect alternative splicing but also partial or full intron retention, we also consider θ as a splicing efficiency metric.

$$\theta_5(D) = \frac{\sum_{A'} n(D, A')}{n(D) + \sum_{A'} n(D, A')} \quad 3$$

and

$$\theta_3(A) = \frac{\sum_{D'} n(D', A)}{n(A) + \sum_{D'} n(D', A)}, \quad 4$$

where $n(D)$ is the number of non-split reads spanning exon-intron boundary of donor D , and $n(A)$ is defined as the number of non-split reads spanning the intron-exon boundary of acceptor A . While we calculate θ for the 5' and 3' splice site separately, we do not distinguish later in the modeling step between θ_5 and θ_3 and hence call it jointly θ in the following.

2 Quick guide to FRASER

Here we quickly show how to do an analysis with *FRASER*, starting from a sample annotation table and the corresponding bam files. First, we create an *FraserDataSet* from the sample annotation and count the relevant reads in the bam files. Then, we compute the ψ/θ values and filter out introns that are just noise. Secondly, we run the full pipeline using the command *FRASER*. In the last step, we extract the results table from the *FraserDataSet* using the *results* function. Additionally, the user can create several analysis plots directly from the fitted *FraserDataSet* object. These plotting functions are described in section 4.4.

```
# load FRASER library
library(FRASER)

# count data
fds <- createTestFraserSettings()
fds <- countRNAData(fds)
fds

## ----- Sample data table -----
## # A tibble: 3 x 5
##   sampleID bamFile condition gene SeqLevelStyle
##   <chr>    <chr>      <int> <chr> <chr>
## 1 sample1 /tmp/Rtmpxhq3B/Rinst366f27e8c763/F~ 1 TIMM~ UCSC
## 2 sample2 /tmp/Rtmpxhq3B/Rinst366f27e8c763/F~ 3 CLPP UCSC
```

FRASER: Find RAre Splicing Events in RNA-seq

```
## 3 sample3 /tmp/Rtmpxhq3B/Rinst366f27e8c763/F~ 2 MCOL~ UCSC
##
## Number of samples:      3
## Number of junctions:    60
## Number of splice sites: 38
## assays(2): rawCountsJ rawCountsSS
##
## ----- Settings -----
## Analysis name:          Data Analysis
## Analysis is strand specific: no
## Working directory:      '/tmp/RtmpXd6AIZ'
##
## ----- BAM parameters -----
## Default used with: bamMapqFilter=0

# compute stats
fds <- calculatePSIValues(fds)

# filtering junction with low expression
fds <- filterExpressionAndVariability(fds, minExpressionInOneSample=20,
  minDeltaPsi=0.0, filter=TRUE)

# use PCA to speed up the tutorial
fds <- FRASER(fds, q=2, implementation="PCA")

##      Min.   1st Qu.   Median     Mean   3rd Qu.     Max.
## 5.509e-08 5.509e-08 5.509e-08 5.509e-08 5.509e-08 5.509e-08
##      Min.   1st Qu.   Median     Mean   3rd Qu.     Max.
## 5.509e-08 5.509e-08 5.509e-08 5.509e-08 5.509e-08 5.509e-08
##      Min.   1st Qu.   Median     Mean   3rd Qu.     Max.
## 5.509e-08 5.509e-08 5.509e-08 5.509e-08 5.509e-08 5.509e-08

# we provide two ways to annotate introns with the corresponding gene symbols:
# the first way uses TxDb-objects provided by the user as shown here
require(TxDb.Hsapiens.UCSC.hg19.knownGene)
txdb <- TxDb.Hsapiens.UCSC.hg19.knownGene
require(org.Hs.eg.db)
orgDb <- org.Hs.eg.db
fds <- annotateRangesWithTxDb(fds, txdb=txdb, orgDb=orgDb)

# alternatively, we also provide a way to use biomaRt for the annotation:
# fds <- annotateRanges(fds)

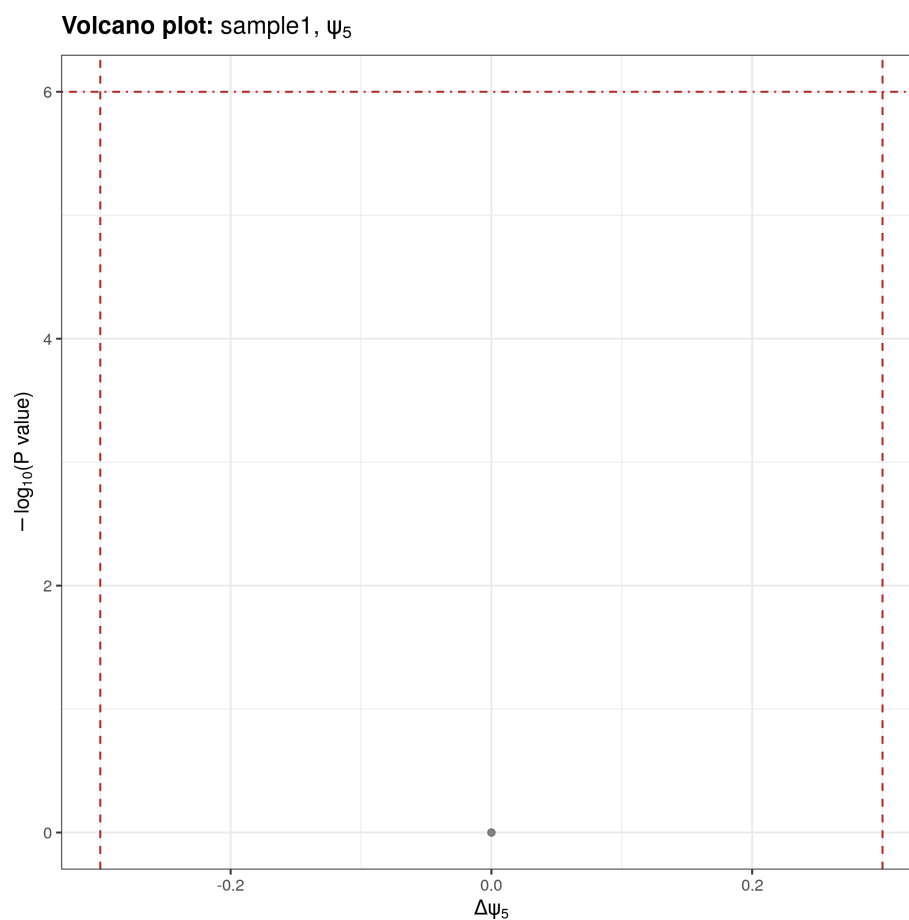
# get results: usually, using p-value cutoff is recommended, but on the small
# example dataset, this would give no results, so we use z-scores here instead
# res <- results(fds, zScoreCutoff=NA, padjCutoff=0.05, deltaPsiCutoff=0.1)
```

FRASER: Find RAre Splicing Events in RNA-seq

```
res <- results(fds, zScoreCutoff=1, padjCutoff=NA, deltaPsiCutoff=0.02)
res

## GRanges object with 0 ranges and 0 metadata columns:
##   seqnames   ranges strand
##   <Rle> <IRanges> <Rle>
##   -----
##   seqinfo: 84 sequences from an unspecified genome

# result visualization
plotVolcano(fds, sampleID="sample1", type="psi5", aggregate=TRUE)
```



3 A detailed *FRASER* analysis

The analysis workflow of *FRASER* for detecting rare aberrant splicing events in RNA-seq data can be divided into the following steps:

1. Data import or Counting reads [3.1](#)

FRASER: Find RAre Splicing Events in RNA-seq

2. Data preprocessing and QC [3.2](#)
3. Correcting for confounders [4.1](#)
4. Calculate P-values [4.2](#)
5. Calculate Z-scores [4.3](#)
6. Visualize the results [4.4](#)

Step 3-5 are wrapped up in one function `FRASER`, but each step can be called individually and parametrized. Either way, data preprocessing should be done before starting the analysis, so that samples failing quality measurements or introns stemming from background noise are discarded.

Detailed explanations of each step are given in the following subsections.

For this tutorial we will use the a small example dataset that is contained in the package.

3.1 Data preparation

3.1.1 Creating a *FraserDataSet* and Counting reads

To start a RNA-seq data analysis with `FRASER` some preparation steps are needed. The first step is the creation of a *FraserDataSet* which derives from a `RangedSummarizedExperiment` object. To create the *FraserDataSet*, sample annotation and two count matrices are needed: one containing counts for the splice junctions, i.e. the split read counts, and one containing the splice site counts, i.e. the counts of non split reads overlapping with the splice sites present in the splice junctions.

You can first create the *FraserDataSet* with only the sample annotation and subsequently count the reads as described in [3.1.1](#). For this, we need a table with basic informations which then can be transformed into a *FraserSettings* object. The minimum of information per sample is an unique sample name, the path to the aligned bam file. Additionally groups can be specified for the P-value calculations later. If a **NA** is assigned no P-values will be calculated. An example sample table is given within the package:

```
sampleTable <- fread(system.file(
  "extdata", "sampleTable.tsv", package="FRASER", mustWork=TRUE))
head(sampleTable)

##   sampleID          bamFile group  gene
## 1: sample1 extdata/bam/sample1.bam    1 TIMMDC1
## 2: sample2 extdata/bam/sample2.bam    3  CLPP
## 3: sample3 extdata/bam/sample3.bam    2 MCOLN1
```

FRASER: Find RAre Splicing Events in RNA-seq

To create a settings object for *FRASER* the constructor `FraserSettings` should be called with at least a `sampleData` table. For an example have a look into the `createTestFraserSettings`. In addition to the `sampleData` you can specify further parameters.

1. The parallel backend (a *BiocParallelParam* object)
2. The read filtering (a *ScanBamParam* object)
3. An output folder for the resulting figures and the cache
4. If the data is strand specific or not

The following shows how to create a example *FraserDataSet* with only the settings options from the sample annotation above:

```
# convert it to a bamFile list
bamFiles <- system.file(sampleTable[,bamFile], package="FRASER", mustWork=TRUE)
sampleTable[,bamFile:=bamFiles]

# create FRASER object
settings <- FraserDataSet(colData=sampleTable,
                          workingDir=tempdir())

# show the FraserSettings object
settings

## ----- Sample data table -----
## # A tibble: 3 x 4
##   sampleID bamFile                                group gene
##   <chr>    <chr>                                <int> <chr>
## 1 sample1 /tmp/Rtmpxhxq3B/Rinst366f27e8c763/FRASER/extdata/bam~    1 TIMMD~
## 2 sample2 /tmp/Rtmpxhxq3B/Rinst366f27e8c763/FRASER/extdata/bam~    3 CLPP
## 3 sample3 /tmp/Rtmpxhxq3B/Rinst366f27e8c763/FRASER/extdata/bam~    2 MCOLN1
##
## ----- Settings -----
## Analysis name:           Data Analysis
## Analysis is strand specific: no
## Working directory:       '/tmp/RtmpXd6AIZ'
##
## ----- BAM parameters -----
## Default used with: bamMapqFilter=0
```

The *FraserDataSet* for this example data can also be generated through the function `createTestFraserSettings`:

```
settings <- createTestFraserSettings()
settings
```


FRASER: Find RAre Splicing Events in RNA-seq

```
## ----- Sample data table -----
## # A tibble: 3 x 4
##   sampleID bamFile          condition gene
##   <chr>    <chr>          <int> <chr>
## 1 sample1 /tmp/Rtmpxhq3B/Rinst366f27e8c763/FRASER/extdata~ 1 TIMMD~
## 2 sample2 /tmp/Rtmpxhq3B/Rinst366f27e8c763/FRASER/extdata~ 3 CLPP
## 3 sample3 /tmp/Rtmpxhq3B/Rinst366f27e8c763/FRASER/extdata~ 2 MCOLN1
##
## ----- Settings -----
## Analysis name:          Data Analysis
## Analysis is strand specific: no
## Working directory:      '/tmp/RtmpXd6AIZ'
##
## ----- BAM parameters -----
## Default used with: bamMapqFilter=0
```

Counting of the reads are straight forward and is done through the `countRNAData` function. The only required parameter is the `FraserSettings` object. First all split reads are extracted from each individual sample and cached if enabled. Then a dataset wide junction map is created (all visible junctions over all samples). After that for each sample the non-spliced reads at each given donor and acceptor site is counted. The resulting *FraserDataSet* object contains two *SummarizedExperiment* objects for each the junctions and the splice sites.

```
# example of how to use parallelization: use 10 cores or the maximal number of
# available cores if fewer than 10 are available and use Snow if on Windows
if(.Platform$OS.type == "unix") {
  register(MulticoreParam(workers=min(10, multicoreWorkers())))
} else {
  register(SnowParam(workers=min(10, multicoreWorkers())))
}

# count reads
fds <- countRNAData(settings)
fds

## ----- Sample data table -----
## # A tibble: 3 x 4
##   sampleID bamFile          condition gene
##   <chr>    <chr>          <int> <chr>
## 1 sample1 /tmp/Rtmpxhq3B/Rinst366f27e8c763/FRASER/extdata~ 1 TIMMD~
## 2 sample2 /tmp/Rtmpxhq3B/Rinst366f27e8c763/FRASER/extdata~ 3 CLPP
## 3 sample3 /tmp/Rtmpxhq3B/Rinst366f27e8c763/FRASER/extdata~ 2 MCOLN1
##
## Number of samples:      3
```

FRASER: Find RAre Splicing Events in RNA-seq

```
## Number of junctions:    60
## Number of splice sites: 38
## assays(2): rawCountsJ rawCountsSS
##
## ----- Settings -----
## Analysis name:           Data Analysis
## Analysis is strand specific: no
## Working directory:       '/tmp/RtmpXd6AIZ'
##
## ----- BAM parameters -----
## Default used with: bamMapqFilter=0
```

3.1.2 Creating a *FraserDataSet* from existing count matrices

If the count matrices already exist, you can use these matrices directly together with the sample annotation from above to create the *FraserDataSet*:

```
# example sample annoation for precalculated count matrices
sampleTable <- fread(system.file("extdata", "sampleTable_countTable.tsv",
                                package="FRASER", mustWork=TRUE))

head(sampleTable)

##   sampleID          bamFile group  gene
## 1: sample1 extdata/bam/sample1.bam    1 TIMMDC1
## 2: sample2 extdata/bam/sample2.bam    1 TIMMDC1
## 3: sample3 extdata/bam/sample3.bam    2 MCOLN1
## 4: sample4 extdata/bam/sample4.bam    3  CLPP
## 5: sample5 extdata/bam/sample5.bam   NA  NHDF
## 6: sample6 extdata/bam/sample6.bam   NA  NHDF

# get raw counts
junctionCts <- fread(system.file("extdata", "raw_junction_counts.tsv.gz",
                                package="FRASER", mustWork=TRUE))

head(junctionCts)

##   seqnames  start    end  width strand sample1 sample2 sample3 sample4
## 1: chr19 7126380 7690902 564523    *      0      1      0      0
## 2: chr19 7413458 7615986 202529    *      0      1      0      0
## 3: chr19 7436801 7703913 267113    *      0      0      0      0
## 4: chr19 7466307 7607189 140883    *      0      0      0      0
## 5: chr19 7471938 7607808 135871    *      1      0      0      0
## 6: chr19 7479042 7625600 146559    *      0      0      0      0
##   sample5 sample6 sample7 sample8 sample9 sample10 sample11 sample12
## 1:      0      0      0      0      0      0      0      0
## 2:      0      0      0      0      0      0      0      0
```

FRASER: Find RAre Splicing Events in RNA-seq

```
## 3:      0      1      0      0      0      0      0      0
## 4:      0      0      1      0      0      0      0      0
## 5:      0      0      0      0      0      0      0      0
## 6:      0      0      0      0      0      0      0      1
##      startID endID
## 1:      1     90
## 2:      2     91
## 3:      3     92
## 4:      4     93
## 5:      5     94
## 6:      6     95

spliceSiteCts <- fread(system.file("extdata", "raw_site_counts.tsv.gz",
  package="FRASER", mustWork=TRUE))
head(spliceSiteCts)

##      seqnames      start      end width strand spliceSiteID  type sample1 sample2
## 1:   chr19 7126379 7126380      2      *           1 Donor      0      0
## 2:   chr19 7413457 7413458      2      *           2 Donor      0      0
## 3:   chr19 7436800 7436801      2      *           3 Donor      0      0
## 4:   chr19 7466306 7466307      2      *           4 Donor      0      0
## 5:   chr19 7471937 7471938      2      *           5 Donor      0      0
## 6:   chr19 7479041 7479042      2      *           6 Donor      0      0
##      sample3 sample4 sample5 sample6 sample7 sample8 sample9 sample10 sample11
## 1:      0      0      0      0      0      0      0      0      0
## 2:      0      0      0      0      0      0      0      0      0
## 3:      0      0      0      0      0      0      0      0      0
## 4:      0      0      0      0      0      0      0      0      0
## 5:      0      0      0      0      0      0      0      0      0
## 6:      0      0      0      0      0      0      0      0      0
##      sample12
## 1:      0
## 2:      0
## 3:      0
## 4:      0
## 5:      0
## 6:      0

# create FRASER object
fds <- FraserDataSet(colData=sampleTable, junctions=junctionCts,
  spliceSites=spliceSiteCts, workingDir=tempdir())
fds

## ----- Sample data table -----
## # A tibble: 12 x 4
##   sampleID bamFile      group gene
##   <chr>    <chr>      <int> <chr>
```

FRASER: Find RAre Splicing Events in RNA-seq

```
## 1 sample1 extdata/bam/sample1.bam 1 TIMMDC1
## 2 sample2 extdata/bam/sample2.bam 1 TIMMDC1
## 3 sample3 extdata/bam/sample3.bam 2 MCOLN1
## 4 sample4 extdata/bam/sample4.bam 3 CLPP
## 5 sample5 extdata/bam/sample5.bam NA NHDF
## 6 sample6 extdata/bam/sample6.bam NA NHDF
## 7 sample7 extdata/bam/sample7.bam NA NHDF
## 8 sample8 extdata/bam/sample8.bam NA NHDF
## 9 sample9 extdata/bam/sample9.bam NA NHDF
## 10 sample10 extdata/bam/sample10.bam NA NHDF
## 11 sample11 extdata/bam/sample11.bam NA NHDF
## 12 sample12 extdata/bam/sample12.bam NA NHDF
##
## Number of samples: 12
## Number of junctions: 123
## Number of splice sites: 165
## assays(2): rawCountsJ rawCountsSS
##
## ----- Settings -----
## Analysis name: Data Analysis
## Analysis is strand specific: no
## Working directory: '/tmp/RtmpXd6AIZ'
##
## ----- BAM parameters -----
## Default used with: bamMapqFilter=0
```

3.2 Data preprocessing and QC

As with gene expression analysis, a good quality control of the raw data is crucial. For some hints please refer to our workshop slides¹.

¹<http://tinyurl.com/RNA-ASHG-presentation>

At the time of writing this vignette, we recommend that the RNA-seq data should be aligned with a splice-aware aligner like STAR^[2] or GEM^[3]. To gain better results, at least 20 samples should be sequenced and they should be processed with the same protocol and origin from the same tissue.

3.2.1 Filtering

Before we can filter the data, we have to compute the main splicing metric: the ψ -value (Percent Spliced In).

```
fds <- calculatePSIValues(fds)
fds
```

FRASER: Find RAre Splicing Events in RNA-seq

```
## ----- Sample data table -----
## # A tibble: 12 x 4
##   sampleID bamFile          group gene
##   <chr>    <chr>          <int> <chr>
## 1 sample1  extdata/bam/sample1.bam      1 TIMMDC1
## 2 sample2  extdata/bam/sample2.bam      1 TIMMDC1
## 3 sample3  extdata/bam/sample3.bam      2 MCOLN1
## 4 sample4  extdata/bam/sample4.bam      3 CLPP
## 5 sample5  extdata/bam/sample5.bam     NA NHDF
## 6 sample6  extdata/bam/sample6.bam     NA NHDF
## 7 sample7  extdata/bam/sample7.bam     NA NHDF
## 8 sample8  extdata/bam/sample8.bam     NA NHDF
## 9 sample9  extdata/bam/sample9.bam     NA NHDF
## 10 sample10 extdata/bam/sample10.bam    NA NHDF
## 11 sample11 extdata/bam/sample11.bam  NA NHDF
## 12 sample12 extdata/bam/sample12.bam  NA NHDF
##
## Number of samples:      12
## Number of junctions:    123
## Number of splice sites: 165
## assays(11): rawCountsJ psi5 ... rawOtherCounts_psiSite delta_psiSite
##
## ----- Settings -----
## Analysis name:          Data Analysis
## Analysis is strand specific: no
## Working directory:      '/tmp/RtmpXd6AIZ'
##
## ----- BAM parameters -----
## Default used with: bamMapqFilter=0
```

Now we can have some cut-offs to filter down the number of junctions we want to test later on.

Currently, we keep only junctions which support the following:

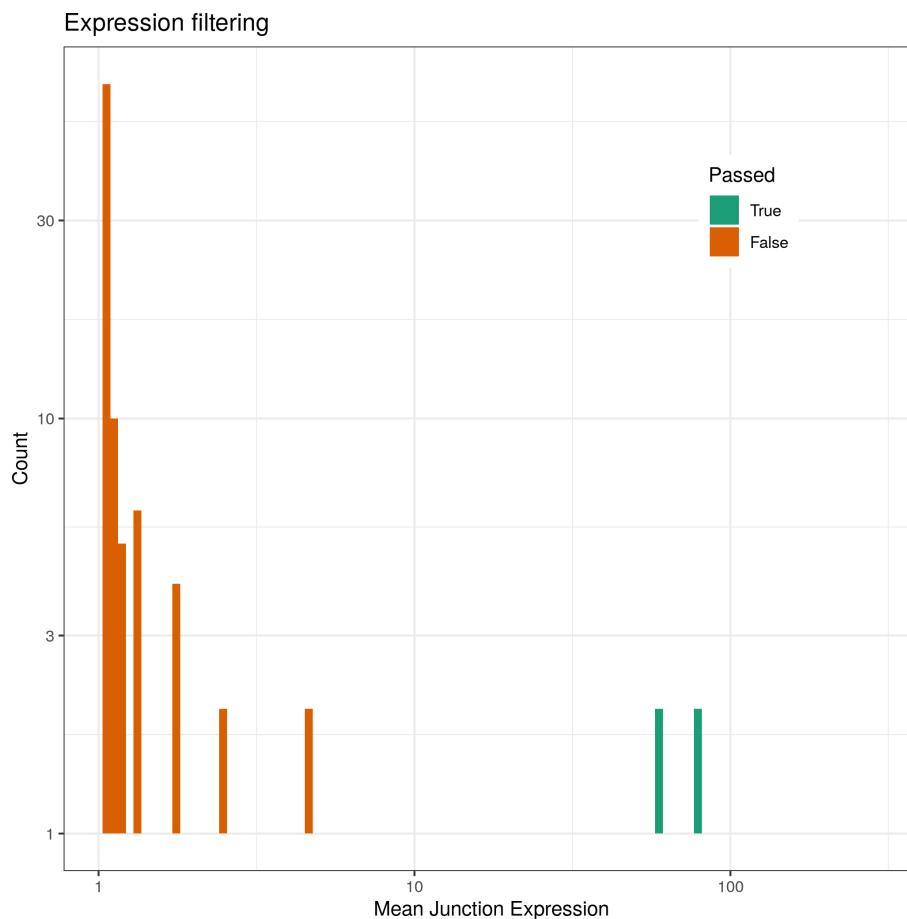
- At least one sample has 20 reads
- 5% of the samples have at least 1 read

Furthermore one could filter for:

- At least one sample has a $|\Delta\psi|$ of 0.1

```
fds <- filterExpressionAndVariability(fds, minDeltaPsi=0.0, filter=FALSE)
plotFilterExpression(fds, bins=100)
```

FRASER: Find RAre Splicing Events in RNA-seq



After looking at the expression distribution between filtered and unfiltered junctions, we can now subset the dataset:

```
fds_filtered <- fds[mcols(fds, type="j"), "passed", ]
fds_filtered

## ----- Sample data table -----
## # A tibble: 12 x 4
##   sampleID bamFile          group gene
##   <chr>    <chr>          <int> <chr>
## 1 sample1  extdata/bam/sample1.bam     1 TIMMDC1
## 2 sample2  extdata/bam/sample2.bam     1 TIMMDC1
## 3 sample3  extdata/bam/sample3.bam     2 MCOLN1
## 4 sample4  extdata/bam/sample4.bam     3 CLPP
## 5 sample5  extdata/bam/sample5.bam    NA NHDF
## 6 sample6  extdata/bam/sample6.bam    NA NHDF
## 7 sample7  extdata/bam/sample7.bam    NA NHDF
## 8 sample8  extdata/bam/sample8.bam    NA NHDF
## 9 sample9  extdata/bam/sample9.bam    NA NHDF
## 10 sample10 extdata/bam/sample10.bam  NA NHDF
```

FRASER: Find RAre Splicing Events in RNA-seq

```
## 11 sample11 extdata/bam/sample11.bam    NA NHDF
## 12 sample12 extdata/bam/sample12.bam    NA NHDF
##
## Number of samples:      12
## Number of junctions:    20
## Number of splice sites: 38
## assays(11): rawCountsJ psi5 ... rawOtherCounts_psiSite delta_psiSite
##
## ----- Settings -----
## Analysis name:          Data Analysis
## Analysis is strand specific: no
## Working directory:      '/tmp/RtmpXd6AIZ'
##
## ----- BAM parameters -----
## Default used with: bamMapqFilter=0

# filtered_fds not further used for this tutorial because the example dataset
# is otherwise too small
```

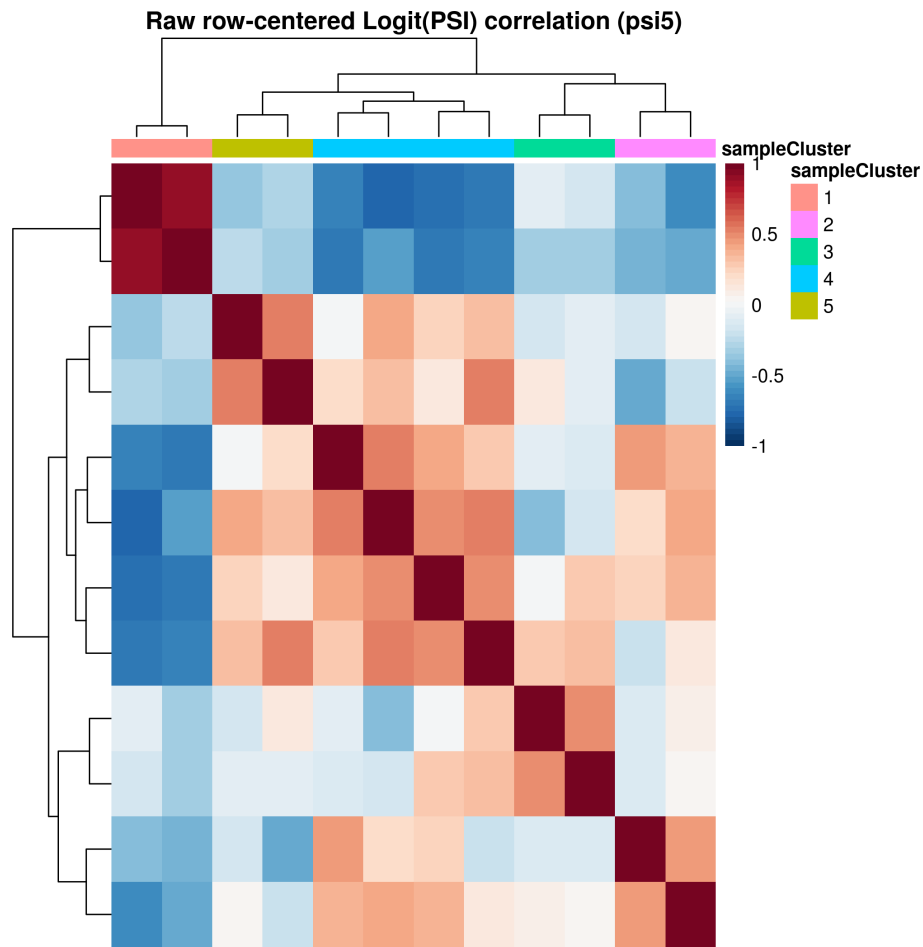
3.2.2 Sample co-variation

Since ψ values are ratios within a sample, one might think that there should not be as much correlation structure as observed in gene expression data within the splicing data.

This is not true as we do see strong sample co-variation across different tissues and cohorts. Let's have a look into our data to see if we do have correlation structure or not. To have a better estimate, we use the logit transformed ψ values to compute the correlation.

```
# Heatmap of the sample correlation
plotCountCorHeatmap(fds, type="psi5", logit=TRUE, normalized=FALSE)
```

FRASER: Find RAre Splicing Events in RNA-seq



It is also possible to visualize the correlation structure of the logit transformed ψ values of the *topJ* most variable introns for all samples:

```
# Heatmap of the intron/sample expression  
plotCountCorHeatmap(fds, type="psi5", logit=TRUE, normalized=FALSE,  
  plotType="junctionSample", topJ=100, minDeltaPsi = 0.01)
```

3.3 Detection of aberrant splicing events

After preprocessing the raw data and visualizing it, we can start our analysis. Let's start with the first step in the aberrant splicing detection: the model fitting.

3.3.1 Fitting the splicing model

During the fitting procedure, we will normalize the data and correct for confounding effects by using a denoising autoencoder. Here we use a predefined latent space with a dimension $q = 10$. Using the correct dimension is crucial to have the best

FRASER: Find RAre Splicing Events in RNA-seq

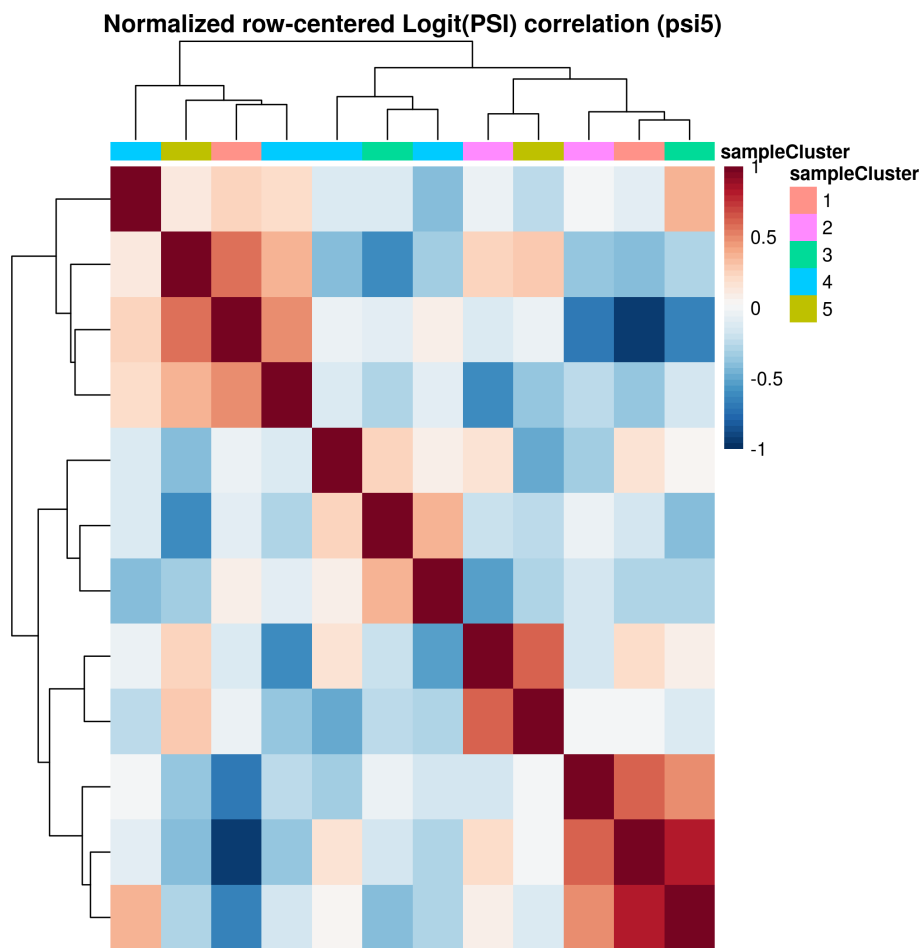
performance (see 4.1.1). Alternatively, one can also use a PCA to correct the data. The wrapper function `FRASER` both fits the model and calculates the p-values and z-scores for all ψ types. For more details see section 4.

```
# This is computational heavy on real size datasets and can take awhile
fds <- FRASER(fds, q=3)

##      Min.   1st Qu.   Median     Mean   3rd Qu.    Max.
## 5.509e-08 5.509e-08 8.842e-08 7.171e-07 8.842e-08 7.922e-05
##      Min.   1st Qu.   Median     Mean   3rd Qu.    Max.
## 5.510e-08 5.510e-08 5.510e-08 2.939e-06 8.840e-08 3.528e-04
##      Min.   1st Qu.   Median     Mean   3rd Qu.    Max.
## 5.500e-08 5.500e-08 8.800e-08 1.795e-04 8.800e-08 5.741e-03
```

To check whether the correction worked, we can have a look at the correlation heatmap using the normalized ψ values from the fit.

```
plotCountCorHeatmap(fds, type="psi5", normalized=TRUE, logit=TRUE)
```



FRASER: Find RAre Splicing Events in RNA-seq

3.3.2 Calling splicing outliers

Before we extract the results, we should add the human readable HGNC symbols. *FRASER* comes already with an annotation function. The function uses *biomaRt* in the background to overlap the genomic ranges with the known HGNC symbols. To have more flexibility on the annotation, one can also provide a custom 'txdb' object to annotate the HGNC symbols.

Here we assume a beta binomial distribution and call outliers based on the significance level. The user can choose between a p value cutoff, a Z score cutoff or a cutoff on the $\Delta\psi$ values between the observed and expected ψ values or both.

```
# annotate introns with the HGNC symbols of the corresponding gene
require(TxDb.Hsapiens.UCSC.hg19.knownGene)
txdb <- TxDb.Hsapiens.UCSC.hg19.knownGene
require(org.Hs.eg.db)
orgDb <- org.Hs.eg.db
fds <- annotateRangesWithTxDb(fds, txdb=txdb, orgDb=orgDb)
# fds <- annotateRanges(fds) # alternative way using biomaRt

# retrieve results with default and recommended cutoffs (padj <= 0.05 and
# |deltaPsi| >= 0.3)
res <- results(fds)
```

3.3.3 Interpreting the result table

The function *results* retrieves significant events based on the specified cutoffs as a *GRanges* object which contains the genomic location of the splice junction or splice site that was found as aberrant and the following additional information:

- sampleID: the sampleID in which this aberrant event occurred
- hgncSymbol: the gene symbol of the gene that contains the splice junction or site if available
- type: the metric for which the aberrant event was detected (either psi5 for ψ_5 , psi3 for ψ_3 or psiSite for θ)
- pValue, padjust, zScore: the p-value, adjusted p-value and z-score of this event
- psiValue: the value of ψ_5 , ψ_3 or θ metric (depending on the type column) of this junction or splice site for the sample in which it is detected as aberrant
- deltaPsi: the $\Delta\psi$ -value of the event in this sample, which is the difference between the actual observed ψ and the expected ψ
- meanCounts: the mean count (k) of reads mapping to this splice junction or site over all samples

FRASER: Find RAre Splicing Events in RNA-seq

- meanTotalCounts: the mean total count (n) of reads mapping to the same donor or acceptor site as this junction or site over all samples
- counts, totalCounts: the count (k) and total count (n) of the splice junction or site for the sample where it is detected as aberrant

Please refer to section 1 for more information about the metrics ψ_5 , ψ_3 and θ and their definition. In general, an aberrant ψ_5 value might indicate aberrant acceptor site usage of the junction where the event is detected; an aberrant ψ_3 value might indicate aberrant donor site usage of the junction where the event is detected; and an aberrant θ value might indicate partial or full intron retention, or exon truncation or elongation. We recommend using a genome browser to investigate interesting detected events in more detail.

```
# to show result visualization functions for this tutorial, zScore cutoff used
res <- results(fds, zScoreCutoff=2, padjCutoff=NA, deltaPsiCutoff=0.1)
res

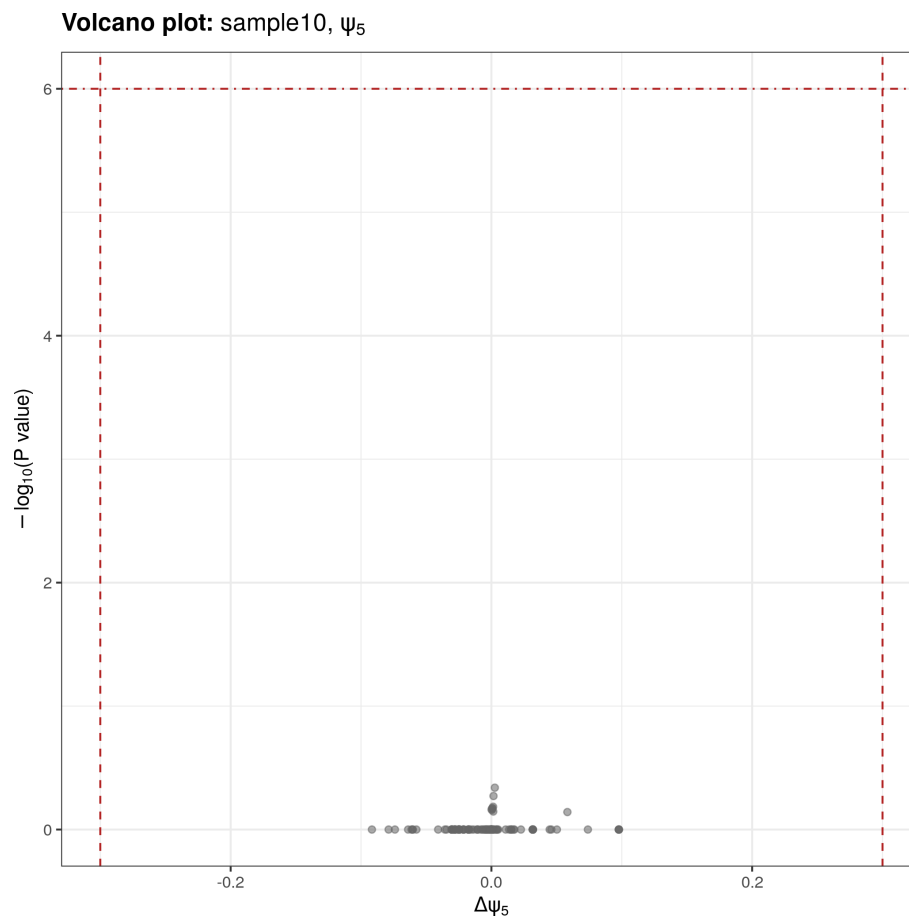
## GRanges object with 2 ranges and 12 metadata columns:
##      seqnames          ranges strand | sampleID hgncSymbol   type
##      <Rle>             <IRanges> <Rle> |   <Rle>      <Rle>   <Rle>
##  119   chr3 119217435-119217436      * | sample9    TIMMDC1 psiSite
##  126   chr3 119217789-119217790      * | sample10   TIMMDC1 psiSite
##      pValue   padjust zScore psiValue deltaPsi meanCounts
##      <numeric> <numeric> <Rle>   <Rle>   <Rle>      <Rle>
##  119  0.011132      1 -2.51   0.84   -0.11      86.58
##  126  0.197310      1 -2.34   0.11   -0.19       2.92
##      meanTotalCounts counts totalCounts
##      <Rle>   <Rle>      <Rle>
##  119      91.58     61      73
##  126       7.25      1       9
##  -----
##  seqinfo: 2 sequences from an unspecified genome; no seqlengths
```

3.4 Finding splicing candidates in patients

Let's have a look at sample 10 and check if we got some splicing candidates for this sample.

```
plotVolcano(fds, type="psi5", "sample10")
```

FRASER: Find RAre Splicing Events in RNA-seq



Which are the splicing events in detail?

```
sampleRes <- res[res$sampleID == "sample10"]
sampleRes

## GRanges object with 1 range and 12 metadata columns:
##      seqnames      ranges strand | sampleID hgncSymbol   type
##      <Rle>         <IRanges> <Rle> |   <Rle>      <Rle>  <Rle>
##  126      chr3 119217789-119217790 * | sample10    TIMMDC1 psiSite
##      pValue  padjust zScore psiValue deltaPsi meanCounts
##      <numeric> <numeric> <Rle>    <Rle>    <Rle>    <Rle>
##  126    0.19731      1  -2.34    0.11    -0.19      2.92
##      meanTotalCounts counts totalCounts
##      <Rle> <Rle>      <Rle>
##  126      7.25      1      9
##  -----
##  seqinfo: 2 sequences from an unspecified genome; no seqlengths
```

To have a closer look at the junction level, use the following functions:

FRASER: Find RAre Splicing Events in RNA-seq

```
plotExpression(fds, type="psi5", result=sampleRes[1])
plotExpectedVsObservedPsi(fds, result=sampleRes[1])
```

3.5 Saving and loading a *FraserDataSet*

A *FraserDataSet* object can be easily saved and reloaded at any time as follows:

```
# saving a fds
workingDir(fds) <- tempdir()
name(fds) <- "ExampleAnalysis"
saveFraserDataSet(fds, dir=workingDir(fds), name=name(fds))

## ----- Sample data table -----
## # A tibble: 12 x 4
##   sampleID bamFile          group gene
##   <chr>    <chr>          <int> <chr>
## 1 sample1  extdata/bam/sample1.bam     1 TIMMDC1
## 2 sample2  extdata/bam/sample2.bam     1 TIMMDC1
## 3 sample3  extdata/bam/sample3.bam     2 MCOLN1
## 4 sample4  extdata/bam/sample4.bam     3 CLPP
## 5 sample5  extdata/bam/sample5.bam    NA NHDF
## 6 sample6  extdata/bam/sample6.bam    NA NHDF
## 7 sample7  extdata/bam/sample7.bam    NA NHDF
## 8 sample8  extdata/bam/sample8.bam    NA NHDF
## 9 sample9  extdata/bam/sample9.bam    NA NHDF
## 10 sample10 extdata/bam/sample10.bam  NA NHDF
## 11 sample11 extdata/bam/sample11.bam  NA NHDF
## 12 sample12 extdata/bam/sample12.bam  NA NHDF
##
## Number of samples:      12
## Number of junctions:    123
## Number of splice sites: 165
## assays(23): rawCountsJ psi5 ... padjBetaBinomial_psiSite
##   zScores_psiSite
##
## ----- Settings -----
## Analysis name:           ExampleAnalysis
## Analysis is strand specific: no
## Working directory:       '/tmp/RtmpXd6AIZ'
##
## ----- BAM parameters -----
## Default used with: bamMapqFilter=0

# two ways of loading a fds by either specifying the directory and anaysis name
```

```
# or directly giving the path the to fds-object.RDS file
fds <- loadFraserDataSet(dir=workingDir(fds), name=name(fds))
fds <- loadFraserDataSet(file=file.path(workingDir(fds),
    "savedObjects", "ExampleAnalysis", "fds-object.RDS"))
```

4 More details on *FRASER*

The function `FRASER` is a convenient wrapper function that takes care of correcting for confounders, fitting the beta binomial distribution and calculating p-values and z-scores for all ψ types. To have more control over the individual steps, the different functions can also be called separately. The following sections give a short explanation of these steps.

4.1 Correction for confounders

The wrapper function `FRASER` and the underlying function `fit` method offer different methods to automatically control for confounders in the data. Currently the following methods are implemented:

- AE: uses a beta-binomial AE
- PCA-BB-Decoder: uses a beta-binomial AE where PCA is used to find the latent space (encoder) due to speed reasons
- PCA: uses PCA for both the encoder and the decoder
- BB: no correction for confounders, fits a beta binomial distribution directly on the raw counts

```
# Using an alternative way to correct splicing ratios
# here: only 2 iteration to speed the calculation up
# for the vignette, the default is 15 iterations
fds <- fit(fds, q=3, type="psi5", implementation="PCA-BB-Decoder",
    iterations=2)

##
## TRUE
## 123
## [1] "Initial PCA loss: 1.88295102819695"
## [1] "Finished with fitting the E matrix. Starting now with the D fit. ..."
## [1] "Tue Jun  9 00:42:00 2020: Iteration: final_1 loss: 1.06958200790331 (mean); 4.1325197650276"
## [1] "Tue Jun  9 00:42:01 2020: Iteration: final_2 loss: 1.03998913773137 (mean); 3.4759683898319"
## [1] "2 Final betabin-AE loss: 1.03998913773137"
```

FRASER: Find RAre Splicing Events in RNA-seq

4.1.1 Finding the dimension of the latent space

For the previous call, the dimension q of the latent space has been fixed to $q = 10$. Since working with the correct q is very important, the *FRASER* package also provides the function `optimHyperParams` that can be used to estimate the dimension q of the latent space of the data. It works by artificially injecting outliers into the data and then comparing the AUC of recalling these outliers for different values of q . Since this hyperparameter optimization step can take some time for the full dataset, we only show it here for a subset of the dataset:

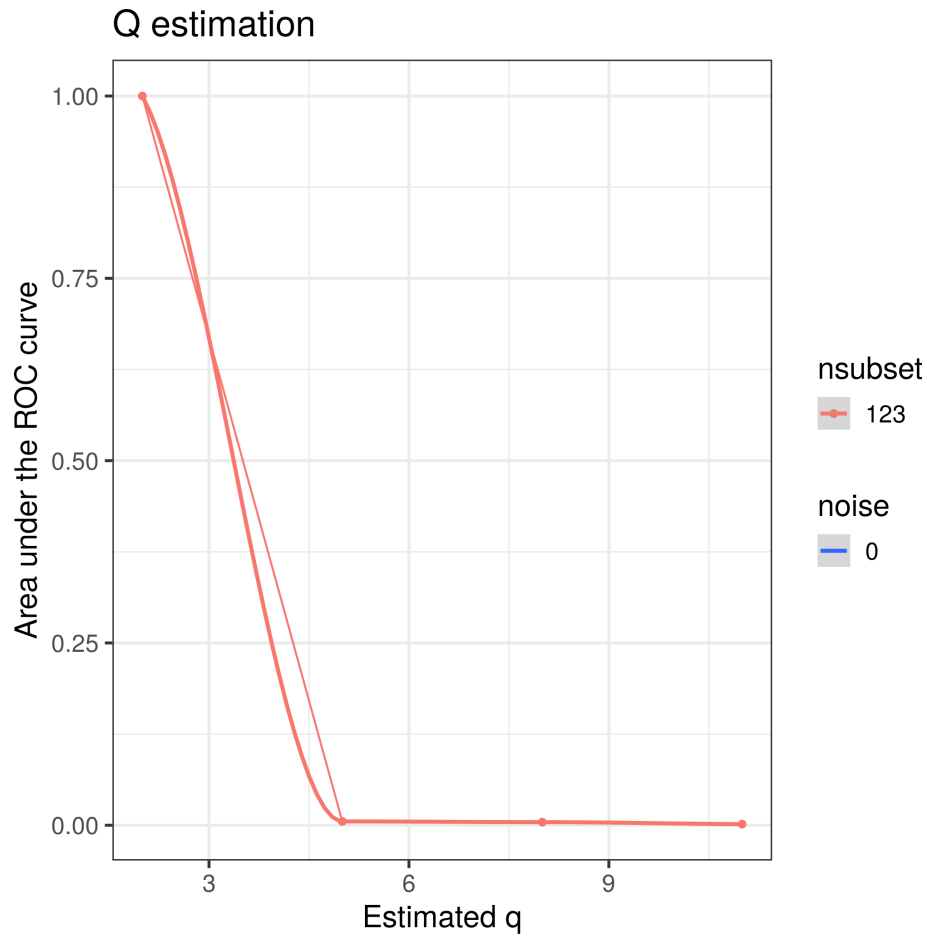
```
set.seed(42)
# hyperparameter optimization
fds <- optimHyperParams(fds, type="psi5", plot=FALSE)

# retrieve the estimated optimal dimension of the latent space
bestQ(fds, type="psi5")

## [1] 2
```

The results from this hyper parameter optimization can be visualized with the function `plotEncDimSearch`.

```
plotEncDimSearch(fds, type="psi5")
```



4.2 P-value calculation

After determining the fit parameters, two-sided beta binomial P-values are computed using the following equation:

$$p_{ij} = 2 \cdot \min \left\{ \frac{1}{2}, \sum_0^{k_{ij}} BB(k_{ij}, n_{ij}, \mu_{ij}, \rho_i), 1 - \sum_0^{k_{ij}-1} BB(k_{ij}, n_{ij}, \mu_{ij}, \rho_i) \right\}, \quad 5$$

where the $\frac{1}{2}$ term handles the case of both terms exceeding 0.5, which can happen due to the discrete nature of counts. Here μ_{ij} are computed as the product of the fitted correction values from the autoencoder and the fitted mean adjustments.

```
fds <- calculatePvalues(fds, type="psi5")
head(pVals(fds, type="psi5"))

## sample1 sample2 sample3 sample4 sample5 sample6 sample7 sample8 sample9
```


FRASER: Find RAre Splicing Events in RNA-seq

```
## 1      1      1      1      1      1      1      1      1      1
## 2      1      1      1      1      1      1      1      1      1
## 3      1      1      1      1      1      1      1      1      1
## 4      1      1      1      1      1      1      1      1      1
## 5      1      1      1      1      1      1      1      1      1
## 6      1      1      1      1      1      1      1      1      1
## sample10 sample11 sample12
## 1      1      1      1
## 2      1      1      1
## 3      1      1      1
## 4      1      1      1
## 5      1      1      1
## 6      1      1      1
```

Afterwards, adjusted p-values can be calculated. Multiple testing correction is done across all junctions in a per-sample fashion using Benjamini-Yekutieli's false discovery rate method[4]. Alternatively, all adjustment methods supported by `p.adjust` can be used via the `method` argument.

```
fds <- calculatePadjValues(fds, type="psi5", method="BY")
head(padjVals(fds, type="psi5"))

## sample1 sample2 sample3 sample4 sample5 sample6 sample7 sample8 sample9
## 1      1      1      1      1      1      1      1      1      1
## 2      1      1      1      1      1      1      1      1      1
## 3      1      1      1      1      1      1      1      1      1
## 4      1      1      1      1      1      1      1      1      1
## 5      1      1      1      1      1      1      1      1      1
## 6      1      1      1      1      1      1      1      1      1
## sample10 sample11 sample12
## 1      1      1      1
## 2      1      1      1
## 3      1      1      1
## 4      1      1      1
## 5      1      1      1
## 6      1      1      1
```

4.3 Z-score calculation

To calculate z-scores on the logit transformed $\Delta\psi$ values and to store them in the *FraserDataSet* object, the function `calculateZScores` can be called. The Z-scores can be used for visualization, filtering, and ranking of samples. The Z-scores are calculated as follows:

$$z_{ij} = \frac{\delta_{ij} - \bar{\delta}_j}{sd(\delta_j)} \quad 6$$

$$\delta_{ij} = \text{logit}\left(\frac{k_{ij} + 1}{n_{ij} + 2}\right) - \text{logit}(\mu_{ij}),$$

where δ_{ij} is the difference on the logit scale between the measured counts and the counts after correction for confounders and $\bar{\delta}_j$ is the mean of intron j .

```
fds <- calculateZscore(fds, type="psi5")
head(zScores(fds, type="psi5"))
```

##	sample1	sample2	sample3	sample4	sample5	sample6
## 1	-2.1755815	2.4067165	-0.1384260	0.17336601	-0.1490333	-0.05973031
## 2	-2.1770128	2.4053175	-0.1385028	0.17401793	-0.1488697	-0.05949617
## 3	-0.2581446	-0.1385868	-0.2818389	-0.52710114	-0.2228608	3.15041270
## 4	-0.2795574	-0.2986968	-0.1564624	-0.17400499	-0.4176445	-0.18609912
## 5	2.6563442	-1.9040835	-0.1603279	-0.05062633	-0.1292613	-0.17316134
## 6	0.3571851	0.2194324	-0.6220061	-0.56612414	-0.3149957	-0.54404760
##	sample7	sample8	sample9	sample10	sample11	sample12
## 1	-0.2336089	-0.40298341	0.01287679	0.1059155	0.04970753	0.41078121
## 2	-0.2332849	-0.40298132	0.01312719	0.1061177	0.05033377	0.41123358
## 3	-0.2514969	-0.08338584	-0.38424823	-0.2005702	-0.36721984	-0.43495950
## 4	3.1348128	-0.62963895	-0.16927603	-0.3435526	-0.44425337	-0.03562676
## 5	-0.2428929	-0.27617132	-0.09323371	0.1507743	-0.04971283	0.27235252
## 6	-0.1055714	-0.02535663	-0.60014295	-0.5307622	-0.26705776	2.99944691

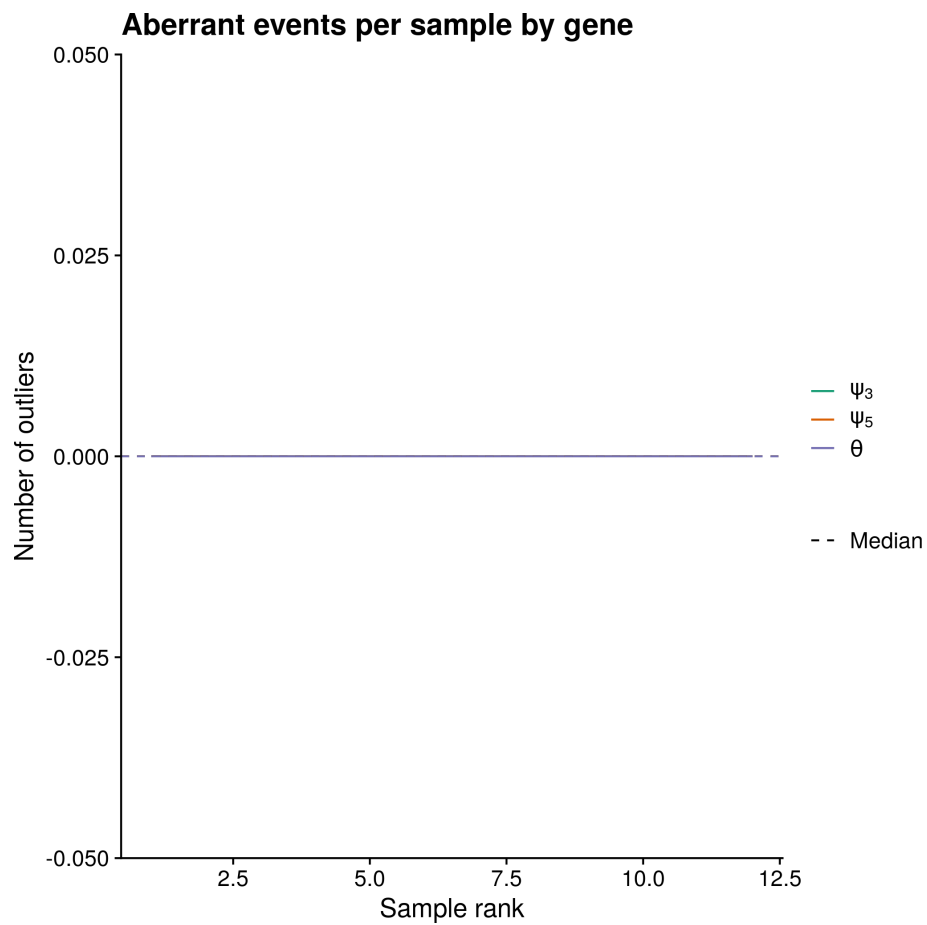
4.4 Result visualization

In addition to the plotting methods `plotVolcano`, `plotExpression`, `plotExpectedVsObservedPsi`, `plotFilterExpression` and `plotEncDimSearch` used above, the *FRASER* package provides two additional functions to visualize the results:

`plotAberrantPerSample` displays the number of aberrant events per sample based on the given cutoff values and `plotQQ` gives a quantile-quantile plot either for a single junction/splice site or globally.

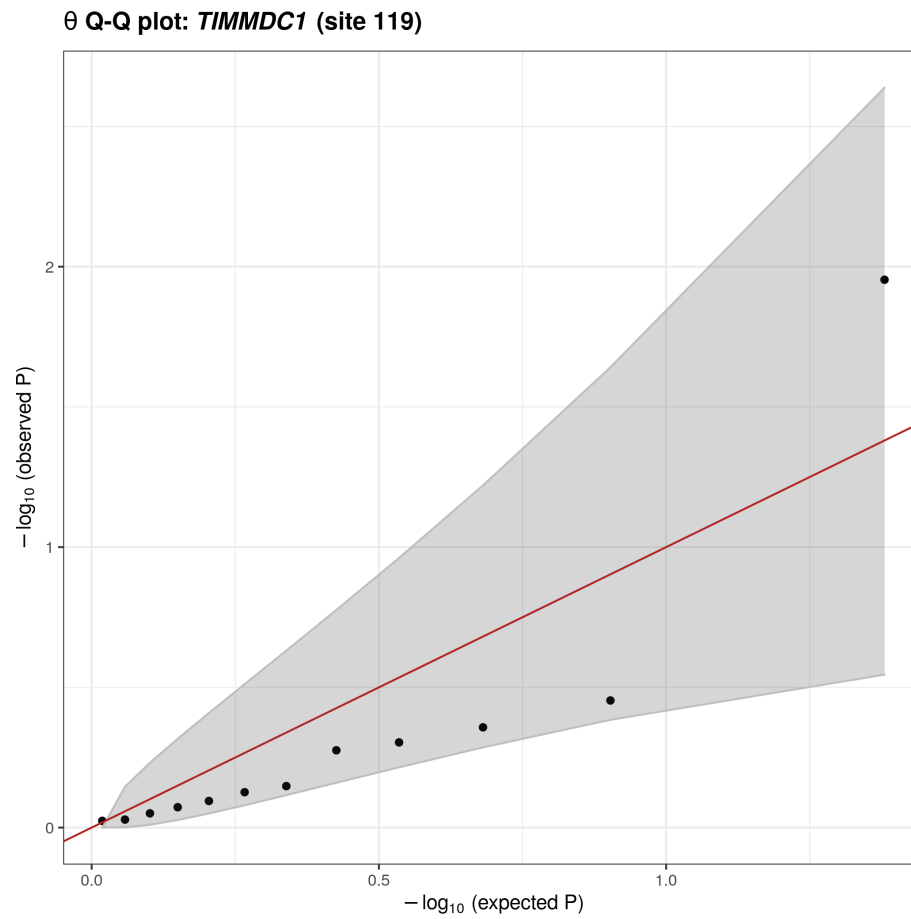
```
plotAberrantPerSample(fds)
```

FRASER: Find RAre Splicing Events in RNA-seq

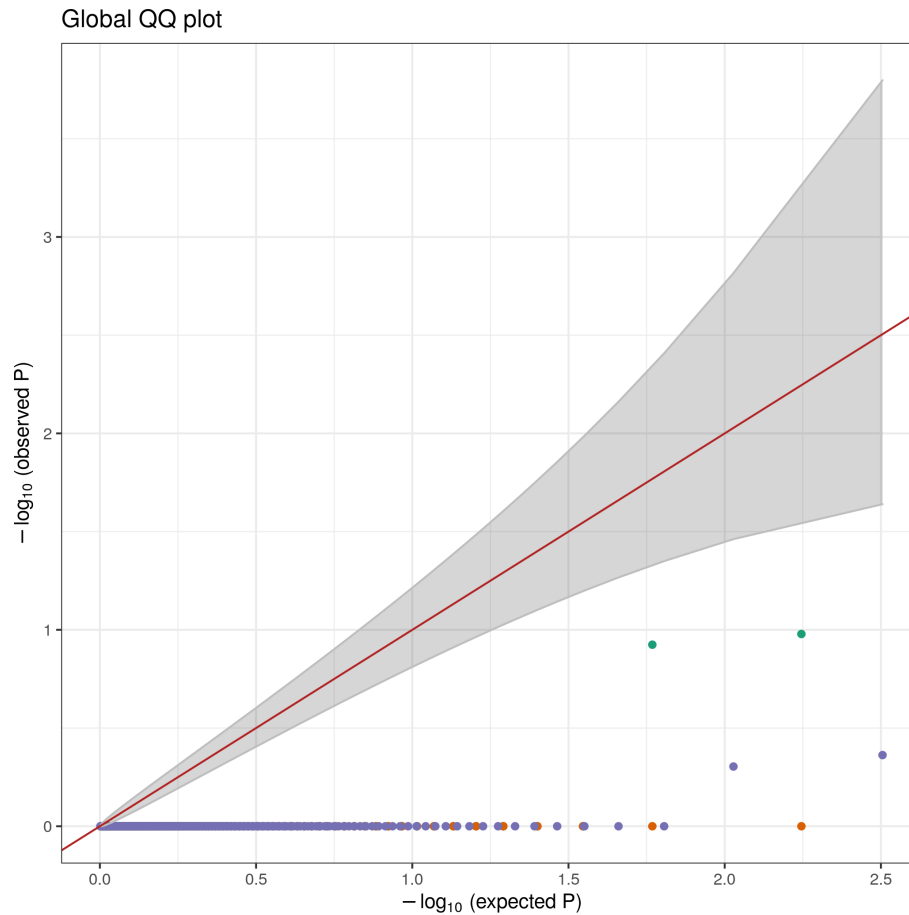


```
# qq-plot for single junction  
plotQQ(fds, result=res[1])
```

FRASER: Find RAre Splicing Events in RNA-seq



```
# global qq-plot (on gene level since aggregate=TRUE)  
plotQQ(fds, aggregate=TRUE, global=TRUE)
```



References

- [1] D. D. Pervouchine, D. G. Knowles, and R. Guigo. Intron-centric estimation of alternative splicing from RNA-seq data. *Bioinformatics*, 29(2):273–274, November 2012. URL: <https://doi.org/10.1093/bioinformatics/bts678>, doi:10.1093/bioinformatics/bts678.
- [2] Alexander Dobin, Carrie A. Davis, Felix Schlesinger, Jorg Drenkow, Chris Zaleski, Sonali Jha, Philippe Batut, Mark Chaisson, and Thomas R. Gingeras. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics*, 29(1):15–21, January 2013. URL: <https://doi.org/10.1093/bioinformatics/bts635>, doi:10.1093/bioinformatics/bts635.
- [3] Santiago Marco-Sola, Michael Sammeth, Roderic Guigó, and Paolo Ribeca. The GEM mapper: fast, accurate and versatile alignment by filtration. *Nature Methods*, 9(12):1185–1188, October 2012. URL: <https://doi.org/10.1038/nmeth.2221>, doi:10.1038/nmeth.2221.

- [4] Yoav Benjamini and Daniel Yekutieli. The control of the false discovery rate in multiple testing under dependency. *Annals of Statistics*, 29(4):1165–1188, 2001. URL: <https://projecteuclid.org/euclid.aos/1013699998>, [arXiv:0801.1095](https://arxiv.org/abs/0801.1095), [doi:10.1214/aos/1013699998](https://doi.org/10.1214/aos/1013699998).

5 Session Info

Here is the output of `sessionInfo()` on the system on which this document was compiled:

```
## R version 4.0.0 (2020-04-24)
## Platform: x86_64-pc-linux-gnu (64-bit)
## Running under: Ubuntu 18.04.4 LTS
##
## Matrix products: default
## BLAS: /home/biocbuild/bbs-3.11-bioc/R/lib/libRblas.so
## LAPACK: /home/biocbuild/bbs-3.11-bioc/R/lib/libRlapack.so
##
## locale:
##  [1] LC_CTYPE=en_US.UTF-8      LC_NUMERIC=C
##  [3] LC_TIME=en_US.UTF-8      LC_COLLATE=C
##  [5] LC_MONETARY=en_US.UTF-8  LC_MESSAGES=en_US.UTF-8
##  [7] LC_PAPER=en_US.UTF-8     LC_NAME=C
##  [9] LC_ADDRESS=C             LC_TELEPHONE=C
## [11] LC_MEASUREMENT=en_US.UTF-8 LC_IDENTIFICATION=C
##
## attached base packages:
## [1] stats4      parallel  stats      graphics  grDevices  utils      datasets
## [8] methods     base
##
## other attached packages:
##  [1] org.Hs.eg.db_3.11.4
##  [2] TxDb.Hsapiens.UCSC.hg19.knownGene_3.2.2
##  [3] GenomicFeatures_1.40.0
##  [4] AnnotationDbi_1.50.0
##  [5] FRASER_1.0.1
##  [6] SummarizedExperiment_1.18.1
##  [7] DelayedArray_0.14.0
##  [8] matrixStats_0.56.0
##  [9] Biobase_2.48.0
## [10] Rsamtools_2.4.0
## [11] Biostrings_2.56.0
## [12] XVector_0.28.0
## [13] GenomicRanges_1.40.0
```

FRASER: Find RAre Splicing Events in RNA-seq

```
## [14] GenomeInfoDb_1.24.0
## [15] IRanges_2.22.2
## [16] S4Vectors_0.26.1
## [17] BiocGenerics_0.34.0
## [18] data.table_1.12.8
## [19] BiocParallel_1.22.0
## [20] knitr_1.28
##
## loaded via a namespace (and not attached):
## [1] nlme_3.1-148          bitops_1.0-6
## [3] bit64_0.9-7           RColorBrewer_1.1-2
## [5] progress_1.2.2        PRR0C_1.3.1
## [7] httr_1.4.1            tools_4.0.0
## [9] backports_1.1.7       utf8_1.1.4
## [11] R6_2.4.1              HDF5Array_1.16.0
## [13] mgcv_1.8-31           lazyeval_0.2.2
## [15] colorspace_1.4-1      DBI_1.1.0
## [17] tidyselect_1.1.0      prettyunits_1.1.1
## [19] bit_1.1-15.2          curl_4.3
## [21] compiler_4.0.0        cli_2.0.2
## [23] plotly_4.9.2.1        labeling_0.3
## [25] rtracklayer_1.48.0    scales_1.1.1
## [27] checkmate_2.0.0       askpass_1.1
## [29] rappdirs_0.3.1        stringr_1.4.0
## [31] digest_0.6.25         rmarkdown_2.2
## [33] R.utils_2.9.2         extraDistr_1.8.11
## [35] Rsubread_2.2.2        pkgconfig_2.0.3
## [37] htmltools_0.4.0       dbplyr_1.4.4
## [39] BSgenome_1.56.0       highr_0.8
## [41] htmlwidgets_1.5.1     rlang_0.4.6
## [43] RSQLite_2.2.0         VGAM_1.1-3
## [45] BBmisc_1.11           DelayedMatrixStats_1.10.0
## [47] farver_2.0.3          generics_0.0.2
## [49] jsonlite_1.6.1        dplyr_1.0.0
## [51] R.oo_1.23.0           RCurl_1.98-1.2
## [53] magrittr_1.5          GenomeInfoDbData_1.2.3
## [55] Matrix_1.2-18         fansi_0.4.1
## [57] munsell_0.5.0         Rcpp_1.0.4.6
## [59] Rhdf5lib_1.10.0       lifecycle_0.2.0
## [61] R.methodsS3_1.8.0     stringi_1.4.6
## [63] yaml_2.2.1            zlibbioc_1.34.0
## [65] rhdf5_2.32.0          BiocFileCache_1.12.0
## [67] grid_4.0.0            blob_1.2.1
## [69] crayon_1.3.4          lattice_0.20-41
## [71] cowplot_1.0.0         splines_4.0.0
```

FRASER: Find RAre Splicing Events in RNA-seq

```
## [73] hms_0.5.3          magick_2.3
## [75] pillar_1.4.4        codetools_0.2-16
## [77] biomaRt_2.44.0      XML_3.99-0.3
## [79] glue_1.4.1          evaluate_0.14
## [81] pcaMethods_1.80.0    BiocManager_1.30.10
## [83] vctrs_0.3.1         tidyr_1.1.0
## [85] gtable_0.3.0        openssl_1.4.1
## [87] purrr_0.3.4         assertthat_0.2.1
## [89] ggplot2_3.3.1       xfun_0.14
## [91] viridisLite_0.3.0   pheatmap_1.0.12
## [93] tibble_3.0.1        GenomicAlignments_1.24.0
## [95] memoise_1.1.0       ellipsis_0.3.1
## [97] BiocStyle_2.16.0
```