

Example workflow for *regsplice* package

Lukas M. Weber

29 April 2016

Package version: regsplice 0.1

Contents

1	Introduction	2
1.1	Example workflow	2
1.2	Data set	2
1.3	SummarizedExperiment class	2
2	Workflow	3
2.1	Complete workflow with wrapper function	3
2.2	Individual steps	3
3	Additional steps for microarray data	3

1 Introduction

The *regsplice* package implements statistical methods for the detection of differential exon usage (differential splicing) in RNA sequencing (RNA-seq) and microarray data sets. The *regsplice* methods are fast and make use of the lasso to improve power compared to standard generalized linear models. The methodology and comparisons to previous approaches are described in our paper:

Title of paper and link to bioRxiv preprint here.

1.1 Example workflow

This vignette demonstrates an example workflow for the *regsplice* package using a small simulated RNA-seq data set.

There are two options for running *regsplice*: you can run a complete workflow in one step using the wrapper function *regsplice*; or you can run the individual functions for each step in sequence, which provides additional insight into the methodology. Both options are demonstrated below.

1.2 Data set

The data set used for the workflow consists of exon-level read counts for a subset of 100 genes from a simulated human RNA-seq data set, consisting of 6 biological samples, with 3 samples in each of 2 conditions.

The original data set is from the paper:

Soneson et al. (2016), *Isoform prefiltering improves performance of count-based methods for analysis of differential transcript usage*, Genome Biology, [available here](#).

Original data files from this paper, containing the simulated RNA-seq reads (FASTQ and BAM files), are available from ArrayExpress at accession code [E-MTAB-3766](#).

Exon bin counts were generated with the Python counting scripts provided with the *DEXSeq* package, using the option to exclude exons from overlapping genes instead of aggregating them into multi-gene complexes (see Soneson et al. 2016, Supplementary Material).

For this workflow, we have selected a subset of the first 100 genes from this simulated data set. The exon-level read counts and the true differential splicing status labels for these 100 genes are saved as tab-delimited TXT files in the *extdata/* directory in the *regsplice* package source code.

1.3 SummarizedExperiment class

The *regsplice* package uses the *SummarizedExperiment* S4 class to enable easier integration into Bioconductor workflows. The *SummarizedExperiment* class stores matrices of data values (such as RNA-seq read counts or microarray expression intensities) along with associated meta-data for rows (genes or exons) and columns (biological samples).

The main benefit of using *SummarizedExperiment* objects is that subsetting operations work intuitively, keeping data and meta-data in sync. The *SummarizedExperiment* class replaces the earlier *ExpressionSet* class. For more information, see the *SummarizedExperiment* Bioconductor vignette, [available here](#).

The rows of *SummarizedExperiment* objects are also naturally grouped into a hierarchy of genes and exons, which is useful for the implementation of the *regsplice* model fitting functions.

2 Workflow

2.1 Complete workflow with wrapper function

How to run a complete workflow in a single step using the *regsplice* wrapper function...

You can provide the data either as filenames / matrices, or as a *SummarizedExperiment* object. If you provide filenames or matrices, you also need to provide gene and exon indices, so that *regsplice* can create the *SummarizedExperiment* object automatically.

...

Below, we show how to run the individual functions for the steps in the *regsplice* workflow in sequence, which provides additional flexibility and insight into the methodology.

2.2 Individual steps

2.2.1 Load data

The first step is to load the data and create a *SummarizedExperiment* object.

```
files_counts <- system.file("extdata", "sample*", package = "regsplice")
file_truth <- system.file("extdata", "truth*", package = "regsplice")
```

2.2.2 Create design matrices

The first step is to create design matrices for each gene

2.2.3 Fit models

regularized model and null model for each gene

2.2.4 Calculate likelihood ratio tests

3 Additional steps for microarray data

If you are using microarray data, you also need to use *limma*/*voom* to convert data...

give example code

Vignette info

Figures

The figure sizes have been customised so that you can easily put two images side-by-side.

```
#plot(1:10)
#plot(10:1)
```

You can enable figure captions by `fig_caption: yes` in YAML:

```
output:
  rmarkdown::html_vignette:
    fig_caption: yes
```

Then you can use the chunk option `fig.cap = "Your figure caption."` in **knitr**.

Math expressions

You can write math expressions, e.g. $Y = X\beta + \epsilon$, footnotes¹, and tables, e.g. using `knitr::kable()`.

¹A footnote here.