

Data preprocessing and creation of the data objects auxiliary for the DEXSeq package

Alejandro Reyes

May 7, 2016

Abstract

This vignette describes the steps that were followed for the generation of the data objects contained in the package *pasilla*.

Contents

1	Downloading the files	1
2	Read alignment and filtering	1
3	Exon count files	2
4	Creation of the <i>DEXSeqDataSet</i> dxd	3

1 Downloading the files

We used the RNA-Seq data from the publication by Brooks et al. [1]. The experiment investigated the effect of siRNA knock-down of *pasilla*, a gene that is known to bind to mRNA in the spliceosome, and which is thought to be involved in the regulation of splicing. The data set contains 3 biological replicates of the knockdown as well as 4 biological replicates for the untreated control. Data files are publicly available in the NCBI Gene Expression Omnibus under the accession GSE18508¹. The read sequences in FASTQ format were extracted from the NCBI short read archive file (.sra files), using the sra toolkit².

2 Read alignment and filtering

The reads in the FASTQ files were aligned using tophat version 1.2.0 with default parameters against the reference *Drosophila melanogaster* genome. Table 1 summarizes the read number and alignment statistics.

¹<http://www.ncbi.nlm.nih.gov/projects/geo/query/acc.cgi?acc=GSE18508>

²http://www.ncbi.nlm.nih.gov/books/NBK47540/#SRA_Download_Guid_B.5_Converting_SRA_for

	file	type	number of lanes	total number of reads	exon counts
1	treated1fb	single-read	5	35158667	15679615
2	treated2fb	paired-end	2	12242535 (x2)	15620018
3	treated3fb	paired-end	2	12443664 (x2)	12733865
4	untreated1fb	single-read	2	17812866	14924838
5	untreated2fb	single-read	6	34284521	20764558
6	untreated3fb	paired-end	2	10542625 (x2)	10283129
7	untreated4fb	paired-end	2	12214974 (x 2)	11653031

Table 1: Read numbers and alignment statistics. The column *exon counts* refers to the number of reads that could be uniquely aligned to an exon.

The reference genome fasta files were obtained from the Ensembl ftp server³. We ran `bowtie-build` to index the fasta file. For more information on this procedure see the bowtie webpage⁴. The indexed form is required by bowtie, and thus tophat.

```
wget ftp://ftp.ensembl.org/pub/release-62/fasta/drosophila_melanogaster/ \
dna/Drosophila_melanogaster.BDGP5.25.62.dna_rm.toplevel.fa.gz
```

```
gunzip Drosophila_melanogaster.BDGP5.25.62.dna_rm.toplevel.fa.gz
bowtie-build Drosophila_melanogaster.BDGP5.25.62.dna_rm.toplevel.fa \
d_melanogaster_BDGP5.25.62
```

We generated the alignment BAM file using tophat. For the single-reads data:

```
tophat bowtie_index reads1.fastq,reads2.fastq,...,readsN.fastq
```

For the paired-end data:

```
tophat -r inner-fragment-size bowtie_index \
reads1_1.fastq,reads2_1.fastq,...,readsN_1.fastq \
reads1_2.fastq,reads2_2.fastq,...,readsN_2.fastq
```

More information on tophat is provided on its webpage⁵. The SAM alignment files from which *pasilla* was generated are available at <http://www-huber.embl.de/pub/DEXSeq/analysis/brooksetal/bam/>.

3 Exon count files

To generate the per-exon read counts, we first needed to define the exonic regions. To this end, we downloaded the file `Drosophila_melanogaster.BDGP5.25.62.gtf.gz` from Ensembl⁶. The script `dexseq_prepare_annotation.py` contained in the *DEXSeq* package was used to extract the exons of the transcripts from the file, define new non-overlapping exonic regions and reformat it to create the file `Dmel.BDGP5.25.62.DEXSeq.chr.gff` contained in `pasilla/extdata`. For example, for this file we ran:

```
wget ftp://ftp.ensembl.org/pub/release-62/gtf/ \
drosophila_melanogaster/Drosophila_melanogaster.BDGP5.25.62.gtf.gz
```

³<http://www.ensembl.org/info/data/ftp/index.html>

⁴<http://bowtie-bio.sourceforge.net/tutorial.shtml>

⁵<http://tophat.cbcb.umd.edu/tutorial.html>

⁶ftp://ftp.ensembl.org/pub/release-62/gtf/drosophila_melanogaster

```
gunzip Drosophila_melanogaster.BDGP5.25.62.gtf.gz
python dexseq_prepare_annotation.py Drosophila_melanogaster.BDGP5.25.62.gtf \
    Dmel.BDGP5.25.62.DEXSeq.chr.gff
```

To count the reads that fell into each non-overlapping exonic part, the script `dexseq_count.py`, which is also contained in the *DEXSeq* package, was used. It took the alignment results in the form of a SAM file (sorted by position in the case of a paired end data) and the `gtf` file `Dmel.BDGP5.25.62.DEXSeq.chr.gff` and returned one file for each biological replicate with the exon counts. For example, for the file `treated1.bam`, which contained single-end alignments, we ran:

```
samtools index treated1.bam
samtools view treated1.bam > treated1.sam
python dexseq_count.py Dmel.BDGP5.25.62.DEXSeq.chr.gff \
    treated1.sam treated1fb.txt
```

For the file `treated2.bam`, which contained paired-end alignments:

```
samtools index treated2.bam
samtools view treated2.bam > treated2.sam
sort -k1,1 -k2,2n treated2.sam > treated2_sorted.sam
python dexseq_count.py -p yes Dmel.BDGP5.25.62.DEXSeq.chr.gff \
    treated2_sorted.sam treated2fb.txt
```

The output of the two HTSeq python scripts is provided in the *pasilla* package:

```
> library("pasilla")
> inDir = system.file("extdata", package="pasilla", mustWork=TRUE)
> dir(inDir)

[1] "Dmel.BDGP5.25.62.DEXSeq.chr.gff" "geneIDsinsubset.txt"
[3] "pasilla_gene_counts.tsv"      "treated1fb.txt"
[5] "treated2fb.txt"              "treated3fb.txt"
[7] "untreated1fb.txt"           "untreated2fb.txt"
[9] "untreated3fb.txt"           "untreated4fb.txt"
```

The Python scripts are built upon the HTSeq library⁷.

4 Creation of the *DEXSeqDataSet* `dxd`

To create an *DEXSeqDataSet* object, we started with a data frame `samples` that contained the sample annotations, as in Table 1.

```
> head(samples)

      condition      type
treated1fb    treated single-read
treated2fb    treated  paired-end
treated3fb    treated  paired-end
```

⁷<http://www-huber.embl.de/users/anders/HTSeq/doc/overview.html>

```
untreated1fb untreated single-read
untreated2fb untreated single-read
untreated3fb untreated paired-end
```

We also needed the annotation file with the per exon annotation.

```
> annotationfile = file.path(inDir, "Dmel.BDGP5.25.62.DEXSeq.chr.gff")
```

With these, we could call the function `DEXSeqDataSet` to construct the object `dxd`.

```
> library("DEXSeq")
> dxd = DEXSeqDataSetFromHTSeq(
+   countfiles = file.path(inDir, paste(rownames(samples), "txt", sep=".")),
+   sampleData=samples,
+   design= ~ sample + exon + condition:exon,
+   flattenedfile = annotationfile)
```

We only wanted to work with data from a subset of genes, which was defined in the following file.

```
> genesforsubset = readLines(file.path(inDir, "geneIDsinsubset.txt"))
> dxd = dxd[genesforsubset,]
```

We save our objects

We saved the objects in the data directory of the package:

```
> save(dxd, file=file.path("../data", "pasillaDEXSeqDataSet.RData"))
```

References

- [1] A. N. Brooks, L. Yang, M. O. Duff, K. D. Hansen, J. W. Park, S. Dudoit, S. E. Brenner, and B. R. Graveley. Conservation of an RNA regulatory map between *Drosophila* and mammals. *Genome Research*, pages 193–202, October 2010.

```
> toLatex(sessionInfo())
```

- R version 3.3.0 (2016-05-03), x86_64-pc-linux-gnu
- Locale: LC_CTYPE=en_US.UTF-8, LC_NUMERIC=C, LC_TIME=en_US.UTF-8, LC_COLLATE=C, LC_MONETARY=en_US.UTF-8, LC_MESSAGES=en_US.UTF-8, LC_PAPER=en_US.UTF-8, LC_NAME=C, LC_ADDRESS=C, LC_TELEPHONE=C, LC_MEASUREMENT=en_US.UTF-8, LC_IDENTIFICATION=C
- Base packages: base, datasets, grDevices, graphics, methods, parallel, stats, stats4, utils
- Other packages: AnnotationDbi 1.34.0, Biobase 2.32.0, BiocGenerics 0.18.0, BiocParallel 1.6.0, DESeq2 1.12.0, DEXSeq 1.18.0, GenomeInfoDb 1.8.0, GenomicRanges 1.24.0, IRanges 2.6.0, RColorBrewer 1.1-2, S4Vectors 0.10.0, SummarizedExperiment 1.2.0, pasilla 0.12.0, xtable 1.8-2
- Loaded via a namespace (and not attached): Biostrings 2.40.0, DBI 0.4, Formula 1.2-1, Hmisc 3.17-4, Matrix 1.2-6, RCurl 1.95-4.8, RSQLite 1.0.0, Rcpp 0.12.4.5, Rsamtools 1.24.0, XML 3.98-1.4, XVector 0.12.0, acepack 1.3-3.3, annotate 1.50.0, biomaRt 2.28.0, bitops 1.0-6, chron 2.3-47, cluster 2.0.4, colorspace 1.2-6, data.table 1.9.6, foreign 0.8-66, genefilter 1.54.0, geneplotter 1.50.0, ggplot2 2.1.0, grid 3.3.0, gridExtra 2.2.1, gtable 0.2.0, hwriter 1.3.2, lattice 0.20-33, latticeExtra 0.6-28, locfit 1.5-9.1, magrittr 1.5, munsell 0.4.3, nnet 7.3-12, plyr 1.8.3, rpart 4.1-10, scales 0.4.0, splines 3.3.0, statmod 1.4.24, stringi 1.0-1, stringr 1.0.0, survival 2.39-2, tools 3.3.0, zlibbioc 1.18.0

Table 2: The output of `sessionInfo` on the build system after running this vignette.