

The **rtracklayer** package

Michael Lawrence

November 23, 2011

1 Introduction

The **rtracklayer** package is an interface (or *layer*) between **R** and genome browsers. Its main purpose is the visualization of genomic annotation *tracks*, whether generated through experimental data analysis performed in **R** or loaded from an external data source. The features of **rtracklayer** may be divided into two categories: 1) the import/export of track data and 2) the control and querying of external genome browser sessions and views.

There are two basic track data structures in Bioconductor: *GRanges*, defined by the *GenomicRanges* package, and *RangedData*, a more general data structure from the *IRanges* package. For most purposes, *GRanges* is preferred in the context of *rtracklayer*; however, *RangedData* has some particular use cases.

rtracklayer supports the import and export of tracks from and to files in various formats, see Section 2.1.4. All positions in a *RangedData* or *GRanges* should be 1-based, as in **R** itself.

The **rtracklayer** package currently interfaces with the **UCSC** web-based genome browser. Other packages may provide drivers for other genome browsers through a plugin system. With **rtracklayer**, the user may start a genome browser session, create and manipulate genomic views, and import/export tracks and sequences to and from a browser. Please note that not all features are necessarily supported by every browser interface.

The rest of this vignette will consist of a number of case studies. First, we consider an experiment investigating microRNA regulation of gene expression, where the microRNA target sites are the primary genomic features of interest.

2 Gene expression and microRNA target sites

This section will demonstrate the features of **rtracklayer** on a microarray dataset from a larger experiment investigating the regulation of human stem cell differentiation by microRNAs. The transcriptome of the cells was measured before and after differentiation by HG-U133plus2 Affymetrix GeneChip arrays. We begin our demonstration by constructing an annotation dataset from the experimental data, and then illustrate the use of the genome browser interface to display interesting genomic regions in the UCSC browser.

2.1 Creating a target site track

For the analysis of the stem cell microarray data, we are interested in the genomic regions corresponding to differentially expressed genes that are known to be targeted by a microRNA. We will represent this information as an annotation track, so that we may view it in the UCSC genome browser.

2.1.1 Constructing the *G*Ranges

In preparation for creating the microRNA target track, we first used **limma** to detect the differentially expressed genes in the microarray experiment. The locations of the microRNA target sites were obtained from MiRBase. The code below stores information about the target sites on differentially expressed genes in the *data.frame* called **targets**, which can also be obtained by entering `data(targets)` when **rtracklayer** is loaded.

```
> library("humanStemCell")
> data(fhesc)
> library("genefilter")
> filtFhesc <- nsFilter(fhesc)[[1]]
> library("limma")
> design <- model.matrix(~filtFhesc$Diff)
> hesclim <- lmFit(filtFhesc, design)
> hesceb <- eBayes(hesclim)
> tab <- topTable(hesceb, coef = 2, adjust.method = "BH", n = 7676)
> tab2 <- tab[(tab$logFC > 1) & (tab$adj.P.Val < 0.01),]
> affyIDs <- tab2$ID
> library("microRNA")
> data(hsTargets)
> library("hgu133plus2.db")
> entrezIDs <- mappedRkeys(hgu133plus2ENTREZID[affyIDs])
> library("org.Hs.eg.db")
> mappedEntrezIDs <- entrezIDs[in% mappedkeys(org.Hs.egENSEMBLTRANS)]
> ensemblIDs <- mappedRkeys(org.Hs.egENSEMBLTRANS[mappedEntrezIDs])
> targetMatches <- match(ensemblIDs, hsTargets$target, 0)
> ## same as data(targets)
> targets <- hsTargets[targetMatches,]
> targets$chrom <- paste("chr", targets$chrom, sep = "")
```

The following code creates the track from the **targets** dataset:

```
> library(rtracklayer)
> library(GenomicRanges)
> ## call data(targets) if skipping first block
> head(targets)
```

	name	target	chrom	start	end
32350	hsa-miR-139-3p	ENST00000372874	chr20	42681711	42681732

```

534942   hsa-miR-135a ENST00000336199   chr1 241718191 241718211
534975   hsa-miR-148a ENST00000366540   chr1 241728912 241728933
534987   hsa-miR-505 ENST00000366539   chr1 241733864 241733885
534991   hsa-miR-505 ENST00000263826   chr1 241733864 241733885
699911   hsa-miR-196b ENST00000321955 chr11 89565074 89565096
strand
32350   -
534942   -
534975   -
534987   -
534991   -
699911   +

```

```

> targetRanges <- IRanges(targets$start, targets$end)
> targetTrack <- with(targets,
+                   GRangesForUCSCGenome("hg18", chrom, targetRanges, strand,
+                   name, target))

```

The `GRangesForUCSCGenome` function constructs a `GRanges` object for the named genome. The stand information, the name of the microRNA and the Ensembl ID of the targeted transcript are stored in the `GRanges`. The chromosome for each site is passed as the `chrom` argument. The chromosome names and lengths for the genome are taken from the UCSC database and stored in the `GRanges` along with the genome identifier. We can retrieve them as follows:

```

> genome(targetTrack)

```

chr1	chr1_random	chr2	chr2_random	chr3
"hg18"	"hg18"	"hg18"	"hg18"	"hg18"
chr3_random	chr4	chr4_random	chr5	chr5_h2_hap1
"hg18"	"hg18"	"hg18"	"hg18"	"hg18"
chr5_random	chr6	chr6_cox_hap1	chr6_qbl_hap2	chr6_random
"hg18"	"hg18"	"hg18"	"hg18"	"hg18"
chr7	chr7_random	chr8	chr8_random	chr9
"hg18"	"hg18"	"hg18"	"hg18"	"hg18"
chr9_random	chr10	chr10_random	chr11	chr11_random
"hg18"	"hg18"	"hg18"	"hg18"	"hg18"
chr12	chr13	chr13_random	chr14	chr15
"hg18"	"hg18"	"hg18"	"hg18"	"hg18"
chr15_random	chr16	chr16_random	chr17	chr17_random
"hg18"	"hg18"	"hg18"	"hg18"	"hg18"
chr18	chr18_random	chr19	chr19_random	chr20
"hg18"	"hg18"	"hg18"	"hg18"	"hg18"
chr21	chr21_random	chr22	chr22_h2_hap1	chr22_random
"hg18"	"hg18"	"hg18"	"hg18"	"hg18"
chrX	chrX_random	chrY	chrM	
"hg18"	"hg18"	"hg18"	"hg18"	

```
> head(seqlengths(targetTrack))
      chr1 chr1_random      chr2 chr2_random      chr3
247249719  1663265  242951149  185571  199501827
chr3_random
  749256
```

While this extra information is not strictly needed to upload data to UCSC, calling `GRangesForUCSCGenome` is an easy way to formally associate interval data to a UCSC genome build. This ensures, for example, that the data will always be uploaded to the correct genome, regardless of browser state. It also immediately validates whether the intervals fall within the bounds of the genome.

For cases where one is not interacting with the UCSC genome browser, and in particular when network access is unavailable, the `GRangesForBSGenome` function behaves the same, except it finds an installed *BSGenome* package and loads it to retrieve the chromosome information.

2.1.2 Accessing track information

The track information is now stored in the R session as a *GRanges* object. It holds the chromosome, start, end and strand for each feature, along with any number of data columns.

The primary feature attributes are the `start`, `end`, `seqnames` and `strand`. There are accessors for each of these, named accordingly. For example, the following code retrieves the chromosome names and then start positions for each feature in the track.

```
> head(seqnames(targetTrack))
'factor' Rle of length 6 with 3 runs
Lengths:  1  4  1
Values : chr20 chr1 chr11
Levels(49): chr1 chr1_random chr2 ... chrX_random chrY chrM

> head(start(targetTrack))
[1] 42681711 241718191 241728912 241733864 241733864 89565074
```

Exercises

1. Get the strand of each feature in the track
2. Calculate the length of each feature
3. Reconstruct (partially) the `targets` *data.frame*

2.1.3 Subsetting a *GRanges*

It is often helpful to extract subsets from *GRanges* instances, especially when uploading to a genome browser. The data can be subset through a matrix-style syntax by feature and column. The conventional `[]` method is employed for subsetting, where the first parameter, *i*, indexes the features and *j* indexes the data columns. Both *i* and *j* may contain numeric, logical and character indices, which behave as expected.

```
> ## get the first 10 targets
> first10 <- targetTrack[1:10]
> ## get pos strand targets
> posTargets <- targetTrack[strand(targetTrack) == "+"]
> ## get the targets on chr1
> chr1Targets <- targetTrack[seqnames(targetTrack) == "chr1"]
```

Exercises

1. Subset the track for all features on the negative strand of chromosome 2.

2.1.4 Exporting and importing tracks

Import and export of *GRanges* and *RangedData* instances is supported in the following formats: Browser Extended Display (BED), versions 1, 2 and 3 of the General Feature Format (GFF), and Wiggle (WIG). Support for additional formats may be provided by other packages through a plugin system.

To save the microRNA target track created above in a format understood by other tools, we could export it as BED. This is done with the `export` function, which accepts a filename or any R connection object as its target. If a target is not given, the serialized string is returned. The desired format is derived, by default, from the extension of the filename. Use the `format` parameter to explicitly specify a format.

```
> export(targetTrack, "targets.bed")
```

To read the data back in a future session, we could use the `import` function. The source of the data may be given as a connection, a filename or a character vector containing the data. Like the `export` function, the format is determined from the filename, by default.

```
> restoredTrack <- import("targets.bed")
```

The `restoredTrack` object is of class *RangedData*. To obtain a *GRanges* object directly upon importing, pass `asRangedData = FALSE` to `import`:

```
> restoredTrack <- import("targets.bed", asRangedData = FALSE)
```

Eventually, `import` will return *GRanges* objects by default; the extra parameter is necessary for backwards compatibility. Thus, if a *RangedData* is desired, new code should specify `asRangedData = TRUE` in anticipation of the change in default value.

Exercises

1. Output the track to a file in the “gff” format.
2. Read the track back into R.
3. Export the track as a character vector.

2.2 Viewing the targets in a genome browser

For the next step in our example, we will load the track into a genome browser for visualization with other genomic annotations. The **rtracklayer** package is capable of interfacing with any genome browser for which a driver exists. In this case, we will interact with the web-based **UCSC** browser, but the same code should work for any browser.

2.2.1 Starting a session

The first step towards interfacing with a browser is to start a browser session, represented in R as a *BrowserSession* object. A *BrowserSession* is primarily a container of tracks and genomic views. The following code creates a *BrowserSession* for the **UCSC** browser:

```
> session <- browserSession("UCSC")
```

Note that the name of any other supported browser could have been given here instead of “UCSC”. To see the names of supported browsers, enter:

```
> genomeBrowsers()
[1] "UCSC"
```

2.2.2 Laying the track

Before a track can be viewed on the genome, it must be loaded into the session using the `track<-` function, as demonstrated below:

```
> track(session, "targets") <- targetTrack
```

The *name* argument should be a character vector that will help identify the track within `session`. Note that the invocation of `track<-` above does not specify an upload format. Thus, the default, “auto”, is used. Since the track does not contain any data values, the track is uploaded as BED. To make this explicit, we could pass “bed” as the *format* parameter.

Exercises

1. Lay a track with the first 100 features of `targetTrack`
Here we use the short-cut `$` syntax for storing the track.

2.2.3 Viewing the track

For **UCSC**, a view roughly corresponds to one tab or window in the web browser. The target sites are distributed throughout the genome, so we will only be able to view a few features at a time. In this case, we will view only the first feature in the track. A convenient way to focus a view on a particular set of features is to subset the track and pass the range of the subtrack to the constructor of the view. Below we take a track subset that contains only the first feature.

```
> subTargetTrack <- targetTrack[1] # get first feature
```

Now we call the `browserView` function to construct the view and pass the subtrack, zoomed out by a factor of 10, as the segment to view. By passing the name of the targets track in the `pack` parameter, we instruct the browser to use the “pack” mode for viewing the track. This results in the name of the microRNA appearing next to the target site glyph.

```
> view <- browserView(session, subTargetTrack * -10, pack = "targets")
```

If multiple ranges are provided, multiple views are launched:

```
> view <- browserView(session, targetTrack[1:5] * -10, pack = "targets")
```

Exercises

1. Create a new view with the same region as `view`, except zoomed out 2X.
2. Create a view with the “targets” track displayed in “full” mode, instead of “packed”.

2.2.4 A shortcut

There is also a shortcut to the above steps. The `browseGenome` function creates a session for a specified browser, loads one or more tracks into the session and creates a view of a given genome segment. In the following code, we create a new **UCSC** session, load the track and view the first two features, all in one call:

```
> browseGenome(targetTrack, range = subTargetTrack * -10)
```

It is even simpler to view the subtrack in **UCSC** by relying on parameter defaults:

```
> browseGenome(subTargetTrack)
```

2.2.5 Downloading tracks

It is possible to query the browser to obtain the names of the loaded tracks and to download the tracks into R. To list the tracks loaded in the browser, enter the following:

```
> loaded_tracks <- trackNames(session)
```

One may download any of the tracks, such as the “targets” track that was loaded previously in this example.

```
> subTargetTrack <- track(session, "targets")
```

The returned object is a *RangedData*, even if the data was originally uploaded as a *GRanges* or other object. To get a *GRanges* instead, pass `asRangedData = FALSE`. By default, the segment of the track downloaded is the current default genome segment associated with the session. One may download track data for any genome segment, such as those on a particular chromosome. Note that this does not distinguish by strand; we are only indicating a position on the genome.

```
> chr1Targets <- track(session, "targets", chr1Targets)
```

Exercises

1. Get the SNP under the first target, displayed in `view`.
2. Get the UCSC gene for the same target.

2.2.6 Accessing view state

The `view` variable is an instance of *BrowserView*, which provides an interface for getting and setting view attributes. Note that for the UCSC browser, changing the view state opens a new view, as a new page must be opened in the web browser.

To programmatically query the segment displayed by a view, use the `range` method for a *BrowserView*.

```
> segment <- range(view)
```

Similarly, one may get and set the names of the visible tracks in the view.

```
> visible_tracks <- trackNames(view)
> trackNames(view) <- visible_tracks
```

The visibility mode (hide, dense, pack, squish, full) of the tracks may be retrieved with the `ucscTrackModes` method.

```
> modes <- ucscTrackModes(view)
```

The returned value, `modes`, is of class *UCSCTrackModes*. The modes may be accessed using the `[]` function. Here, we set the mode of our “targets” track to “full” visibility.


```

> modes["targets"]
> modes["targets"] <- "full"
> ucscTrackModes(view) <- modes

```

Existing browser views for a session may be retrieved by calling the `browserViews` method on the `browserSession` instance.

```

> views <- browserViews(session)
> length(views)

```

Exercises

1. Retrieve target currently visible in the view.
2. Limit the view to display only the SNP, UCSC gene and target track.
3. Hide the UCSC gene track.

3 CPNE1 expression and HapMap SNPs

Included with the `rtracklayer` package is a track object (created by the `GGtools` package) with features from a subset of the SNPs on chromosome 20 from 60 HapMap founders in the CEU cohort. Each SNP has an associated data value indicating its association with the expression of the CPNE1 gene according to a Cochran-Armitage 1df test. The top 5000 scoring SNPs were selected for the track.

We load the track presently.

```

> library(rtracklayer)
> data(cpneTrack)

```

3.1 Loading and manipulating the track

The data values for a track are stored in the columns on the `RangedData` instance. Often, a track contains a single column of numeric values, conventionally known as the `score`. The `score` function retrieves the column named `score` or, if one does not exist, the first column in the `RangedData`, as long as it is numeric. Otherwise, `NULL` is returned.

```

> head(score(cpneTrack))

rs4814683 rs6076506 rs6139074 rs1418258 rs7274499 rs6116610
0.16261691 0.02170423 0.47098379 0.16261691 0.05944578 0.18101862

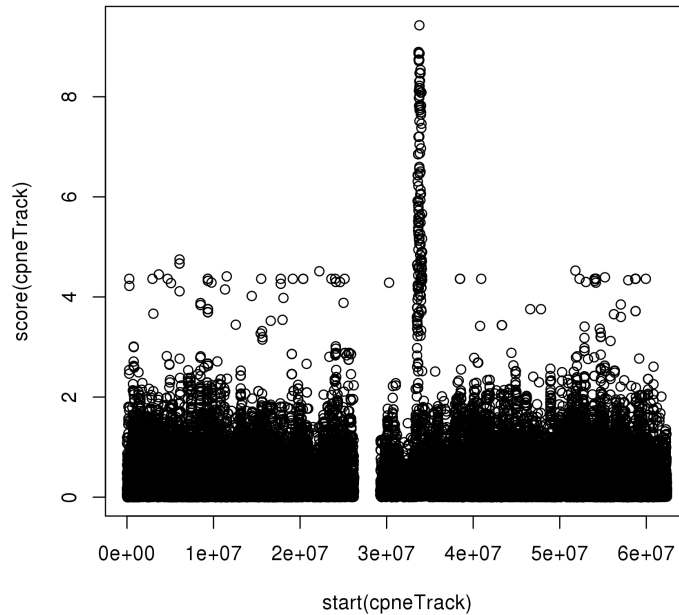
```

One use of extracting the data values is to plot the data.

```

> plot(start(cpneTrack), score(cpneTrack))

```



3.2 Browsing the SNPs

We now aim to view some of the SNPs in the UCSC browser. Unlike the microRNA target site example above, this track has quantitative information, which requires special consideration for visualization.

3.2.1 Laying a WIG track

To view the SNP locations as a track in a genome browser, we first need to upload the track to a fresh session. In the code below, we use the `[[<-` alias of `track<-`.

```
> session <- browserSession()
> session$cpne <- cpneTrack
```

Note that because `cpneTrack` contains data values and its features do not overlap, it is uploaded to the browser in the WIG format. One limitation of the WIG format is that it is not possible to encode strand information. Thus, each strand needs to have its own track, and `rtracklayer` does this automatically, unless only one strand is represented in the track (as in this case). One could pass “bed” to the `format` parameter of `track<-` to prevent the split, but tracks uploaded as BED are much more limited compared to WIG tracks in terms of visualization options.

To form the labels for the WIG subtracks, “p” is concatenated onto the plus track and “m” onto the minus track. Features with missing track information are placed in a track named with the “na” postfix. It is important to note that the subtracks must be identified individually when, for example, downloading the track or changing track visibility.

3.2.2 Plotting the SNP track

To plot the data values for the SNP’s in a track, we need to create a *browserView*. We will view the region spanning the first 5 SNPs in the track, which will be displayed in the “full” mode.

```
> view <- browserView(session, range(cpneTrack[1:5,]), full = "cpne")
```

The UCSC browser will plot the data values as bars. There are several options available for tweaking the plot, as described in the help for the *GraphTrackLine* class. These need to be specified laying the track, so we will lay a new track named “cpne2”. First, we will turn the *autoScale* option off, so that the bars will be scaled globally, rather than locally to the current view. Then we could turn on the *yLineOnOff* option to add horizontal line that could represent some sort of cut-off. The position of the line is specified by *yLineMark*. We set it arbitrarily to the 25% quantile.

```
> track(session, "cpne2", autoScale = FALSE, yLineOnOff = TRUE,
+       yLineMark = quantile(score(cpneTrack), .25)) <- cpneTrack
> view <- browserView(session, range(cpneTrack[1:5,]), full = "cpne2")
```

4 Binding sites for NRSF

Another common type of genomic feature is transcription factor binding sites. Here we will use the **Biostrings** package to search for matches to the binding motif for NRSF, convert the result to a track, and display a portion of it in the UCSC browser.

4.1 Creating the binding site track

We will use the **Biostrings** package to search human chromosome 1 for NRSF binding sites. The binding sequence motif is assumed to be *TCAGCACCATG-GACAG*, though in reality it is more variable. To perform the search, we run *matchPattern* on the positive strand of chromosome 1.

```
> library(BSgenome.Hsapiens.UCSC.hg19)
> nrsfHits <- matchPattern("TCAGCACCATGGACAG", Hsapiens[["chr1"]])
> length(nrsfHits) # number of hits
```

```
[1] 2
```

We then convert the hits, stored as a *Views* object (a particular type of *IRanges* object), to a *RangedData* instance.

```
> nrsfTrack <- GenomicData(ranges(nrsfHits), strand="+", chrom="chr1",
+                           genome = "hg19")
```

GenomicData is a convenience function that by default constructs a *RangedData* object. A *GRanges* may be obtained by passing `asRangedData = FALSE`; however, a *RangedData* may be preferred in this case, as the `nrsfHits` object is stored as a *Views* object in `nrsfTrack` without any loss of information, in particular the sequence of `chr1`.

4.2 Browsing the binding sites

Now that the NRSF binding sites are stored as a track, we can upload them to the UCSC browser and view them. Below, load the track and we view the region around the first hit in a single call to `browseGenome`.

```
> session <- browseGenome(nrsfTrack, range = range(nrsfTrack[1,]) * -10)
```

We observe significant conservation across mammal species in the region of the motif.

5 Conclusion

These case studies have demonstrated a few of the most important features of **rtracklayer**. Please see the package documentation for more details.

The following is the session info that generated this vignette:

```
> sessionInfo()

R version 2.14.0 (2011-10-31)
Platform: x86_64-unknown-linux-gnu (64-bit)

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=C               LC_NAME=C
 [9] LC_ADDRESS=C            LC_TELEPHONE=C
[11] LC_MEASUREMENT=en_US.UTF-8 LC_IDENTIFICATION=C

attached base packages:
[1] stats      graphics  grDevices  utils      datasets  methods
[7] base

other attached packages:
```

[1] BSgenome.Hsapiens.UCSC.hg19_1.3.17
[2] BSgenome_1.22.0
[3] Biostrings_2.22.0
[4] GenomicRanges_1.6.4
[5] IRanges_1.12.3
[6] rtracklayer_1.14.4
[7] RCurl_1.7-0
[8] bitops_1.0-4.1
[9] microRNA_1.12.0
[10] limma_3.10.0
[11] genefilter_1.36.0
[12] humanStemCell_0.2.5
[13] hgu133plus2.db_2.6.3
[14] org.Hs.eg.db_2.6.4
[15] RSQLite_0.10.0
[16] DBI_0.2-5
[17] AnnotationDbi_1.16.5
[18] Biobase_2.14.0

loaded via a namespace (and not attached):

[1] XML_3.4-3 annotate_1.32.0 splines_2.14.0
[4] survival_2.36-10 tools_2.14.0 xtable_1.6-0
[7] zlibbioc_1.0.0