

Gene Expression Assessment

Stephen Blatti

November 23, 2017

We have shown how Bioconductor provides resources for studying the human genome sequence as well as SNPs. Bioconductor also provides resources that permit us to obtain information about genes. We will see how these databases can be quite complex. But before learning about these here we present some experimental data.

The driving force behind the formation of the Bioconductor project was the emergence of high throughput measurement of gene expression data. Unlike genome sequence and SNP data, gene expression data varies from cell to cell, from tissue to tissue and from individual to individual. Statistical techniques such as those implemented in R were natural tools to parse out variability and perform statistical inference. Furthermore, the ambitious use of newly invented technologies added issues of measurement error and bias to already difficult challenge.

Note that in the previous assessments we focused on static database information: genome sequence and SNPs. In previous courses we have seen the ‘tissuesGeneExpression’ data which is experimental data measured with microarrays. If you have not installed it you can do it like this:

```
library(devtools)
```

```
## Warning: package 'devtools' was built under R version 3.4.2
```

```
#install_github("genomicsclass/tissuesGeneExpression")
```

You can load the data:

```
library(tissuesGeneExpression)
data(tissuesGeneExpression)
head(e[,1:5])
```

```
##          GSM11805.CEL.gz GSM11814.CEL.gz GSM11823.CEL.gz GSM11830.CEL.gz
## 1007_s_at      10.191267      10.509167      10.272027      10.252952
## 1053_at        6.040463        6.696075        6.144663        6.575153
## 117_at         7.447409        7.775354        7.696235        8.478135
## 121_at         12.025042       12.007817       11.633279       11.075286
## 1255_g_at      5.269269        5.180389        5.301714        5.372235
## 1294_at        8.535176        8.587241        8.277414        8.603650
##          GSM12067.CEL.gz
## 1007_s_at      10.157605
## 1053_at        6.606701
## 117_at         8.116336
## 121_at         10.832528
## 1255_g_at      5.334905
## 1294_at        8.303227
```

```
table(tissue)
```

```
## tissue
## cerebellum      colon endometrium hippocampus      kidney      liver
##          38          34          15          31          39          26
## placenta
##          6
```

The rows of the matrix `e` are the features, in this case representing genes, and the columns are the samples. The entries of the matrix are gene expression measurements (in log scale) obtained using a microarray technology.

Once the `tissuesGeneExpression` package is loaded, and `data(tissuesGeneExpression)` is run, you have object `e`, `tab`, and `tissue` in your workspace. You can work with them separately but it is preferable in Bioconductor to unify them in an object. In 2017, the preferred object type is `SummarizedExperiment`. Let's perform the computations and then discuss their utility.

```
library(SummarizedExperiment)

## Loading required package: GenomicRanges
## Warning: package 'GenomicRanges' was built under R version 3.4.2
## Loading required package: stats4
## Loading required package: BiocGenerics
## Loading required package: parallel
##
## Attaching package: 'BiocGenerics'
## The following objects are masked from 'package:parallel':
##
##   clusterApply, clusterApplyLB, clusterCall, clusterEvalQ,
##   clusterExport, clusterMap, parApply, parCapply, parLapply,
##   parLapplyLB, parRapply, parSapply, parSapplyLB
## The following objects are masked from 'package:stats':
##
##   IQR, mad, sd, var, xtabs
## The following objects are masked from 'package:base':
##
##   anyDuplicated, append, as.data.frame, cbind, colMeans,
##   colnames, colSums, do.call, duplicated, eval, evalq, Filter,
##   Find, get, grep, grepl, intersect, is.unsorted, lapply,
##   lengths, Map, mapply, match, mget, order, paste, pmax,
##   pmax.int, pmin, pmin.int, Position, rank, rbind, Reduce,
##   rowMeans, rownames, rowSums, sapply, setdiff, sort, table,
##   tapply, union, unique, unsplit, which, which.max, which.min
## Loading required package: S4Vectors
##
## Attaching package: 'S4Vectors'
## The following object is masked from 'package:base':
##
##   expand.grid
## Loading required package: IRanges
## Loading required package: GenomeInfoDb
## Warning: package 'GenomeInfoDb' was built under R version 3.4.2
## Loading required package: Biobase
## Welcome to Bioconductor
##
```

```
## Vignettes contain introductory material; view with
## 'browseVignettes()'. To cite Bioconductor, see
## 'citation("Biobase")', and for packages 'citation("pkgname)".

## Loading required package: DelayedArray
## Loading required package: matrixStats

##
## Attaching package: 'matrixStats'

## The following objects are masked from 'package:Biobase':
##
## anyMissing, rowMedians

##
## Attaching package: 'DelayedArray'

## The following objects are masked from 'package:matrixStats':
##
## colMaxs, colMins, colRanges, rowMaxs, rowMins, rowRanges

## The following object is masked from 'package:base':
##
## apply

tissSE = SummarizedExperiment(list(rma=e))
colData(tissSE) = DataFrame(tab)
```

The storage concept for the assay results is that we have a matrix with rows corresponding to features (in this case, genes) and columns corresponding to samples (in this case, extracts from particular tissues). The basic idea is that we unite metadata about experimental samples (colData) with metadata about features (rowData) and the numerical assay data in a single object. We also define methods such as “[” so that the R expression `X[G, S]` for SummarizedExperiment `X` defines a new SummarizedExperiment with features restricted to those specified in `G`, and samples restricted to those specified in `S`. When we want the numerical values for all features and samples in a SummarizedExperiment `X`, we use `assay(X)`. This approach reduces the risk of mismatches between sample characteristic data and assay data, as selections are coordinated through the underlying code for `[`. *### Localization of expression to tissues*

```
#Look at the data for the feature with ID "209169_at". You can index the rows of

#assay(tissSE)

#directly with this character string. For example,

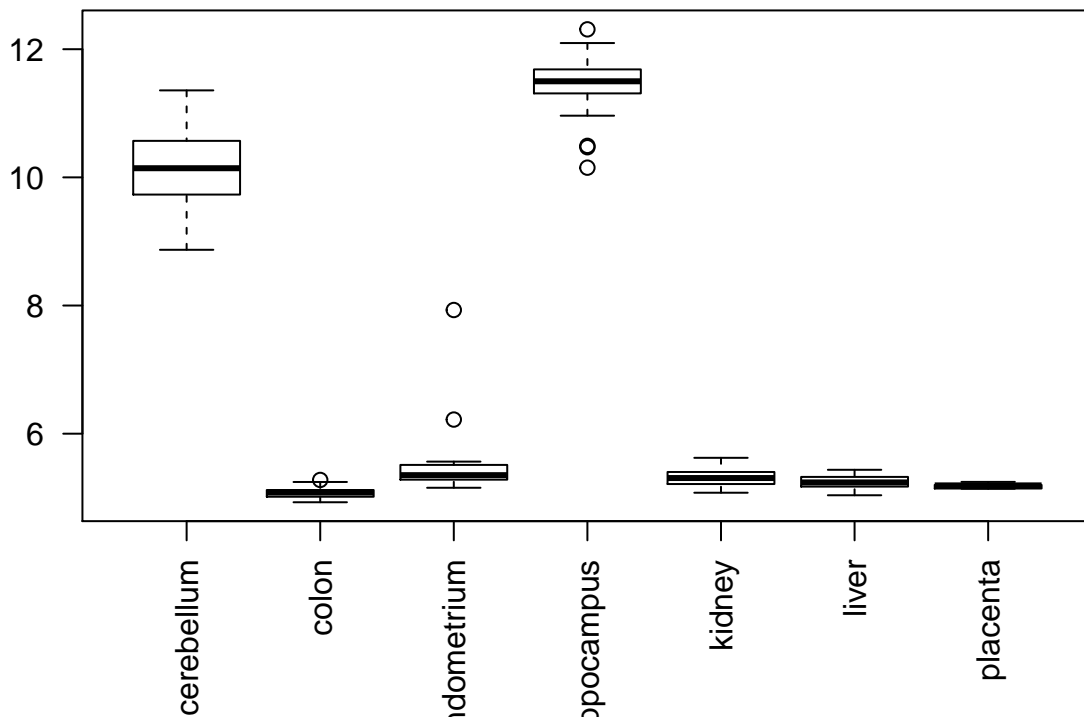
mean(assay(tissSE["209169_at",]))

## [1] 7.26365
#is about 7.26
```

Which of the following best describes the data? (Hint: stratify assay data for the feature by tissue and create boxplots)

This is human data and this gene has the same sequence across all tissues thus there is no difference in expression across tissues
This gene is expressed in the brain but not the other tissues
This gene is differentially expressed between all tissues
The individual to individual variability is much larger than the difference between tissues

```
boxplot(e["209169_at",]~tissue,las=2)
```



```
# boxplot(assay(tissSE)["209169_at",]~tissSE$Tissue, las=2, mar=c(6,4,2,2) )
```

Comparing genes for tissue-specificity

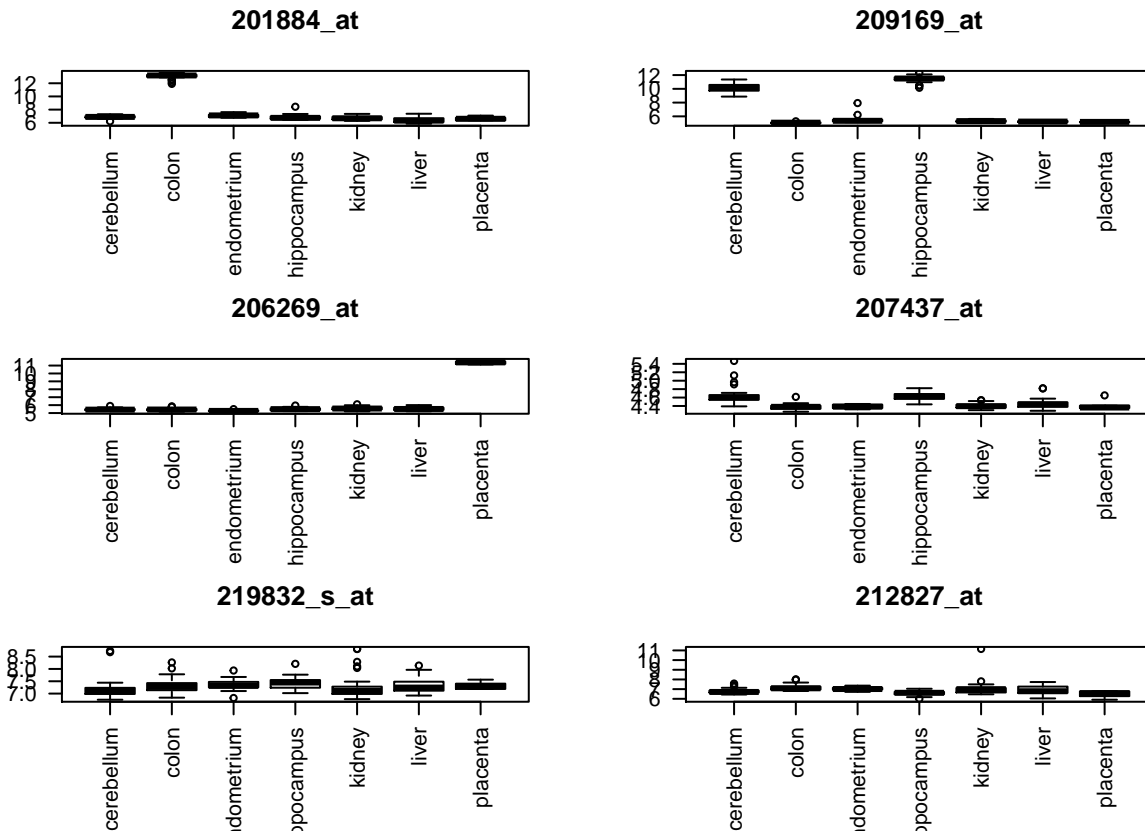
Below is a vector of 6 IDs which index features of 'tissSE':

```
IDs = c("201884_at", "209169_at", "206269_at", "207437_at", "219832_s_at", "212827_at")
```

Which of the following ID(s) appear to represent a gene specific to placenta? Be careful when you are picking, to pick the correct name or names. Names often look similar.

```
"201884_at"
"209169_at"
"206269_at"
"207437_at"
"219832_s_at"
"212827_at"
```

```
par(mfrow=c(3,2))
sapply(IDs,function(x){boxplot(e[x,]~tissue,las=2,main=x)})
```



```
##      201884_at  209169_at  206269_at  207437_at  219832_s_at
## stats Numeric,35 Numeric,35 Numeric,35 Numeric,35 Numeric,35
## n      Numeric,7  Numeric,7  Numeric,7  Numeric,7  Numeric,7
## conf   Numeric,14 Numeric,14 Numeric,14 Numeric,14 Numeric,14
## out    Numeric,6  Numeric,7  Numeric,6  Numeric,10 Numeric,13
## group  Numeric,6  Numeric,7  Numeric,6  Numeric,10 Numeric,13
## names  Character,7 Character,7 Character,7 Character,7 Character,7
##      212827_at
## stats Numeric,35
## n      Numeric,7
## conf   Numeric,14
## out    Numeric,8
## group  Numeric,8
## names  Character,7
```

```
# IDs = c("201884_at", "209169_at", "206269_at", "207437_at", "219832_s_at", "212827_at")
# par(mar=c(6,4,2,2), mfrow=c(3,2))
# for(i in IDs){
#   boxplot(assay(tissSE)[i,]~tissSE$Tissue, las=2, main=i)
# }
```

Discovery of microarray annotation in Bioconductor

Note that there is much existing work on gene function and all we have here are identifiers provided by the manufacturer of the machine that makes the measurements. How would we go about finding more information about gene “206269_at” for example? Does it have a known function? Where is it on the genome? What is

its sequence? One of the strengths of Bioconductor is that it connects R, an existing comprehensive toolbox for data analysis, with the existing comprehensive databases annotating the genome. We will learn about these powerful resources in this class.

The microarray product used to make the measurements described here is the Affymetrix Human GeneChip HG133A. Search the Bioconductor website and determine which of the following packages provides a connection to gene information:

```
Biobase
simpleaffy
hgu133a2cdf
hgu133a.db correct
affy
```

Oligo sequences on affymetrix arrays

The affymetrix technology for mRNA abundance measurement is based on hybridization of highly processed biological sample material to synthetic oligonucleotide “probes” that are on fixed positions of the microarray surface. Bioconductor provides detailed information on the probe and array structure as published by affymetrix.

Install and attach the hgu133aprobe package.

```
library(BiocInstaller)

## Bioconductor version 3.5 (BiocInstaller 1.26.1), ?biocLite for help
## A newer version of Bioconductor is available for this version of R,
##   ?BiocUpgrade for help

#biocLite("hgu133aprobe")
library(hgu133aprobe)

## Loading required package: AnnotationDbi

head(hgu133aprobe)

##           sequence      x      y Probe.Set.Name
## 1 CACCCAGCTGGTCCTGTGGATGGGA 467 181      1007_s_at
## 2 GCCCCACTGGACAACACTGATTCCT 531 299      1007_s_at
## 3 TGGACCCCACTGGCTGAGAATCTGG  86 557      1007_s_at
## 4 AAATGTTTCCTTGTGCCTGCTCCTG 365 115      1007_s_at
## 5 TCCTTGTGCCTGCTCCTGTACTTGT 207 605      1007_s_at
## 6 TGCCTGCTCCTGTACTTGTCTCAG 593 599      1007_s_at
##  Probe.Interrogation.Position Target.Strandedness
## 1                      3330      Antisense
## 2                      3443      Antisense
## 3                      3512      Antisense
## 4                      3563      Antisense
## 5                      3570      Antisense
## 6                      3576      Antisense
```

The field “sequence” gives 25 base-pair sequences of oligonucleotides that are in the 3’ UTR region of the gene associated with the array “probe set”.

You will learn how to use this information to check for accuracy of annotation, to assess risk of cross-hybridization, etc. This table is essentially a large data.frame. How many oligos are used to interrogate

samples for gene GCM1, annotated to probe 206269_at? You will need to work with the Probe.Set.Name field of the hgu133aprobe data.frame.

```
sum(hgu133aprobe$Probe.Set.Name == "206269_at")
```

```
## [1] 11
```

We'll conclude this series with a quick illustration of annotation enhancement of a SummarizedExperiment.

```
library(BiocInstaller)
#biocLite("hgu133a.db")
library(hgu133a.db)
```

```
## Loading required package: org.Hs.eg.db
```

```
##
```

```
##
```

```
sym = mapIds(hgu133a.db, keys=rownames(tissSE), column="SYMBOL", keytype="PROBEID")
```

```
## 'select()' returned 1:many mapping between keys and columns
```

```
nm = mapIds(hgu133a.db, keys=rownames(tissSE), column="GENENAME", keytype="PROBEID")
```

```
## 'select()' returned 1:many mapping between keys and columns
```

```
rowData(tissSE) = DataFrame(symbol=sym, genename=nm)
```

To restrict attention to genes with 'phosphatase' in their names, use code like:

```
tissSE[ grep("phosphatase", rowData(tissSE)$genename), ] ### Counting features in a SummarizedExperiment
```

Set up the rowData for tissSE as noted above. How many features are annotated to genes with 'kinase' in their name?

```
tissSE[ grep("kinase", rowData(tissSE)$genename), ]
```

```
## class: SummarizedExperiment
```

```
## dim: 1065 189
```

```
## metadata(0):
```

```
## assays(1): rma
```

```
## rownames(1065): 1007_s_at 200075_s_at ... 59644_at 632_at
```

```
## rowData names(2): symbol genename
```

```
## colnames: NULL
```

```
## colData names(6): filename DB_ID ... SubType ClinicalGroup
```