# Class12

Dennis Kim

Already submitted the first part but lost the r project file, time to catch up on some old data lost

```
library(BiocManager)
library(DESeq2)
```

Loading required package: S4Vectors

Loading required package: stats4

Loading required package: BiocGenerics


Attaching package: 'BiocGenerics'

The following objects are masked from 'package:stats':

    IQR, mad, sd, var, xtabs

The following objects are masked from 'package:base':

    anyDuplicated, aperm, append, as.data.frame, basename, cbind,
    colnames, dirname, do.call, duplicated, eval, evalq, Filter, Find,
    get, grep, grepl, intersect, is.unsorted, lapply, Map, mapply,
    match, mget, order, paste, pmax, pmax.int, pmin, pmin.int,
    Position, rank, rbind, Reduce, rownames, sapply, setdiff, sort,
    table, tapply, union, unique, unsplit, which.max, which.min


Attaching package: 'S4Vectors'

The following objects are masked from 'package:base':

    expand.grid, I, unname

Loading required package: IRanges


Attaching package: 'IRanges'

The following object is masked from 'package:grDevices':

    windows

Loading required package: GenomicRanges

Loading required package: GenomeInfoDb

Loading required package: SummarizedExperiment

Loading required package: MatrixGenerics

Loading required package: matrixStats


Attaching package: 'MatrixGenerics'

The following objects are masked from 'package:matrixStats':

    colAlls, colAnyNAs, colAnys, colAvgsPerRowSet, colCollapse,
    colCounts, colCummaxs, colCummins, colCumprods, colCumsums,
    colDiffs, colIQRDiffs, colIQRs, colLogSumExps, colMadDiffs,
    colMads, colMaxs, colMeans2, colMedians, colMins, colOrderStats,
    colProds, colQuantiles, colRanges, colRanks, colSdDiffs, colSds,
    colSums2, colTabulates, colVarDiffs, colVars, colWeightedMads,
    colWeightedMeans, colWeightedMedians, colWeightedSds,
    colWeightedVars, rowAlls, rowAnyNAs, rowAnys, rowAvgsPerColSet,
    rowCollapse, rowCounts, rowCummaxs, rowCummins, rowCumprods,
    rowCumsums, rowDiffs, rowIQRDiffs, rowIQRs, rowLogSumExps,
    rowMadDiffs, rowMads, rowMaxs, rowMeans2, rowMedians, rowMins,

```
    rowOrderStats, rowProds, rowQuantiles, rowRanges, rowRanks,
    rowSdDiffs, rowSds, rowSums2, rowTabulates, rowVarDiffs, rowVars,
    rowWeightedMads, rowWeightedMeans, rowWeightedMedians,
    rowWeightedSds, rowWeightedVars
```

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")'.


Attaching package: 'Biobase'

The following object is masked from 'package:MatrixGenerics':

    rowMedians

The following objects are masked from 'package:matrixStats':

    anyMissing, rowMedians

```r
counts <- read.csv("airway_scaledcounts.csv", row.names=1)
metadata <-  read.csv("airway_metadata.csv")
```

```r
library(DESeq2)
```

```r
dds <- DESeqDataSetFromMatrix(countData=counts,
                              colData=metadata,
                              design=~dex)
```

converting counts to integer mode

Warning in DESeqDataSet(se, design = design, ignoreRank): some variables in
design formula are characters, converting to factors
```

```
dds
```

```
class: DESeqDataSet
dim: 38694 8
metadata(1): version
assays(1): counts
rownames(38694): ENSG00000000003 ENSG00000000005 ... ENSG00000283120
  ENSG00000283123
rowData names(0):
colnames(8): SRR1039508 SRR1039509 ... SRR1039520 SRR1039521
colData names(4): id dex celltype geo_id
```

```
dds <- DESeq(dds)
```

```
estimating size factors
```

```
estimating dispersions
```

```
gene-wise dispersion estimates
```

```
mean-dispersion relationship
```

```
final dispersion estimates
```

```
fitting model and testing
```

```
res <- results(dds)
res
```

```
log2 fold change (MLE): dex treated vs control
Wald test p-value: dex treated vs control
DataFrame with 38694 rows and 6 columns
                  baseMean log2FoldChange      lfcSE      stat     pvalue
                 <numeric>      <numeric> <numeric> <numeric> <numeric>
ENSG00000000003   747.1942     -0.3507030  0.168246 -2.084470 0.0371175
ENSG00000000005     0.0000             NA        NA        NA         NA
ENSG00000000419   520.1342      0.2061078  0.101059  2.039475 0.0414026
```

```
ENSG00000000457   322.6648       0.0245269   0.145145   0.168982 0.8658106
ENSG00000000460    87.6826      -0.1471420   0.257007  -0.572521 0.5669691
...                     ...             ...        ...        ...       ...
ENSG00000283115   0.000000              NA         NA         NA        NA
ENSG00000283116   0.000000              NA         NA         NA        NA
ENSG00000283119   0.000000              NA         NA         NA        NA
ENSG00000283120   0.974916       -0.668258    1.69456  -0.394354  0.693319
ENSG00000283123   0.000000              NA         NA         NA        NA
                       padj
                  <numeric>
ENSG00000000003   0.163035
ENSG00000000005         NA
ENSG00000000419   0.176032
ENSG00000000457   0.961694
ENSG00000000460   0.815849
...                     ...
ENSG00000283115         NA
ENSG00000283116         NA
ENSG00000283119         NA
ENSG00000283120         NA
ENSG00000283123         NA
```

```r
summary(res)
```

```
out of 25258 with nonzero total read count
adjusted p-value < 0.1
LFC > 0 (up)       : 1563, 6.2%
LFC < 0 (down)     : 1188, 4.7%
outliers [1]       : 142, 0.56%
low counts [2]     : 9971, 39%
(mean count < 10)
[1] see 'cooksCutoff' argument of ?results
[2] see 'independentFiltering' argument of ?results
```

```r
res05 <- results(dds, alpha=0.05)
summary(res05)
```

```
out of 25258 with nonzero total read count
```

```
adjusted p-value < 0.05
LFC > 0 (up)        : 1236, 4.9%
LFC < 0 (down)      : 933, 3.7%
outliers [1]        : 142, 0.56%
low counts [2]      : 9033, 36%
(mean count < 6)
[1] see 'cooksCutoff' argument of ?results
[2] see 'independentFiltering' argument of ?results
```

done catching up, time to work on the new lab

```r
library("AnnotationDbi")
library("org.Hs.eg.db")
```

```r
columns(org.Hs.eg.db)
```

```
 [1] "ACCNUM"      "ALIAS"       "ENSEMBL"       "ENSEMBLPROT"  "ENSEMBLTRANS"
 [6] "ENTREZID"    "ENZYME"      "EVIDENCE"      "EVIDENCEALL"  "GENENAME"
[11] "GENETYPE"    "GO"          "GOALL"         "IPI"          "MAP"
[16] "OMIM"        "ONTOLOGY"    "ONTOLOGYALL"   "PATH"         "PFAM"
[21] "PMID"        "PROSITE"     "REFSEQ"        "SYMBOL"       "UCSCKG"
[26] "UNIPROT"
```

We can use the mapIds() function to add individual columns to our results table. We provide the row names of our results table as a key, and specify that keytype=ENSEMBL. The column argument tells the mapIds() function which information we want, and the multiVals argument tells the function what to do if there are multiple possible values for a single input value. Here we ask to just give us back the first one that occurs in the database.

> Q11. Run the mapIds() function two more times to add the Entrez ID and UniProt accession and GENENAME as new columns called res$entrez, res$uniprot and res$genename.

```r
res$symbol <- mapIds(org.Hs.eg.db,
                     keys=row.names(res), # Our genenames
                     keytype="ENSEMBL",        # The format of our genenames
                     column="SYMBOL",          # The new format we want to add
                     multiVals="first")
```

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

```r
res$entrez <- mapIds(org.Hs.eg.db,
                     keys=row.names(res), # Our genenames
                     keytype="ENSEMBL",        # The format of our genenames
                     column="ENTREZID",         # The new format we want to add
                     multiVals="first")
```

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

```r
res$uniprot <- mapIds(org.Hs.eg.db,
                      keys=row.names(res), # Our genenames
                      keytype="ENSEMBL",        # The format of our genenames
                      column="UNIPROT",          # The new format we want to add
                      multiVals="first")
```

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

```r
res$genename <- mapIds(org.Hs.eg.db,
                       keys=row.names(res), # Our genenames
                       keytype="ENSEMBL",        # The format of our genenames
                       column="GENENAME",         # The new format we want to add
                       multiVals="first")
```

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

```r
head(res)
```

```
log2 fold change (MLE): dex treated vs control
Wald test p-value: dex treated vs control
DataFrame with 6 rows and 10 columns
                 baseMean log2FoldChange    lfcSE      stat    pvalue
                <numeric>      <numeric> <numeric> <numeric> <numeric>
ENSG00000000003 747.194195     -0.3507030  0.168246 -2.084470 0.0371175
ENSG00000000005   0.000000             NA        NA        NA        NA
ENSG00000000419 520.134160      0.2061078  0.101059  2.039475 0.0414026
ENSG00000000457 322.664844      0.0245269  0.145145  0.168982 0.8658106
```

```
ENSG00000000460   87.682625      -0.1471420   0.257007 -0.572521 0.5669691
ENSG00000000938    0.319167      -1.7322890   3.493601 -0.495846 0.6200029
                        padj      symbol      entrez      uniprot
                   <numeric> <character> <character> <character>
ENSG00000000003   0.163035      TSPAN6        7105  A0A024RCI0
ENSG00000000005         NA        TNMD       64102      Q9H2S6
ENSG00000000419   0.176032        DPM1        8813      O60762
ENSG00000000457   0.961694       SCYL3       57147      Q8IZE3
ENSG00000000460   0.815849     C1orf112       55732  A0A024R922
ENSG00000000938         NA         FGR        2268      P09769
                        genename
                       <character>
ENSG00000000003         tetraspanin 6
ENSG00000000005           tenomodulin
ENSG00000000419 dolichyl-phosphate m..
ENSG00000000457 SCY1 like pseudokina..
ENSG00000000460 chromosome 1 open re..
ENSG00000000938 FGR proto-oncogene, ..
```

```r
ord <- order( res$padj )
#View(res[ord,])
head(res[ord,])
```

```
log2 fold change (MLE): dex treated vs control
Wald test p-value: dex treated vs control
DataFrame with 6 rows and 10 columns
                  baseMean log2FoldChange      lfcSE      stat      pvalue
                 <numeric>      <numeric> <numeric> <numeric>   <numeric>
ENSG00000152583    954.771        4.36836 0.2371268   18.4220 8.74490e-76
ENSG00000179094    743.253        2.86389 0.1755693   16.3120 8.10784e-60
ENSG00000116584   2277.913       -1.03470 0.0650984  -15.8944 6.92855e-57
ENSG00000189221   2383.754        3.34154 0.2124058   15.7319 9.14433e-56
ENSG00000120129   3440.704        2.96521 0.2036951   14.5571 5.26424e-48
ENSG00000148175  13493.920        1.42717 0.1003890   14.2164 7.25128e-46
                       padj      symbol      entrez      uniprot
                  <numeric> <character> <character> <character>
ENSG00000152583 1.32441e-71     SPARCL1        8404  A0A024RDE1
ENSG00000179094 6.13966e-56        PER1        5187      O15534
ENSG00000116584 3.49776e-53     ARHGEF2        9181      Q92974
ENSG00000189221 3.46227e-52        MAOA        4128      P21397
ENSG00000120129 1.59454e-44       DUSP1        1843      B4DU40
```

```
ENSG00000148175 1.83034e-42          STOM          2040       F8VSL7
                                      genename
                                    <character>
ENSG00000152583             SPARC like 1
ENSG00000179094 period circadian reg..
ENSG00000116584 Rho/Rac guanine nucl..
ENSG00000189221    monoamine oxidase A
ENSG00000120129 dual specificity pho..
ENSG00000148175                stomatin
```

```r
write.csv(res[ord,], "deseq_results.csv")
```

Now we can load the packages and setup the KEGG data-sets we need. The gageData package has pre-compiled databases mapping genes to KEGG pathways and GO terms for common organisms. kegg.sets.hs is a named list of 229 elements. Each element is a character vector of member gene Entrez IDs for a single KEGG pathway.

```r
library(pathview)
```

```
################################################################################
Pathview is an open source software package distributed under GNU General
Public License version 3 (GPLv3). Details of GPLv3 is available at
http://www.gnu.org/licenses/gpl-3.0.html. Particullary, users are required to
formally cite the original Pathview paper (not just mention it) in publications
or products. For details, do citation("pathview") within R.

The pathview downloads and uses KEGG data. Non-academic uses may require a KEGG
license agreement (details at http://www.kegg.jp/kegg/legal.html).
################################################################################
```

```r
library(gage)
```

```r
library(gageData)
```

```r
data(kegg.sets.hs)
```

```r
# Examine the first 2 pathways in this kegg set for humans
```

```
head(kegg.sets.hs, 2)
```

```
$`hsa00232 Caffeine metabolism`
[1] "10"    "1544" "1548" "1549" "1553" "7498" "9"

$`hsa00983 Drug metabolism - other enzymes`
 [1] "10"     "1066"  "10720" "10941" "151531" "1548"  "1549"  "1551"
 [9] "1553"   "1576"  "1577"  "1806"  "1807"   "1890"  "221223" "2990"
[17] "3251"   "3614"  "3615"  "3704"  "51733"  "54490" "54575" "54576"
[25] "54577"  "54578" "54579" "54600" "54657"  "54658" "54659" "54963"
[33] "574537" "64816" "7083"  "7084"  "7172"   "7363"  "7364"  "7365"
[41] "7366"   "7367"  "7371"  "7372"  "7378"   "7498"  "79799" "83549"
[49] "8824"   "8833"  "9"     "978"
```

```
foldchanges = res$log2FoldChange
names(foldchanges) = res$entrez
head(foldchanges)
```

```
      7105         64102         8813        57147        55732         2268
-0.35070302          NA   0.20610777   0.02452695  -0.14714205  -1.73228897
```

Run gage pathway analysis

```
keggres = gage(foldchanges, gsets=kegg.sets.hs)
attributes(keggres)
```

```
$names
[1] "greater" "less"     "stats"
```

```
# Look at the first three down (less) pathways
head(keggres$less, 3)
```

```
                                  p.geomean stat.mean       p.val
hsa05332 Graft-versus-host disease 0.0004250461 -3.473346 0.0004250461
hsa04940 Type I diabetes mellitus  0.0017820293 -3.002352 0.0017820293
hsa05310 Asthma                    0.0020045888 -3.009050 0.0020045888
                                       q.val set.size       exp1
```

```
hsa05332 Graft-versus-host disease 0.09053483        40 0.0004250461
hsa04940 Type I diabetes mellitus  0.14232581        42 0.0017820293
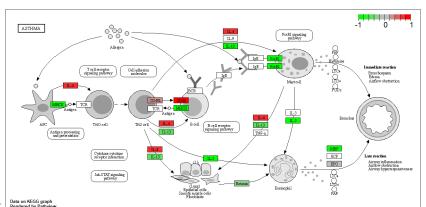hsa05310 Asthma                     0.14232581        29 0.0020045888
```

Let's pull up the highlighted pathways and show our differentially expressed genes on the pathway. I will use the "hsa" KEGG id to get the pathway from KEGGand my 'foldchange' vector to show my genes.

```
pathview(gene.data=foldchanges, pathway.id="hsa05310")
```

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

```
Info: Working in directory C:/Users/denjk/OneDrive/Desktop/BioInformatics Lab/Class12/Class1
```

```
Info: Writing image file hsa05310.pathview.png
```



Put this image into my document