

erccdashboard Package Vignette

Sarah A. Munro

October 26, 2021

This vignette describes the use of the `erccdashboard` R package to analyze External RNA Controls Consortium (ERCC) spike-in control ratio mixtures in gene expression experiments. If you use this package for method validation of your gene expression experiments please cite our manuscript that describes this R package using citation("erccdashboard").

In this vignette we demonstrate analysis of two types of gene expression experiments from the SEQC project that used ERCC control ratio mixture spike-ins:

- Rat toxicogenomics methimazole-treated and control samples
- Human reference RNA samples from the MAQC I project, Universal Human Reference RNA (UHRR) and Human Brain Reference RNA (HBRR)

A subset of the large data set produced in the SEQC study are provided here as examples. The three sets of example data are:

1. Rat toxicogenomics RNA-Seq gene expression count data
2. UHRR/HBRR RNA-Seq gene expression count data
3. UHRR/HBRR Microarray gene expression fluorescent intensity data

1 Rat Toxicogenomics Example: MET (methimazole treatment) and CTL (control) Experiment

1.1 Load data and define input parameters

Load the package gene expression data.

```
> data(SEQC.Example)
```

The R workspace should now contain 5 objects Three of these objects are gene expression experiment expression measures:

- `UHRR.HBRR.arrayDat` - Fluorescent signal data from an Illumina beadarray microarray experiment with UHRR and HBRR in the SEQC interlaboratory study
- `MET.CTL.countDat` - RNA-Seq count data from a rat toxicogenomics experiment
- `UHRR.HBRR.countDat` - RNA-Seq count data from Lab 5 in the SEQC interlaboratory study with UHRR and HBRR

The other two objects are vectors of total reads for the 2 sequencing experiments

- `MET.CTL.totalReads` - total sequenced reads factors for each column in the corresponding rat experiment count table
- `UHRR.HBRR.totalReads` - total sequenced reads factors for each column in the corresponding UHRR/HBRR count table

1.2 Quick analysis: runDashboard

To run the default analysis function `runDashboard` on the MET-CTL rat toxicogenomics RNA-Seq experiment, the following input arguments are required:

```
> datType = "count" # "count" for RNA-Seq data, "array" for microarray data
> isNorm = FALSE # flag to indicate if input expression measures are already
> # normalized, default is FALSE
> exTable = MET.CTL.countDat # the expression measure table
> filenameRoot = "RatTox" # user defined filename prefix for results files
> sample1Name = "MET" # name for sample 1 in the experiment
> sample2Name = "CTL" # name for sample 2 in the experiment
> erccmix = "RatioPair" # name of ERCC mixture design, "RatioPair" is default
> erccdilution = 1/100 # dilution factor used for Ambion spike-in mixtures
> spikeVol = 1 # volume (in microliters) of diluted spike-in mixture added to
> # total RNA mass
> totalRNAmass = 0.500 # mass (in micrograms) of total RNA
> choseFDR = 0.05 # user defined false discovery rate (FDR), default is 0.05
```

The first input argument, `datType`, indicates whether that data is integer count data from an RNA-Seq experiment (“count”) or data from a microarray experiment (“array”). The `isNorm` argument indicates if the input expression measures are already normalized, the default value is `FALSE`. If you want to use normalized RNA-Seq or microarray data in your analysis, the `isNorm` argument must be set to `TRUE`. If `isNorm` is `TRUE`, then the software asks if the input data is length normalized. Type Y at the command line in the R console if the data is length normalized (e.g. FPKM or RPKM data) otherwise type N.

If the data is normalized, then `limma` will be used for array data differential expression (DE) testing, but for normalized RNA-Seq data, DE testing results must be generated outside of the `erccdashboard` pipeline and the DE testing results should be provided in the working directory in a file named ‘filenameRoot ERCC Pvals.csv’ (for details on how to do this see section **3.1 Flexibility in Differential Expression Testing**).

The third argument, `exTable`, is the expression measure table. Take a look at the expression measure table from the RatTox experiment, to see an example `exTable` argument:

```
> head(MET.CTL.countDat)
      Feature MET_1 MET_2 MET_3 CTL_1 CTL_2 CTL_3
16499 ERCC-00002 16629 18798 26568 36600 45436 25163
16500 ERCC-00003  1347  1565  1983  3048  3447  2195
16501 ERCC-00004  4569  5570  6755  1240  1484   902
16502 ERCC-00009   811   869  1123   909  1073   537
16503 ERCC-00012     0     0     0     0     0     0
16504 ERCC-00013     3     1     2     1     5     1
```

The first column of the expression measure table, `Feature`, contains unique names for all the transcripts that were quantified in this experiment. The remaining columns represent replicates of the pair of samples, in this expression measure table the control sample is labeled CTL and the treatment sample is labeled MET. An underscore is included to separate the sample names from the replicate numbers during analysis. This column name format `Sample_Rep` is required for the columns of any input expression measure table. Only one underscore (–) should be used in the column names.

The default differential expression testing of RNA-Seq experiments in the `erccdashboard` is done with the `QuasiSeq` package, which requires the use of integer count data. If you are trying to use either a different approach for DE testing of RNA-Seq count data (e.g. `edgeR` or `DESeq`) or RNA-Seq data that is not integer count data (e.g. FPKM data from `Cufflinks`) please see section **3.1 Flexibility in Differential Expression Testing**.

The `erccdashboard` default normalization for RNA-Seq count data is 75th percentile (also known as upper quartile) normalization. It is optional to provide a vector of per replicate normalization factors through the input argument `repNormFactor`, such as a vector of total reads for each replicate. The example total reads vectors we provide here were derived from the FASTQ files associated with each column in the RNA-Seq experiment count tables. Any `repNormFactor` vector will be used as a library size normalization factor for each column of `exTable`. This will be adjusted to be a per million reads factor.

For any experiment the sample spiked with ERCC Mix 1 is `sample1Name` and the sample spiked with ERCC Mix 2 is `sample2Name`. In this experiment `sample1Name` = MET and `sample2Name` = CTL. For a more robust experimental design the reverse spike-in design could be created using additional replicates of the treatment and control samples. ERCC Mix 2 would be spiked into MET samples and ERCC Mix 1 would be spiked into CTL control replicates.

The dilution factor of the pure Ambion ERCC mixes prior to spiking into total RNA samples is `erccdilution`. The amount of diluted ERCC mix spiked into the total RNA sample is `spikeVol` (units are μL). The mass of total RNA spiked with the diluted ERCC mix is `totalRNAmass` (units are μg).

The final required input parameter, `choseFDR`, is the False Discovery Rate (FDR) for differential expression testing. A typical choice would be 0.05 (5% FDR), so this is the default `choseFDR` value. For the rat data since most genes are not differentially expressed a less conservative FDR is chosen ($\text{FDR} = 0.1$) and for the UHRR and HBRR reference RNA samples $\text{FDR} = 0.01$ is chosen, because there is a large number of differentially expressed genes for this pair of samples.

The function `runDashboard.R` is provided for convenient default `erccdashboard` analysis. Execution of the `runDashboard` function calls the default functions for `erccdashboard` analysis and reports parameters and progress to the R console. The functions called within `runDashboard.R` are also available to the user (details provided in **Section 4**).

All data and analysis results are stored in the list object `exDat`. For convenience the main diagnostic figures are saved to a pdf file and the `exDat` object is saved to an `.RData` object named using the `filenameRoot` provided by the user.

Use the following command to run the default `runDashboard` script:

```
> exDat <- runDashboard(datType="count", isNorm = FALSE,
                        exTable=MET.CTL.countDat,
                        filenameRoot="RatTox", sample1Name="MET",
                        sample2Name="CTL", erccmix="RatioPair",
                        erccdilution=1/100, spikeVol=1,
                        totalRNAmass=0.500, choseFDR=0.1)
```

Initializing the `exDat` list structure...

`choseFDR` = 0.1

`repNormFactor` is NULL

Filename root is: RatTox.MET.CTL

Transcripts were removed with a mean count < 1 or more than 2 replicates with 0 counts.

Original data contained 16590 transcripts.

After filtering 11570 transcripts remain for analysis.

A total of 29 out of 92

ERCC controls were filtered from the data set

The excluded ERCCs are:

ERCC-00012 ERCC-00014 ERCC-00016 ERCC-00017 ERCC-00024

ERCC-00041 ERCC-00048 ERCC-00057 ERCC-00061 ERCC-00073

ERCC-00075 ERCC-00081 ERCC-00083 ERCC-00086 ERCC-00097

ERCC-00098 ERCC-00104 ERCC-00117 ERCC-00120 ERCC-00123

ERCC-00126 ERCC-00134 ERCC-00137 ERCC-00138 ERCC-00142

ERCC-00147 ERCC-00150 ERCC-00156 ERCC-00164

```

repNormFactor is NULL,
  Using Default Upper Quartile Normalization Method - 75th percentile

normVec:
438 517 473 397 546 389
Check for sample mRNA fraction differences(r_m)...

Number of ERCC Controls Used in r_m estimate
63

Outlier ERCCs for GLM r_m Estimate:
ERCC-00062 ERCC-00060 ERCC-00022 ERCC-00043 ERCC-00003
ERCC-00004 ERCC-00046

GLM log(r_m) estimate:
0.2340727

GLM log(r_m) estimate weighted s.e.:
0.3347723

Number of ERCCs in Mix 1 dyn range: 63

Number of ERCCs in Mix 2 dyn range: 63
These ERCCs were not included in the signal-abundance plot,
because not enough non-zero replicate measurements of these
controls were obtained for both samples:

ERCC-00058 ERCC-00067 ERCC-00077 ERCC-00168 ERCC-00028
ERCC-00033 ERCC-00040 ERCC-00109 ERCC-00154 ERCC-00158

Saving dynRangePlot to exDat

Starting differential expression tests

Show log.offset
6.082219 6.248043 6.159095 5.983936 6.302619 5.963579
[1] 1
[1] 2
[1] 3
[1] 4
Disp = 0.06252 , BCV = 0.25
Disp = 0.06249 , BCV = 0.25
Finished DE testing
Finished examining dispersions

Threshold P-value
0.006370529

Generating ROC curve and AUC statistics...

```

Area Under the Curve (AUC) Results:

Ratio	AUC	Detected	Spiked
4:1	0.983	16	23
1:1.5	0.408	16	23
1:2	0.708	16	23

Estimating ERCC LODR

```
.....
Ratio LODR Estimate 90% CI Lower Bound 90% CI Upper Bound
4:1          22          16          24
1:1.5        Inf          <NA>          <NA>
1:2          Inf          <NA>          <NA>
```

LODR estimates are available to code ratio-abundance plot

Saving main dashboard plots to pdf file...

Saving exDat list to .RData file...

Analysis completed.

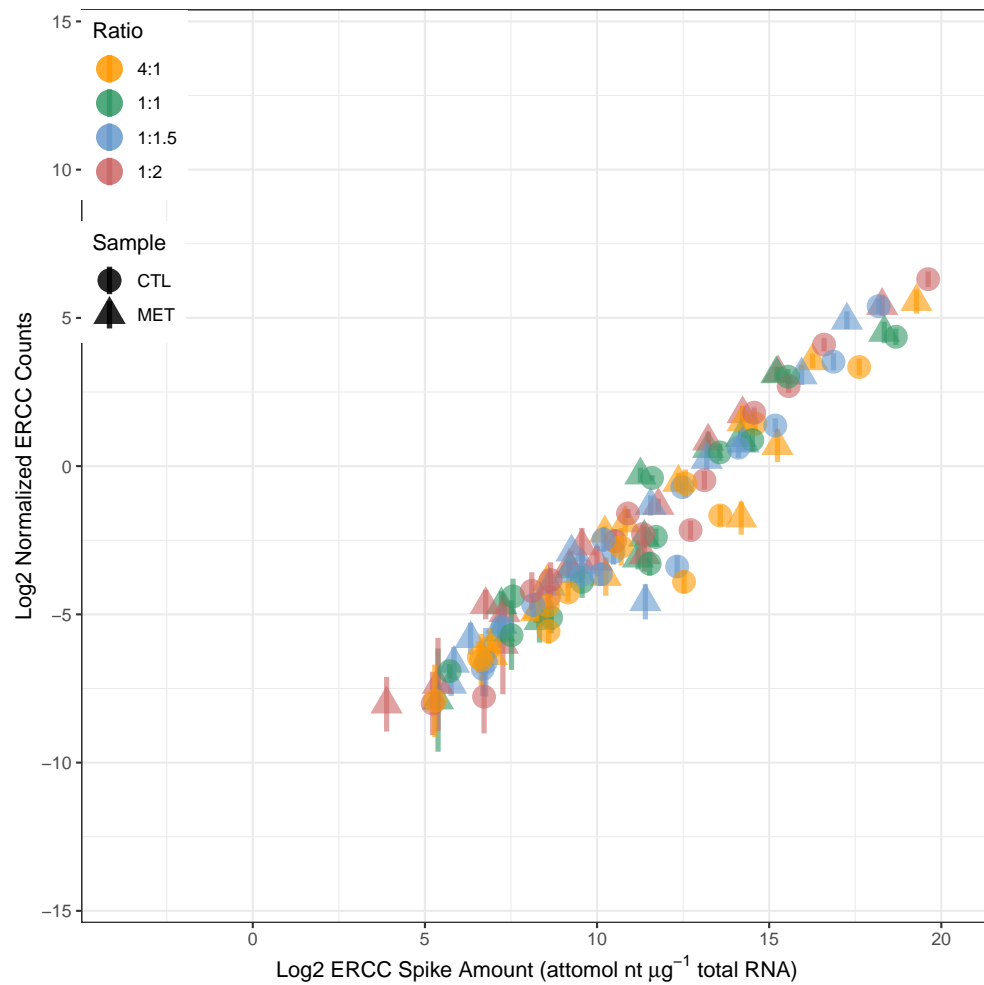
1.3 Results of dashboard analysis

The summary function will give a top level view of the exDat list structure. The str function will give more detail. It is a good idea to set the max.level argument in the str function, because by the end of the analysis the exDat structure is quite large.

```
> summary(exDat)
      Length Class      Mode
sampleInfo    11  -none-    list
plotInfo       9  -none-    list
erccInfo       4  -none-    list
Transcripts    7  data.frame list
designMat       3  data.frame list
sampleNames    2  -none-   character
idCols        6  data.frame list
normERCCDat    7  data.frame list
normFactor     6  -none-    numeric
mnLibeFactor   1  -none-    numeric
spikeFraction  1  -none-    numeric
idColsAdj      6  data.frame list
Results       12  -none-    list
Figures        7  -none-    list
```

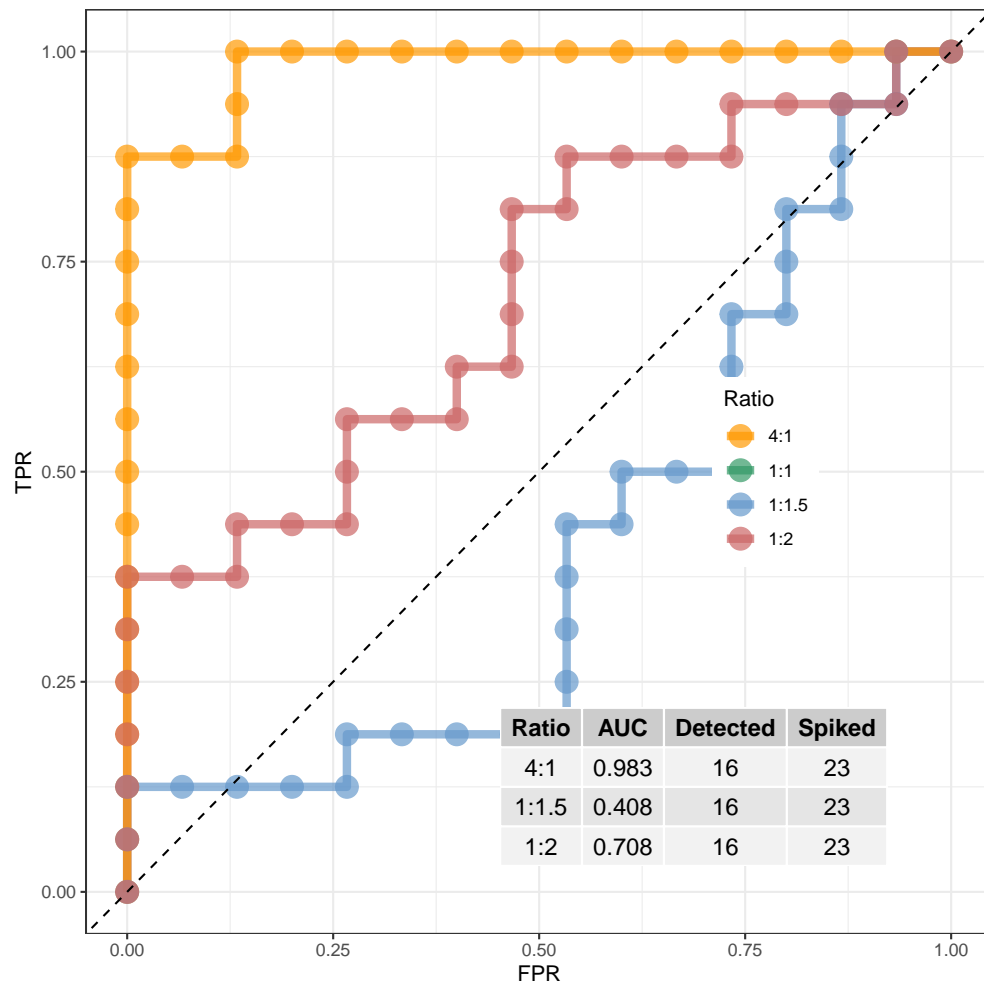
The figures from the analysis are stored in `exDat$Figures`. The four main diagnostic figures that are saved to a pdf file are the `dynRangePlot`, `rocPlot`, `lodrERCCPlot`, and `maPlot`.

```
> grid.arrange(exDat$Figures$dynRangePlot)
```



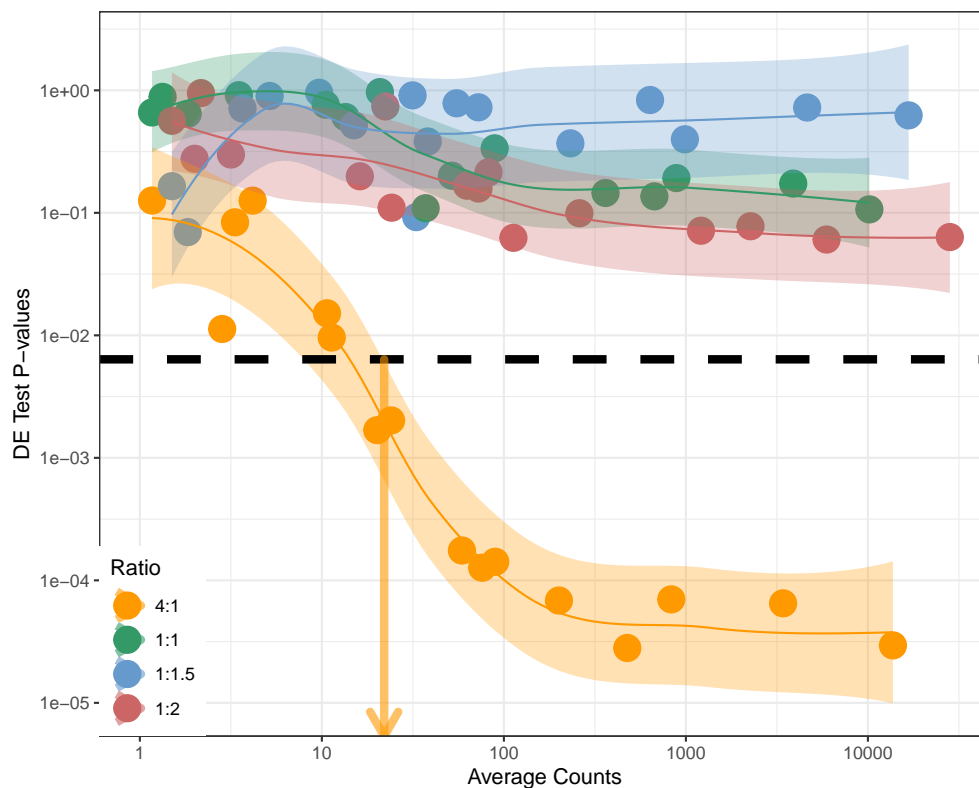
For this particular experiment the relationship between abundance and signal for the ERCC controls shows that the measurement results span a 2^{15} dynamic range. These ERCC mixtures were designed to span a 2^{20} dynamic range, but there was insufficient evidence to reliably quantify ERCC transcripts at low abundances.

```
> grid.arrange(exDat$Figures$rocPlot)
```



The receiver operator characteristic (ROC) curve and the Area Under the Curve (AUC) statistic provide evidence of the diagnostic power for detecting differential expression in this rat toxicogenomics experiment. As expected with increased fold change, diagnostic power increases. The AUC summary statistic for different experiments can be used to compare diagnostic performance.

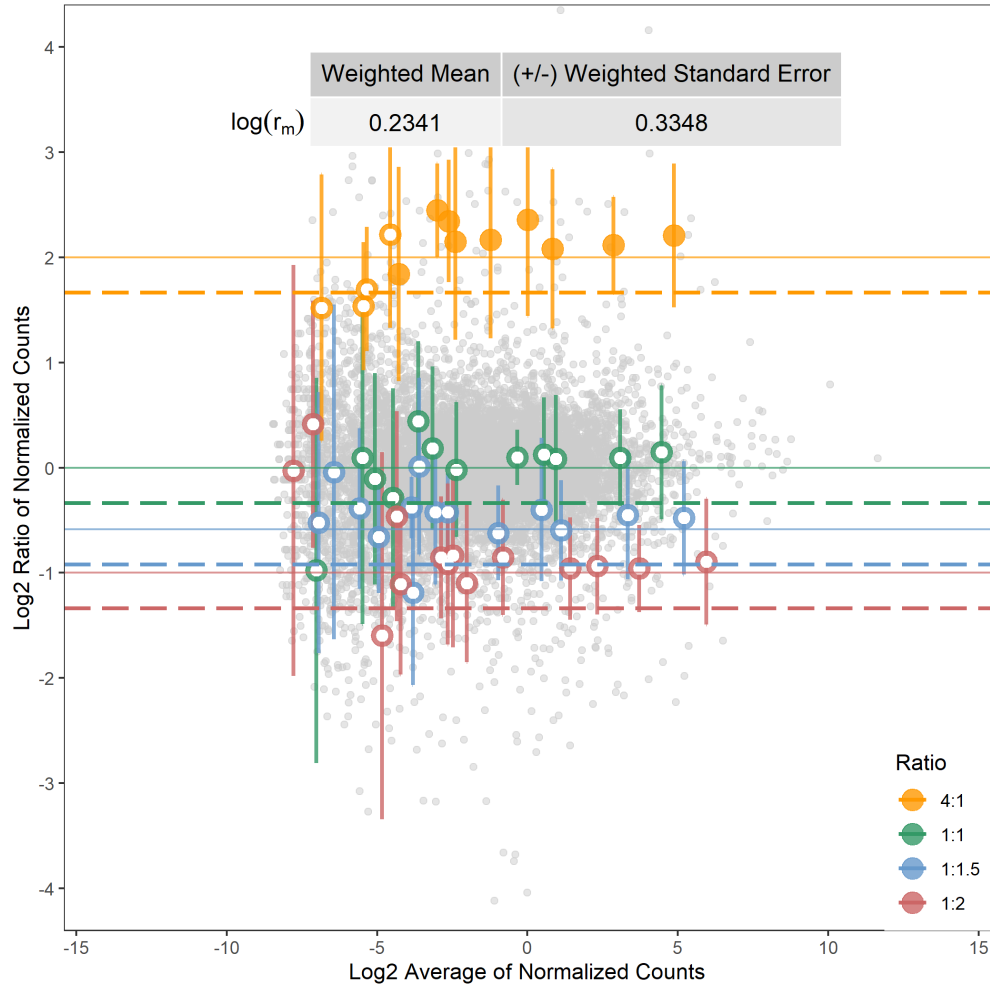
```
> grid.arrange(exDat$Figures$lodrERCCPlot)
```



Ratio	LODR Estimate	90% CI Lower Bound	90% CI Upper Bound
4:1	22	16	24
1:1.5	Inf	NA	NA
1:2	Inf	NA	NA

By modeling the relationship between average signal and p-values we can obtain Limit of Detection of Ratios (LODR) estimates for each differential fold change (or Ratio, indicated by color) and a threshold p-value, `p.thresh`, indicated by the dotted black line. LODR values can be compared between experiments to evaluate detection of differences between samples as a function of transcript abundance. If the input data used in `erccdashboard` analysis are unnormalized RNA-Seq counts, then the LODR estimates from the `lodrERCCPlot` are unnormalized counts. If normalized RNA-Seq data or microarray data are used for analysis then the LODR estimates will have corresponding units.


```
> grid.arrange(exDat$Figures$maPlot)
```



2 Comparison of Performance Between Experiments

The performance metrics provided here derived from measurements of ERCC ratios in gene expression experiments (AUC, LODR, r_m , and the standard deviations of the ER

```

exDat <- initDat(datType = datType, isNorm = isNorm, exTable = exTable,
  repNormFactor = repNormFactor, filenameRoot = filenameRoot,
  sample1Name = sample1Name, sample2Name = sample2Name,
  erccmix = erccmix, erccd
```

3.2 Options for LODR Estimation

The default behavior of `runDashboard` is to use the `estLODR` function to

attached base packages:

```
[1] stats      graphics  grDevices  utils      datasets  
[6] methods    base
```