

Package ‘edgeR’

April 9, 2015

Version 3.8.6

Date 2015/03/09

Title Empirical analysis of digital gene expression data in R

Author Yunshun Chen <yuchen@wehi.edu.au>, Davis McCarthy <dmccarthy@wehi.edu.au>, Aaron Lun <alun@wehi.edu.au>, Xiaobei Zhou <xiaobei.zhou@uzh.ch>, Mark Robinson <mark.robinson@imls.uzh.ch>, Gordon Smyth <smyth@wehi.edu.au>

Maintainer Yunshun Chen <yuchen@wehi.edu.au>, Aaron Lun <alun@wehi.edu.au>, Mark Robinson <mark.robinson@imls.uzh.ch>, Davis McCarthy <dmccarthy@wehi.edu.au>, Gordon Smyth <smyth@wehi.edu.au>

License GPL (>=2)

Depends R (>= 2.15.0), limma

Imports methods

Suggests MASS, statmod, splines, locfit, KernSmooth

URL <http://bioinf.wehi.edu.au/edgeR>

biocViews GeneExpression, Transcription, AlternativeSplicing, Coverage, DifferentialExpression, DifferentialSplicing, GeneSetEnrichment, Genetics, Bayesian, Clustering, Regression, TimeCourse, SAGE, Sequencing, ChIPSeq, RNASeq, BatchEffect, MultipleComparison, Normalization, QualityControl

Description Differential expression analysis of RNA-seq and digital gene expression profiles with biological replication. Uses empirical Bayes estimation and exact tests based on the negative binomial distribution. Also useful for differential signal analysis with other types of genome-scale count data.

R topics documented:

edgeR-package	3
adjustedProfileLik	4
as.data.frame	6
as.matrix	7
aveLogCPM	7

binomTest	9
calcNormFactors	10
camera.DGEList	12
commonCondLogLikDerDelta	14
condLogLikDerSize	15
cpm	16
cutWithMinN	17
decideTestsDGE	18
DGEEexact-class	19
DGEGLM-class	20
DGEList	21
DGEList-class	22
DGELRT-class	23
dglmStdResid	24
diffSpliceDGE	27
dim	29
dimnames	30
dispBinTrend	31
dispCoxReid	33
dispCoxReidInterpolateTagwise	35
dispCoxReidSplineTrend	37
edgeRUsersGuide	38
equalizeLibSizes	39
estimateCommonDisp	41
estimateDisp	42
estimateExonGenewiseDisp	44
estimateGLMCommonDisp	45
estimateGLMRobustDisp	47
estimateGLMTagwiseDisp	48
estimateGLMTrendedDisp	50
estimateTagwiseDisp	52
estimateTrendedDisp	54
exactTest	55
expandAsMatrix	58
getCounts	58
getPriorN	59
glmFit	61
glmQLFit	64
goana.DGELRT	66
gof	68
goodTuring	70
loessByCol	71
maPlot	72
maximizeInterpolant	74
maximizeQuadratic	75
meanvar	76
mglm	78
movingAverageByCol	81

nbinomDeviance	82
normalizeChIPtoInput	83
plotBCV	84
plotExonUsage	85
plotMDS.DGEList	86
plotQLDisp	88
plotSmear	90
plotSpliceDGE	91
predFC	92
processAmplicons	94
q2qnbinom	96
readDGE	97
roast.DGEList	99
spliceVariants	100
splitIntoGroups	102
subsetting	103
sumTechReps	104
systematicSubset	105
thinCounts	106
topSpliceDGE	107
topTags	108
treatDGE	109
validDGEList	111
weightedCondLogLikDerDelta	112
WLEB	113
zscoreNBinom	114

Index**116**

edgeR-package

*Empirical analysis of digital gene expression data in R***Description**

edgeR is a package for the analysis of digital gene expression data arising from RNA sequencing technologies such as SAGE, CAGE, Tag-seq or RNA-seq, with emphasis on testing for differential expression.

Particular strengths of the package include the ability to estimate biological variation between replicate libraries, and to conduct exact tests of significance which are suitable for small counts. The package is able to make use of even minimal numbers of replicates.

An extensive User's Guide is available, and can be opened by typing `edgeRUsersGuide()` at the R prompt. Detailed help pages are also provided for each individual function.

The edgeR package implements original statistical methodology described in the publications below.

Author(s)

Mark Robinson <mrobinson@wehi.edu.au>, Davis McCarthy <dmccarthy@wehi.edu.au>, Yunshun Chen <yuchen@wehi.edu.au>, Aaron Lun <alun@wehi.edu.au>, Gordon Smyth

References

- Robinson MD and Smyth GK (2007). Moderated statistical tests for assessing differences in tag abundance. *Bioinformatics* 23, 2881-2887
- Robinson MD and Smyth GK (2008). Small-sample estimation of negative binomial dispersion, with applications to SAGE data. *Biostatistics*, 9, 321-332
- Robinson MD, McCarthy DJ and Smyth GK (2010). edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics* 26, 139-140
- McCarthy, DJ, Chen, Y, Smyth, GK (2012). Differential expression analysis of multifactor RNA-Seq experiments with respect to biological variation. *Nucleic Acids Research* 40, 4288-4297.
- Lund, SP, Nettleton, D, McCarthy, DJ, Smyth, GK (2012). Detecting differential expression in RNA-sequence data using quasi-likelihood with shrunken dispersion estimates. *Statistical Applications in Genetics and Molecular Biology*. (Accepted 31 July 2012)

adjustedProfileLik	<i>Adjusted Profile Likelihood for the Negative Binomial Dispersion Parameter</i>
--------------------	---

Description

Compute adjusted profile-likelihoods for estimating the dispersion parameters of genewise negative binomial GLMs.

Usage

```
adjustedProfileLik(dispersion, y, design, offset, weights=NULL, adjust=TRUE,
                  start=NULL, get.coef=FALSE)
```

Arguments

dispersion	numeric scalar or vector of dispersions.
y	numeric matrix of counts.
design	numeric matrix giving the design matrix.
offset	numeric matrix of same size as y giving offsets for the log-linear models. Can be a scalar or a vector of length ncol(y), in which case it is expanded out to a matrix.
weights	optional numeric matrix giving observation weights.
adjust	logical, if TRUE then Cox-Reid adjustment is made to the log-likelihood, if FALSE then the log-likelihood is returned without adjustment.
start	numeric matrix of starting values for the GLM coefficients, to be passed to glmFit .
get.coef	logical, specifying whether fitted GLM coefficients should be returned.

Details

For each row of data, compute the adjusted profile-likelihood for estimating the dispersion parameter of the negative binomial glm. The adjusted profile likelihood is described by McCarthy et al (2012), and is based on the method of Cox and Reid (1987).

The adjusted profile likelihood is an approximate log-likelihood for the dispersion parameter, conditional on the estimated values of the coefficients in the NB log-linear models. The conditional likelihood approach is a technique for adjusting the likelihood function to allow for the fact that nuisance parameters have to be estimated in order to evaluate the likelihood. When estimating the dispersion, the nuisance parameters are the coefficients in the linear model.

This implementation calls the LAPACK library to perform the Cholesky decomposition during adjustment estimation.

The purpose of `start` and `get.coef` is to allow hot-starting for multiple calls to `adjustedProfileLik`, when only the dispersion is altered. Specifically, the returned GLM coefficients from one call with `get.coef==TRUE` can be used as the `start` values for the next call.

Value

If `get.coef==FALSE`, a vector of adjusted profile log-likelihood values is returned containing one element for each row of `y`.

Otherwise, a list is returned containing `apl`, the aforementioned vector of adjusted profile likelihoods; and `beta`, a numeric matrix of fitted GLM coefficients.

Author(s)

Yunshun Chen, Gordon Smyth, Aaron Lun

References

Cox, DR, and Reid, N (1987). Parameter orthogonality and approximate conditional inference. *Journal of the Royal Statistical Society Series B* 49, 1-39.

McCarthy, DJ, Chen, Y, Smyth, GK (2012). Differential expression analysis of multifactor RNA-Seq experiments with respect to biological variation. *Nucleic Acids Research* 40, 4288-4297. <http://nar.oxfordjournals.org/content/40/10/4288>

See Also

[glmFit](#)

Examples

```
y <- matrix(rnbinom(1000, mu=10, size=2), ncol=4)
design <- matrix(1, 4, 1)
dispersion <- 0.5
apl <- adjustedProfileLik(dispersion, y, design, offset=0)
apl
```

as.data.frame *Turn a TopTags Object into a Dataframe*

Description

Turn a TopTags object into a data.frame.

Usage

```
## S3 method for class TopTags
as.data.frame(x, row.names = NULL, optional = FALSE, ...)
```

Arguments

x	an object of class TopTags
row.names	NULL or a character vector giving the row names for the data frame. Missing values are not allowed.
optional	logical. If TRUE, setting row names and converting column names (to syntactic names) is optional.
...	additional arguments to be passed to or from methods.

Details

This method combines all the components of x which have a row for each tag (transcript) into a data.frame.

Value

A data.frame.

Author(s)

Gordon Smyth

See Also

[as.data.frame](#) in the base package.

as.matrix	<i>Turn a DGEList Object into a Matrix</i>
-----------	--

Description

Coerce a digital gene expression object into a numeric matrix by extracting the count values.

Usage

```
## S3 method for class DGEList  
as.matrix(x,...)
```

Arguments

x	an object of class DGEList.
...	additional arguments, not used for these methods.

Details

This method extracts the matrix of counts.

This involves loss of information, so the original data object is not recoverable.

Value

A numeric matrix.

Author(s)

Gordon Smyth

See Also

[as.matrix](#) in the base package or [as.matrix.RGList](#) in the limma package.

aveLogCPM	<i>Average Log Counts Per Million</i>
-----------	---------------------------------------

Description

Compute average log₂ counts-per-million for each row of counts.

Usage

```
## S3 method for class DGEList
aveLogCPM(y, normalized.lib.sizes=TRUE, prior.count=2, dispersion=NULL, ...)
## Default S3 method:
aveLogCPM(y, lib.size=NULL, offset=NULL, prior.count=2, dispersion=NULL,
           weights=NULL, ...)
```

Arguments

<code>y</code>	numeric matrix containing counts. Rows for tags and columns for libraries.
<code>normalized.lib.sizes</code>	logical, use normalized library sizes?
<code>prior.count</code>	average value to be added to each count, to avoid infinite values on the log-scale.
<code>dispersion</code>	numeric scalar or vector of negative-binomial dispersions. Defaults to 0.05.
<code>lib.size</code>	numeric vector of library sizes. Defaults to <code>colSums(y)</code> . Ignored if <code>offset</code> is not NULL.
<code>offset</code>	numeric matrix of offsets for the log-linear models.
<code>weights</code>	optional numeric matrix of observation weights.
<code>...</code>	other arguments are not currently used.

Details

This function uses `mg1mOneGroup` to compute average counts-per-million (AveCPM) for each row of counts, and returns `log2(AveCPM)`. An average value of `prior.count` is added to the counts before running `mg1mOneGroup`.

This function is similar to

```
log2(rowMeans(cpm(y, ...))),
```

but with the refinement that larger library sizes are given more weight in the average. The two versions will agree for large values of the dispersion.

Value

Numeric vector giving `log2(AveCPM)` for each row of `y`.

Author(s)

Gordon Smyth

See Also

See [cpm](#) for individual logCPM values, rather than tagwise averages.

The computations for `aveLogCPM` are done by [mg1mOneGroup](#).

Examples

```

y <- matrix(c(0,100,30,40),2,2)
lib.size <- c(1000,10000)

# With disp large, the function is equivalent to row-wise averages of individual cpms:
aveLogCPM(y, dispersion=1e4)
cpm(y, log=TRUE, prior.count=2)

# With disp=0, the function is equivalent to pooling the counts before dividing by lib.size:
aveLogCPM(y,prior.count=0,dispersion=0)
cpms <- rowSums(y)/sum(lib.size)*1e6
log2(cpms)

```

binomTest

*Exact Binomial Tests for Comparing Two Digital Libraries***Description**

Computes p-values for differential abundance for each tag between two digital libraries, conditioning on the total count for each tag. The counts in each group as a proportion of the whole are assumed to follow a binomial distribution.

Usage

```
binomTest(y1, y2, n1=sum(y1), n2=sum(y2), p=n1/(n1+n2))
```

Arguments

y1	integer vector giving counts in first library. Non-integer values are rounded to the nearest integer.
y2	integer vector giving counts in second library. Of same length as x. Non-integer values are rounded to the nearest integer.
n1	total number of tags in first library. Non-integer values are rounded to the nearest integer. Not required if p is supplied.
n2	total number of tags in second library. Non-integer values are rounded to the nearest integer. Not required if p is supplied.
p	expected proportion of y1 to the total under the null hypothesis.

Details

This function can be used to compare two libraries from SAGE, RNA-Seq, ChIP-Seq or other sequencing technologies with respect to technical variation.

An exact two-sided binomial test is computed for each tag. This test is closely related to Fisher's exact test for 2x2 contingency tables but, unlike Fisher's test, it conditions on the total number of counts for each tag. The null hypothesis is that the expected counts are in the same proportions as the library sizes, i.e., that the binomial probability for the first library is $n1/(n1+n2)$.

The two-sided rejection region is chosen analogously to Fisher's test. Specifically, the rejection region consists of those values with smallest probabilities under the null hypothesis.

When the counts are reasonably large, the binomial test, Fisher's test and Pearson's chisquare all give the same results. When the counts are smaller, the binomial test is usually to be preferred in this context.

This function replaces the earlier `sage.test` functions in the `statmod` and `sagenhaft` packages. It produces the same results as `binom.test` in the `stats` package, but is much faster.

Value

Numeric vector of p-values.

Author(s)

Gordon Smyth

References

http://en.wikipedia.org/wiki/Binomial_test

http://en.wikipedia.org/wiki/Fisher's_exact_test

http://en.wikipedia.org/wiki/Serial_analysis_of_gene_expression

<http://en.wikipedia.org/wiki/RNA-Seq>

See Also

[sage.test](#) (statmod package), [binom.test](#) (stats package)

Examples

```
binomTest(c(0,5,10),c(0,30,50),n1=10000,n2=15000)
# Univariate equivalents:
binom.test(5,5+30,p=10000/(10000+15000))$p.value
binom.test(10,10+50,p=10000/(10000+15000))$p.value
```

calcNormFactors

Calculate Normalization Factors to Align Columns of a Count Matrix

Description

Calculate normalization factors to scale the raw library sizes.

Usage

```
## S3 method for class DGEList
calcNormFactors(object, method=c("TMM","RLE","upperquartile","none"),
                refColumn=NULL, logratioTrim=.3, sumTrim=0.05, doWeighting=TRUE,
                Acutoff=-1e10, p=0.75, ...)

## Default S3 method:
calcNormFactors(object, lib.size=NULL, method=c("TMM","RLE",
        "upperquartile","none"), refColumn=NULL, logratioTrim=.3,
        sumTrim=0.05, doWeighting=TRUE, Acutoff=-1e10, p=0.75, ...)
```

Arguments

object	either a matrix of raw (read) counts or a DGEList object
lib.size	numeric vector of library sizes of the object.
method	normalization method to be used
refColumn	column to use as reference for method="TMM". Can be a column number or a numeric vector of length nrow(object).
logratioTrim	amount of trim to use on log-ratios ("M" values) for method="TMM"
sumTrim	amount of trim to use on the combined absolute levels ("A" values) for method="TMM"
doWeighting	logical, whether to compute (asymptotic binomial precision) weights for method="TMM"
Acutoff	cutoff on "A" values to use before trimming for method="TMM"
p	percentile (between 0 and 1) of the counts that is aligned when method="upperquartile"
...	further arguments that are not currently used.

Details

method="TMM" is the weighted trimmed mean of M-values (to the reference) proposed by Robinson and Oshlack (2010), where the weights are from the delta method on Binomial data. If refColumn is unspecified, the library whose upper quartile is closest to the mean upper quartile is used.

method="RLE" is the scaling factor method proposed by Anders and Huber (2010). We call it "relative log expression", as median library is calculated from the geometric mean of all columns and the median ratio of each sample to the median library is taken as the scale factor.

method="upperquartile" is the upper-quartile normalization method of Bullard et al (2010), in which the scale factors are calculated from the 75% quantile of the counts for each library, after removing transcripts which are zero in all libraries. This idea is generalized here to allow scaling by any quantile of the distributions.

If method="none", then the normalization factors are set to 1.

For symmetry, normalization factors are adjusted to multiply to 1. The effective library size is then the original library size multiplied by the scaling factor.

Note that rows that have zero counts for all columns are trimmed before normalization factors are computed. Therefore rows with all zero counts do not affect the estimated factors.

Value

If `object` is a matrix, the output is a vector with length `ncol(object)` giving the relative normalization factors. If `object` is a `DGEList`, then it is returned as output with the relative normalization factors in `object$samples$norm.factors`.

Author(s)

Mark Robinson, Gordon Smyth

References

Anders, S, Huber, W (2010). Differential expression analysis for sequence count data *Genome Biology* 11, R106.

Bullard JH, Purdom E, Hansen KD, Dudoit S. (2010) Evaluation of statistical methods for normalization and differential expression in mRNA-Seq experiments. *BMC Bioinformatics* 11, 94.

Robinson MD, Oshlack A (2010). A scaling normalization method for differential expression analysis of RNA-seq data. *Genome Biology* 11, R25.

Examples

```
y <- matrix( rpois(1000, lambda=5), nrow=200 )
calcNormFactors(y)
```

camera.DGEList	<i>Competitive Gene Set Test for Digital Gene Expression Data Accounting for Inter-gene Correlation</i>
----------------	---

Description

Test whether a set of genes is highly ranked relative to other genes in terms of differential expression, accounting for inter-gene correlation.

Usage

```
## S3 method for class DGEList
camera(y, index, design=NULL, contrast=ncol(design), ...)
```

Arguments

<code>y</code>	a <code>DGEList</code> object containing dispersion estimates.
<code>index</code>	an index vector or a list of index vectors. Can be any vector such that <code>y[indices,]</code> selects the rows corresponding to the test set.
<code>design</code>	the design matrix.
<code>contrast</code>	the contrast of the linear model coefficients for which the test is required. Can be an integer specifying a column of design, or else a numeric vector of same length as the number of columns of design.
<code>...</code>	other arguments are passed to <code>camera.default</code> .

Details

The camera gene set test was proposed by Wu and Smyth (2012) for microarray data. This function makes the camera test available for digital gene expression data. The negative binomial count data is converted to approximate normal deviates by computing mid-p quantile residuals (Dunn and Smyth, 1996; Routledge, 1994) under the null hypothesis that the contrast is zero. See [camera](#) for more description of the test and for a complete list of possible arguments.

The design matrix defaults to the `model.matrix(~y$samples$group)`.

Value

A data.frame. See [camera](#) for details.

Author(s)

Yunshun Chen, Gordon Smyth

References

Dunn, K. P., and Smyth, G. K. (1996). Randomized quantile residuals. *J. Comput. Graph. Statist.*, 5, 236-244. <http://www.statsci.org/smyth/pubs/residual.html>

Routledge, RD (1994). Practicing safe statistics with the mid-p. *Canadian Journal of Statistics* 22, 103-110.

Wu, D, and Smyth, GK (2012). Camera: a competitive gene set test accounting for inter-gene correlation. *Nucleic Acids Research* 40, e133. <http://nar.oxfordjournals.org/content/40/17/e133>

See Also

[roast.DGEList](#), [camera](#).

Examples

```
mu <- matrix(10, 100, 4)
group <- factor(c(0,0,1,1))
design <- model.matrix(~group)

# First set of 10 genes that are genuinely differentially expressed
iset1 <- 1:10
mu[iset1,3:4] <- mu[iset1,3:4]+10

# Second set of 10 genes are not DE
iset2 <- 11:20

# Generate counts and create a DGEList object
y <- matrix(rnbinom(100*4, mu=mu, size=10),100,4)
y <- DGEList(counts=y, group=group)

# Estimate dispersions
y <- estimateDisp(y, design)
```

```
camera(y, iset1, design)
camera(y, iset2, design)

camera(y, list(set1=iset1,set2=iset2), design)
```

commonCondLogLikDerDelta

Conditional Log-Likelihoods in Terms of Delta

Description

Common conditional log-likelihood parameterized in terms of delta ($\phi / (\phi+1)$)

Usage

```
commonCondLogLikDerDelta(y, delta, der = 0)
```

Arguments

y	list with elements comprising the matrices of count data (or pseudocounts) for the different groups
delta	delta ($\phi / (\phi+1)$) parameter of negative binomial
der	derivative, either 0 (the function), 1 (first derivative) or 2 (second derivative)

Details

The common conditional log-likelihood is constructed by summing over all of the individual tag conditional log-likelihoods. The common conditional log-likelihood is taken as a function of the dispersion parameter (ϕ), and here parameterized in terms of delta ($\phi / (\phi+1)$). The value of delta that maximizes the common conditional log-likelihood is converted back to the ϕ scale, and this value is the estimate of the common dispersion parameter used by all tags.

Value

numeric scalar of function/derivative evaluated at given delta

Author(s)

Davis McCarthy

See Also

[estimateCommonDisp](#) is the user-level function for estimating the common dispersion parameter.

Examples

```
counts<-matrix(rnbinom(20,size=1,mu=10),nrow=5)
d<-DGEList(counts=counts,group=rep(1:2,each=2),lib.size=rep(c(1000:1001),2))
y<-splitIntoGroups(d)
l11<-commonCondLogLikDerDelta(y,delta=0.5,der=0)
l12<-commonCondLogLikDerDelta(y,delta=0.5,der=1)
```

condLogLikDerSize	<i>Conditional Log-Likelihood of the Dispersion for a Single Group of Replicate Libraries</i>
-------------------	---

Description

Derivatives of the negative-binomial log-likelihood with respect to the dispersion parameter for each tag/transcript, conditional on the mean count, for a single group of replicate libraries of the same size.

Usage

```
condLogLikDerSize(y, r, der=1L)
condLogLikDerDelta(y, delta, der=1L)
```

Arguments

y	matrix of counts, all counts in each row having the same population mean
r	numeric vector or scalar, size parameter of negative binomial distribution, equal to 1/dispersion
delta	numeric vector or scalar, delta parameter of negative binomial, equal to dispersion/(1+dispersion)
der	integer specifying derivative required, either 0 (the function), 1 (first derivative) or 2 (second derivative)

Details

The library sizes must be equalized before running this function. This function carries out the actual mathematical computations for the conditional log-likelihood and its derivatives, calculating the conditional log-likelihood for each tag/transcript. Derivatives are with respect to either the size (*r*) or the delta parametrization (*delta*) of the dispersion.

Value

vector of row-wise derivatives with respect to *r* or *delta*

Author(s)

Mark Robinson, Davis McCarthy, Gordon Smyth

Examples

```
y <- matrix(rnbinom(10,size=1,mu=10),nrow=5)
condLogLikDerSize(y,r=1,der=1)
condLogLikDerDelta(y,delta=0.5,der=1)
```

cpm

*Counts per Million or Reads per Kilobase per Million***Description**

Computes counts per million (CPM) or reads per kilobase per million (RPKM) values.

Usage

```
## S3 method for class DGEList
cpm(x, normalized.lib.sizes=TRUE, log=FALSE, prior.count=0.25, ...)
## Default S3 method:
cpm(x, lib.size=NULL, log=FALSE, prior.count=0.25, ...)
## S3 method for class DGEList
rpkm(x, gene.length=NULL, normalized.lib.sizes=TRUE, log=FALSE, prior.count=0.25, ...)
## Default S3 method:
rpkm(x, gene.length, lib.size=NULL, log=FALSE, prior.count=0.25, ...)
```

Arguments

x	matrix of counts or a DGEList object
normalized.lib.sizes	logical, use normalized library sizes?
lib.size	library size, defaults to colSums(x).
log	logical, if TRUE then log2 values are returned.
prior.count	average count to be added to each observation to avoid taking log of zero. Used only if log=TRUE.
gene.length	vector of length nrow(x) giving gene length in bases, or the name of the column x\$genes containing the gene lengths.
...	other arguments are not currently used

Details

CPM or RPKM values are useful descriptive measures for the expression level of a gene or transcript. By default, the normalized library sizes are used in the computation for DGEList objects but simple column sums for matrices.

If log-values are computed, then a small count, given by prior.count but scaled to be proportional to the library size, is added to x to avoid taking the log of zero.

The rpkm method for DGEList objects will try to find the gene lengths in a column of x\$genes called Length or length. Failing that, it will look for any column name containing "length" in any capitalization.

Value

numeric matrix of CPM or RPKM values.

Note

aveLogCPM(x), rowMeans(cpm(x, log=TRUE)) and log2(rowMeans(cpm(x))) all give slightly different results.

Author(s)

Davis McCarthy, Gordon Smyth

See Also

[aveLogCPM](#)

Examples

```
y <- matrix(rnbinom(20, size=1, mu=10), 5, 4)
cpm(y)

d <- DGEList(counts=y, lib.size=1001:1004)
cpm(d)
cpm(d, log=TRUE)

d$genes$Length <- c(1000, 2000, 500, 1500, 3000)
rpkm(d)
```

cutWithMinN

Cut numeric vector into non-empty intervals

Description

Discretizes a numeric vector. Divides the range of x into intervals, so that each interval contains a minimum number of values, and codes the values in x according to which interval they fall into.

Usage

```
cutWithMinN(x, intervals=2, min.n=1)
```

Arguments

x	numeric vector.
intervals	number of intervals required.
min.n	minimum number of values in any interval. Must be greater than length(x)/intervals.

Details

This function strikes a compromise between the base functions `cut`, which by default cuts a vector into equal length intervals, and `quantile`, which is suited to finding equally populated intervals. It finds a partition of the `x` values that is as close as possible to equal length intervals while keeping at least `min.n` values in each interval.

Tied values of `x` are broken by random jittering, so the partition may vary slightly from run to run if there are many tied values.

Value

A list with components:

<code>group</code>	integer vector of same length as <code>x</code> indicating which interval each value belongs to.
<code>breaks</code>	numeric vector of length <code>intervals+1</code> giving the left and right limits of each interval.

Author(s)

Gordon Smyth

See Also

[cut](#), [quantile](#).

Examples

```
x <- c(1,2,3,4,5,6,7,100)
cutWithMinN(x,intervals=3,min.n=1)
```

decideTestsDGE

Multiple Testing Across Genes and Contrasts

Description

Classify a series of related differential expression statistics as up, down or not significant. A number of different multiple testing schemes are offered which adjust for multiple testing down the genes as well as across contrasts for each gene.

Usage

```
decideTestsDGE(object, adjust.method="BH", p.value=0.05, lfc=0)
```

Arguments

<code>object</code>	deDGEList object, output from <code>exactTest</code> , or DGELRT object, output from <code>DGELRT</code> , from which p-values for differential expression and log-fold change values may be extracted.
<code>adjust.method</code>	character string specifying p-value adjustment method. Possible values are "none", "BH", "fdr" (equivalent to "BH"), "BY" and "holm". See p.adjust for details.
<code>p.value</code>	numeric value between 0 and 1 giving the desired size of the test
<code>lfc</code>	numeric value giving the desired absolute minimum log-fold-change

Details

These functions implement multiple testing procedures for determining whether each log-fold change in a matrix of log-fold changes should be considered significantly different from zero.

Value

An object of class `TestResults` (see [TestResults](#)). This is essentially a numeric matrix with elements -1, 0 or 1 depending on whether each DE p-value is classified as significant with negative log-fold change, not significant or significant with positive log-fold change, respectively.

Author(s)

Davis McCarthy, Gordon Smyth

See Also

Adapted from [decideTests](#) in the `limma` package.

DGEEexact-class

differential expression of Digital Gene Expression data - class

Description

A list-based S4 class for for storing results of a differential expression analysis for DGE data.

List Components

For objects of this class, rows correspond to genomic features and columns to statistics associated with the differential expression analysis. The genomic features are called genes, but in reality might correspond to transcripts, tags, exons etc.

Objects of this class contain the following list components:

`table`: data frame containing columns for the log₂-fold-change, logFC, the average log₂-counts-per-million, logCPM, and the two-sided p-value PValue.

`comparison`: vector giving the two experimental groups/conditions being compared.

`genes`: a data frame containing information about each gene (can be NULL).

Methods

This class inherits directly from class `list`, so `DGEEExact` objects can be manipulated as if they were ordinary lists. However they can also be treated as if they were matrices for the purposes of subsetting.

The dimensions, row names and column names of a `DGEEExact` object are defined by those of `table`, see `dim.DGEEExact` or `dimnames.DGEEExact`.

`DGEEExact` objects can be subsetted, see [subsetting](#).

`DGEEExact` objects also have a `show` method so that printing produces a compact summary of their contents.

Author(s)

edgeR team. First created by Mark Robinson and Davis McCarthy

See Also

Other classes defined in edgeR are [DGEList-class](#), [DGEGLM-class](#), [DGEGLRT-class](#), [TopTags-class](#)

DGEGLM-class

Digital Gene Expression Generalized Linear Model results - class

Description

A list-based S4 class for storing results of a GLM fit to each gene in a DGE dataset.

List Components

For objects of this class, rows correspond to genomic features and columns to coefficients in the linear model. The genomic features are called genes, but in reality might correspond to transcripts, tags, exons etc.

Objects of this class contain the following list components:

`coefficients`: matrix containing the coefficients computed from fitting the model defined by the design matrix to each gene in the dataset.

`df.residual`: vector containing the residual degrees of freedom for the model fit to each gene in the dataset.

`deviance`: vector giving the deviance from the model fit to each gene.

`design`: design matrix for the full model from the likelihood ratio test.

`offset`: scalar, vector or matrix of offset values to be included in the GLMs for each gene.

`samples`: data frame containing information about the samples comprising the dataset.

`genes`: data frame containing information about the genes or tags for which we have DGE data (can be `NULL` if there is no information available).

`dispersion`: scalar or vector providing the value of the dispersion parameter used in the negative binomial GLM for each gene.

`lib.size`: vector providing the effective library size for each sample in the dataset.

`weights`: matrix of weights used in the GLM fitting for each gene.

`fitted.values`: the fitted (expected) values from the GLM for each gene.

`AveLogCPM`: numeric vector giving average log2 counts per million for each gene.

Methods

This class inherits directly from class `list` so any operation appropriate for lists will work on objects of this class.

The dimensions, row names and column names of a `DGEGLM` object are defined by those of the dataset, see `dim.DGEGLM` or `dimnames.DGEGLM`.

`DGEGLM` objects can be subsetted, see [subsetting](#).

`DGEGLM` objects also have a `show` method so that printing produces a compact summary of their contents.

Author(s)

edgeR team. First created by Davis McCarthy.

See Also

Other classes defined in edgeR are [DGEList-class](#), [DGEExact-class](#), [DGELRT-class](#), [TopTags-class](#)

DGEList

DGEList Constructor

Description

Creates a `DGEList` object from a table of counts (rows=features, columns=samples), group indicator for each column, library size (optional) and a table of feature annotation (optional).

Usage

```
DGEList(counts = matrix(0, 0, 0), lib.size = colSums(counts),
        norm.factors = rep(1, ncol(counts)), group = rep(1, ncol(counts)), genes = NULL,
        remove.zeros = FALSE)
```

Arguments

<code>counts</code>	numeric matrix of read counts.
<code>lib.size</code>	numeric vector giving the total count (sequence depth) for each library.
<code>norm.factors</code>	numeric vector of normalization factors that modify the library sizes.
<code>group</code>	vector or factor giving the experimental group/condition for each sample/library.
<code>genes</code>	data frame containing annotation information for the tags/transcripts/genes.
<code>remove.zeros</code>	logical, whether to remove rows that have 0 total count.

Details

To facilitate programming pipelines, NULL values can be input for `lib.size`, `norm.factors` or `group`, in which case the default value is used as if the argument had been missing.

Value

a `DGEList` object

Author(s)

edgeR team. First created by Mark Robinson.

See Also

[DGEList-class](#)

Examples

```
y <- matrix(rnbinom(10000,mu=5,size=2),ncol=4)
d <- DGEList(counts=y, group=rep(1:2,each=2))
dim(d)
colnames(d)
d$samples
```

DGEList-class

Digital Gene Expression data - class

Description

A list-based S4 class for storing read counts and associated information from digital gene expression or sequencing technologies.

List Components

For objects of this class, rows correspond to genomic features and columns to samples. The genomic features are called genes, but in reality might correspond to transcripts, tags, exons etc. Objects of this class contain the following essential list components:

`counts`: numeric matrix of read counts, one row for each gene and one column for each sample.

`samples`: `data.frame` with a row for each sample and columns `group`, `lib.size` and `norm.factors` containing the group labels, library sizes and normalization factors. Other columns can be optionally added to give more detailed sample information.

Optional components include:

`genes`: `data.frame` giving annotation information for each gene. Same number of rows as `counts`.

`AveLogCPM`: numeric vector giving average log₂ counts per million for each gene.

`common.dispersion`: numeric scalar giving the overall dispersion estimate.

trended.dispersion: numeric vector giving trended dispersion estimates for each gene.
 tagwise.dispersion: numeric vector giving tagwise dispersion estimates for each gene.
 offset: numeric matrix of same size as counts giving offsets for use in log-linear models.

Methods

This class inherits directly from class `list`, so `DGEList` objects can be manipulated as if they were ordinary lists. However they can also be treated as if they were matrices for the purposes of subsetting.

The dimensions, row names and column names of a `DGEList` object are defined by those of counts, see `dim.DGEList` or `dimnames.DGEList`.

`DGEList` objects can be subsetted, see `subsetting`.

`DGEList` objects also have a `show` method so that printing produces a compact summary of their contents.

Author(s)

edgeR team. First created by Mark Robinson.

See Also

`DGEList` constructs `DGEList` objects. Other classes defined in edgeR are `DGEEexact-class`, `DGEGLM-class`, `DGELRT-class`, `TopTags-class`

DGELRT-class

Digital Gene Expression Likelihood Ratio Test data and results - class

Description

A list-based S4 class for storing results of a GLM-based differential expression analysis for DGE data.

List Components

For objects of this class, rows correspond to genomic features and columns to statistics associated with the differential expression analysis. The genomic features are called genes, but in reality might correspond to transcripts, tags, exons etc.

Objects of this class contain the following list components:

`table`: data frame containing the log-concentration (i.e. expression level), the log-fold change in expression between the two groups/conditions and the exact p-value for differential expression, for each gene.

`coefficients.full`: matrix containing the coefficients computed from fitting the full model (fit using `glmFit` and a given design matrix) to each gene in the dataset.

`coefficients.null`: matrix containing the coefficients computed from fitting the null model to each gene in the dataset. The null model is the model to which the full model is compared, and is fit using `glmFit` and dropping selected column(s) (i.e. coefficient(s)) from the design matrix for the full model.

`design`: design matrix for the full model from the likelihood ratio test.

`...`: if the argument `y` to `glmLRT` (which produces the `DGELRT` object) was itself a `DGEList` object, then the `DGELRT` will contain all of the elements of `y`, except for the table of counts and the table of pseudocounts.

Methods

This class inherits directly from class `list`, so `DGELRT` objects can be manipulated as if they were ordinary lists. However they can also be treated as if they were matrices for the purposes of subsetting.

The dimensions, row names and column names of a `DGELRT` object are defined by those of `table`, see `dim.DGELRT` or `dimnames.DGELRT`.

`DGELRT` objects can be subsetted, see [subsetting](#).

`DGELRT` objects also have a `show` method so that printing produces a compact summary of their contents.

Author(s)

edgeR team. First created by Davis McCarthy

See Also

Other classes defined in edgeR are [DGEList-class](#), [DGEExact-class](#), [DGEGLM-class](#), [TopTags-class](#)

dglmStdResid

Visualize the mean-variance relationship in DGE data using standardized residuals

Description

Appropriate modelling of the mean-variance relationship in DGE data is important for making inferences about differential expression. However, the standard approach to visualizing the mean-variance relationship is not appropriate for general, complicated experimental designs that require generalized linear models (GLMs) for analysis. Here are functions to compute standardized residuals from a Poisson GLM and plot them for bins based on overall expression level of tags as a way to visualize the mean-variance relationship. A rough estimate of the dispersion parameter can also be obtained from the standardized residuals.

Usage

```
dglmStdResid(y, design, dispersion=0, offset=0, nbins=100, make.plot=TRUE,
             xlab="Mean", ylab="Ave. binned standardized residual", ...)
getDispersions(binned.object)
```


Arguments

<code>y</code>	numeric matrix of counts, each row represents one tag, each column represents one DGE library.
<code>design</code>	numeric matrix giving the design matrix of the GLM. Assumed to be full column rank.
<code>dispersion</code>	numeric scalar or vector giving the dispersion parameter for each GLM. Can be a scalar giving one value for all tags, or a vector of length equal to the number of tags giving tag-wise dispersions.
<code>offset</code>	numeric vector or matrix giving the offset that is to be included in the log-linear model predictor. Can be a vector of length equal to the number of libraries, or a matrix of the same size as <code>y</code> .
<code>nbins</code>	scalar giving the number of bins (formed by using the quantiles of the genewise mean expression levels) for which to compute average means and variances for exploring the mean-variance relationship. Default is 100 bins
<code>make.plot</code>	logical, whether or not to plot the mean standardized residual for binned data (binned on expression level). Provides a visualization of the mean-variance relationship. Default is TRUE.
<code>xlab</code>	character string giving the label for the x-axis. Standard graphical parameter. If left as the default, then the x-axis label will be set to "Mean".
<code>ylab</code>	character string giving the label for the y-axis. Standard graphical parameter. If left as the default, then the y-axis label will be set to "Ave. binned standardized residual".
<code>...</code>	further arguments passed on to <code>plot</code>
<code>binned.object</code>	list object, which is the output of <code>dglmStdResid</code> .

Details

This function is useful for exploring the mean-variance relationship in the data. Raw or pooled variances cannot be used for complex experimental designs, so instead we can fit a Poisson model using the appropriate design matrix to each tag and use the standardized residuals in place of the pooled variance (as in `plotMeanVar`) to visualize the mean-variance relationship in the data. The function will plot the average standardized residual for observations split into `nbins` bins by overall expression level. This provides a useful summary of how the variance of the counts change with respect to average expression level (abundance). A line showing the Poisson mean-variance relationship (mean equals variance) is always shown to illustrate how the genewise variances may differ from a Poisson mean-variance relationship. A log-log scale is used for the plot.

The function `mg1mLS` is used to fit the Poisson models to the data. This code is fast for fitting models, but does not compute the value for the leverage, technically required to compute the standardized residuals. Here, we approximate the standardized residuals by replacing the usual denominator of $(1 - \text{leverage})$ by $(1 - p/n)$, where n is the number of observations per tag (i.e. number of libraries) and p is the number of parameters in the model (i.e. number of columns in the full-rank design matrix).

Value

`dglmStdResid` produces a mean-variance plot based on standardized residuals from a Poisson model fit for each tag for the DGE data. `dglmStdResid` returns a list with the following elements:

<code>ave.means</code>	vector of the average expression level within each bin of observations
<code>ave.std.resid</code>	vector of the average standardized Poisson residual within each bin of tags
<code>bin.means</code>	list containing the average (mean) expression level (given by the fitted value from the given Poisson model) for observations divided into bins based on amount of expression
<code>bin.std.resid</code>	list containing the standardized residual from the given Poisson model for observations divided into bins based on amount of expression
<code>means</code>	vector giving the fitted value for each observed count
<code>standardized.residuals</code>	vector giving approximate standardized residual for each observed count
<code>bins</code>	list containing the indices for the observations, assigning them to bins
<code>nbins</code>	scalar giving the number of bins used to split up the observed counts
<code>ngenes</code>	scalar giving the number of genes/tags in the dataset
<code>nlibs</code>	scalar giving the number of libraries in the dataset

`getDispersions` computes the dispersion from the standardized residuals and returns a list with the following components:

<code>bin.dispersion</code>	vector giving the estimated dispersion value for each bin of observed counts, computed using the average standardized residual for the bin
<code>bin.dispersion.used</code>	vector giving the actual estimated dispersion value to be used. Some computed dispersions using the method in this function can be negative, which is not allowed. We use the dispersion value from the nearest bin of higher expression level with positive dispersion value in place of any negative dispersions.
<code>dispersion</code>	vector giving the estimated dispersion for each observation, using the binned dispersion estimates from above, so that all of the observations in a given bin get the same dispersion value.

Author(s)

Davis McCarthy

See Also

[plotMeanVar](#), [plotMDS.DGEList](#), [plotSmear](#) and [maPlot](#) provide more ways of visualizing DGE data.

Examples

```

y <- matrix(rnbinom(1000,mu=10,size=2),ncol=4)
design <- model.matrix(~c(0,0,1,1)+c(0,1,0,1))
binned <- dglmStdResid(y, design, dispersion=0.5)

getDispersions(binned)$bin.dispersion.used # Look at the estimated dispersions for the bins

```

diffSpliceDGE *Test for Differential Exon Usage*

Description

Given a negative binomial generalized log-linear model fit at the exon level, test for differential exon usage between experimental conditions.

Usage

```
diffSpliceDGE(fit.exon, coef=ncol(fit.exon$design), geneid, exonid=NULL, verbose=TRUE)
```

Arguments

fit.exon	an DGEGLM fitted model object produced by glmFit. Rows should correspond to exons.
coef	integer indicating which coefficient of the generalized linear model is to be tested for differential exon usage. Defaults to the last coefficient.
geneid	gene identifiers. Either a vector of length nrow(fit.exon) or the name of the column of fit.exon\$genes containing the gene identifiers. Rows with the same ID are assumed to belong to the same gene.
exonid	exon identifiers. Either a vector of length nrow(fit.exon) or the name of the column of fit.exon\$genes containing the exon identifiers.
verbose	logical, if TRUE some diagnostic information about the number of genes and exons is output.

Details

This function tests for differential exon usage for each gene for a given coefficient of the generalized linear model.

Testing for differential exon usage is equivalent to testing whether the exons in each gene have the same log-fold changes as the other exons in the same gene. At exon-level, each exon is compared to the average of all other exons for the same gene using quasi-likelihood F-tests. At gene-level, two different tests are provided. The first is converting exon-level p-values to gene-level p-values by Simes method. The other is an F-test for differences between the exon log-fold-changes within each gene.

Value

diffSpliceDGE produces an object of class DGELRT containing the component design from `fit.exon` plus the following new components:

<code>comparison</code>	character string describing the coefficient being tested.
<code>coefficients</code>	numeric vector of coefficients on the natural log scale. Each coefficient is the difference between the log-fold-change for that exon versus the average log-fold-change for the rest exons within the same gene.
<code>exon.F</code>	numeric vector of F-statistics for exons.
<code>exon.df.test</code>	numeric vector of testing degrees of freedom for exons.
<code>exon.df.prior</code>	numeric vector of prior degrees of freedom for exons.
<code>exon.df.residual</code>	numeric vector of residual degrees of freedom for exons.
<code>exon.p.value</code>	numeric vector of p-values for exons.
<code>genes</code>	data.frame of exon annotation
<code>genecolname</code>	character string giving the name of the column of genes containing gene IDs.
<code>exoncolname</code>	character string giving the name of the column of genes containing exon IDs.
<code>gene.df.test</code>	numeric vector of testing degrees of freedom for genes.
<code>gene.df.prior</code>	numeric vector of prior degrees of freedom for genes.
<code>gene.df.residual</code>	numeric vector of residual degrees of freedom for genes.
<code>gene.Simes.p.value</code>	numeric vector of Simes' p-values for genes.
<code>gene.F</code>	numeric vector of F-statistics for gene-level test.
<code>gene.F.p.value</code>	numeric vector of F-test p-values for genes.
<code>gene.genes</code>	data.frame of gene annotation.

The information and testing results for both exons and genes are sorted by `geneid` and by `exonid` within gene.

Author(s)

Yunshun Chen and Gordon Smyth

Examples

```
# Gene exon annotation
Gene <- paste("G", 1:10, sep="")
Gene <- rep(Gene, each=10)
Exon <- paste("Ex", 1:10, sep="")
Gene.Exon <- paste(Gene, Exon, sep=".")
genes <- data.frame(GeneID=Gene, Gene.Exon=Gene.Exon)

design <- model.matrix(~c(0,0,0,1,1,1))
mu <- matrix(20, 100, 6)
```

```

mu[1,4:6] <- 200
counts <- matrix(rnbinom(600,mu=mu,size=20),100,6)

y <- DGEList(counts=counts, lib.size=rep(1e6,6), genes=genes)
gfit <- glmFit(y, design, dispersion=0.05)

ds <- diffSpliceDGE(gfit, geneid="GeneID")
topSpliceDGE(ds)
plotSpliceDGE(ds)

```

dim	<i>Retrieve the Dimensions of a DGEList, DGEEexact, DGEGLM, DGELRT or TopTags Object</i>
-----	--

Description

Retrieve the number of rows (transcripts) and columns (libraries) for an DGEList, DGEEexact or TopTags Object.

Usage

```

## S3 method for class DGEList
dim(x)
## S3 method for class DGEList
length(x)

```

Arguments

x an object of class DGEList, DGEEexact, TopTags, DGEGLM or DGELRT

Details

Digital gene expression data objects share many analogies with ordinary matrices in which the rows correspond to transcripts or genes and the columns to arrays. These methods allow one to extract the size of microarray data objects in the same way that one would do for ordinary matrices.

A consequence is that row and column commands `nrow(x)`, `ncol(x)` and so on also work.

Value

Numeric vector of length 2. The first element is the number of rows (genes) and the second is the number of columns (arrays).

Author(s)

Gordon Smyth, Davis McCarthy

See Also

[dim](#) in the base package.

[02.Classes](#) gives an overview of data classes used in LIMMA.

Examples

```
M <- A <- matrix(11:14,4,2)
rownames(M) <- rownames(A) <- c("a","b","c","d")
colnames(M) <- colnames(A) <- c("A1","A2")
MA <- new("MAList",list(M=M,A=A))
dim(M)
ncol(M)
nrow(M)
length(M)
```

 dimnames

Retrieve the Dimension Names of a DGE Object

Description

Retrieve the dimension names of a digital gene expression data object.

Usage

```
## S3 method for class DGEList
dimnames(x)
## S3 replacement method for class DGEList
dimnames(x) <- value
```

Arguments

`x` an object of class `DGEList`, `DGEEexact`, `DGEGLM`, `DGELRT` or `TopTags`
`value` a possible value for `dimnames(x)`, see [dimnames](#)

Details

The dimension names of a DGE data object are the same as those of the most important component of that object.

Setting dimension names is currently only permitted for `DGEList` or `DGEGLM` objects.

A consequence is that `rownames` and `colnames` will work as expected.

Value

Either `NULL` or a list of length 2. If a list, its components are either `NULL` or a character vector the length of the appropriate dimension of `x`.

Author(s)

Gordon Smyth

See Also[dimnames](#) in the base package.

dispBinTrend

*Estimate Dispersion Trend by Binning for NB GLMs***Description**

Estimate the abundance-dispersion trend by computing the common dispersion for bins of genes of similar AveLogCPM and then fitting a smooth curve.

Usage

```
dispBinTrend(y, design=NULL, offset=NULL, df = 5, span=0.3, min.n=400,
             method.bin="CoxReid", method.trend="spline", AveLogCPM=NULL,
             weights=NULL, ...)
```

Arguments

y	numeric matrix of counts
design	numeric matrix giving the design matrix for the GLM that is to be fit.
offset	numeric scalar, vector or matrix giving the offset (in addition to the log of the effective library size) that is to be included in the NB GLM for the transcripts. If a scalar, then this value will be used as an offset for all transcripts and libraries. If a vector, it should be have length equal to the number of libraries, and the same vector of offsets will be used for each transcript. If a matrix, then each library for each transcript can have a unique offset, if desired. In <code>adjustedProfileLik</code> the <code>offset</code> must be a matrix with the same dimension as the table of counts.
df	degrees of freedom for spline curve.
span	span used for loess curve.
min.n	minimum number of genes in a bins.
method.bin	method used to estimate the dispersion in each bin. Possible values are "CoxReid", "Pearson" or "deviance".
method.trend	type of curve to smooth the bins. Possible values are "spline" for a natural cubic regression spline or "loess" for a linear lowess curve.
AveLogCPM	numeric vector giving average log ₂ counts per million for each gene
weights	optional numeric matrix giving observation weights
...	other arguments are passed to <code>estimateGLMCommonDisp</code>

Details

Estimate a dispersion parameter for each of many negative binomial generalized linear models by computing the common dispersion for genes sorted into bins based on overall AveLogCPM. A regression natural cubic splines or a linear loess curve is used to smooth the trend and extrapolate a value to each gene.

If there are fewer than `min.n` rows of `y` with at least one positive count, then one bin is used. The number of bins is limited to 1000.

Value

list with the following components:

AveLogCPM	numeric vector containing the overall AveLogCPM for each gene
dispersion	numeric vector giving the trended dispersion estimate for each gene
bin.AveLogCPM	numeric vector of length equal to <code>nbins</code> giving the average (mean) AveLogCPM for each bin
bin.dispersion	numeric vector of length equal to <code>nbins</code> giving the estimated common dispersion for each bin

Author(s)

Davis McCarthy and Gordon Smyth

References

McCarthy, DJ, Chen, Y, Smyth, GK (2012). Differential expression analysis of multifactor RNA-Seq experiments with respect to biological variation. *Nucleic Acids Research* 40, 4288-4297. <http://nar.oxfordjournals.org/content/40/10/4288>

See Also

[estimateGLMTrendedDisp](#)

Examples

```
ntags <- 1000
nlibs <- 4
means <- seq(5,10000,length.out=ntags)
y <- matrix(rnbinom(ntags*nlibs,mu=rep(means,nlibs),size=0.1*means),nrow=ntags,ncol=nlibs)
keep <- rowSums(y) > 0
y <- y[keep,]
group <- factor(c(1,1,2,2))
design <- model.matrix(~group) # Define the design matrix for the full model
out <- dispBinTrend(y, design, min.n=100, span=0.3)
with(out, plot(AveLogCPM, sqrt(dispersion)))
```


dispCoxReid

*Estimate Common Dispersion for Negative Binomial GLMs***Description**

Estimate a common dispersion parameter across multiple negative binomial generalized linear models.

Usage

```
dispCoxReid(y, design=NULL, offset=NULL, weights=NULL, AveLogCPM=NULL, interval=c(0,4),
            tol=1e-5, min.row.sum=5, subset=10000)
dispDeviance(y, design=NULL, offset=NULL, interval=c(0,4), tol=1e-5, min.row.sum=5,
             subset=10000, AveLogCPM=NULL, robust=FALSE, trace=FALSE)
dispPearson(y, design=NULL, offset=NULL, min.row.sum=5, subset=10000,
            AveLogCPM=NULL, tol=1e-6, trace=FALSE, initial.dispersion=0.1)
```

Arguments

y	numeric matrix of counts. A glm is fitted to each row.
design	numeric design matrix, as for glmFit .
offset	numeric vector or matrix of offsets for the log-linear models, as for glmFit . Defaults to <code>log(colSums(y))</code> .
weights	optional numeric matrix giving observation weights
AveLogCPM	numeric vector giving average log2 counts per million.
interval	numeric vector of length 2 giving minimum and maximum allowable values for the dispersion, passed to <code>optimize</code> .
tol	the desired accuracy, see <code>optimize</code> or <code>uniroot</code> .
min.row.sum	integer. Only rows with at least this number of counts are used.
subset	integer, number of rows to use in the calculation. Rows used are chosen evenly spaced by AveLogCPM.
trace	logical, should iteration information be output?
robust	logical, should a robust estimator be used?
initial.dispersion	starting value for the dispersion

Details

These are low-level (non-object-orientated) functions called by `estimateGLMCommonDisp`.

`dispCoxReid` maximizes the Cox-Reid adjusted profile likelihood (Cox and Reid, 1987). `dispPearson` sets the average Pearson goodness of fit statistics to its (asymptotic) expected value. This is also known as the *pseudo-likelihood* estimator. `dispDeviance` sets the average residual deviance statistic to its (asymptotic) expected values. This is also known as the *quasi-likelihood* estimator.

Robinson and Smyth (2008) and McCarthy et al (2011) showed that the Pearson (pseudo-likelihood) estimator typically under-estimates the true dispersion. It can be seriously biased when the number of libraries (`ncol(y)`) is small. On the other hand, the deviance (quasi-likelihood) estimator typically over-estimates the true dispersion when the number of libraries is small. Robinson and Smyth (2008) and McCarthy et al (2011) showed the Cox-Reid estimator to be the least biased of the three options.

`dispCoxReid` uses `optimize` to maximize the adjusted profile likelihood. `dispDeviance` uses `uniroot` to solve the estimating equation. The robust options use an order statistic instead the mean statistic, and have the effect that a minority of tags with very large (outlier) dispersions should have limited influence on the estimated value. `dispPearson` uses a globally convergent Newton iteration.

Value

Numeric vector of length one giving the estimated common dispersion.

Author(s)

Gordon Smyth

References

Cox, DR, and Reid, N (1987). Parameter orthogonality and approximate conditional inference. *Journal of the Royal Statistical Society Series B* 49, 1-39.

Robinson MD and Smyth GK (2008). Small-sample estimation of negative binomial dispersion, with applications to SAGE data. *Biostatistics*, 9, 321-332

McCarthy, DJ, Chen, Y, Smyth, GK (2012). Differential expression analysis of multifactor RNA-Seq experiments with respect to biological variation. *Nucleic Acids Research*. <http://nar.oxfordjournals.org/content/early/2012/02/06/nar.gks042> (Published online 28 January 2012)

See Also

[estimateGLMCommonDisp](#), [optimize](#), [uniroot](#)

Examples

```
ntags <- 100
nlibs <- 4
y <- matrix(rnbinom(ntags*nlibs,mu=10,size=10),nrow=ntags,ncol=nlibs)
group <- factor(c(1,1,2,2))
lib.size <- rowSums(y)
design <- model.matrix(~group)
disp <- dispCoxReid(y, design, offset=log(lib.size), subset=100)
```

 dispCoxReidInterpolateTagwise

Estimate Tagwise Dispersion for Negative Binomial GLMs by Cox-Reid Adjusted Profile Likelihood

Description

Estimate tagwise dispersion parameters across multiple negative binomial generalized linear models using weighted Cox-Reid Adjusted Profile-likelihood and cubic spline interpolation over a tagwise grid.

Usage

```
dispCoxReidInterpolateTagwise(y, design, offset=NULL, dispersion, trend=TRUE,
                              AveLogCPM=NULL, min.row.sum=5, prior.df=10,
                              span=0.3, grid.npts=11, grid.range=c(-6,6),
                              weights=NULL)
```

Arguments

y	numeric matrix of counts
design	numeric matrix giving the design matrix for the GLM that is to be fit.
offset	numeric scalar, vector or matrix giving the offset (in addition to the log of the effective library size) that is to be included in the NB GLM for the transcripts. If a scalar, then this value will be used as an offset for all transcripts and libraries. If a vector, it should be have length equal to the number of libraries, and the same vector of offsets will be used for each transcript. If a matrix, then each library for each transcript can have a unique offset, if desired. In <code>adjustedProfileLik</code> the offset must be a matrix with the same dimension as the table of counts.
dispersion	numeric scalar or vector giving the dispersion(s) towards which the tagwise dispersion parameters are shrunk.
trend	logical, whether abundance-dispersion trend is used for smoothing.
AveLogCPM	numeric vector giving average log ₂ counts per million for each tag.
min.row.sum	numeric scalar giving a value for the filtering out of low abundance tags. Only tags with total sum of counts above this value are used. Low abundance tags can adversely affect the estimation of the common dispersion, so this argument allows the user to select an appropriate filter threshold for the tag abundance.
prior.df	numeric scalar, prior degsmoothing parameter that indicates the weight to give to the common likelihood compared to the individual tag's likelihood; default <code>getPriorN(object)</code> gives a value for <code>prior.n</code> that is equivalent to giving the common likelihood 20 prior degrees of freedom in the estimation of the tag/genewise dispersion.
span	numeric parameter between 0 and 1 specifying proportion of data to be used in the local regression moving window. Larger numbers give smoother fits.

grid.npts	numeric scalar, the number of points at which to place knots for the spline-based estimation of the tagwise dispersion estimates.
grid.range	numeric vector of length 2, giving relative range, in terms of $\log_2(\text{dispersion})$, on either side of trendline for each tag for spline grid points.
weights	optional numeric matrix giving observation weights

Details

In the edgeR context, `dispCoxReidInterpolateTagwise` is a low-level function called by `estimateGLMTagwiseDisp`.

`dispCoxReidInterpolateTagwise` calls the function `maximizeInterpolant` to fit cubic spline interpolation over a tagwise grid.

Value

`dispCoxReidInterpolateTagwise` produces a vector of tagwise dispersions having the same length as the number of genes in the count data.

Author(s)

Yunshun Chen, Gordon Smyth

References

Cox, DR, and Reid, N (1987). Parameter orthogonality and approximate conditional inference. *Journal of the Royal Statistical Society Series B* 49, 1-39.

McCarthy, DJ, Chen, Y, Smyth, GK (2012). Differential expression analysis of multifactor RNA-Seq experiments with respect to biological variation. *Nucleic Acids Research* 40, 4288-4297. <http://nar.oxfordjournals.org/content/40/10/4288>

See Also

[estimateGLMTagwiseDisp](#), [maximizeInterpolant](#)

Examples

```
y <- matrix(rnbinom(1000, mu=10, size=2), ncol=4)
design <- matrix(1, 4, 1)
dispersion <- 0.5
d <- dispCoxReidInterpolateTagwise(y, design, dispersion=dispersion)
d
```

 dispCoxReidSplineTrend

Estimate Dispersion Trend for Negative Binomial GLMs

Description

Estimate trended dispersion parameters across multiple negative binomial generalized linear models using Cox-Reid adjusted profile likelihood.

Usage

```
dispCoxReidSplineTrend(y, design, offset=NULL, df = 5, subset=10000, AveLogCPM=NULL,
                       method.optim="Nelder-Mead", trace=0)
dispCoxReidPowerTrend(y, design, offset=NULL, subset=10000, AveLogCPM=NULL,
                      method.optim="Nelder-Mead", trace=0)
```

Arguments

y	numeric matrix of counts
design	numeric matrix giving the design matrix for the GLM that is to be fit.
offset	numeric scalar, vector or matrix giving the offset (in addition to the log of the effective library size) that is to be included in the NB GLM for the transcripts. If a scalar, then this value will be used as an offset for all transcripts and libraries. If a vector, it should be have length equal to the number of libraries, and the same vector of offsets will be used for each transcript. If a matrix, then each library for each transcript can have a unique offset, if desired. In <code>adjustedProfileLik</code> the offset must be a matrix with the same dimension as the table of counts.
df	integer giving the degrees of freedom of the spline function, see <code>ns</code> in the <code>splines</code> package.
subset	integer, number of rows to use in the calculation. Rows used are chosen evenly spaced by <code>AveLogCPM</code> using <code>cutWithMinN</code> .
AveLogCPM	numeric vector giving average log2 counts per million for each gene
method.optim	the method to be used in <code>optim</code> . See <code>optim</code> for more detail.
trace	logical, should iteration information be output?

Details

In the edgeR context, these are low-level functions called by `estimateGLMTrendedDisp`.

`dispCoxReidSplineTrend` and `dispCoxReidPowerTrend` fit abundance trends to the tagwise dispersions. `dispCoxReidSplineTrend` fits a regression spline whereas `dispCoxReidPowerTrend` fits a log-linear trend of the form $a \times \exp(\text{abundance})^b + c$. In either case, `optim` is used to maximize the adjusted profile likelihood (Cox and Reid, 1987).

Value

List containing numeric vectors dispersion and abundance containing the estimated dispersion and abundance for each transcript. The vectors are of the same length as `nrow(y)`.

Author(s)

Yunshun Chen, Davis McCarthy, Gordon Smyth

References

Cox, DR, and Reid, N (1987). Parameter orthogonality and approximate conditional inference. *Journal of the Royal Statistical Society Series B* 49, 1-39.

See Also

[estimateGLMTrendedDisp](#)

Examples

```
design <- matrix(1,4,1)
y <- matrix((rnbinom(400,mu=100,size=5)),100,4)
d1 <- dispCoxReidSplineTrend(y, design, df=3)
d2 <- dispCoxReidPowerTrend(y, design)
with(d2,plot(AveLogCPM,sqrt(dispersion)))
```

edgeRUsersGuide

View edgeR User's Guide

Description

Finds the location of the edgeR User's Guide and optionally opens it.

Usage

```
edgeRUsersGuide(view=TRUE)
```

Arguments

`view` logical, should the document be opened using the default PDF document reader?

Details

The function `vignette("edgeR")` will find the short edgeR Vignette which describes how to obtain the edgeR User's Guide. The User's Guide is not itself a true vignette because it is not automatically generated using [Sweave](#) during the package build process. This means that it cannot be found using `vignette`, hence the need for this special function.

If the operating system is other than Windows, then the PDF viewer used is that given by `Sys.getenv("R_PDFVIEWER")`. The PDF viewer can be changed using `Sys.putenv(R_PDFVIEWER=)`.

Value

Character string giving the file location. If `view=TRUE`, the PDF document reader is started and the User's Guide is opened, as a side effect.

Author(s)

Gordon Smyth

See Also

[system](#)

Examples

```
# To get the location:
edgeRUsersGuide(view=FALSE)
# To open in pdf viewer:
## Not run: edgeRUsersGuide()
```

equalizeLibSizes *Equalize Library Sizes by Quantile-to-Quantile Normalization*

Description

Adjusts counts so that the effective library sizes are equal, preserving fold-changes between groups and preserving biological variability within each group.

Usage

```
equalizeLibSizes(object, dispersion=NULL, common.lib.size)
```

Arguments

<code>object</code>	DGEList object
<code>dispersion</code>	numeric vector of dispersion parameters. By default, is extracted from <code>object</code> or, if <code>object</code> contains no dispersion information, is set to 0.05.
<code>common.lib.size</code>	numeric scalar, the library size to normalize to; default is the geometric mean of the original effective library sizes

Details

This function implements the quantile-quantile normalization method of Robinson and Smyth (2008). It computes normalized counts, or pseudo-counts, used by `exactTest` and `estimateCommonDisp`.

The output pseudo-counts are the counts that would have theoretically arisen had the effective library sizes been equal for all samples. The pseudo-counts are computed in such a way as to preserve fold-change differences between the groups defined by `object$samples$group` as well as biological variability within each group. Consequently, the results will depend on how the groups are defined.

Note that the column sums of the `pseudo.counts` matrix will not generally be equal, because the effective library sizes are not necessarily the same as actual library sizes and because the normalized pseudo counts are not equal to expected counts.

Value

A list with components

`pseudo.counts` numeric matrix of normalized pseudo-counts
`common.lib.size` normalized library size

Note

This function is intended mainly for internal edgeR use. It is not normally called directly by users.

Author(s)

Mark Robinson, Davis McCarthy, Gordon Smyth

References

Robinson MD and Smyth GK (2008). Small-sample estimation of negative binomial dispersion, with applications to SAGE data. *Biostatistics*, 9, 321-332. <http://biostatistics.oxfordjournals.org/content/9/2/321>

See Also

[q2qnbinom](#)

Examples

```
ngenes <- 1000
nlibs <- 2
counts <- matrix(0, ngenes, nlibs)
colnames(counts) <- c("Sample1", "Sample2")
counts[,1] <- rpois(ngenes, lambda=10)
counts[,2] <- rpois(ngenes, lambda=20)
summary(counts)
y <- DGEList(counts=counts)
out <- equalizeLibSizes(y)
summary(out$pseudo.counts)
```

estimateCommonDisp	<i>Estimate Common Negative Binomial Dispersion by Conditional Maximum Likelihood</i>
--------------------	---

Description

Maximizes the negative binomial conditional common likelihood to give the estimate of the common dispersion across all tags.

Usage

```
estimateCommonDisp(object, tol=1e-06, rowsum.filter=5, verbose=FALSE)
```

Arguments

object	DGEList object
tol	the desired accuracy, passed to optimize
rowsum.filter	numeric scalar giving a value for the filtering out of low abundance tags in the estimation of the common dispersion. Only tags with total sum of counts above this value are used in the estimation of the common dispersion.
verbose	logical, if TRUE estimated dispersion and BCV will be printed to standard output.

Details

Implements the method of Robinson and Smyth (2008) for estimating a common dispersion parameter by conditional maximum likelihood. The method of conditional maximum likelihood assumes that library sizes are equal, which is not true in general, so pseudocounts (counts adjusted so that the library sizes are equal) need to be calculated. The function `equalizeLibSizes` is called to adjust the counts using a quantile-to-quantile method, but this requires a fixed value for the common dispersion parameter. To obtain a good estimate for the common dispersion, pseudocounts are calculated under the Poisson model (dispersion is zero) and these pseudocounts are used to give an estimate of the common dispersion. This estimate of the common dispersion is then used to recalculate the pseudocounts, which are used to provide a final estimate of the common dispersion.

Value

Returns object with the following added components:

common.dispersion	estimate of the common dispersion.
pseudo.counts	numeric matrix of quantile-quantile normalized counts. These are counts adjusted so that the library sizes are equal, while preserving differences between groups and variability within each group.
pseudo.lib.size	the common library size to which the counts have been adjusted

Author(s)

Mark Robinson, Davis McCarthy, Gordon Smyth

References

Robinson MD and Smyth GK (2008). Small-sample estimation of negative binomial dispersion, with applications to SAGE data. *Biostatistics*, 9, 321-332

See Also

[equalizeLibSizes](#)

Examples

```
# True dispersion is 1/5=0.2
y <- matrix(rnbinom(1000,mu=10,size=5),ncol=4)
d <- DGEList(counts=y,group=c(1,1,2,2),lib.size=c(1000:1003))
d <- estimateCommonDisp(d, verbose=TRUE)
```

estimateDisp	<i>Estimate Common, Trended and Tagwise Negative Binomial dispersions by weighted likelihood empirical Bayes</i>
--------------	--

Description

Maximizes the negative binomial likelihood to give the estimate of the common, trended and tag-wise dispersions across all tags.

Usage

```
estimateDisp(y, design=NULL, prior.df=NULL, trend.method="locfit", span=NULL,
             min.row.sum=5, grid.length=21, grid.range=c(-10,10), robust=FALSE,
             winsor.tail.p=c(0.05,0.1), tol=1e-06)
```

Arguments

y	DGEList object
design	numeric design matrix
prior.df	prior degrees of freedom. It is used in calculating prior.n.
trend.method	method for estimating dispersion trend. Possible values are "none", "movingave", "loess" and "locfit".
span	width of the smoothing window, as a proportion of the data set.
min.row.sum	numeric scalar giving a value for the filtering out of low abundance tags. Only tags with total sum of counts above this value are used. Low abundance tags can adversely affect the dispersion estimation, so this argument allows the user to select an appropriate filter threshold for the tag abundance.

grid.length	the number of points on which the interpolation is applied for each tag.
grid.range	the range of the grid points around the trend on a log2 scale.
robust	logical, should the estimation of prior.df be robustified against outliers?
winsor.tail.p	numeric vector of length 1 or 2, giving left and right tail proportions of the deviances to Winsorize when estimating prior.df.
tol	the desired accuracy, passed to <code>optimize</code>

Details

This function calculates a matrix of likelihoods for each gene at a set of dispersion grid points, and then applies weighted likelihood empirical Bayes method to obtain posterior dispersion estimates. If there is no design matrix, it calculates the quantile conditional likelihood for each gene (tag) and then maximize it. The method is same as in the function `estimateCommonDisp` and `estimateTagwiseDisp`. If a design matrix is given, it then calculates the adjusted profile log-likelihood for each gene (tag) and then maximize it. It is similar to the functions `estimateGLMCommonDisp`, `estimateGLMTrendedDisp` and `estimateGLMTagwiseDisp`.

Value

Returns object with the following added components:

common.dispersion	estimate of the common dispersion.
trended.dispersion	estimates of the trended dispersions.
tagwise.dispersion	tag- or gene-wise estimates of the dispersion parameter.
logCPM	the tag abundance in log average counts per million.
prior.df	prior degrees of freedom. It is a vector when robust method is used.
prior.n	estimate of the prior weight, i.e. the smoothing parameter that indicates the weight to put on the common likelihood compared to the individual tag's likelihood.
span	width of the smoothing window used in estimating dispersions.

Author(s)

Yunshun Chen, Gordon Smyth

References

Chen, Y, Lun, ATL, and Smyth, GK (2014). Differential expression analysis of complex RNA-seq experiments using edgeR. In: *Statistical Analysis of Next Generation Sequence Data*, Somnath Datta and Daniel S Nettleton (eds), Springer, New York. <http://www.statsci.org/smyth/pubs/edgeRChapterPreprint.pdf>

See Also

[estimateCommonDisp](#), [estimateTagwiseDisp](#), [estimateGLMCommonDisp](#), [estimateGLMTrendedDisp](#), [estimateGLMTagwiseDisp](#)

Examples

```
# True dispersion is 1/5=0.2
y <- matrix(rnbinom(1000, mu=10, size=5), ncol=4)
group <- c(1,1,2,2)
design <- model.matrix(~group)
d <- DGEList(counts=y, group=group)
d1 <- estimateDisp(d)
d2 <- estimateDisp(d, design)
```

estimateExonGenewiseDisp

Estimate Genewise Dispersions from Exon-Level Count Data

Description

Estimate a dispersion value for each gene from exon-level count data by collapsing exons into the genes to which they belong.

Usage

```
estimateExonGenewiseDisp(y, geneID, group=NULL)
```

Arguments

y	either a matrix of exon-level counts or a <code>DGEList</code> object with (at least) elements <code>counts</code> (table of counts summarized at the exon level) and <code>samples</code> (data frame containing information about experimental group, library size and normalization factor for the library size). Each row of y should represent one exon.
geneID	vector of length equal to the number of rows of y, which provides the gene identifier for each exon in y. These identifiers are used to group the relevant exons into genes for the gene-level analysis of splice variation.
group	factor supplying the experimental group/condition to which each sample (column of y) belongs. If <code>NULL</code> (default) the function will try to extract it from y, which only works if y is a <code>DGEList</code> object.

Details

This function can be used to compute genewise dispersion estimates (for an experiment with a one-way, or multiple group, layout) from exon-level count data. `estimateCommonDisp` and `estimateTagwiseDisp` are used to do the computation and estimation, and the default arguments for those functions are used.

Value

estimateExonGenewiseDisp returns a vector of genewise dispersion estimates, one for each unique geneID.

Author(s)

Davis McCarthy, Gordon Smyth

See Also

[estimateCommonDisp](#) and related functions for estimating the dispersion parameter for the negative binomial model.

Examples

```
# generate exon counts from NB, create list object
y<-matrix(rnbinom(40,size=1,mu=10),nrow=10)
d<-DGEList(counts=y,group=rep(1:2,each=2))
genes <- rep(c("gene.1","gene.2"), each=5)
estimateExonGenewiseDisp(d, genes)
```

estimateGLMCommonDisp *Estimate Common Dispersion for Negative Binomial GLMs*

Description

Estimates a common negative binomial dispersion parameter for a DGE dataset with a general experimental design.

Usage

```
## S3 method for class DGEList
estimateGLMCommonDisp(y, design=NULL, method="CoxReid",
                      subset=10000, verbose=FALSE, ...)

## Default S3 method:
estimateGLMCommonDisp(y, design=NULL, offset=NULL,
                      method="CoxReid", subset=10000, AveLogCPM=NULL,
                      verbose=FALSE, weights=NULL,...)
```

Arguments

y	object containing read counts, as for glmFit .
design	numeric design matrix, as for glmFit .
offset	numeric vector or matrix of offsets for the log-linear models, as for glmFit .
method	method for estimating the dispersion. Possible values are "CoxReid", "Pearson" or "deviance".

subset	maximum number of rows of <code>y</code> to use in the calculation. Rows used are chosen evenly spaced by AveLogCPM using <code>systematicSubset</code> .
AveLogCPM	numeric vector giving average log2 counts per million for each gene
verbose	logical, if TRUE estimated dispersion and BCV will be printed to standard output.
weights	optional numeric matrix giving observation weights
...	other arguments are passed to lower-level functions. See <code>dispCoxReid</code> , <code>dispPearson</code> and <code>dispDeviance</code> for details.

Details

This function calls `dispCoxReid`, `dispPearson` or `dispDeviance` depending on the method specified. See `dispCoxReid` for details of the three methods and a discussion of their relative performance.

Value

The default method returns a numeric vector of length 1 containing the estimated common dispersion.

The `DGEList` method returns the same `DGEList` `y` as input but with `common.dispersion` as an added component. The output object will also contain a component `AveLogCPM` if it was not already present in `y`.

Author(s)

Gordon Smyth, Davis McCarthy, Yunshun Chen

References

McCarthy, DJ, Chen, Y, Smyth, GK (2012). Differential expression analysis of multifactor RNA-Seq experiments with respect to biological variation. *Nucleic Acids Research* 40, 4288-4297. <http://nar.oxfordjournals.org/content/40/10/4288>

See Also

`dispCoxReid`, `dispPearson`, `dispDeviance`

`estimateGLMTrendedDisp` for trended dispersion and `estimateGLMTagwiseDisp` for tagwise dispersions in the context of a generalized linear model.

`estimateCommonDisp` for common dispersion or `estimateTagwiseDisp` for tagwise dispersion in the context of a multiple group experiment (one-way layout).

Examples

```
# True dispersion is 1/size=0.1
y <- matrix(rnbinom(1000,mu=10,size=10),ncol=4)
d <- DGEList(counts=y,group=c(1,1,2,2))
design <- model.matrix(~group, data=d$samples)
d1 <- estimateGLMCommonDisp(d, design, verbose=TRUE)
```

```
# Compare with classic CML estimator:
d2 <- estimateCommonDisp(d, verbose=TRUE)

# See example(glmFit) for a different example
```

estimateGLMRobustDisp *Empirical Robust Bayes Tagwise Dispersions for Negative Binomial GLMs using Observation Weights*

Description

Compute a robust estimate of the negative binomial dispersion parameter for each tag or transcript, with expression levels specified by a log-linear model, using observation weights. These observation weights will be stored and used later for estimating regression parameters.

Usage

```
estimateGLMRobustDisp(y, design = NULL, prior.df = 10, update.trend = TRUE,
  trend.method = "bin.loess", maxit = 6, k = 1.345,
  residual.type = "pearson", verbose = FALSE,
  record = FALSE)
```

Arguments

y	a DGEList object.
design	numeric design matrix, as for glmFit .
prior.df	prior degrees of freedom.
update.trend	logical. Should the trended dispersion be re-estimated at each iteration?
trend.method	method (low-level function) used to estimate the trended dispersions. estimateGLMTrendedDisp
maxit	maximum number of iterations for weighted estimateGLMTagwiseDisp .
k	the tuning constant for Huber estimator. If the absolute value of residual (r) is less than k, its observation weight is 1, otherwise k/abs(r).
residual.type	type of residual (r) used for estimation observation weight
verbose	logical. Should verbose comments be printed?
record	logical. Should information for each iteration be recorded (and returned as a list)?

Details

At times, because of the moderation of dispersion estimates towards a trended values, features (typically, genes) can be sensitive to outliers and causing false positives. That is, since the dispersion estimates are moderated downwards toward the trend and because the regression parameter estimates may be affected by the outliers, genes are deemed significantly differential expressed. The function uses an iterative procedure where weights are calculated from residuals and estimates are made after re-weighting.

Note: it is not necessary to first calculate the common, trended and tagwise dispersion estimates. If these are not available, the function will first calculate this (in an unweighted) fashion.

Value

estimateGLMRobustDisp produces a DGEList object, which contains the (robust) tagwise dispersion parameter estimate for each tag for the negative binomial model that maximizes the weighted Cox-Reid adjusted profile likelihood, as well as the observation weights. The observation weights are calculated using residuals and the Huber function.

Note that when record=TRUE, a simple list of DGEList objects is returned, one for each iteration (this is for debugging or tracking purposes).

Author(s)

Xiaobei Zhou, Mark D. Robinson

References

Zhou X, Lindsay H, Robinson MD (2014). Robustly detecting differential expression in RNA sequencing data using observation weights. *Nucleic Acids Research*, 42(11), e91.

See Also

This function calls [estimateGLMTrendedDisp](#) and [estimateGLMTagwiseDisp](#).

Examples

```
y <- matrix(rnbinom(100*6,mu=10,size=1/0.1),ncol=6)
d <- DGEList(counts=y,group=c(1,1,1,2,2,2),lib.size=c(1000:1005))
d <- calcNormFactors(d)
design <- model.matrix(~group, data=d$samples) # Define the design matrix for the full model
d <- estimateGLMRobustDisp(d, design)
summary(d$tagwise.dispersion)
```

estimateGLMTagwiseDisp

Empirical Bayes Tagwise Dispersions for Negative Binomial GLMs

Description

Compute an empirical Bayes estimate of the negative binomial dispersion parameter for each tag or transcript, with expression levels specified by a log-linear model.

Usage

```
## S3 method for class DGEList
estimateGLMTagwiseDisp(y, design=NULL, prior.df=10,
                       trend=!is.null(y$trended.dispersion), span=NULL, ...)
## Default S3 method:
estimateGLMTagwiseDisp(y, design=NULL, offset=NULL, dispersion,
                       prior.df=10, trend=TRUE, span=NULL, AveLogCPM=NULL,
                       weights=NULL, ...)
```


Arguments

<code>y</code>	matrix of counts or a <code>DGEList</code> object.
<code>design</code>	numeric design matrix, as for <code>glmFit</code> .
<code>trend</code>	logical. Should the prior be the trended dispersion (TRUE) or the common dispersion (FALSE)?
<code>offset</code>	offset matrix for the log-linear model, as for <code>glmFit</code> . Defaults to the log-effective library sizes.
<code>dispersion</code>	common or trended dispersion estimates, used as an initial estimate for the tagwise estimates. By default uses values stored in the <code>DGEList</code> object.
<code>prior.df</code>	prior degrees of freedom.
<code>span</code>	width of the smoothing window, in terms of proportion of the data set. Default value decreases with the number of tags.
<code>AveLogCPM</code>	numeric vector giving average log ₂ counts per million for each gene
<code>weights</code>	optional numeric matrix giving observation weights
<code>...</code>	other arguments are passed to <code>dispCoxReidInterpolateTagwise</code> .

Details

This function implements the empirical Bayes strategy proposed by McCarthy et al (2012) for estimating the tagwise negative binomial dispersions. The experimental conditions are specified by design matrix allowing for multiple explanatory factors. The empirical Bayes posterior is implemented as a conditional likelihood with tag-specific weights, and the conditional likelihood is computed using Cox-Reid approximate conditional likelihood (Cox and Reid, 1987).

The prior degrees of freedom determines the weight given to the global dispersion trend. The larger the prior degrees of freedom, the more the tagwise dispersions are squeezed towards the global trend.

This function calls the lower-level function `dispCoxReidInterpolateTagwise`.

Value

`estimateGLMTagwiseDisp.DGEList` produces a `DGEList` object, which contains the tagwise dispersion parameter estimate for each tag for the negative binomial model that maximizes the Cox-Reid adjusted profile likelihood. The tagwise dispersions are simply added to the `DGEList` object provided as the argument to the function.

`estimateGLMTagwiseDisp.default` returns a vector of the tagwise dispersion estimates.

Author(s)

Gordon Smyth, Davis McCarthy

References

Cox, DR, and Reid, N (1987). Parameter orthogonality and approximate conditional inference. *Journal of the Royal Statistical Society Series B* 49, 1-39.

McCarthy, DJ, Chen, Y, Smyth, GK (2012). Differential expression analysis of multifactor RNA-Seq experiments with respect to biological variation. *Nucleic Acids Research* 40, 4288-4297. <http://nar.oxfordjournals.org/content/40/10/4288>

See Also

[estimateGLMCommonDisp](#) for common dispersion and [estimateGLMTrendedDisp](#) for trended dispersion in the context of a generalized linear model.

[estimateCommonDisp](#) for common dispersion or [estimateTagwiseDisp](#) for tagwise dispersion in the context of a multiple group experiment (one-way layout).

Examples

```
y <- matrix(rnbinom(1000,mu=10,size=10),ncol=4)
d <- DGEList(counts=y,group=c(1,1,2,2),lib.size=c(1000:1003))
design <- model.matrix(~group, data=d$samples) # Define the design matrix for the full model
d <- estimateGLMTrendedDisp(d, design, min.n=10)
d <- estimateGLMTagwiseDisp(d, design)
summary(d$tagwise.dispersion)
```

estimateGLMTrendedDisp

Estimate Trended Dispersion for Negative Binomial GLMs

Description

Estimates the abundance-dispersion trend by Cox-Reid approximate profile likelihood.

Usage

```
## S3 method for class DGEList
estimateGLMTrendedDisp(y, design=NULL, method="auto", ...)
## Default S3 method:
estimateGLMTrendedDisp(y, design=NULL, offset=NULL, AveLogCPM=NULL,
                        method="auto", weights=NULL, ...)
```

Arguments

`y` a matrix of counts or a `DGEList` object.)
`design` numeric design matrix, as for [glmFit](#).

method	method (low-level function) used to estimate the trended dispersions. Possible values are "auto" (default, switch to "bin.spline" method if the number of tags is greater than 200 and "power" method otherwise), "bin.spline", "bin.loess" (which both result in a call to dispBinTrend), "power" (call to dispCoxReidPowerTrend), or "spline" (call to dispCoxReidSplineTrend).
offset	numeric scalar, vector or matrix giving the linear model offsets, as for <code>glmFit</code> .
AveLogCPM	numeric vector giving average log ₂ counts per million for each gene.
weights	optional numeric matrix giving observation weights
...	other arguments are passed to lower-level functions <code>dispBinTrend</code> , <code>dispCoxReidPowerTrend</code> or <code>dispCoxReidSplineTrend</code> .

Details

Estimates the dispersion parameter for each transcript (tag) with a trend that depends on the overall level of expression for the transcript for a DGE dataset for general experimental designs by using Cox-Reid approximate conditional inference for a negative binomial generalized linear model for each transcript (tag) with the unadjusted counts and design matrix provided.

The function provides an object-orientated interface to lower-level functions.

Value

When the input object is a `DGEList`, `estimateGLMTrendedDisp` produces a `DGEList` object, which contains the estimates of the trended dispersion parameter for the negative binomial model according to the method applied.

When the input object is a numeric matrix, the output of one of the lower-level functions `dispBinTrend`, `dispCoxReidPowerTrend` or `dispCoxReidSplineTrend` is returned.

Author(s)

Gordon Smyth, Davis McCarthy, Yunshun Chen

References

Cox, DR, and Reid, N (1987). Parameter orthogonality and approximate conditional inference. *Journal of the Royal Statistical Society Series B* 49, 1-39.

McCarthy, DJ, Chen, Y, Smyth, GK (2012). Differential expression analysis of multifactor RNA-Seq experiments with respect to biological variation. *Nucleic Acids Research* 40, 4288-4297. <http://nar.oxfordjournals.org/content/40/10/4288>

See Also

`dispBinTrend`, `dispCoxReidPowerTrend` and `dispCoxReidSplineTrend` for details on how the calculations are done.

Examples

```

ntags <- 250
nlibs <- 4
y <- matrix(rnbinom(ntags*nlibs,mu=10,size=10),ntags,nlibs)
d <- DGEList(counts=y,group=c(1,1,2,2),lib.size=c(1000:1003))
design <- model.matrix(~group, data=d$samples)
disp <- estimateGLMTrendedDisp(d, design, min.n=25, df=3)
plotBCV(disp)

```

estimateTagwiseDisp *Estimate Empirical Bayes Tagwise Dispersion Values*

Description

Estimates tagwise dispersion values by an empirical Bayes method based on weighted conditional maximum likelihood.

Usage

```

estimateTagwiseDisp(object, prior.df=10, trend="movingave", span=NULL, method="grid",
                    grid.length=11, grid.range=c(-6,6), tol=1e-06, verbose=FALSE)

```

Arguments

object	object of class DGEList containing (at least) the elements counts (table of raw counts), group (factor indicating group), lib.size (numeric vector of library sizes) and pseudo.alt (numeric matrix of quantile-adjusted pseudocounts calculated under the alternative hypothesis of a true difference between groups; recommended to use the DGEList object provided as the output of estimateCommonDisp)
prior.df	prior degrees of freedom.
trend	method for estimating dispersion trend. Possible values are "none", "movingave" and "loess".
span	width of the smoothing window, as a proportion of the data set.
method	method for maximizing the posterior likelihood. Possible values are "grid" for interpolation on grid points or "optimize" to call the function of the same name.
grid.length	for method="grid", the number of points on which the interpolation is applied for each tag.
grid.range	for method="grid", the range of the grid points around the trend on a log2 scale.
tol	for method="optimize", the tolerance for Newton-Rhapson iterations.
verbose	logical, if TRUE then diagnostic output is produced during the estimation process.

Details

This function implements the empirical Bayes strategy proposed by Robinson and Smyth (2007) for estimating the tagwise negative binomial dispersions. The experimental design is assumed to be a oneway layout with one or more experimental groups. The empirical Bayes posterior is implemented as a conditional likelihood with tag-specific weights.

The prior values for the dispersions are determined by a global trend. The individual tagwise dispersions are then squeezed towards this trend. The prior degrees of freedom determines the weight given to the prior. The larger the prior degrees of freedom, the more the tagwise dispersions are squeezed towards the global trend. If the number of libraries is large, the prior becomes less important and the tagwise dispersion are determined more by the individual tagwise data.

If trend="none", then the prior dispersion is just a constant, the common dispersion. Otherwise, the trend is determined by a moving average (trend="movingave") or loess smoother applied to the tagwise conditional log-likelihood. method="loess" applies a loess curve of degree 0 as implemented in [loessByCol](#).

method="optimize" is not recommended for routine use as it is very slow. It is included for testing purposes.

Value

An object of class DGEList with the same components as for [estimateCommonDisp](#) plus the following:

prior.n	estimate of the prior weight, i.e. the smoothing parameter that indicates the weight to put on the common likelihood compared to the individual tag's likelihood; prior.n of 10 means that the common likelihood is given 10 times the weight of the individual tag/gene's likelihood in the estimation of the tag/genewise dispersion
tagwise.dispersion	tag- or gene-wise estimates of the dispersion parameter

Author(s)

Mark Robinson, Davis McCarthy, Yunshun Chen and Gordon Smyth

References

Robinson, MD, and Smyth, GK (2007). Moderated statistical tests for assessing differences in tag abundance. *Bioinformatics* 23, 2881-2887. <http://bioinformatics.oxfordjournals.org/content/23/21/2881>

See Also

[estimateCommonDisp](#) is usually run before estimateTagwiseDisp.
[movingAverageByCol](#) and [loessByCol](#) implement the moving average or loess smoothers.

Examples

```
# See ?exactTest or ?estimateTrendedDisp for examples
```

estimateTrendedDisp *Estimate Empirical Bayes Trended Dispersion Values*

Description

Estimates trended dispersion values by an empirical Bayes method.

Usage

```
estimateTrendedDisp(object, method="bin.spline", df=5, span=2/3)
```

Arguments

object	object of class <code>DGEList</code> containing (at least) the elements counts (table of raw counts), <code>group</code> (factor indicating group), <code>lib.size</code> (numeric vector of library sizes) and <code>pseudo.alt</code> (numeric matrix of quantile-adjusted pseudocounts calculated under the alternative hypothesis of a true difference between groups; recommended to use the <code>DGEList</code> object provided as the output of <code>estimateCommonDisp</code>)
method	method used to estimate the trended dispersions. Possible values are "spline", and "loess".
df	integer giving the degrees of freedom of the spline function if "spline" method is used, see <code>ns</code> in the <code>splines</code> package. Default is 5.
span	scalar, passed to <code>loess</code> to determine the amount of smoothing for the loess fit when "loess" method is used. Default is 2/3.

Details

This function takes the binned common dispersion and abundance, and fits a smooth curve through these binned values using either natural cubic splines or loess. From this smooth curve it predicts the dispersion value for each gene based on the gene's overall abundance. This results in estimates for the NB dispersion parameter which have a dependence on the overall expression level of the gene, and thus have an abundance-dependent trend.

Value

An object of class `DGEList` with the same components as for `estimateCommonDisp` plus the trended dispersion estimates for each gene or tag.

Author(s)

Yunshun Chen and Gordon Smyth

See Also

`estimateCommonDisp` estimates a common value for the dispersion parameter for all tags/genes - should generally be run before `estimateTrendedDisp`.

Examples

```

ngenes <- 1000
nlib <- 4
log2cpm <- seq(from=0,to=16,length=ngenes)
lib.size <- 1e7
mu <- 2^log2cpm * lib.size * 1e-6
dispersion <- 1/sqrt(mu) + 0.1
counts <- rnbinom(ngenes*nlib, mu=mu, size=1/dispersion)
counts <- matrix(counts,ngenes,nlib)
y <- DGEList(counts,lib.size=rep(lib.size,nlib))
y <- estimateCommonDisp(y)
y <- estimateTrendedDisp(y)
y <- estimateTagwiseDisp(y)
plotBCV(y)

```

exactTest	<i>Exact Tests for Differences between Two Groups of Negative-Binomial Counts</i>
-----------	---

Description

Compute genewise exact tests for differences in the means between two groups of negative-binomially distributed counts.

Usage

```

exactTest(object, pair=1:2, dispersion="auto", rejection.region="doubletail",
          big.count=900, prior.count=0.125)
exactTestDoubleTail(y1, y2, dispersion=0, big.count=900)
exactTestBySmallP(y1, y2, dispersion=0)
exactTestByDeviance(y1, y2, dispersion=0)
exactTestBetaApprox(y1, y2, dispersion=0)

```

Arguments

object	an object of class <code>DGEList</code> .
pair	vector of length two, either numeric or character, providing the pair of groups to be compared; if a character vector, then should be the names of two groups (e.g. two levels of <code>object\$samples\$group</code>); if numeric, then groups to be compared are chosen by finding the levels of <code>object\$samples\$group</code> corresponding to those numeric values and using those levels as the groups to be compared; if <code>NULL</code> , then first two levels of <code>object\$samples\$group</code> (a factor) are used. Note that the first group listed in the pair is the baseline for the comparison—so if the pair is <code>c("A", "B")</code> then the comparison is $B - A$, so genes with positive log-fold change are up-regulated in group B compared with group A (and vice versa for genes with negative log-fold change).

dispersion	either a numeric vector of dispersions or a character string indicating that dispersions should be taken from the data object. If a numeric vector, then can be either of length one or of length equal to the number of tags. Allowable character values are "common", "trended", "tagwise" or "auto". Default behavior ("auto" is to use most complex dispersions found in data object.
rejection.region	type of rejection region for two-sided exact test. Possible values are "doubletail", "smallp" or "deviance".
big.count	count size above which asymptotic beta approximation will be used.
prior.count	average prior count used to shrink log-fold-changes. Larger values produce more shrinkage.
y1	numeric matrix of counts for the first the two experimental groups to be tested for differences. Rows correspond to genes or transcripts and columns to libraries. Libraries are assumed to be equal in size - e.g. adjusted pseudocounts from the output of equalizeLibSizes .
y2	numeric matrix of counts for the second of the two experimental groups to be tested for differences. Rows correspond to genes or transcripts and columns to libraries. Libraries are assumed to be equal in size - e.g. adjusted pseudocounts from the output of equalizeLibSizes . Must have the same number of rows as y1.

Details

The functions test for differential expression between two groups of count libraries. They implement the exact test proposed by Robinson and Smyth (2008) for a difference in mean between two groups of negative binomial random variables. The functions accept two groups of count libraries, and a test is performed for each row of data. For each row, the test is conditional on the sum of counts for that row. The test can be viewed as a generalization of the well-known exact binomial test (implemented in `binomTest`) but generalized to overdispersed counts.

`exactTest` is the main user-level function, and produces an object containing all the necessary components for downstream analysis. `exactTest` calls one of the low level functions `exactTestDoubleTail`, `exactTestBetaApprox`, `exactTestBySmallP` or `exactTestByDeviance` to do the p-value computation. The low level functions all assume that the libraries have been normalized to have the same size, i.e., to have the same expected column sum under the null hypothesis. `exactTest` equalizes the library sizes using [equalizeLibSizes](#) before calling the low level functions.

The functions `exactTestDoubleTail`, `exactTestBySmallP` and `exactTestByDeviance` correspond to different ways to define the two-sided rejection region when the two groups have different numbers of samples. `exactTestBySmallP` implements the method of small probabilities as proposed by Robinson and Smyth (2008). This method corresponds exactly to `binomTest` as the dispersion approaches zero, but gives poor results when the dispersion is very large. `exactTestDoubleTail` computes two-sided p-values by doubling the smaller tail probability. `exactTestByDeviance` uses the deviance goodness of fit statistics to define the rejection region, and is therefore equivalent to a conditional likelihood ratio test.

Note that `rejection.region="smallp"` is no longer recommended. It is preserved as an option only for backward compatibility with early versions of `edgeR`. `rejection.region="deviance"` has good theoretical statistical properties but is relatively slow to compute. `rejection.region="doubletail"`

is just slightly more conservative than `rejection.region="deviance"`, but is recommended because of its much greater speed. For general remarks on different types of rejection regions for exact tests see Gibbons and Pratt (1975).

`exactTestBetaApprox` implements an asymptotic beta distribution approximation to the conditional count distribution. It is called by the other functions for rows with both group counts greater than `big.count`.

Value

`exactTest` produces an object of class `DGEEExact` containing the following components:

<code>table</code>	data frame containing columns for the log2-fold-change, <code>logFC</code> , the average log2-counts-per-million, <code>logCPM</code> , and the two-sided p-value <code>PValue</code>
<code>comparison</code>	character vector giving the names of the two groups being compared
<code>genes</code>	optional data frame containing annotation for transcript; taken from object

The low-level functions, `exactTestDoubleTail` etc, produce a numeric vector of genewise p-values, one for each row of `y1` and `y2`.

Author(s)

Mark Robinson, Davis McCarthy, Gordon Smyth

References

Robinson MD and Smyth GK (2008). Small-sample estimation of negative binomial dispersion, with applications to SAGE data. *Biostatistics*, 9, 321-332.

Gibbons, JD and Pratt, JW (1975). P-values: interpretation and methodology. *The American Statistician* 29, 20-25.

See Also

[equalizeLibSizes](#), [binomTest](#)

Examples

```
# generate raw counts from NB, create list object
y <- matrix(rnbinom(80,size=1/0.2,mu=10),nrow=20,ncol=4)
d <- DGEList(counts=y, group=c(1,1,2,2), lib.size=rep(1000,4))

de <- exactTest(d, dispersion=0.2)
topTags(de)

# same p-values using low-level function directly
p.value <- exactTestDoubleTail(y[,1:2], y[,3:4], dispersion=0.2)
sort(p.value)[1:10]
```

expandAsMatrix	<i>expandAsMatrix</i>
----------------	-----------------------

Description

Expand scalar or vector to a matrix.

Usage

```
expandAsMatrix(x, dim)
```

Arguments

x	scalar, vector or matrix. If a vector, length must match one of the output dimensions.
dim	required dimension for the output matrix.

Details

This function expands a row or column vector to be a matrix. It is used internally in edgeR to convert offsets to a matrix.

Value

Numeric matrix of dimension dim.

Author(s)

Gordon Smyth

Examples

```
expandAsMatrix(1:3,c(4,3))  
expandAsMatrix(1:4,c(4,3))
```

getCounts	<i>Extract Specified Component of a DGEList Object</i>
-----------	--

Description

getCounts(y) returns the matrix of read counts y\$counts.

getOffset(y) returns offsets for the log-linear predictor account for sequencing depth and possibly other normalization factors. Specifically it returns the matrix y\$offset if it is non-null, otherwise it returns the log product of lib.size and norm.factors from y\$samples.

getDispersion(y) returns the most complex dispersion estimates (common, trended or tagwise) found in y.

Usage

```
getCounts(y)
getOffset(y)
getDispersion(y)
```

Arguments

`y` DGEList object containing (at least) the elements counts (table of raw counts), `group` (factor indicating group) and `lib.size` (numeric vector of library sizes)

Value

`getCounts` returns the matrix of counts. `getOffset` returns a numeric matrix or vector. `getDispersion` returns vector of dispersion values.

Author(s)

Mark Robinson, Davis McCarthy, Gordon Smyth

See Also

[DGEList-class](#)

Examples

```
# generate raw counts from NB, create list object
y <- matrix(rnbinom(20,size=5,mu=10),5,4)
d <- DGEList(counts=y, group=c(1,1,2,2), lib.size=1001:1004)
getCounts(d)
getOffset(d)
d <- estimateCommonDisp(d)
getDispersion(d)
```

getPriorN

Get a Recommended Value for Prior N from DGEList Object

Description

Returns the `lib.size` component of the `samples` component of DGEList object multiplied by the `norm.factors` component

Usage

```
getPriorN(y, design=NULL, prior.df=20)
```

Arguments

<code>y</code>	a <code>DGEList</code> object with (at least) <code>elements</code> counts (table of unadjusted counts) and <code>samples</code> (data frame containing information about experimental group, library size and normalization factor for the library size)
<code>design</code>	numeric matrix (optional argument) giving the design matrix for the GLM that is to be fit. Must be of full column rank. If provided <code>design</code> is used to determine the number of parameters to be fit in the statistical model and therefore the residual degrees of freedom. If left as the default (<code>NULL</code>) then the <code>y\$samples\$group</code> element of the <code>DGEList</code> object is used to determine the residual degrees of freedom.
<code>prior.df</code>	numeric scalar giving the weight, in terms of prior degrees of freedom, to be given to the common parameter likelihood when estimating tagwise dispersion estimates.

Details

When estimating tagwise dispersion values using [estimateTagwiseDisp](#) or [estimateGLMtagwiseDisp](#) we need to decide how much weight to give to the common parameter likelihood in order to smooth (or stabilize) the dispersion estimates. The best choice of value for the `prior.n` parameter varies between datasets depending on the number of samples in the dataset and the complexity of the model to be fit. The value of `prior.n` should be inversely proportional to the residual degrees of freedom. We have found that choosing a value for `prior.n` that is equivalent to giving the common parameter likelihood 20 degrees of freedom generally gives a good amount of smoothing for the tagwise dispersion estimates. This function simply recommends an appropriate value for `prior.n`—to be used as an argument for [estimateTagwiseDisp](#) or [estimateGLMtagwiseDisp](#)—given the experimental design at hand and the chosen prior degrees of freedom.

Value

`getPriorN` returns a numeric scalar

Author(s)

Davis McCarthy, Gordon Smyth

See Also

[DGEList](#) for more information about the `DGEList` class. [as.matrix.DGEList](#).

Examples

```
# generate raw counts from NB, create list object
y<-matrix(rnbinom(20,size=1,mu=10),nrow=5)
d<-DGEList(counts=y,group=rep(1:2,each=2),lib.size=rep(c(1000:1001),2))
getPriorN(d)
```

Description

Fit a negative binomial generalized log-linear model to the read counts for each gene or transcript. Conduct genewise statistical tests for a given coefficient or coefficient contrast.

Usage

```
## S3 method for class DGEList
glmFit(y, design=NULL, dispersion=NULL, prior.count=0.125, start=NULL, ...)
## Default S3 method:
glmFit(y, design=NULL, dispersion=NULL, offset=NULL, lib.size=NULL, weights=NULL,
       prior.count=0.125, start=NULL, ...)
glmLRT(glmfit, coef=ncol(glmfit$design), contrast=NULL, test="chisq")
```

Arguments

<code>y</code>	an object that contains the raw counts for each library (the measure of expression level); alternatively, a matrix of counts, or a <code>DGEList</code> object with (at least) elements <code>counts</code> (table of unadjusted counts) and <code>samples</code> (data frame containing information about experimental group, library size and normalization factor for the library size)
<code>design</code>	numeric matrix giving the design matrix for the tagwise linear models. Must be of full column rank. Defaults to a single column of ones, equivalent to treating the columns as replicate libraries.
<code>dispersion</code>	numeric scalar or vector of negative binomial dispersions. Can be a common value for all tags, or a vector of values can provide a unique dispersion value for each tag. If <code>NULL</code> will be extracted from <code>y</code> , with order of precedence: tagwise dispersion, trended dispersions, common dispersion.
<code>offset</code>	numeric matrix of same size as <code>y</code> giving offsets for the log-linear models. Can be a scalar or a vector of length <code>ncol{y}</code> , in which case it is expanded out to a matrix.
<code>weights</code>	optional numeric matrix giving prior weights for the observations (for each library and transcript) to be used in the GLM calculations. Not supported by methods <code>"linesearch"</code> or <code>"levenberg"</code> .
<code>lib.size</code>	numeric vector of length <code>ncol(y)</code> giving library sizes. Only used if <code>offset=NULL</code> , in which case <code>offset</code> is set to <code>log(lib.size)</code> . Defaults to <code>colSums(y)</code> .
<code>prior.count</code>	average prior count to be added to observation to shrink the estimated log-fold-changes towards zero.
<code>start</code>	optional numeric matrix of initial estimates for the linear model coefficients.
<code>...</code>	other arguments are passed to lower level fitting functions.
<code>glmfit</code>	a <code>DGEGLM</code> object, usually output from <code>glmFit</code> .

<code>coef</code>	integer or character vector indicating which coefficients of the linear model are to be tested equal to zero. Values must be columns or column names of design. Defaults to the last coefficient. Ignored if <code>contrast</code> is specified.
<code>contrast</code>	numeric vector or matrix specifying one or more contrasts of the linear model coefficients to be tested equal to zero. Number of rows must equal to the number of columns of design. If specified, then takes precedence over <code>coef</code> .
<code>test</code>	which test (distribution) to use in calculating the p-values. Possible values are "F" or "chisq".

Details

`glmFit` and `glmLRT` implement generalized linear model (glm) methods developed by McCarthy et al (2012).

`glmFit` fits genewise negative binomial glms, all with the same design matrix but possibly different dispersions, offsets and weights. When the design matrix defines a one-way layout, or can be re-parametrized to a one-way layout, the glms are fitting very quickly using `mg1mOneGroup`. Otherwise the default fitting method, implemented in `mg1mLevenberg` a Fisher scoring algorithm with Levenberg-style damping.

Positive `prior.count` cause the returned coefficients to be shrunk in such a way that fold-changes between the treatment conditions are decreased. In particular, infinite fold-changes are avoided. Larger values cause more shrinkage. The returned coefficients are affected but not the likelihood ratio tests or p-values.

`glmLRT` conducts likelihood ratio tests for one or more coefficients in the linear model. If `coef` is used, the null hypothesis is that all the coefficients indicated by `coef` are equal to zero. If `contrast` is non-null, then the null hypothesis is that the specified contrasts of the coefficients are equal to zero. For example, a contrast of $c(0, 1, -1)$, assuming there are three coefficients, would test the hypothesis that the second and third coefficients are equal.

Value

`glmFit` produces an object of class `DGEGLM` containing components `counts`, `samples`, `genes` and `abundance` from `y` plus the following new components:

<code>design</code>	design matrix as input.
<code>weights</code>	matrix of weights as input.
<code>df.residual</code>	numeric vector of residual degrees of freedom, one for each tag.
<code>offset</code>	numeric matrix of linear model offsets.
<code>dispersion</code>	vector of dispersions used for the fit.
<code>coefficients</code>	numeric matrix of estimated coefficients from the glm fits, on the natural log scale, of size <code>nrow(y)</code> by <code>ncol(design)</code> .
<code>fitted.values</code>	matrix of fitted values from glm fits, same number of rows and columns as <code>y</code> .
<code>deviance</code>	numeric vector of deviances, one for each tag.

`glmLRT` produces objects of class `DGELRT` with the same components as for `glmfit` plus the following:

table	data frame with the same rows as y containing the log2-fold changes, likelihood ratio statistics and p-values, ready to be displayed by topTags..
comparison	character string describing the coefficient or the contrast being tested.

The data frame table contains the following columns:

logFC	log2-fold change of expression between conditions being tested.
logCPM	average log2-counts per million, the average taken over all libraries in y.
LR	likelihood ratio statistics.
PValue	p-values.

Author(s)

Davis McCarthy and Gordon Smyth

References

McCarthy, DJ, Chen, Y, Smyth, GK (2012). Differential expression analysis of multifactor RNA-Seq experiments with respect to biological variation. *Nucleic Acids Research* 40, 4288-4297. <http://nar.oxfordjournals.org/content/40/10/4288>

See Also

Low-level computations are done by [mg1mOneGroup](#) or [mg1mLevenberg](#).
[topTags](#) displays results from glmLRT.

Examples

```
nlibs <- 3
ntags <- 100
dispersion.true <- 0.1

# Make first transcript respond to covariate x
x <- 0:2
design <- model.matrix(~x)
beta.true <- cbind(Beta1=2,Beta2=c(2,rep(0,ntags-1)))
mu.true <- 2^(beta.true %*% t(design))

# Generate count data
y <- rnbinom(ntags*nlibs,mu=mu.true,size=1/dispersion.true)
y <- matrix(y,ntags,nlibs)
colnames(y) <- c("x0","x1","x2")
rownames(y) <- paste("Gene",1:ntags,sep="")
d <- DGEList(y)

# Normalize
d <- calcNormFactors(d)

# Fit the NB GLMs
fit <- glmFit(d, design, dispersion=dispersion.true)
```

```
# Likelihood ratio tests for trend
results <- glmLRT(fit, coef=2)
topTags(results)

# Estimate the dispersion (may be unreliable with so few tags)
d <- estimateGLMCommonDisp(d, design, verbose=TRUE)
```

glmQLFit

Quasi-likelihood methods with empirical Bayes shrinkage

Description

Fit a quasi-likelihood negative binomial generalized log-linear model to count data. Conduct gene-wise statistical tests for a given coefficient or coefficient contrast.

Usage

```
glmQLFit(y, design=NULL, dispersion=NULL, abundance.trend=TRUE, robust=FALSE, winsor.tail.p=c(0.05, 0), 0)
glmQLFTest(glmfit, coef=ncol(glmfit$design), contrast=NULL)
```

Arguments

<code>y</code>	a <code>DGEList</code> object containing count and sample data.
<code>design</code>	numeric matrix giving the design matrix for the tagwise linear models.
<code>dispersion</code>	numeric scalar or vector of negative binomial dispersions. Defaults to the trended dispersion, or the common dispersion (if no trend is available), or a value of 0.05 (if no common value is available).
<code>abundance.trend</code>	logical, whether to allow an abundance-dependent trend when estimating the prior values for the quasi-likelihood multiplicative dispersion parameter.
<code>robust</code>	logical, whether to estimate the prior degrees of freedom robustly.
<code>winsor.tail.p</code>	numeric vector of length 2 giving proportion to trim (Winsorize) from lower and upper tail of the distribution of genewise deviances when estimating the hyperparameters. Positive values produce robust empirical Bayes ignoring outlier small or large deviances. Only used when <code>robust=TRUE</code> .
<code>...</code>	other arguments are passed to <code>glmFit</code> .
<code>glmfit</code>	a <code>DGEGLM</code> object, usually output from <code>qlmQLFit</code> .
<code>coef</code>	integer or character vector indicating which coefficients of the linear model are to be tested equal to zero.
<code>contrast</code>	numeric vector or matrix specifying one or more contrasts of the linear model coefficients to be tested equal to zero.

Details

glmQLFTest implements the quasi-likelihood method of Lund et al (2012). It behaves the same as glmLRT except that it replaces likelihood ratio tests with quasi-likelihood F-tests for coefficients in the linear model. This function calls the limma function [squeezeVar](#) to conduct empirical Bayes smoothing of the genewise multiplicative dispersions. Note that the QuasiSeq package provides an alternative implementation of Lund et al (2012), with slightly different glm, trend and FDR methods.

There are a number of subtleties involved in the use of QL models. The first is that the negative binomial dispersions *must* be trended or common values. This is because the function assumes that the supplied values are the true values. For the trended/common values, the assumption is reasonable as information from many genes improves precision. This is not the case for the tagwise dispersions due to the limited information for each gene.

Another subtlety involves the handling of zero counts. Observations with fitted values of zero provide no residual degrees of freedom. This must be considered when computing the value of the quasi-likelihood dispersion for genes with many zeros. Finally, a lower bound is defined for the p-value of each gene, based on the likelihood ratio test. This avoids spurious results involving weak shrinkage with very low quasi-likelihood dispersions.

Value

glmQLFit produces an object of class DGEGLM with the same components as that produced by [glmFit](#), plus:

df.residual	a numeric vector containing the number of residual degrees of freedom for the GLM fit of each gene.
s2.fit	a list containing df.prior, the prior degrees of freedom; and var.prior, the location of the prior distribution. Both are numeric vectors if abundance.trend=TRUE and scalars otherwise. var.post is a numeric vector containing the shrunk quasi-likelihood dispersion for each gene.
df.prior	a numeric vector or scalar containing the prior degrees of freedom, same as that in s2.fit.

glmQFTest produces objects of class DGELRT with the same components as for [glmfit](#) plus the following:

table	data frame with the same rows as y containing the log2-fold changes, F-statistics and p-values, ready to be displayed by topTags..
comparison	character string describing the coefficient or the contrast being tested.

The data frame table contains the following columns:

logFC	log2-fold change of expression between conditions being tested.
logCPM	average log2-counts per million, the average taken over all libraries in y.
F	F-statistics.
PValue	p-values.

Author(s)

Davis McCarthy and Gordon Smyth, with modifications by Aaron Lun

References

Lund, SP, Nettleton, D, McCarthy, DJ, and Smyth, GK (2012). Detecting differential expression in RNA-sequence data using quasi-likelihood with shrunken dispersion estimates. *Statistical Applications in Genetics and Molecular Biology* Volume 11, Issue 5, Article 8. <http://www.statsci.org/smyth/pubs/QuasiSeqPreprint.pdf>

See Also

`topTags` displays results from `glmQLFTest`.

`plotQLDisp` can be used to visualize the distribution of QL dispersions after EB shrinkage from `glmQLFit`.

The `QuasiSeq` package gives an alternative implementation of `glmQLFTest` based on the same statistical ideas.

Examples

```
nlibs <- 4
ntags <- 1000
dispersion.true <- 1/rchisq(ntags, df=10)
design <- model.matrix(~factor(c(1,1,2,2)))

# Generate count data
y <- rnbinom(ntags*nlibs,mu=20,size=1/dispersion.true)
y <- matrix(y,ntags,nlibs)
d <- DGEList(y)
d <- calcNormFactors(d)

# Fit the NB GLMs with QL methods
d <- estimateDisp(d, design)
fit <- glmQLFit(d, design)
results <- glmQLFTest(fit)
topTags(results)
fit <- glmQLFit(d, design, robust=TRUE)
results <- glmQLFTest(fit)
topTags(results)
fit <- glmQLFit(d, design, abundance.trend=FALSE)
results <- glmQLFTest(fit)
topTags(results)
```

Description

Test for over-representation of gene ontology (GO) terms in the up and down differentially expressed genes from a linear model fit.

Usage

```
## S3 method for class DGELRT
goana(de, geneid = rownames(de), FDR = 0.05, species = "Hs",
      trend = FALSE, ...)
```

Arguments

de	an DGELRT object.
geneid	Entrez Gene identifiers. Either a vector of length nrow(de) or the name of the column of de\$genes containing the Entrez Gene IDs.
FDR	false discovery rate cutoff for differentially expressed genes. Numeric value between 0 and 1.
species	species identifier. Possible values are "Hs", "Mm", "Rn" or "Dm".
trend	adjust analysis for gene length or abundance? Can be logical, or a numeric vector of covariate values, or the name of the column of de\$genes containing the covariate values. If TRUE, then de\$AveLogCPM is used as the covariate.
...	any other arguments are passed to goana.default.

Details

Performs Gene Ontology enrichment analyses for the up and down differentially expressed genes from a linear model analysis. The Entrez Gene ID must be supplied for each gene.

If trend=FALSE, the function computes one-sided hypergeometric tests equivalent to Fisher's exact test.

If trend=TRUE or a covariate is supplied, then a trend is fitted to the differential expression results and the method of Young et al (2010) is used to adjust for this trend. The adjusted test uses Wallenius' noncentral hypergeometric distribution.

Value

A data frame with a row for each GO term and the following columns:

Term	GO term.
Ont	ontology that the GO term belongs to. Possible values are "BP", "CC" and "MF".
N	Number of genes in the GO term.
Up	number of up-regulated differentially expressed genes.
Down	number of down-regulated differentially expressed genes.
P.Up	p-value for over-representation of GO term in up-regulated genes.
P.Down	p-value for over-representation of GO term in down-regulated genes.

The row names of the data frame give the GO term IDs.

Author(s)

Gordon Smyth and Yifang Hu

References

Young, M. D., Wakefield, M. J., Smyth, G. K., Oshlack, A. (2010). Gene ontology analysis for RNA-seq: accounting for selection bias. *Genome Biology* 11, R14. <http://genomebiology.com/2010/11/2/R14>

See Also

[goana](#), [topGO](#)

Examples

```
## Not run:

fit <- glmFit(y, design)
lrt <- glmLRT(fit)
go <- goana(lrt)
topGO(go, ont="BP", sort = "up")
topGO(go, ont="BP", sort = "down")

## End(Not run)
```

gof

Goodness of Fit Tests for Multiple GLM Fits

Description

Conducts deviance goodness of fit tests for each fit in a DGEGLM object

Usage

```
gof(glmfit, pcutoff=0.1, adjust="holm", plot=FALSE,
     main="qq-plot of genewise goodness of fit", ...)
```

Arguments

glmfit	DGEGLM object containing results from fitting NB GLMs to genes in a DGE dataset. Output from glmFit.
pcutoff	scalar giving the cut-off value for the Holm-adjusted p-value. Genes with Holm-adjusted p-values lower than this cutoff value are flagged as ‘dispersion outlier’ genes.
adjust	method used to adjust goodness of fit p-values for multiple testing.
plot	logical, if TRUE a qq-plot is produced.
main	character, title for the plot.
...	other arguments are passed to qqnorm.

Details

If `plot=TRUE`, produces a plot similar to Figure 2 of McCarthy et al (2012).

Value

This function returns a list with the following components:

<code>gof.statistics</code>	numeric vector of deviance statistics, which are the statistics used for the goodness of fit test
<code>gof.pvalues</code>	numeric vector of p-values providing evidence of poor fit; computed from the chi-square distribution on the residual degrees of freedom from the GLM fits.
<code>outlier</code>	logical vector indicating whether or not each gene is a ‘dispersion outlier’ (i.e., the model fit is poor for that gene indicating that the dispersion estimate is not good for that gene).
<code>df</code>	scalar, the residual degrees of freedom from the GLM fit for which the goodness of fit statistics have been computed. Also the degrees of freedom for the goodness of fit statistics for the LR (chi-square) test for significance.

Author(s)

Davis McCarthy and Gordon Smyth

References

McCarthy, DJ, Chen, Y, Smyth, GK (2012). Differential expression analysis of multifactor RNA-Seq experiments with respect to biological variation. *Nucleic Acids Research* 40, 4288-4297 <http://nar.oxfordjournals.org/content/40/10/4288>

See Also

[qqnorm](#).

[glmFit](#) for more information on fitting NB GLMs to DGE data.

Examples

```
nlibs <- 3
ntags <- 100
dispersion.true <- 0.1

# Make first transcript respond to covariate x
x <- 0:2
design <- model.matrix(~x)
beta.true <- cbind(Beta1=2,Beta2=c(2,rep(0,ntags-1)))
mu.true <- 2^(beta.true %*% t(design))

# Generate count data
y <- rnbinom(ntags*nlibs,mu=mu.true,size=1/dispersion.true)
y <- matrix(y,ntags,nlibs)
colnames(y) <- c("x0","x1","x2")
```

```
rownames(y) <- paste("Gene",1:ntags,sep="")
d <- DGEList(y)

# Normalize
d <- calcNormFactors(d)

# Fit the NB GLMs
fit <- glmFit(d, design, dispersion=dispersion.true)
# Check how good the fit is for each gene
gof(fit)
```

goodTuring

Good-Turing Frequency Estimation

Description

Non-parametric empirical Bayes estimates of the frequencies of observed (and unobserved) species.

Usage

```
goodTuring(x, conf=1.96)
goodTuringPlot(x)
goodTuringProportions(counts)
```

Arguments

x	numeric vector of non-negative integers, representing the observed frequency of each species.
conf	confidence factor, as a quantile of the standard normal distribution, used to decide for what values the log-linear relationship between frequencies and frequencies of frequencies is acceptable.
counts	matrix of counts

Details

Observed counts are assumed to be Poisson distributed. Using a non-parametric empirical Bayes strategy, the algorithm evaluates the posterior expectation of each species mean given its observed count. The posterior means are then converted to proportions. In the empirical Bayes step, the counts are smoothed by assuming a log-linear relationship between frequencies and frequencies of frequencies. The fundamentals of the algorithm are from Good (1953). Gale and Sampson (1995) proposed a simplified algorithm with a rule for switching between the observed and smoothed frequencies, and it is Gale and Sampson's simplified algorithm that is implemented here. The number of zero values in x are not used in the algorithm, but is returned by this function.

Sampson gives a C code version on his webpage at <http://www.grsampson.net/RGoodTur.html> which gives identical results to this function.

goodTuringPlot plots log-probability (i.e., log frequencies of frequencies) versus log-frequency.

goodTuringProportions runs goodTuring on each column of data, then uses the results to predict the proportion of each tag in each library.

Value

goodTuring returns a list with components

count	observed frequencies, i.e., the unique positive values of x
n	frequencies of frequencies
n0	frequency of zero, i.e., number of zeros found in x
proportion	estimated proportion of each species given its count
P0	estimated combined proportion of all undetected species

goodTuringProportions returns a matrix of proportions of the same size as counts.

Author(s)

Aaron Lun and Gordon Smyth, adapted from Sampson's C code from <http://www.grsampson.net/RGoodTur.html>

References

Gale, WA, and Sampson, G (1995). Good-Turing frequency estimation without tears. *Journal of Quantitative Linguistics* 2, 217-237.

Examples

```
# True means of observed species
lambda <- rnbinom(10000,mu=2,size=1/10)
lambda <- lambda[lambda>1]

# Observed frequencies
Ntrue <- length(lambda)
x <- rpois(Ntrue, lambda=lambda)
freq <- goodTuring(x)
goodTuringPlot(x)
```

loessByCol

Locally Weighted Mean By Column

Description

Smooth columns of matrix by non-robust loess curves of degree 0.

Usage

```
loessByCol(y, x=NULL, span=0.5)
locfitByCol(y, x=NULL, weights=1, span=0.5, degree=0)
```

Arguments

y	numeric matrix of response variables.
x	numeric covariate vector of length <code>nrow(y)</code> , defaults to equally spaced.
span	width of the smoothing window, in terms of proportion of the data set. Larger values produce smoother curves.
weights	relative weights of each observation, one for each covariate value.
degree	degree of local polynomial fit

Details

Fits a loess curve with degree 0 to each column of the response matrix, using the same covariate vector for each column. The smoothed column values are tricube-weighted means of the original values.

`locfitByCol` uses the `locfit.raw` function of the `locfit` package.

Value

A list containing a numeric matrix with smoothed columns and a vector of leverages for each covariate value.

`locfitByCol` returns a numeric matrix.

Author(s)

Aaron Lun for `loessByCol`, replacing earlier R code by Davis McCarthy. Gordon Smyth for `locfitByCol`.

See Also

[loess](#)

Examples

```
y <- matrix(rnorm(100*3), nrow=100, ncol=3)
head(y)
out <- loessByCol(y)
head(out$fitted.values)
```

maPlot

Plots Log-Fold Change versus Log-Concentration (or, M versus A) for Count Data

Description

To represent counts that were low (e.g. zero in 1 library and non-zero in the other) in one of the two conditions, a 'smear' of points at low A value is presented.

Usage

```
maPlot(x, y, logAbundance=NULL, logFC=NULL, normalize=FALSE, plot.it=TRUE,
       smearWidth=1, col=NULL, allCol="black", lowCol="orange", deCol="red",
       de.tags=NULL, smooth.scatter=FALSE, lowess=FALSE, ...)
```

Arguments

x	vector of counts or concentrations (group 1)
y	vector of counts or concentrations (group 2)
logAbundance	vector providing the abundance of each tag on the log ₂ scale. Purely optional (default is NULL), but in combination with logFC provides a more direct way to create an MA-plot if the log-abundance and log-fold change are available.
logFC	vector providing the log-fold change for each tag for a given experimental contrast. Default is NULL, only to be used together with logAbundance as both need to be non-null for their values to be used.
normalize	logical, whether to divide x and y vectors by their sum
plot.it	logical, whether to produce a plot
smearWidth	scalar, width of the smear
col	vector of colours for the points (if NULL, uses allCol and lowCol)
allCol	colour of the non-smearred points
lowCol	colour of the smearred points
deCol	colour of the DE (differentially expressed) points
de.tags	indices for tags identified as being differentially expressed; use exactTest to identify DE genes
smooth.scatter	logical, whether to produce a 'smooth scatter' plot using the KernSmooth::smoothScatter function or just a regular scatter plot; default is FALSE, i.e. produce a regular scatter plot
lowess	logical, indicating whether or not to add a lowess curve to the MA-plot to give an indication of any trend in the log-fold change with log-concentration
...	further arguments passed on to plot

Details

The points to be smeared are identified as being equal to the minimum in one of the two groups. The smear is created by using random uniform numbers of width smearWidth to the left of the minimum A value.

Value

a plot to the current device (if plot.it=TRUE), and invisibly returns the M (logFC) and A (logConc) values used for the plot, plus identifiers w and v of genes for which M and A values, or just M values, respectively, were adjusted to make a nicer looking plot.

Author(s)

Mark Robinson, Davis McCarthy

See Also

[plotSmear](#)

Examples

```
y <- matrix(rnbinom(10000,mu=5,size=2),ncol=4)
maPlot(y[,1], y[,2])
```

maximizeInterpolant *Maximize a function given a table of values by spline interpolation.*

Description

Maximize a function given a table of values by spline interpolation.

Usage

```
maximizeInterpolant(x, y)
```

Arguments

x numeric vector of the inputs of the function.
y numeric matrix of function values at the values of x. Columns correspond to x values and each row corresponds to a different function to be maximized.

Details

Calculates the cubic spline interpolant for each row the method of Forsythe et al (1977) using the function `fmm_spline` from `splines.c` in the `stats` package). Then calculates the derivatives of the spline segments adjacent to the input with the maximum function value. This allows identification of the maximum of the interpolating spline.

Value

numeric vector of input values at which the function maximums occur.

Author(s)

Aaron Lun, improving on earlier code by Gordon Smyth

References

Forsythe, G. E., Malcolm, M. A. and Moler, C. B. (1977). *Computer Methods for Mathematical Computations*, Prentice-Hall.

Examples

```
x <- seq(0,1,length=10)
y <- rnorm(10,1,1)
maximizeInterpolant(x,y)
```

maximizeQuadratic	<i>Maximize a function given a table of values by quadratic interpolation.</i>
-------------------	--

Description

Maximize a function given a table of values by quadratic interpolation.

Usage

```
maximizeQuadratic(y, x=1:ncol(y))
```

Arguments

y	numeric matrix of response values.
x	numeric matrix of inputs of the function of same dimension as y. If a vector, must be a row vector of length equal to ncol(y).

Details

For each row of y, finds the three x values bracketing the maximum of y, interpolates a quadratic polynomial through these y for these three values and solves for the location of the maximum of the polynomial.

Value

numeric vector of length equal to nrow(y) giving the x-value at which y is maximized.

Author(s)

Yunshun Chen and Gordon Smyth

See Also

[maximizeInterpolant](#)

Examples

```
y <- matrix(rnorm(5*9),5,9)
maximizeQuadratic(y)
```

 meanvar

Explore the mean-variance relationship for DGE data

Description

Appropriate modelling of the mean-variance relationship in DGE data is important for making inferences about differential expression. Here are functions to compute tag/gene means and variances, as well as looking at these quantities when data is binned based on overall expression level.

Usage

```
plotMeanVar(object, meanvar=NULL, show.raw.vars=FALSE, show.tagwise.vars=FALSE,
            show.binned.common.disp.vars=FALSE, show.ave.raw.vars=TRUE,
            scalar=NULL, NBlines=FALSE, nbins=100, log.axes="xy", xlab=NULL,
            ylab=NULL, ...)
binMeanVar(x, group, nbins=100, common.dispersion=FALSE, object=NULL)
```

Arguments

object	DGEList object containing the raw data and dispersion value. According to the method desired for computing the dispersion, either <code>estimateCommonDisp</code> and (possibly) <code>estimateTagwiseDisp</code> should be run on the DGEList object before using <code>plotMeanVar</code> . The argument <code>object</code> must be supplied in the function <code>binMeanVar</code> if common dispersion values are to be computed for each bin.
meanvar	list (optional) containing the output from <code>binMeanVar</code> or the returned value of <code>plotMeanVar</code> . Providing this object as an argument will save time in computing the tag/gene means and variances when producing a mean-variance plot.
show.raw.vars	logical, whether or not to display the raw (pooled) gene/tag variances on the mean-variance plot. Default is FALSE.
show.tagwise.vars	logical, whether or not to display the estimated genewise/tagwise variances on the mean-variance plot. Default is FALSE.
show.binned.common.disp.vars	logical, whether or not to compute the common dispersion for each bin of tags and show the variances computed from those binned common dispersions and the mean expression level of the respective bin of tags. Default is FALSE.
show.ave.raw.vars	logical, whether or not to show the average of the raw variances for each bin of tags plotted against the average expression level of the tags in the bin. Averages are taken on the square root scale as regular arithmetic means are likely to be upwardly biased for count data, whereas averaging on the square scale gives a better summary of the mean-variance relationship in the data. The default is TRUE.
scalar	vector (optional) of scaling values to divide counts by. Would expect to have this the same length as the number of columns in the count matrix (i.e. the number of libraries).

NBline	logical, whether or not to add a line on the graph showing the mean-variance relationship for a NB model with common dispersion.
nbins	scalar giving the number of bins (formed by using the quantiles of the genewise mean expression levels) for which to compute average means and variances for exploring the mean-variance relationship. Default is 100 bins
log.axes	character vector indicating if any of the axes should use a log scale. Default is "xy", which makes both y and x axes on the log scale. Other valid options are "x" (log scale on x-axis only), "y" (log scale on y-axis only) and "" (linear scale on x- and y-axis).
xlab	character string giving the label for the x-axis. Standard graphical parameter. If left as the default NULL, then the x-axis label will be set to "logConc".
ylab	character string giving the label for the y-axis. Standard graphical parameter. If left as the default NULL, then the x-axis label will be set to "logConc".
...	further arguments passed on to plot
x	matrix of count data, with rows representing tags/genes and columns representing samples
group	factor giving the experimental group or condition to which each sample (i.e. column of x or element of y) belongs
common.dispersion	logical, whether or not to compute the common dispersion for each bin of tags.

Details

This function is useful for exploring the mean-variance relationship in the data. Raw variances are, for each gene, the pooled variance of the counts from each sample, divided by a scaling factor (by default the effective library size). The function will plot the average raw variance for tags split into nbins bins by overall expression level. The averages are taken on the square-root scale as for count data the arithmetic mean is upwardly biased. Taking averages on the square-root scale provides a useful summary of how the variance of the gene counts change with respect to expression level (abundance). A line showing the Poisson mean-variance relationship (mean equals variance) is always shown to illustrate how the genewise variances may differ from a Poisson mean-variance relationship. Optionally, the raw variances and estimated tagwise variances can also be plotted. Estimated tagwise variances can be calculated using either qCML estimates of the tagwise dispersions (`estimateTagwiseDisp`) or Cox-Reid conditional inference estimates (`CRDisp`). A log-log scale is used for the plot.

Value

`plotMeanVar` produces a mean-variance plot for the DGE data using the options described above. `plotMeanVar` and `binMeanVar` both return a list with the following components:

avemeans	vector of the average expression level within each bin of genes, with the average taken on the square-root scale
avevars	vector of the average raw pooled gene-wise variance within each bin of genes, with the average taken on the square-root scale
bin.means	list containing the average (mean) expression level for genes divided into bins based on amount of expression

bin.vars	list containing the pooled variance for genes divided into bins based on amount of expression
means	vector giving the mean expression level for each gene
vars	vector giving the pooled variance for each gene
bins	list giving the indices of the tags in each bin, ordered from lowest expression bin to highest

Author(s)

Davis McCarthy

See Also

[plotMDS.DGEList](#), [plotSmear](#) and [maPlot](#) provide more ways of visualizing DGE data.

Examples

```

y <- matrix(rnbinom(1000,mu=10,size=2),ncol=4)
d <- DGEList(counts=y,group=c(1,1,2,2),lib.size=c(1000:1003))
plotMeanVar(d) # Produce a straight-forward mean-variance plot
# Produce a mean-variance plot with the raw variances shown and save the means
# and variances for later use
meanvar <- plotMeanVar(d, show.raw.vars=TRUE)
## If we want to show estimated tagwise variances on the plot, we must first estimate them!
d <- estimateCommonDisp(d) # Obtain an estimate of the dispersion parameter
d <- estimateTagwiseDisp(d) # Obtain tagwise dispersion estimates
# Use previously saved object to speed up plotting
plotMeanVar(d, meanvar=meanvar, show.tagwise.vars=TRUE, NBlind=TRUE)
## We could also estimate common/tagwise dispersions using the Cox-Reid methods with an
## appropriate design matrix

```

mglm

Fit Negative Binomial Generalized Linear Model to Multiple Response Vectors: Low Level Functions

Description

Fit the same log-link negative binomial or Poisson generalized linear model (GLM) to each row of a matrix of counts.

Usage

```

mglmOneGroup(y, dispersion=0, offset=0, weights=NULL, maxit=50, tol=1e-10,
             verbose=FALSE, coef.start=NULL)
mglmOneWay(y, design=NULL, dispersion=0, offset=0, weights=NULL, maxit=50,
           tol=1e-10, coef.start=NULL)
mglmLevenberg(y, design, dispersion=0, offset=0, weights=NULL,
             coef.start=NULL, start.method="null", maxit=200, tol=1e-06)
designAsFactor(design)

```

Arguments

<code>y</code>	numeric matrix containing the negative binomial counts. Rows for tags and columns for libraries.
<code>design</code>	numeric matrix giving the design matrix of the GLM. Assumed to be full column rank.
<code>dispersion</code>	numeric scalar or vector giving the dispersion parameter for each GLM. Can be a scalar giving one value for all tags, or a vector of length equal to the number of tags giving tag-wise dispersions.
<code>offset</code>	numeric vector or matrix giving the offset that is to be included in the log-linear model predictor. Can be a scalar, a vector of length equal to the number of libraries, or a matrix of the same size as <code>y</code> .
<code>weights</code>	numeric vector or matrix of non-negative quantitative weights. Can be a vector of length equal to the number of libraries, or a matrix of the same size as <code>y</code> .
<code>coef.start</code>	numeric matrix of starting values for the linear model coefficients. Number of rows should agree with <code>y</code> and number of columns should agree with <code>design</code> .
<code>start.method</code>	method used to generate starting values when <code>coef.stat=NULL</code> . Possible values are "null" to start from the null model of equal expression levels or "y" to use the data as starting value for the mean.
<code>tol</code>	numeric scalar giving the convergence tolerance. For <code>mglmOneGroup</code> , convergence is judged successful when the step size falls below <code>tol</code> in absolute size.
<code>maxit</code>	scalar giving the maximum number of iterations for the Fisher scoring algorithm.
<code>verbose</code>	logical. If TRUE, warnings will be issued when <code>maxit</code> iterations are exceeded before convergence is achieved.

Details

The functions `mglmOneGroup`, `mglmOneWay` and `mglmLevenberg` all fit negative binomial generalized linear models, with the same design matrix but possibly different dispersions, offsets and weights, to a series of response vectors. The functions are all low-level functions in that they operate on atomic objects such as matrices. They are used as work-horses by higher-level functions in the edgeR package, especially by `glmFit`.

`mglmOneGroup` fit the null model, with intercept term only, to each response vector. In other words, it treats the libraries as belonging to one group. It implements Fisher scoring with a score-statistic stopping criterion for each tag. Excellent starting values are available for the null model, so this function seldom has any problems with convergence. It is used by other edgeR functions to compute the overall abundance for each tag.

`mglmLevenberg` fits an arbitrary log-linear model to each response vector. It implements a Levenberg-Marquardt modification of the glm scoring algorithm to prevent divergence. The main computation is implemented in C++.

All these functions treat the dispersion parameter of the negative binomial distribution as a known input.

`deviances.function` chooses the appropriate deviance function to use given a scalar or vector of dispersion parameters. If the dispersion values are zero, then the Poisson deviance function is returned; if the dispersion values are positive, then the negative binomial deviance function is returned.

Value

`mglimOneGroup` produces a vector of length equal to the number of tags/genes (number of rows of `y`) providing the single coefficient from the GLM fit for each tag/gene. This can be interpreted as a measure of the 'average expression' level of the tag/gene.

`mglimLevenberg` produces a list with the following components:

<code>coefficients</code>	matrix of estimated coefficients for the linear models
<code>fitted.values</code>	matrix of fitted values
<code>deviance</code>	residual deviances
<code>iter</code>	number of iterations used
<code>fail</code>	logical vector indicating tags for which the maximum damping was exceeded before convergence was achieved

`deviances.function` returns a function to calculate the deviance as appropriate for the given values of the dispersion.

`designAsFactor` returns a factor of length equal to `nrow(design)`.

Author(s)

Gordon Smyth, Yunshun Chen, Davis McCarthy, Aaron Lun. C++ code by Aaron Lun.

References

McCarthy, DJ, Chen, Y, Smyth, GK (2012). Differential expression analysis of multifactor RNA-Seq experiments with respect to biological variation. *Nucleic Acids Research* 40, 4288-4297. <http://nar.oxfordjournals.org/content/40/10/4288>

See Also

[glmFit](#), for more object-orientated GLM modelling for DGE data.

Examples

```
y <- matrix(rnbinom(1000,mu=10,size=2),ncol=4)
lib.size <- colSums(y)
dispersion <- 0.1

abundance <- mglimOneGroup(y, dispersion=dispersion, offset=log(lib.size))
AveLogCPM <- log1p(exp(1e6*abundance))/log(2)
summary(AveLogCPM)

## Same as above:
AveLogCPM <- aveLogCPM(y, dispersion, offset=log(lib.size))

## Fit the NB GLM to the counts with a given design matrix
f1 <- factor(c(1,1,2,2))
f2 <- factor(c(1,2,1,2))
x <- model.matrix(~f1+f2)
fit <- mglimLevenberg(y, x, dispersion=dispersion, offset=log(lib.size))
head(fit$coefficients)
```

movingAverageByCol *Moving Average Smoother of Matrix Columns*

Description

Apply a moving average smoother to the columns of a matrix.

Usage

```
movingAverageByCol(x, width=5, full.length=TRUE)
```

Arguments

x	numeric matrix
width	integer, width of window of rows to be averaged
full.length	logical value, should output have same number of rows as input?

Details

If `full.length=TRUE`, narrower windows are used at the start and end of each column to make a column of the same length as input. If `FALSE`, all values are average of `width` input values, so the number of rows is less than input.

Value

Numeric matrix containing smoothed values. If `full.length=TRUE`, of same dimension as `x`. If `full.length=FALSE`, has `width-1` fewer rows than `x`.

Author(s)

Gordon Smyth

Examples

```
x <- matrix(rpois(20,lambda=5),10,2)
movingAverageByCol(x,3)
```

nbinomDeviance	<i>Negative Binomial Deviance</i>
----------------	-----------------------------------

Description

Fit the same log-link negative binomial or Poisson generalized linear model (GLM) to each row of a matrix of counts.

Usage

```
nbinomUnitDeviance(y, mean, dispersion=0)
nbinomDeviance(y, mean, dispersion=0, weights=NULL)
```

Arguments

y	numeric vector or matrix containing the negative binomial counts. If a matrix, then rows for tags and columns for libraries. nbinomDeviance treats a vector as a matrix with one row.
mean	numeric vector matrix of expected values, of same dimension as y.
dispersion	numeric vector or matrix of negative binomial dispersions. Can be a scalar, or a vector of length equal to the number of tags, or a matrix of same dimensions as y.
weights	numeric vector or matrix of non-negative weights, as for glmFit.

Details

nbinomUnitDeviance computes the unit deviance for each y observation. nbinomDeviance computes the total residual deviance for each row of y observation, i.e., weighted row sums of the unit deviances.

Care is taken to ensure accurate computation for small dispersion values.

Value

nbinomUnitDeviance returns a numeric vector or matrix of the same size as y.

nbinomDeviance returns a numeric vector of length equal to the number of rows of y.

Author(s)

Gordon Smyth, Yunshun Chen, Aaron Lun. C++ code by Aaron Lun.

References

Jorgensen, B. (2006). Generalized linear models. Encyclopedia of Environmetrics, Wiley. <http://onlinelibrary.wiley.com/doi/10.1002/9780470057339.vag010/full>.

McCarthy, DJ, Chen, Y, Smyth, GK (2012). Differential expression analysis of multifactor RNA-Seq experiments with respect to biological variation. *Nucleic Acids Research* 40, 4288-4297. <http://nar.oxfordjournals.org/content/40/10/4288>

Examples

```

y <- matrix(1:6,3,2)
mu <- matrix(3,3,2)
nbinomUnitDeviance(y,mu,dispersion=0.2)
nbinomDeviance(y,mu,dispersion=0.2)

```

normalizeChIPtoInput *Normalize ChIP-Seq Read Counts to Input and Test for Enrichment*

Description

Normalize ChIP-Seq read counts to input control values, then test for significant enrichment relative to the control.

Usage

```

normalizeChIPtoInput(input, response, dispersion=0.01, niter=6, loss="p", plot=FALSE,
                    verbose=FALSE, ...)
calcNormOffsetsforChIP(input, response, dispersion=0.01, niter=6, loss="p", plot=FALSE,
                    verbose=FALSE, ...)

```

Arguments

input	numeric vector of non-negative input values, not necessarily integer.
response	vector of non-negative integer counts of some ChIP-Seq mark for each gene or other genomic feature.
dispersion	negative binomial dispersion, must be positive.
niter	number of iterations.
loss	loss function to be used when fitting the response counts to the input: "p" for cumulative probabilities or "z" for z-value.
plot	if TRUE, a plot of the fit is produced.
verbose	if TRUE, working estimates from each iteration are output.
...	other arguments are passed to the plot function.

Details

normalizeChIPtoInput identifies significant enrichment for a ChIP-Seq mark relative to input values. The ChIP-Seq mark might be for example transcriptional factor binding or an epigenetic mark. The function works on the data from one sample. Replicate libraries are not explicitly accounted for, and would normally be pooled before using this function.

ChIP-Seq counts are assumed to be summarized by gene or similar genomic feature of interest.

This function makes the assumption that a non-negligible proportion of the genes, say 25% or more, are not truly marked by the ChIP-Seq feature of interest. Unmarked genes are further assumed to have counts at a background level proportional to the input. The function aligns the counts to

the input so that the counts for the unmarked genes behave like a random sample. The function estimates the proportion of marked genes, and removes marked genes from the fitting process. For this purpose, marked genes are those with a Holm-adjusted mid-p-value less than 0.5.

The read counts are treated as negative binomial. The dispersion parameter is not estimated from the data; instead a reasonable value is assumed to be given.

calcNormOffsetsforChIP returns a numeric matrix of offsets, ready for linear modelling.

Value

normalizeChIPtoInput returns a list with components

p.value numeric vector of p-values for enrichment.

scaling.factor factor by which input is scaled to align with response counts for unmarked genes.

prop.enriched proportion of marked genes, as internally estimated

calcNormOffsetsforChIP returns a numeric matrix of offsets.

Author(s)

Gordon Smyth

plotBCV

Plot Biological Coefficient of Variation

Description

Plot genewise biological coefficient of variation (BCV) against gene abundance (in log₂ counts per million).

Usage

```
plotBCV(y, xlab="Average log CPM", ylab="Biological coefficient of variation",
        pch=16, cex=0.2, col.common="red", col.trend="blue", col.tagwise="black", ...)
```

Arguments

y a DGEList object.

xlab label for the x-axis.

ylab label for the y-axis.

pch the plotting symbol. See [points](#) for more details.

cex plot symbol expansion factor. See [points](#) for more details.

col.common color of line showing common dispersion

col.trend color of line showing dispersion trend

col.tagwise color of points showing tagwise dispersions

... any other arguments are passed to plot.

Details

The BCV is the square root of the negative binomial dispersion. This function displays the common, trended and tagwise BCV estimates.

Value

A plot is created on the current graphics device.

Author(s)

Davis McCarthy, Yunshun Chen, Gordon Smyth

Examples

```
BCV.true <- 0.1
y <- DGEList(matrix(rnbinom(6000, size = 1/BCV.true^2, mu = 10),1000,6))
y <- estimateCommonDisp(y)
y <- estimateTrendedDisp(y)
y <- estimateTagwiseDisp(y)
plotBCV(y)
```

plotExonUsage

Create a Plot of Exon Usage from Exon-Level Count Data

Description

Create a plot of exon usage for a given gene by plotting the (un)transformed counts for each exon, coloured by experimental group.

Usage

```
plotExonUsage(y, geneID, group=NULL, transform="none", counts.per.million=TRUE,
              legend.coords=NULL, ...)
```

Arguments

y	either a matrix of exon-level counts, a list containing a matrix of counts for each exon or a DGEList object with (at least) elements counts (table of counts summarized at the exon level) and samples (data frame containing information about experimental group, library size and normalization factor for the library size). Each row of y should represent one exon.
geneID	character string giving the name of the gene for which exon usage is to be plotted.
group	factor supplying the experimental group/condition to which each sample (column of y) belongs. If NULL (default) the function will try to extract it from y, which only works if y is a DGEList object.

transform character, supplying the method of transformation to be applied to the exon counts, if any. Options are "none" (original counts are preserved), "sqrt" (square-root transformation) and "log2" (log2 transformation). Default is "none".

counts.per.million logical, if TRUE then counts per million (as determined from total library sizes) will be plotted for each exon, if FALSE the raw read counts will be plotted. Using counts per million effectively normalizes for different read depth among the different samples, which can make the exon usage plots easier to interpret.

legend.coords optional vector of length 2 giving the x- and y-coordinates of the legend on the plot. If NULL (default), the legend will be automatically placed near the top right corner of the plot.

... optional further arguments to be passed on to plot.

Details

This function produces a simple plot for comparing exon usage between different experimental conditions for a given gene.

Value

plotExonUsage (invisibly) returns the transformed matrix of counts for the gene being plotted and produces a plot to the current device.

Author(s)

Davis McCarthy, Gordon Smyth

See Also

[spliceVariants](#) for methods to detect genes with evidence for alternative exon usage.

Examples

```
# generate exon counts from NB, create list object
y<-matrix(rnbinom(40,size=1,mu=10),nrow=10)
rownames(y) <- rep(c("gene.1","gene.2"), each=5)
d<-DGEList(counts=y,group=rep(1:2,each=2))
plotExonUsage(d, "gene.1")
```

plotMDS.DGEList	<i>Multidimensional scaling plot of distances between digital gene expression profiles</i>
-----------------	--

Description

Plot samples on a two-dimensional scatterplot so that distances on the plot approximate the expression differences between the samples.

Usage

```
## S3 method for class DGEList
plotMDS(x, top = 500, labels = NULL, pch = NULL, cex = 1,
        dim.plot = c(1,2), ndim = max(dim.plot), gene.selection = "pairwise",
        xlab = NULL, ylab = NULL, method = "logFC", prior.count = 2,
        ...)
```

Arguments

x	a DGEList object.
top	number of top genes used to calculate pairwise distances.
labels	character vector of sample names or labels. If x has no column names, then defaults the index of the samples.
pch	plotting symbol or symbols. See points for possible values. Ignored if labels is non-NULL.
cex	numeric vector of plot symbol expansions. See text for possible values.
dim.plot	which two dimensions should be plotted, numeric vector of length two.
ndim	number of dimensions in which data is to be represented
gene.selection	character, "pairwise" to choose the top genes separately for each pairwise comparison between the samples or "common" to select the same genes for all comparisons. Only used when method="logFC".
xlab	x-axis label
ylab	y-axis label
method	method used to compute distances. Possible values are "logFC" or "bcv".
prior.count	average prior count to be added to observation to shrink the estimated log-fold-changes towards zero. Only used when method="logFC".
...	any other arguments are passed to plot.

Details

The default method (method="logFC") is to convert the counts to log-counts-per-million using `cpm` and to pass these to the `limma plotMDS` function. This method calculates distances between samples based on log₂ fold changes. See the [plotMDS help page](#) for details.

The alternative method (method="bcv") calculates distances based on biological coefficient of variation. A set of top genes are chosen that have largest biological variation between the libraries (those with largest tagwise dispersion treating all libraries as one group). Then the distance between each pair of libraries (columns) is the biological coefficient of variation (square root of the common dispersion) between those two libraries alone, using the top genes.

The number of genes (top) chosen for this exercise should roughly correspond to the number of differentially expressed genes with materially large fold-changes. The default setting of 500 genes is widely effective and suitable for routine use, but a smaller value might be chosen for when the samples are distinguished by a specific focused molecular pathway. Very large values (greater than 1000) are not usually so effective.

Note that the "bcv" method is slower than the "logFC" method when there are many libraries.

Value

An object of class `MDS` is invisibly returned and a plot is created on the current graphics device.

Author(s)

Yunshun Chen, Mark Robinson and Gordon Smyth

See Also

[plotMDS](#), [cmdscale](#), [as.dist](#)

Examples

```
# Simulate DGE data for 1000 genes(tags) and 6 samples.
# Samples are in two groups
# First 200 genes are differentially expressed in second group

ngenes <- 1000
nlib <- 6
counts <- matrix(rnbinom(ngenes*nlib, size=1/10, mu=20),ngenes,nlib)
rownames(counts) <- paste("Gene",1:ngenes)
group <- gl(2,3,labels=c("Grp1","Grp2"))
counts[1:200,group=="Grp2"] <- counts[1:200,group=="Grp2"] + 10
y <- DGEList(counts,group=group)
y <- calcNormFactors(y)

# without labels, indexes of samples are plotted.
col <- as.numeric(group)
mds <- plotMDS(y, top=200, col=col)

# or labels can be provided, here group indicators:
plotMDS(mds, col=col, labels=group)
```

plotQLDisp

Plot the quasi-likelihood dispersion

Description

Plot the genewise quasi-likelihood dispersion against the gene abundance (in log₂ counts per million).

Usage

```
plotQLDisp(glmfit, xlab="Average Log2 CPM", ylab="Quarter-Root Mean Deviance", pch=16, cex=0.2,
  col.shrunk="red", col.trend="blue", col.raw="black", ...)
```


Arguments

<code>glmfit</code>	a DGEGLM object produced by <code>glmQLFit</code> .
<code>xlab</code>	label for the x-axis.
<code>ylab</code>	label for the y-axis.
<code>pch</code>	the plotting symbol. See points for more details.
<code>cex</code>	plot symbol expansion factor. See points for more details.
<code>col.shrunk</code>	color of the points representing the shrunk quasi-likelihood dispersions.
<code>col.trend</code>	color of line showing dispersion trend.
<code>col.raw</code>	color of points showing the unshrunk dispersions.
<code>...</code>	any other arguments are passed to <code>plot</code> .

Details

This function displays the quarter-root of the quasi-likelihood dispersions for all genes, before and after shrinkage towards a trend. If `glmfit` was constructed without an abundance trend, the function instead plots a horizontal line (of colour `col.trend`) at the common value towards which dispersions are shrunk. The quarter-root transformation is applied to improve visibility for dispersions around unity.

Value

A plot is created on the current graphics device.

Author(s)

Aaron Lun, based on code by Davis McCarthy and Gordon Smyth

Examples

```
nbdisp <- 1/rchisq(1000, df=10)
y <- DGEList(matrix(rnbinom(6000, size = 1/nbdisp, mu = 10),1000,6))
design <- model.matrix(~factor(c(1,1,1,2,2,2)))
y <- estimateDisp(y, design)

fit <- glmQLFit(y, design)
plotQLDisp(fit)

fit <- glmQLFit(y, design, abundance.trend=FALSE)
plotQLDisp(fit)
```

plotSmear	<i>Plots log-Fold Change versus log-Concentration (or, M versus A) for Count Data</i>
-----------	---

Description

Both of these functions plot the log-fold change (i.e. the log of the ratio of expression levels for each tag between two experimental groups) against the log-concentration (i.e. the overall average expression level for each tag across the two groups). To represent counts that were low (e.g. zero in 1 library and non-zero in the other) in one of the two conditions, a 'smear' of points at low A value is presented in plotSmear.

Usage

```
plotSmear(object, pair=NULL, de.tags=NULL, xlab="Average logCPM", ylab="logFC", pch=19,
          cex=0.2, smearWidth=0.5, panel.first=grid(), smooth.scatter=FALSE, lowess=FALSE, ...)
```

Arguments

object	DGEList, DGEEexact or DGELRT object containing data to produce an MA-plot.
pair	pair of experimental conditions to plot (if NULL, the first two conditions are used). Ignored if object is a DGELRT object.
de.tags	rownames for tags identified as being differentially expressed; use exactTest to identify DE genes
xlab	x-label of plot
ylab	y-label of plot
pch	scalar or vector giving the character(s) to be used in the plot; default value of 19 gives a round point.
cex	character expansion factor, numerical value giving the amount by which plotting text and symbols should be magnified relative to the default; default cex=0.2 to make the plotted points smaller
smearWidth	width of the smear
panel.first	an expression to be evaluated after the plot axes are set up but before any plotting takes place; the default grid() draws a background grid to aid interpretation of the plot
smooth.scatter	logical, whether to produce a 'smooth scatter' plot using the KernSmooth::smoothScatter function or just a regular scatter plot; default is FALSE, i.e. produce a regular scatter plot
lowess	logical, indicating whether or not to add a lowess curve to the MA-plot to give an indication of any trend in the log-fold change with log-concentration
...	further arguments passed on to plot

Details

plotSmear is a more sophisticated and superior way to produce an 'MA plot'. plotSmear resolves the problem of plotting tags that have a total count of zero for one of the groups by adding the 'smear' of points at low A value. The points to be smeared are identified as being equal to the minimum estimated concentration in one of the two groups. The smear is created by using random uniform numbers of width smearWidth to the left of the minimum A. plotSmear also allows easy highlighting of differentially expressed (DE) tags.

Value

A plot to the current device

Author(s)

Mark Robinson, Davis McCarthy

See Also

[maPlot](#)

Examples

```
y <- matrix(rnbinom(10000,mu=5,size=2),ncol=4)
d <- DGEList(counts=y, group=rep(1:2,each=2), lib.size=colSums(y))
rownames(d$counts) <- paste("tag",1:nrow(d$counts),sep=".")
d <- estimateCommonDisp(d)
plotSmear(d)

# find differential expression
de <- exactTest(d)

# highlighting the top 500 most DE tags
de.tags <- rownames(topTags(de, n=500)$table)
plotSmear(d, de.tags=de.tags)
```

plotSpliceDGE

Plot exons on differentially spliced gene

Description

Plot exons of differentially spliced gene.

Usage

```
plotSpliceDGE(lrt, geneid=NULL, rank=1L, FDR = 0.05)
```

Arguments

lrt	GLMLRT object produced by diffSpliceDGE.
geneid	character string, ID of the gene to plot.
rank	integer, if geneid=NULL then this ranked gene will be plotted.
FDR	numeric, mark exons with false discovery rate less than this cutoff.

Details

Plots interaction log-fold-change by exon for the specified gene.

Value

A plot is created on the current graphics device.

Author(s)

Yunshun Chen, Yifang Hu and Gordon Smyth

See Also

[diffSpliceDGE](#)

Examples

```
# See \link{diffSpliceDGE}
```

predFC	<i>Predictive log-fold changes</i>
--------	------------------------------------

Description

Computes estimated coefficients for a NB glm in such a way that the log-fold-changes are shrunk towards zero.

Usage

```
## S3 method for class DGEList
predFC(y, design=NULL, prior.count=0.125, offset=NULL, dispersion=NULL, weights=NULL, ...)
## Default S3 method:
predFC(y, design=NULL, prior.count=0.125, offset=NULL, dispersion=0, weights=NULL, ...)
```

Arguments

<code>y</code>	a matrix of counts or a <code>DGEList</code> object
<code>design</code>	the design matrix for the experiment
<code>prior.count</code>	the average prior count to be added to each observation. Larger values produce more shrinkage.
<code>offset</code>	numeric vector or matrix giving the offset in the log-linear model predictor, as for <code>glmFit</code> . Usually equal to log library sizes.
<code>dispersion</code>	numeric vector of negative binomial dispersions.
<code>weights</code>	optional numeric matrix giving observation weights
<code>...</code>	other arguments are passed to <code>glmFit</code> .

Details

This function computes predictive log-fold changes (pfc) for a NB glm. The pfc are posterior Bayesian estimators of the true log-fold-changes. They are predictive of values that might be replicated in a future experiment.

Specifically the function adds a small prior count to each observation before estimating the glm. The actual prior count that is added is proportion to the library size. This has the effect that any log-fold-change that was zero prior to augmentation remains zero and non-zero log-fold-changes are shrunk towards zero.

The prior counts can be viewed as equivalent to a prior belief that the log-fold changes are small, and the output can be viewed as posterior log-fold-changes from this Bayesian viewpoint. The output coefficients are called *predictive* log fold-changes because, depending on the prior, they may be a better prediction of the true log fold-changes than the raw estimates.

Log-fold changes for transcripts with low counts are shrunk more than transcript with high counts. In particular, infinite log-fold-changes arising from zero counts are avoided. The exact degree to which this is done depends on the negative binomial dispersion.

If `design=NULL`, then the function returns a matrix of the same size as `y` containing log₂ counts-per-million, with zero values for the counts avoided. This equivalent to choosing `design` to be the identity matrix with the same number of columns as `y`.

Value

Numeric matrix of linear model coefficients (if `design` is given) or logCPM (if `design=NULL`) on the log₂ scale.

Author(s)

Belinda Phipson and Gordon Smyth

References

Phipson, B. (2013). *Empirical Bayes modelling of expression profiles and their associations*. PhD Thesis. University of Melbourne, Australia. <http://repository.unimelb.edu.au/10187/17614>

See Also

[glmFit](#), [exactTest](#)

Examples

```
# generate counts for a two group experiment with n=2 in each group and 100 genes
dispersion <- 0.1
y <- matrix(rnbinom(400,size=1/dispersion,mu=4),nrow=100)
y <- DGEList(y,group=c(1,1,2,2))
design <- model.matrix(~group, data=y$samples)

#estimate the predictive log fold changes
predlfc<-predFC(y,design,dispersion=dispersion,prior.count=1)
logfc <- predFC(y,design,dispersion=dispersion,prior.count=0)
logfc.truncated <- pmax(pmin(logfc,100),-100)

#plot predFCs vs logFCs
plot(predlfc[,2],logfc.truncated[,2],xlab="Predictive log fold changes",ylab="Raw log fold changes")
abline(a=0,b=1)
```

processAmplicons

Process raw data from pooled genetic sequencing screens

Description

Given a list of sample-specific index (barcode) sequences and hairpin/sgRNA-specific sequences from an amplicon sequencing screen, generate a DGEList of counts from the raw fastq file/(s) containing the sequence reads.

Usage

```
processAmplicons(readfile, readfile2=NULL, barcodefile, hairpinfile,
                 barcodeStart=1, barcodeEnd=5,
                 barcodeStartRev=NULL, barcodeEndRev=NULL,
                 hairpinStart=37, hairpinEnd=57,
                 allowShifting=FALSE, shiftingBase=3,
                 allowMismatch=FALSE, barcodeMismatchBase=1,
                 hairpinMismatchBase=2, allowShiftedMismatch=FALSE,
                 verbose=FALSE)
```

Arguments

readfile	character vector giving one or more fastq filenames
readfile2	character vector giving one or more fastq filenames for reverse read, default to NULL
barcodefile	filename containing sample-specific barcode ids and sequences
hairpinfile	filename containing hairpin/sgRNA-specific ids and sequences

barcodeStart	numeric value, starting position (inclusive) of barcode sequence in reads
barcodeEnd	numeric value, ending position (inclusive) of barcode sequence in reads
barcodeStartRev	numeric value, starting position (inclusive) of barcode sequence in reverse reads, default to NULL
barcodeEndRev	numeric value, ending position (inclusive) of barcode sequence in reverse reads, default to NULL
hairpinStart	numeric value, starting position (inclusive) of hairpin/sgRNA sequence in reads
hairpinEnd	numeric value, ending position (inclusive) of hairpin/sgRNA sequence in reads
allowShifting	logical, indicates whether a given hairpin/sgRNA can be matched to a neighbouring position
shiftingBase	numeric value of maximum number of shifted bases from input hairpinStart and hairpinEnd should the program check for a hairpin/sgRNA match when allowShifting is TRUE
allowMismatch	logical, indicates whether sequence mismatch is allowed
barcodeMismatchBase	numeric value of maximum number of base sequence mismatches allowed in a barcode sequence when allowShifting is TRUE
hairpinMismatchBase	numeric value of maximum number of base sequence mismatches allowed in a hairpin/sgRNA sequence when allowShifting is TRUE
allowShiftedMismatch	logical, effective when allowShifting and allowMismatch are both TRUE. It indicates whether we check for sequence mismatches at a shifted position.
verbose	if TRUE, output program progress

Details

The input barcode file and hairpin/sgRNA files are tab-separated text files with at least two columns (named 'ID' and 'Sequences') containing the sample or hairpin/sgRNA ids and a second column indicating the sample index or hairpin/sgRNA sequences to be matched. If readfile2, barcodeStartRev and barcodeEndRev are specified, a third column 'SequencesReverse' is expected in the barcode file. The barcode file may also contain a 'group' column that indicates which experimental group a sample belongs to. Additional columns in each file will be included in the respective \$samples or \$genes data.frames of the final codeDGEList object. These files, along with the fastq file/(s) are assumed to be in the current working directory.

To compute the count matrix, matching to the given barcodes and hairpins/sgRNAs is conducted in two rounds. The first round looks for an exact sequence match for the given barcode sequences and hairpin/sgRNA sequences at the locations specified. If allowShifting is set to TRUE, the program also checks if a given hairpin/sgRNA sequence can be found at a neighbouring position in the read. For hairpins/sgRNAs without a match, the program performs a second round of matching which allows for sequence mismatches. The program checks parameter allowShifting to see if matches can be found at shifted positions in the read and allowShiftedMismatch accommodates mismatches at the shifted positions. The maximum number of mismatch bases in barcode and hairpin/sgRNA are specified by the parameters barcodeMismatchBase and hairpinMismatchBase.

The program outputs a `DGEList` object, with a count matrix indicating the number of times each barcode and hairpin/sgRNA combination could be matched in reads from input fastq file/(s).

For further examples and data, refer to the Case studies available from <http://bioinf.wehi.edu.au/shRNAseq/>.

Value

Returns a `DGEList` object with following components:

<code>counts</code>	read count matrix tallying up the number of reads with particular barcode and hairpin/sgRNA matches. Each row is a hairpin and each column is a sample
<code>genes</code>	In this case, hairpin/sgRNA-specific information (ID, sequences, corresponding target gene) may be recorded in this <code>data.frame</code>
<code>lib.size</code>	auto-calculated column sum of the counts matrix

Author(s)

Zhiyin Dai and Matthew Ritchie

References

Dai Z, Sheridan JM, et al. (2014). shRNA-seq data analysis with edgeR. F1000Research, <http://f1000research.com/articles/10>

q2qnbinom	<i>Quantile to Quantile Mapping between Negative-Binomial Distributions</i>
-----------	---

Description

Interpolated quantile to quantile mapping between negative-binomial distributions with the same dispersion but different means. The Poisson distribution is a special case.

Usage

```
q2qpois(x, input.mean, output.mean)
q2qnbinom(x, input.mean, output.mean, dispersion=0)
```

Arguments

<code>x</code>	numeric matrix of counts.
<code>input.mean</code>	numeric matrix of population means for <code>x</code> . If a vector, then of the same length as <code>nrow(x)</code> .
<code>output.mean</code>	numeric matrix of population means for the output values. If a vector, then of the same length as <code>nrow(x)</code> .
<code>dispersion</code>	numeric scalar, vector or matrix giving negative binomial dispersion values.

Details

This function finds the quantile with the same left and right tail probabilities relative to the output mean as x has relative to the input mean. `q2qpois` is equivalent to `q2qnbinom` with `dispersion=0`.

In principle, `q2qnbinom` gives similar results to calling `pnbinom` followed by `qnbinom` as in the example below. However this function avoids infinite values arising from rounding errors and does appropriate interpolation to return continuous values.

`q2qnbinom` is called by [equalizeLibSizes](#) to perform quantile-to-quantile normalization.

Value

numeric matrix of same dimensions as x , with output `.mean` as the new nominal population mean.

Author(s)

Gordon Smyth

See Also

[equalizeLibSizes](#)

Examples

```
x <- 15
input.mean <- 10
output.mean <- 20
dispersion <- 0.1
q2qnbinom(x, input.mean, output.mean, dispersion)

# Similar in principle:
qnbinom(pnbinom(x, mu=input.mean, size=1/dispersion), mu=output.mean, size=1/dispersion)
```

readDGE

Read and Merge a Set of Files Containing DGE Data

Description

Reads and merges a set of text files containing digital gene expression data.

Usage

```
readDGE(files, path=NULL, columns=c(1,2), group=NULL, labels=NULL, ...)
```

Arguments

files	character vector of filenames, or alternatively a data.frame with a column containing the file names of the files containing the libraries of counts and, optionally, columns containing the group to which each library belongs, descriptions of the other samples and other information.
path	character string giving the directory containing the files. The default is the current working directory.
columns	numeric vector stating which two columns contain the tag names and counts, respectively
group	vector, or preferably a factor, indicating the experimental group to which each library belongs. If group is not NULL, then this argument overrides any group information included in the files argument.
labels	character vector giving short names to associate with the libraries. Defaults to the file names.
...	other are passed to read.delim

Details

Each file is assumed to contained digital gene expression data for one sample (or library), with transcript identifiers in the first column and counts in the second column. Transcript identifiers are assumed to be unique and not repeated in any one file. By default, the files are assumed to be tab-delimited and to contain column headings. The function forms the union of all transcripts and creates one big table with zeros where necessary.

Value

DGEList object

Author(s)

Mark Robinson and Gordon Smyth

See Also

[DGEList](#) provides more information about the DGEList class and the function DGEList, which can also be used to construct a DGEList object, if readDGE is not required to read in and construct a table of counts from separate files.

Examples

```
# Read all .txt files from current working directory

## Not run: files <- dir(pattern="*\\.txt$")
RG <- readDGE(files)
## End(Not run)
```

`roast.DGEList`*Rotation Gene Set Tests for Digital Gene Expression Data*

Description

Rotation gene set testing for Negative Binomial generalized linear models.

Usage

```
## S3 method for class DGEList
roast(y, index=NULL, design=NULL, contrast=ncol(design), ...)
## S3 method for class DGEList
mroast(y, index=NULL, design=NULL, contrast=ncol(design), ...)
```

Arguments

<code>y</code>	DGEList object.
<code>index</code>	index vector specifying which rows (genes) of <code>y</code> are in the test set. This can be a vector of indices, or a logical vector of the same length as <code>statistics</code> , or any vector such as <code>y[,iset,]</code> contains the values for the gene set to be tested. Defaults to all genes. For <code>mroast</code> a list of index vectors.
<code>design</code>	design matrix
<code>contrast</code>	contrast for which the test is required. Can be an integer specifying a column of <code>design</code> , or else a contrast vector of length equal to the number of columns of <code>design</code> .
<code>...</code>	other arguments are passed to <code>link{roast.default}</code> or <code>link{mroast.default}</code> .

Details

The roast gene set test was proposed by Wu et al (2010) for microarray data. This function makes the roast test available for digital gene expression data. The negative binomial count data is converted to approximate normal deviates by computing mid-p quantile residuals (Dunn and Smyth, 1996; Routledge, 1994) under the null hypothesis that the contrast is zero. See [roast](#) for more description of the test and for a complete list of possible arguments.

The design matrix defaults to the `model.matrix(~y$samples$group)`.

`mroast` performs roast tests for a multiple of gene sets.

Value

`roast` produces an object of class [Roast](#). See [roast](#) for details.

`mroast` produces a `data.frame`. See [mroast](#) for details.

Author(s)

Yunshun Chen and Gordon Smyth

References

- Dunn, K. P., and Smyth, G. K. (1996). Randomized quantile residuals. *J. Comput. Graph. Statist.*, 5, 236-244. <http://www.statsci.org/smyth/pubs/residual.html>
- Routledge, RD (1994). Practicing safe statistics with the mid-p. *Canadian Journal of Statistics* 22, 103-110.
- Wu, D, Lim, E, Francois Vaillant, F, Asselin-Labat, M-L, Visvader, JE, and Smyth, GK (2010). ROAST: rotation gene set tests for complex microarray experiments. *Bioinformatics* 26, 2176-2182. <http://bioinformatics.oxfordjournals.org/content/26/17/2176>

See Also

[roast](#), [camera](#), [DGEList](#)

Examples

```
mu <- matrix(10, 100, 4)
group <- factor(c(0,0,1,1))
design <- model.matrix(~group)

# First set of 10 genes that are genuinely differentially expressed
iset1 <- 1:10
mu[iset1,3:4] <- mu[iset1,3:4]+10

# Second set of 10 genes are not DE
iset2 <- 11:20

# Generate counts and create a DGEList object
y <- matrix(rnbinom(100*4, mu=mu, size=10),100,4)
y <- DGEList(counts=y, group=group)

# Estimate dispersions
y <- estimateDisp(y, design)

roast(y, iset1, design, contrast=2)
mroast(y, iset1, design, contrast=2)
mroast(y, list(set1=iset1, set2=iset2), design, contrast=2)
```

spliceVariants

Identify Genes with Splice Variants

Description

Identify genes exhibiting evidence for splice variants (alternative exon usage/transcript isoforms) from exon-level count data using negative binomial generalized linear models.

Usage

```
spliceVariants(y, geneID, dispersion=NULL, group=NULL, estimate.genewise.disp=TRUE,
               trace=FALSE)
```

Arguments

<code>y</code>	either a matrix of exon-level counts or a <code>DGEList</code> object with (at least) <code>elements</code> counts (table of counts summarized at the exon level) and <code>samples</code> (data frame containing information about experimental group, library size and normalization factor for the library size). Each row of <code>y</code> should represent one exon.
<code>geneID</code>	vector of length equal to the number of rows of <code>y</code> , which provides the gene identifier for each exon in <code>y</code> . These identifiers are used to group the relevant exons into genes for the gene-level analysis of splice variation.
<code>dispersion</code>	scalar (in future a vector will also be allowed) supplying the negative binomial dispersion parameter to be used in the negative binomial generalized linear model.
<code>group</code>	factor supplying the experimental group/condition to which each sample (column of <code>y</code>) belongs. If <code>NULL</code> (default) the function will try to extract it from <code>y</code> , which only works if <code>y</code> is a <code>DGEList</code> object.
<code>estimate.genewise.disp</code>	logical, should genewise dispersions (as opposed to a common dispersion value) be computed if the dispersion argument is <code>NULL</code> ?
<code>trace</code>	logical, whether or not verbose comments should be printed as function is run. Default is <code>FALSE</code> .

Details

This function can be used to identify genes showing evidence of splice variation (i.e. alternative splicing, alternative exon usage, transcript isoforms). A negative binomial generalized linear model is used to assess evidence, for each gene, given the counts for the exons for each gene, by fitting a model with an interaction between exon and experimental group and comparing this model (using a likelihood ratio test) to a null model which does not contain the interaction. Genes that show significant evidence for an interaction between exon and experimental group by definition show evidence for splice variation, as this indicates that the observed differences between the exon counts between the different experimental groups cannot be explained by consistent differential expression of the gene across all exons. The function `topTags` can be used to display the results of `spliceVariants` with genes ranked by evidence for splice variation.

Value

`spliceVariants` returns a `DGEEExact` object, which contains a table of results for the test of differential splicing between experimental groups (alternative exon usage), a data frame containing the gene identifiers for which results were obtained and the dispersion estimate(s) used in the statistical models and testing.

Author(s)

Davis McCarthy, Gordon Smyth

See Also

[estimateExonGenewiseDisp](#) for more information about estimating genewise dispersion values from exon-level counts. [DGEList](#) for more information about the `DGEList` class. [topTags](#) for more

information on displaying ranked results from spliceVariants. [estimateCommonDisp](#) and related functions for estimating the dispersion parameter for the negative binomial model.

Examples

```
# generate exon counts from NB, create list object
y<-matrix(rnbinom(40,size=1,mu=10),nrow=10)
d<-DGEList(counts=y,group=rep(1:2,each=2))
genes <- rep(c("gene.1","gene.2"), each=5)
disp <- 0.2
spliceVariants(d, genes, disp)
```

splitIntoGroups	<i>Split the Counts or Pseudocounts from a DGEList Object According To Group</i>
-----------------	--

Description

Split the counts from a DGEList object according to group, creating a list where each element consists of a numeric matrix of counts for a particular experimental group. Given a pair of groups, split pseudocounts for these groups, creating a list where each element is a matrix of pseudocounts for a particular group.

Usage

```
splitIntoGroups(object)
splitIntoGroupsPseudo(pseudo, group, pair)
```

Arguments

object	DGEList, object containing (at least) the elements counts (table of raw counts), group (factor indicating group) and lib.size (numeric vector of library sizes)
pseudo	numeric matrix of quantile-adjusted pseudocounts to be split
group	factor indicating group to which libraries/samples (i.e. columns of pseudo belong; must be same length as ncol(pseudo))
pair	vector of length two stating pair of groups to be split for the pseudocounts

Value

splitIntoGroups outputs a list in which each element is a matrix of count counts for an individual group. splitIntoGroupsPseudo outputs a list with two elements, in which each element is a numeric matrix of (pseudo-)count data for one of the groups specified.

Author(s)

Davis McCarthy

Examples

```
# generate raw counts from NB, create list object
y<-matrix(rnbinom(80,size=1,mu=10),nrow=20)
d<-DGEList(counts=y,group=rep(1:2,each=2),lib.size=rep(c(1000:1001),2))
rownames(d$counts)<-paste("tagno",1:nrow(d$counts),sep=".")
z1<-splitIntoGroups(d)

z2<-splitIntoGroupsPseudo(d$counts,d$group,pair=c(1,2))
```

subsetting

*Subset DGEList, DGEGLM, DGEEExact and DGELRT Objects***Description**

Extract a subset of a DGEList, DGEGLM, DGEEExact or DGELRT object.

Usage

```
## S3 method for class DGEList
object[i, j, keep.lib.sizes=TRUE]
## S3 method for class DGEGLM
object[i, j]
## S3 method for class DGEEExact
object[i, j]
## S3 method for class DGELRT
object[i, j]
## S3 method for class TopTags
object[i, j]
```

Arguments

object object of class DGEList, DGEGLM, DGEEExact or DGELRT. For subsetListOfArrays, any list of conformal matrices and vectors.

i, j elements to extract. *i* subsets the tags or genes while *j* subsets the libraries. Note that columns of DGEGLM, DGEEExact and DGELRT objects cannot be subsetted.

keep.lib.sizes logical, if TRUE the lib.sizes will be kept unchanged on output, otherwise they will be recomputed as the column sums of the counts of the remaining rows.

Details

i, j may take any values acceptable for the matrix components of object of class DGEList. See the [Extract](#) help entry for more details on subsetting matrices. For DGEGLM, DGEEExact and DGELRT objects, only rows (i.e. *i*) may be subsetted.

Value

An object of the same class as object holding data from the specified subset of rows and columns.

Author(s)

Davis McCarthy, Gordon Smyth

See Also

[Extract](#) in the base package.

Examples

```
d <- matrix(rnbinom(16,size=1,mu=10),4,4)
rownames(d) <- c("a","b","c","d")
colnames(d) <- c("A1","A2","B1","B2")
d <- DGEList(counts=d,group=factor(c("A","A","B","B")))
d[1:2,]
d[1:2,2]
d[,2]
d <- estimateCommonDisp(d)
results <- exactTest(d)
results[1:2,]
# NB: cannot subset columns for DGEEExact objects
```

sumTechReps

Sum Over Replicate Samples

Description

Condense the columns of a matrix or DGEList object so that counts are summed over technical replicate samples.

Usage

```
## Default S3 method:
sumTechReps(x, ID=colnames(x), ...)
## S3 method for class DGEList
sumTechReps(x, ID=colnames(x), ...)
```

Arguments

x	a numeric matrix or DGEList object.
ID	sample identifier.
...	other arguments are not currently used.

Details

A new matrix or DGEList object is computed in which the counts for technical replicate samples are replaced by their sums.

Value

A data object of the same class as `x` with a column for each unique value of ID. Columns are in the same order as the ID values first occur in the ID vector.

Author(s)

Gordon Smyth and Yifang Hu

See Also

[rowsum](#).

Examples

```
x <- matrix(rpois(8*3,lambda=5),8,3)
colnames(x) <- c("a","a","b")
sumTechReps(x)
```

systematicSubset	<i>Take a systematic subset of indices.</i>
------------------	---

Description

Take a systematic subset of indices stratified by a ranking variable.

Usage

```
systematicSubset(n, order.by)
```

Arguments

<code>n</code>	integer giving the size of the subset.
<code>order.by</code>	numeric vector of the values by which the indices are ordered.

Value

`systematicSubset` returns a vector of size `n`.

Author(s)

Gordon Smyth

See Also

[order](#)

Examples

```
y <- rnorm(100, 1, 1)
systematicSubset(20, y)
```

`thinCounts`*Binomial or Multinomial Thinning of Counts*

Description

Reduce the size of Poisson-like counts by binomial thinning.

Usage

```
thinCounts(x, prob=NULL, target.size=min(colSums(x)))
```

Arguments

<code>x</code>	numeric vector or array of non-negative integers.
<code>prob</code>	numeric scalar or vector of same length as <code>x</code> , the expected proportion of the events to keep.
<code>target.size</code>	integer scale or vector of the same length as <code>NCOL{x}</code> , the desired total column counts. Must be not greater than column sum of <code>x</code> . Ignored if <code>prob</code> is not <code>NULL</code> .

Details

If `prob` is not `NULL`, then this function calls `rbinom` with `size=x` and `prob=prob` to generate the new counts. This is classic binomial thinning. The new column sums are random, with expected values determined by `prob`.

If `prob` is `NULL`, then this function does multinomial thinning of the counts to achieve specified column totals. The default behavior is to thin the columns to have the same column sum, equal to the smallest column sum of `x`.

If the elements of `x` are Poisson, then binomial thinning produces new Poisson random variables with expected values reduced by factor `prob`. If the elements of each column of `x` are multinomial, then multinomial thinning produces a new multinomial observation with a reduced sum.

Value

A vector or array of the same dimensions as `x`, with thinned counts.

Author(s)

Gordon Smyth

Examples

```
x <- rpois(10,lambda=10)
thinCounts(x,prob=0.5)
```

topSpliceDGE	<i>Top table of differentially spliced genes or exons</i>
--------------	---

Description

Top table ranking the most differentially spliced genes or exons.

Usage

```
topSpliceDGE(lrt, level="gene", gene.test="Simes", number=10, FDR=1)
```

Arguments

lrt	DGELRT object produced by diffSpliceDGE.
level	character string, should the table be by "exon" or by "gene".
gene.test	character string, choice for the gene-level p-values. Possible values are "Simes" and "F".
number	integer, maximum number of rows to output.
FDR	numeric, only show exons or genes with false discovery rate less than this cutoff.

Details

Ranks exons or genes by p-values.

Value

A data.frame with any annotation columns found in `fit` plus the following columns

NExons	number of exons if level="gene"
Gene.Exon	exon annotation if level="exon"
logFC	log-fold change of one exon vs all the exons for the same gene (if level="exon")
F	F-statistics for exons if level="exon"
P.Value	p-value
FDR	false discovery rate

Author(s)

Yunshun Chen and Gordon Smyth

Examples

```
# See \link{diffSpliceDGE}
```

topTags	<i>Table of the Top Differentially Expressed Tags</i>
---------	---

Description

Extracts the top DE tags in a data frame for a given pair of groups, ranked by p-value or absolute log-fold change.

Usage

```
topTags(object, n=10, adjust.method="BH", sort.by="PValue")
```

Arguments

object	a DGEEexact object (output from <code>exactTest</code>) or a DGELRT object (output from <code>glmLRT</code>), containing the (at least) the elements <code>table</code> : a data frame containing the log-concentration (i.e. expression level), the log-fold change in expression between the two groups/conditions and the p-value for differential expression, for each tag. If it is a <code>DGEEexact</code> object, then <code>topTags</code> will also use the <code>comparison</code> element, which is a vector giving the two experimental groups/conditions being compared. The object may contain other elements that are not used by <code>topTags</code> .
n	scalar, number of tags to display/return
adjust.method	character string stating the method used to adjust p-values for multiple testing, passed on to <code>p.adjust</code>
sort.by	character string, should the top tags be sorted by p-value (" <code>PValue</code> "), by absolute log-fold change (" <code>logFC</code> "), or not sorted (" <code>none</code> ").

Value

an object of class `TopTags` containing the following elements for the top `n` most differentially expressed tags as determined by `sort.by`:

table	a data frame containing the elements <code>logFC</code> , the log-abundance ratio, i.e. fold change, for each tag in the two groups being compared, <code>logCPM</code> , the log-average concentration/abundance for each tag in the two groups being compared, <code>PValue</code> , exact p-value for differential expression using the NB model, <code>FDR</code> , the p-value adjusted for multiple testing as found using <code>p.adjust</code> using the method specified.
adjust.method	character string stating the method used to adjust p-values for multiple testing.
comparison	a vector giving the names of the two groups being compared.
test	character string stating the name of the test.

The dimensions, row names and column names of a `TopTags` object are defined by those of `table`, see `dim.TopTags` or `dimnames.TopTags`.

`TopTags` objects also have a `show` method so that printing produces a compact summary of their contents.

Author(s)

Mark Robinson, Davis McCarthy, Gordon Smyth

References

Robinson MD, Smyth GK (2008). Small-sample estimation of negative binomial dispersion, with applications to SAGE data. *Biostatistics* 9, 321-332.

Robinson MD, Smyth GK (2007). Moderated statistical tests for assessing differences in tag abundance. *Bioinformatics* 23, 2881-2887.

See Also

[exactTest](#), [glmLRT](#), [p.adjust](#).

Analogous to [topTable](#) in the limma package.

Examples

```
# generate raw counts from NB, create list object
y <- matrix(rnbinom(80,size=1,mu=10),nrow=20)
d <- DGEList(counts=y,group=rep(1:2,each=2),lib.size=rep(c(1000:1001),2))
rownames(d$counts) <- paste("tag",1:nrow(d$counts),sep=".")

# estimate common dispersion and find differences in expression
# here we demonstrate the exact methods, but the use of topTags is
# the same for a GLM analysis
d <- estimateCommonDisp(d)
de <- exactTest(d)

# look at top 10
topTags(de)
# Can specify how many tags to view
tp <- topTags(de, n=15)
# Here we view top 15
tp
# Or order by fold change instead
topTags(de,sort.by="logFC")
```

treatDGE

Testing for Differential Expression Relative to a Threshold

Description

Conduct genewise statistical tests for a given coefficient or coefficient contrast relative to a specified threshold.

Usage

```
treatDGE(glmfit, coef=ncol(glmfit$design), contrast=NULL, lfc=0)
```

Arguments

<code>glmfit</code>	a DGEGLM object, usually output from <code>glmFit</code> .
<code>coef</code>	integer or character vector indicating which coefficients of the linear model are to be tested equal to zero. Values must be columns or column names of <code>design</code> . Defaults to the last coefficient. Ignored if <code>contrast</code> is specified.
<code>contrast</code>	numeric vector specifying the contrast of the linear model coefficients to be tested against the log ₂ -fold change threshold. Length must equal to the number of columns of <code>design</code> . If specified, then takes precedence over <code>coef</code> .
<code>lfc</code>	numeric scalar specifying the absolute value of the log ₂ -fold change threshold above which differential expression is to be considered.

Details

`treatDGE` implements a two-sided modified likelihood ratio test.

Value

`treatDGE` produces an object of class `DGELRT` with the same components as for `glmfit` plus the following:

<code>lfc</code>	absolute value of the specified log ₂ -fold change threshold.
<code>table</code>	data frame with the same rows as <code>glmfit</code> containing the log ₂ -fold changes, average log ₂ -counts per million and p-values, ready to be displayed by <code>topTags</code> ..
<code>comparison</code>	character string describing the coefficient or the contrast being tested.

The data frame `table` contains the following columns:

<code>logFC</code>	log ₂ -fold change of expression between conditions being tested.
<code>logCPM</code>	average log ₂ -counts per million, the average taken over all libraries.
<code>PValue</code>	p-values.

Author(s)

Yunshun Chen and Gordon Smyth

Examples

```
ngenes <- 100
n1 <- 3
n2 <- 3
nlibs <- n1+n2
mu <- 100
phi <- 0.1
group <- c(rep(1,n1), rep(2,n2))
design <- model.matrix(~as.factor(group))

### 4-fold change for the first 5 genes
i <- 1:5
fc <- 4
```

```
mu <- matrix(mu, ngenes, nlibs)
mu[i, 1:n1] <- mu[i, 1:n1]*fc

counts <- matrix(rnbinom(ngenes*nlibs, mu=mu, size=1/phi), ngenes, nlibs)
d <- DGEList(counts=counts, lib.size=rep(1e6, nlibs), group=group)

gfit <- glmFit(d, design, dispersion=phi)
tr <- treatDGE(gfit, coef=2, lfc=1)
topTags(tr)
```

validDGEList*Check for Valid DGEList object*

Description

Check for existence of standard components of DGEList object.

Usage

```
validDGEList(y)
```

Arguments

y DGEList object.

Details

This function checks that the standard counts and samples components of a DGEList object are present.

Value

DGEList with missing components added.

Author(s)

Gordon Smyth

See Also

[DGEList](#)

Examples

```
counts <- matrix(rpois(4*2, lambda=5), 4, 2)
dge <- new("DGEList", list(counts=counts))
validDGEList(dge)
```

 weightedCondLogLikDerDelta

Weighted Conditional Log-Likelihood in Terms of Delta

Description

Weighted conditional log-likelihood parameterized in terms of delta ($\phi / (\phi+1)$) for a given tag/gene - maximized to find the smoothed (moderated) estimate of the dispersion parameter

Usage

```
weightedCondLogLikDerDelta(y, delta, tag, prior.n=10, ntags=nrow(y[[1]]), der=0)
```

Arguments

y	list with elements comprising the matrices of count data (or pseudocounts) for the different groups
delta	delta ($\phi / (\phi+1)$) parameter of negative binomial
tag	tag/gene at which the weighted conditional log-likelihood is evaluated
prior.n	smoothing parameter that indicates the weight to put on the common likelihood compared to the individual tag's likelihood; default 10 means that the common likelihood is given 10 times the weight of the individual tag/gene's likelihood in the estimation of the tag/genewise dispersion
ntags	numeric scalar number of tags/genes in the dataset to be analysed
der	derivative, either 0 (the function), 1 (first derivative) or 2 (second derivative)

Details

This function computes the weighted conditional log-likelihood for a given tag, parameterized in terms of delta. The value of delta that maximizes the weighted conditional log-likelihood is converted back to the phi scale, and this value is the estimate of the smoothed (moderated) dispersion parameter for that particular tag. The delta scale for convenience (delta is bounded between 0 and 1).

Value

numeric scalar of function/derivative evaluated for the given tag/gene and delta

Author(s)

Mark Robinson, Davis McCarthy

Examples

```
counts<-matrix(rnbinom(20,size=1,mu=10),nrow=5)
d<-DGEList(counts=counts,group=rep(1:2,each=2),lib.size=rep(c(1000:1001),2))
y<-splitIntoGroups(d)
l11<-weightedCondLogLikDerDelta(y,delta=0.5,tag=1,prior.n=10,der=0)
l12<-weightedCondLogLikDerDelta(y,delta=0.5,tag=1,prior.n=10,der=1)
```

WLEB

*Calculate Weighted Likelihood Empirical Bayes Estimates***Description**

Estimates the parameters which maximize the given log-likelihood matrix using empirical Bayes method.

Usage

```
WLEB(theta, loglik, prior.n, covariate, trend.method="locfit", span=NULL, overall=TRUE,
      trend=TRUE, individual=TRUE, m0=NULL, m0.out=FALSE)
```

Arguments

theta	numeric vector of values of the parameter at which the log-likelihoods are calculated.
loglik	numeric matrix of log-likelihood of all the candidates at those values of parameter.
prior.n	numeric scaler, estimate of the prior weight, i.e. the smoothing parameter that indicates the weight to put on the common likelihood compared to the individual's likelihood.
covariate	numeric vector of values across which a parameter trend is fitted
trend.method	method for estimating the parameter trend. Possible values are "none", "movingave" and "loess".
span	width of the smoothing window, as a proportion of the data set.
overall	logical, should a single value of the parameter which maximizes the sum of all the log-likelihoods be estimated?
trend	logical, should a parameter trend (against the covariate) which maximizes the local shared log-likelihoods be estimated?
individual	logical, should individual estimates of all the candidates after applying empirical Bayes method along the trend be estimated?
m0	numeric matrix of local shared log-likelihoods. If Null, they will be calculated using the method selected by trend.method.
m0.out	logical, should local shared log-likelihoods be included in the output?

Details

This function is a generic function that calculates an overall estimate, trend estimates and individual estimates for each candidate given the values of the log-likelihood of all the candidates at some specified parameter values.

Value

A list with the following:

overall	the parameter estimate that maximizes the sum of all the log-likelihoods.
trend	the estimated trended parameters against the covariate.
individual	the individual estimates of all the candidates after applying empirical Bayes method along the trend.
shared.loglik	the estimated numeric matrix of local shared log-likelihoods

Author(s)

Yunshun Chen, Gordon Smyth

See Also

[locfitByCol](#), [movingAverageByCol](#) and [loessByCol](#) implement the local fit, moving average or loess smoothers.

Examples

```
y <- matrix(rpois(100, lambda=10), ncol=4)
theta <- 7:14
loglik <- matrix(0, nrow=nrow(y), ncol=length(theta))
for(i in 1:nrow(y))
  for(j in 1:length(theta))
    loglik[i,j] <- sum(dpois(y[i,], theta[j], log=TRUE))
covariate <- log(rowSums(y))
out <- WLEB(theta, loglik, prior.n=3, covariate)
out
```

zscoreNBinom

Z-score Equivalents of Negative Binomial Deviate

Description

Compute z-score equivalents of negative binomial random deviates.

Usage

```
zscoreNBinom(q, size, mu)
```

Arguments

q	numeric vector or matrix giving negative binomial random values.
size	negative binomial size parameter (>0).
mu	mean of negative binomial distribution (>0).

Details

This function computes the mid-p value of q, then converts to the standard normal deviate with the same cumulative probability distribution value.

Care is taken to do the computations accurately in both tails of the distributions.

Value

Numeric vector or matrix giving equivalent deviates from a standard normal distribution.

Author(s)

Gordon Smyth

See Also

[pnbinom](#), [qnorm](#) in the stats package.

Examples

```
zscoreNBinom(c(0,10,100), mu=10, size=1/10)
```

Index

- *Topic **algebra**
 - dglmStdResid, 24
 - dispCoxReidInterpolateTagwise, 35
 - estimateTagwiseDisp, 52
 - exactTest, 55
 - gof, 68
 - meanvar, 76
 - splitIntoGroups, 102
 - topTags, 108
 - WLEB, 113
- *Topic **array**
 - as.data.frame, 6
 - as.matrix, 7
 - dim, 29
- *Topic **category**
 - cutWithMinN, 17
- *Topic **classes**
 - DGEEexact-class, 19
 - DGEGLM-class, 20
 - DGEList-class, 22
 - DGELRT-class, 23
- *Topic **distribution**
 - zscoreNBinom, 114
- *Topic **documentation**
 - edgeRUsersGuide, 38
- *Topic **file**
 - commonCondLogLikDerDelta, 14
 - getPriorN, 59
 - readDGE, 97
 - weightedCondLogLikDerDelta, 112
- *Topic **gene set test**
 - goana.DGELRT, 66
- *Topic **hplot**
 - expandAsMatrix, 58
 - plotExonUsage, 85
 - plotMDS.DGEList, 86
- *Topic **htest**
 - binomTest, 9
 - decideTestsDGE, 18
 - spliceVariants, 100
- *Topic **interpolation**
 - maximizeInterpolant, 74
 - maximizeQuadratic, 75
- *Topic **models**
 - dispCoxReidSplineTrend, 37
 - estimateExonGenewiseDisp, 44
 - estimateGLMCommonDisp, 45
 - glmFit, 61
 - glmQLFit, 64
 - goodTuring, 70
 - thinCounts, 106
- *Topic **package**
 - edgeR-package, 3
- *Topic **plot**
 - plotBCV, 84
 - plotQLDisp, 88
- *Topic **smooth**
 - movingAverageByCol, 81
- *Topic **subset**
 - systematicSubset, 105
 - [.DGEEexact (subsetting), 103
 - [.DGEGLM (subsetting), 103
 - [.DGELRT (subsetting), 103
 - [.DGEList (subsetting), 103
 - [.TopTags (subsetting), 103
 - 02.Classes, 30
 - adjustedProfileLik, 4
 - as.data.frame, 6, 6
 - as.dist, 88
 - as.matrix, 7, 7
 - as.matrix.DGEList, 60
 - as.matrix.RGList, 7
 - aveLogCPM, 7, 17
 - binMeanVar (meanvar), 76
 - binom.test, 10
 - binomTest, 9, 57

- calcNormFactors, 10
- calcNormOffsetsforChIP
 - (normalizeChIPtoInput), 83
- camera, 13
- camera.default, 12
- camera.DGEList, 12, 100
- cmdscale, 88
- commonCondLogLikDerDelta, 14
- condLogLikDerDelta (condLogLikDerSize), 15
- condLogLikDerSize, 15
- cpm, 8, 16
- cut, 18
- cutWithMinN, 17, 37

- decideTests, 19
- decideTestsDGE, 18
- designAsFactor (mglim), 78
- DGEEExact, 108
- DGEEExact-class, 19
- DGEGLM-class, 20
- DGEList, 21, 22, 23, 39, 55, 60, 95, 96, 98, 101, 111
- DGEList-class, 22
- DGELRT, 108
- DGELRT-class, 23
- dglmStdResid, 24
- diffSpliceDGE, 27, 92
- dim, 29, 30
- dim.DGEEExact, 20
- dim.DGEGLM, 21
- dim.DGEList, 23
- dim.DGELRT, 24
- dim.TopTags, 108
- dimnames, 30, 30, 31
- dimnames.DGEEExact, 20
- dimnames.DGEGLM, 21
- dimnames.DGEList, 23
- dimnames.DGELRT, 24
- dimnames.TopTags, 108
- dimnames<- .DGEGLM (dimnames), 30
- dimnames<- .DGEList (dimnames), 30
- dispBinTrend, 31, 51
- dispCoxReid, 33, 46
- dispCoxReidInterpolateTagwise, 35, 49
- dispCoxReidPowerTrend, 51
- dispCoxReidPowerTrend
 - (dispCoxReidSplineTrend), 37
- dispCoxReidSplineTrend, 37, 51

- dispDeviance, 46
- dispDeviance (dispCoxReid), 33
- dispPearson, 46
- dispPearson (dispCoxReid), 33

- edgeR (edgeR-package), 3
- edgeR-package, 3
- edgeRUsersGuide, 38
- equalizeLibSizes, 39, 42, 56, 57, 97
- estimateCommonDisp, 14, 41, 44–46, 50, 53, 54, 102
- estimateDisp, 42
- estimateExonGenewiseDisp, 44, 101
- estimateGLMCommonDisp, 34, 44, 45, 50
- estimateGLMRobustDisp, 47
- estimateGLMTagwiseDisp, 36, 44, 46, 47, 48, 48, 60
- estimateGLMTrendedDisp, 32, 38, 44, 46–48, 50, 50
- estimateTagwiseDisp, 44, 46, 50, 52, 60
- estimateTrendedDisp, 54
- exactTest, 55, 94, 109
- exactTestBetaApprox (exactTest), 55
- exactTestByDeviance (exactTest), 55
- exactTestBySmallP (exactTest), 55
- exactTestDoubleTail (exactTest), 55
- expandAsMatrix, 58
- Extract, 103, 104

- getCounts, 58
- getDispersion (getCounts), 58
- getDispersions (dglmStdResid), 24
- getOffset (getCounts), 58
- getPriorN, 59
- glmFit, 4, 5, 33, 45, 47, 49–51, 61, 64, 65, 69, 80, 93, 94
- glmLRT, 109
- glmLRT (glmFit), 61
- glmQLFit, 64, 89
- glmQLFTest (glmQLFit), 64
- goana, 68
- goana.DGEEExact (goana.DGELRT), 66
- goana.DGELRT, 66
- gof, 68
- goodTuring, 70
- goodTuringPlot (goodTuring), 70
- goodTuringProportions (goodTuring), 70

- length.DGEEExact (dim), 29

- length.DGEGLM (dim), 29
- length.DGEList (dim), 29
- length.DGELRT (dim), 29
- length.TopTags (dim), 29
- locfitByCol, 114
- locfitByCol (loessByCol), 71
- loess, 72
- loessByCol, 53, 71, 114

- maPlot, 26, 72, 78, 91
- maximizeInterpolant, 36, 74, 75
- maximizeQuadratic, 75
- MDS, 88
- meanvar, 76
- mglm, 78
- mglmLevenberg, 62, 63
- mglmLevenberg (mglm), 78
- mglmOneGroup, 8, 62, 63
- mglmOneGroup (mglm), 78
- mglmOneWay (mglm), 78
- movingAverageByCol, 53, 81, 114
- mroast, 99
- mroast.DGEList (roast.DGEList), 99

- nbinomDeviance, 82
- nbinomUnitDeviance (nbinomDeviance), 82
- normalizeChIPtoInput, 83

- optim, 37
- optimize, 34, 41, 43
- order, 105

- p.adjust, 19, 109
- plotBCV, 84
- plotExonUsage, 85
- plotMDS, 88
- plotMDS.DGEList, 26, 78, 86
- plotMeanVar, 26
- plotMeanVar (meanvar), 76
- plotQLDisp, 66, 88
- plotSmear, 26, 74, 78, 90
- plotSpliceDGE, 91
- pnbinom, 115
- points, 84, 87, 89
- predFC, 92
- processAmplicons, 94

- q2qnbinom, 40, 96
- q2qpois (q2qnbinom), 96

- qnorm, 115
- qqnorm, 69
- quantile, 18

- readDGE, 97
- Roast, 99
- roast, 99, 100
- roast.DGEList, 13, 99
- rowsum, 105
- rpkm (cpm), 16

- sage.test, 10
- show, DGEEExact-method (DGEEExact-class), 19
- show, DGEGLM-method (DGEGLM-class), 20
- show, DGELRT-method (DGELRT-class), 23
- show, TopTags-method (topTags), 108
- spliceVariants, 86, 100
- splitIntoGroups, 102
- splitIntoGroupsPseudo (splitIntoGroups), 102
- squeezeVar, 65
- subsetting, 20, 21, 23, 24, 103
- sumTechReps, 104
- Sweave, 38
- system, 39
- systematicSubset, 46, 105

- TestResults, 19
- text, 87
- thinCounts, 106
- topGO, 68
- topSpliceDGE, 107
- topTable, 109
- topTags, 63, 66, 101, 108
- TopTags-class (topTags), 108
- treatDGE, 109

- uniroot, 34

- validDGEList, 111

- weightedCondLogLikDerDelta, 112
- WLEB, 113

- zscoreNBinom, 114