

Package ‘variancePartition’

November 28, 2021

Type Package

Title Quantify and interpret divers of variation in multilevel gene expression experiments

Version 1.24.0

Date 2021-10-27

Maintainer Gabriel E. Hoffman <gabriel.hoffman@mssm.edu>

Description Quantify and interpret multiple sources of biological and technical variation in gene expression experiments. Uses a linear mixed model to quantify variation in gene expression attributable to individual, tissue, time point, or technical variables. Includes dream differential expression analysis for repeated measures.

VignetteBuilder knitr

License GPL (>= 2)

URL <http://bioconductor.org/packages/variancePartition>,
<https://github.com/GabrielHoffman/variancePartition>

BugReports <https://github.com/GabrielHoffman/variancePartition/issues>

Suggests BiocStyle, knitr, pander, rmarkdown, edgeR, dendextend, tximport, tximportData, ballgown, DESeq2, RUnit, BiocGenerics, r2glmm, readr

biocViews RNASeq, GeneExpression, GeneSetEnrichment, DifferentialExpression, BatchEffect, QualityControl, Regression, Epigenetics, FunctionalGenomics, Transcriptomics, Normalization, Preprocessing, Microarray, ImmunoOncology, Software

Depends R (>= 4.0.0), ggplot2, limma, BiocParallel

Imports MASS, pbkrtest (>= 0.4-4), lmerTest, Matrix, iterators, foreach, doParallel, gplots, RhpcBLASctl, progress, reshape2, aod, scales, Rdpack, rlang, lme4 (>= 1.1-10), grDevices, graphics, Biobase, methods, utils, stats

RoxygenNote 7.1.2

RdMacros Rdpack

git_url <https://git.bioconductor.org/packages/variancePartition>

git_branch RELEASE_3_14

git_last_commit 509f650

git_last_commit_date 2021-10-27

Date/Publication 2021-11-28

Author Gabriel Hoffman [aut, cre]

R topics documented:

as.data.frame.varPartResults	3
as.matrix,varPartResults-method	4
calcVarPart	5
canCorPairs	7
classifyTestsF	8
classifyTestsF,MArrayLM2-method	9
colinearityScore	10
dream	11
dscchisq	14
eBayes,MArrayLM2-method	14
ESS	15
extractVarPart	17
fitExtractVarPartModel	18
fitVarPartModel	22
getContrast	26
get_prediction	27
ggColorHue	28
isRunnableFormula	29
makeContrastsDream	29
MArrayLM2-class	31
plotCompareP	31
plotContrasts	32
plotCorrMatrix	33
plotCorrStructure	34
plotPercentBars	35
plotStratify	37
plotStratifyBy	38
plotVarianceEstimates	40
plotVarPart	40
rdf.merMod	42
rdf_from_matrices	43
reOnly	44
residuals,MArrayLM-method	44
residuals,MArrayLM2-method	45
residuals,VarParFitList-method	45
residuals.MArrayLM2	46
shrinkageMetric	47

sortCols	47
VarParCIList-class	49
VarParFitList-class	49
varParFrac-class	49
varPartConInf	50
varPartData	52
varPartDEdata	53
varPartResults-class	53
voomWithDreamWeights	54

Index **56**

as.data.frame.varPartResults
Convert to data.frame

Description

Convert varPartResults to data.frame

Usage

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

Arguments

x	varPartResults
row.names	pass thru to generic
optional	pass thru to generic
...	other arguments.

Value

data.frame

Examples

```
# load library
# library(variancePartition)

# load simulated data:
# geneExpr: matrix of gene expression values
# info: information/metadata about each sample
data(varPartData)

# Specify variables to consider
# Age is continuous so we model it as a fixed effect
```

```
# Individual and Tissue are both categorical, so we model them as random effects
form <- ~ Age + (1|Individual) + (1|Tissue)

# Fit model
varPart <- fitExtractVarPartModel( geneExpr[1:5,], form, info )

# convert to matrix
as.data.frame(varPart)
```

as.matrix,varPartResults-method

Convert to matrix

Description

Convert varPartResults to matrix

Usage

```
## S4 method for signature 'varPartResults'
as.matrix(x, ...)
```

Arguments

x	varPartResults
...	other arguments.

Value

matrix

Examples

```
# load library
# library(variancePartition)

# load simulated data:
# geneExpr: matrix of gene expression values
# info: information/metadata about each sample
data(varPartData)

# Specify variables to consider
# Age is continuous so we model it as a fixed effect
# Individual and Tissue are both categorical, so we model them as random effects
form <- ~ Age + (1|Individual) + (1|Tissue)

# Fit model
varPart <- fitExtractVarPartModel( geneExpr[1:5,], form, info )
```

```
# convert to matrix
as.matrix(varPart)
```

calcVarPart *Compute variance statistics*

Description

Compute fraction of variation attributable to each variable in regression model. Also interpretable as the intra-class correlation after correcting for all other variables in the model.

Usage

```
calcVarPart(fit, showWarnings = TRUE, returnFractions = TRUE, ...)

## S4 method for signature 'lm'
calcVarPart(fit, showWarnings = TRUE, returnFractions = TRUE, ...)

## S4 method for signature 'lmerMod'
calcVarPart(fit, showWarnings = TRUE, returnFractions = TRUE, ...)

## S4 method for signature 'glm'
calcVarPart(fit, showWarnings = TRUE, returnFractions = TRUE, ...)

## S4 method for signature 'negbin'
calcVarPart(fit, showWarnings = TRUE, returnFractions = TRUE, ...)

## S4 method for signature 'glmerMod'
calcVarPart(fit, showWarnings = TRUE, returnFractions = TRUE, ...)
```

Arguments

fit	model fit from lm() or lmer()
showWarnings	show warnings about model fit (default TRUE)
returnFractions	default: TRUE. If TRUE return fractions that sum to 1. Else return unscaled variance components.
...	additional arguments (not currently used)

Details

For linear model, variance fractions are computed based on the sum of squares explained by each component. For the linear mixed model, the variance fractions are computed by variance component estimates for random effects and sum of squares for fixed effects.

For a generalized linear model, the variance fraction also includes the contribution of the link function so that fractions are reported on the linear (i.e. link) scale rather than the observed (i.e. response) scale. For linear regression with an identity link, fractions are the same on both scales. But for logit or probit links, the fractions are not well defined on the observed scale due to the transformation imposed by the link function.

The variance implied by the link function is the variance of the corresponding distribution:

logit -> logistic distribution -> variance is $\pi^2/3$

probit -> standard normal distribution -> variance is 1

For the Poisson distribution with rate λ , the variance is $\log(1 + 1/\lambda)$.

For the negative binomial distribution with rate λ and shape θ , the variance is $\log(1 + 1/\lambda + 1/\theta)$.

Variance decomposition is reviewed by Nakagawa and Schielzeth (2012), and expanded to other GLMs by Nakagawa, Johnson and Schielzeth (2017). See McKelvey and Zavoina (1975) for early work on applying to GLMs. Also see DeMaris (2002)

We note that Nagelkerke's pseudo R^2 evaluates the variance explained by the full model. Instead, a variance partitioning approach evaluates the variance explained by each term in the model, so that the sum of each systematic plus random term sums to 1 (Hoffman and Schadt, 2016; Nakagawa and Schielzeth, 2012).

Value

fraction of variance explained / ICC for each variable in the regression model

References

Nakagawa S, Johnson PC, Schielzeth H (2017). "The coefficient of determination R^2 and intra-class correlation coefficient from generalized linear mixed-effects models revisited and expanded." *Journal of the Royal Society Interface*, **14**(134), 20170213.

Nakagawa S, Schielzeth H (2013). "A general and simple method for obtaining R^2 from generalized linear mixed-effects models." *Methods in ecology and evolution*, **4**(2), 133–142.

McKelvey RD, Zavoina W (1975). "A statistical model for the analysis of ordinal level dependent variables." *Journal of mathematical sociology*, **4**(1), 103–120.

DeMaris A (2002). "Explained variance in logistic regression: A Monte Carlo study of proposed measures." *Sociological Methods & Research*, **31**(1), 27–74.

Hoffman GE, Schadt EE (2016). "variancePartition: interpreting drivers of variation in complex gene expression studies." *BMC bioinformatics*, **17**(1), 1–13.

Examples

```
library(lme4)
data(varPartData)

# Linear mixed model
fit <- lmer( geneExpr[1,] ~ (1|Tissue) + Age, info)
calcVarPart( fit )

# Linear model
```

```
# Note that the two models produce slightly different results
# This is expected: they are different statistical estimates
# of the same underlying value
fit <- lm( geneExpr[1,] ~ Tissue + Age, info)
calcVarPart( fit )
```

canCorPairs	<i>canCorPairs</i>
-------------	--------------------

Description

Assess correlation between all pairs of variables in a formula

Usage

```
canCorPairs(formula, data, showWarnings = TRUE)
```

Arguments

formula	standard additive linear model formula (doesn't support random effects currently, so just change the syntax)
data	data.frame with the data for the variables in the formula
showWarnings	default to true

Details

Canonical Correlation Analysis (CCA) is similar to correlation between two vectors, except that CCA can accommodate matrices as well. For a pair of variables, canCorPairs assesses the degree to which they co-vary and contain the same information. Variables in the formula can be a continuous variable or a discrete variable expanded to a matrix (which is done in the backend of a regression model). For a pair of variables, canCorPairs uses CCA to compute the correlation between these variables and returns the pairwise correlation matrix.

Statistically, let rho be the array of correlation values returned by the standard R function `canCor` to compute CCA. `canCorPairs` returns $\sqrt{\text{mean}(\text{rho}^2)}$, which is the fraction of the maximum possible correlation. When comparing a two vectors, or a vector and a matrix, this gives the same value as the absolute correlation. When comparing two sets of categorical variables (i.e. expanded to two matrices), this is equivalent to Cramer's V statistic.

Note that CCA returns correlation values between 0 and 1.

Value

Matrix of correlation values between all pairs of variables.

Examples

```
# load library
# library(variancePartition)

# load simulated data:
data(varPartData)

# specify formula
form <- ~ Individual + Tissue + Batch + Age + Height

# Compute Canonical Correlation Analysis (CCA)
# between all pairs of variables
# returns absolute correlation value
C = canCorPairs( form, info)

# Plot correlation matrix
plotCorrMatrix( C )
```

`classifyTestsF`*Multiple Testing Genewise Across Contrasts*

Description

For each gene, classify a series of related t-statistics as up, down or not significant.

Usage

```
classifyTestsF(object, ...)
```

Arguments

<code>object</code>	numeric matrix of t-statistics or an 'MArrayLM2' object from which the t-statistics may be extracted.
<code>...</code>	additional arguments

Details

Works like `limma::classifyTestsF`, except `object` can have a list of covariance matrices `object$cov.coefficients.list`, instead of just one in `object$cov.coefficients`

See Also

`limma::classifyTestsF`

 classifyTestsF,MArrayLM2-method

Multiple Testing Genewise Across Contrasts

Description

For each gene, classify a series of related t-statistics as up, down or not significant.

Usage

```
## S4 method for signature 'MArrayLM2'
classifyTestsF(
  object,
  cor.matrix = NULL,
  df = Inf,
  p.value = 0.01,
  fstat.only = FALSE
)
```

Arguments

object	numeric matrix of t-statistics or an 'MArrayLM2' object from which the t-statistics may be extracted.
cor.matrix	covariance matrix of each row of t-statistics. Defaults to the identity matrix.
df	numeric vector giving the degrees of freedom for the t-statistics. May have length 1 or length equal to the number of rows of tstat.
p.value	numeric value between 0 and 1 giving the desired size of the test
fstat.only	logical, if 'TRUE' then return the overall F-statistic as for 'FStat' instead of classifying the test results

Details

Works like `limma::classifyTestsF`, except object can have a list of covariance matrices `object$cov.coefficients.list`, instead of just one in `object$cov.coefficients`

See Also

`limma::classifyTestsF`

colinearityScore	<i>Collinearity score</i>
------------------	---------------------------

Description

Collinearity score for a regression model indicating if variables are too highly correlated to give meaningful results

Usage

```
colinearityScore(fit)
```

Arguments

fit regression model fit from lm() or lmer()

Value

Returns the collinearity score between 0 and 1, where a score > 0.999 means the degree of collinearity is too high. This function reports the correlation matrix between coefficient estimates for fixed effects. The collinearity score is the maximum absolute correlation value of this matrix. Note that the values are the correlation between the parameter estimates, and not between the variables themselves.

Examples

```
# load library
# library(variancePartition)

# load simulated data:
data(varPartData)
form <- ~ Age + (1|Individual) + (1|Tissue)

res <- fitVarPartModel( geneExpr[1:10,], form, info )

# evaluate the collinearity score on the first model fit
# this reports the correlation matrix between coefficients estimates
# for fixed effects
# the collinearity score is the maximum absolute correlation value
# If the collinearity score > .999 then the variance partition
# estimates may be problematic
# In that case, a least one variable should be omitted
colinearityScore(res[[1]])
```

dream *Differential expression with linear mixed model*

Description

Fit linear mixed model for differential expression and perform hypothesis test on fixed effects as specified in the contrast matrix L

Usage

```
dream(
  exprObj,
  formula,
  data,
  L,
  ddf = c("Satterthwaite", "Kenward-Roger"),
  useWeights = TRUE,
  weightsMatrix = NULL,
  control = lme4::lmerControl(calc.derivs = FALSE, check.rankX = "stop.deficient"),
  suppressWarnings = FALSE,
  quiet = FALSE,
  BPPARAM = bpparam(),
  computeResiduals = TRUE,
  REML = TRUE,
  ...
)
```

Arguments

exprObj	matrix of expression data (g genes x n samples), or ExpressionSet, or EList returned by voom() from the limma package
formula	specifies variables for the linear (mixed) model. Must only specify covariates, since the rows of exprObj are automatically used as a response. e.g.: ~ a + b + (1 c) Formulas with only fixed effects also work, and lmFit() followed by contrasts.fit() are run.
data	data.frame with columns corresponding to formula
L	contrast matrix specifying a linear combination of fixed effects to test
ddf	Specify "Satterthwaite" or "Kenward-Roger" method to estimate effective degrees of freedom for hypothesis testing in the linear mixed model. Note that Kenward-Roger is more accurate, but is <i>*much*</i> slower. Satterthwaite is a good enough approximation for most datasets.
useWeights	if TRUE, analysis uses heteroskedastic error estimates from voom(). Value is ignored unless exprObj is an EList() from voom() or weightsMatrix is specified
weightsMatrix	matrix the same dimension as exprObj with observation-level weights from voom(). Used only if useWeights is TRUE

control	control settings for <code>lmer()</code>
suppressWarnings	if TRUE, do not stop because of warnings or errors in model fit
quiet	suppress message, default FALSE
BPPARAM	parameters for parallel evaluation
computeResiduals	if TRUE, compute residuals and extract with <code>residuals(fit)</code> . Setting to FALSE saves memory
REML	use restricted maximum likelihood to fit linear mixed model. default is TRUE. See Details.
...	Additional arguments for <code>lmer()</code> or <code>lm()</code>

Details

A linear (mixed) model is fit for each gene in `exprObj`, using formula to specify variables in the regression (Hoffman and Roussos, 2021). If categorical variables are modeled as random effects (as is recommended), then a linear mixed model is used. For example if formula is $\sim a + b + (1|c)$, then the model is

```
fit <- lmer( exprObj[j,] ~ a + b + (1|c), data=data)
```

`useWeights=TRUE` causes `weightsMatrix[j,]` to be included as weights in the regression model.

Note: Fitting the model for 20,000 genes can be computationally intensive. To accelerate computation, models can be fit in parallel using `BiocParallel` to run code in parallel. Parallel processing must be enabled before calling this function. See below.

The regression model is fit for each gene separately. Samples with missing values in either gene expression or metadata are omitted by the underlying call to `lmer`.

Hypothesis tests and degrees of freedom are produced by `lmerTest` and `pbkrtest` packages

While `REML=TRUE` is required by `lmerTest` when `ddf='Kenward-Roger'`, `ddf='Satterthwaite'` can be used with `REML` as `TRUE` or `FALSE`. Since the Kenward-Roger method gave the best power with an accurate control of false positive rate in our simulations, and since the Satterthwaite method with `REML=TRUE` gives p-values that are slightly closer to the Kenward-Roger p-values, `REML=TRUE` is the default. See Vignette "3) Theory and practice of random effects and REML"

Value

`MArrayLM2` object (just like `MArrayLM` from `limma`), and the directly estimated p-value (without `eBayes`)

References

Hoffman GE, Roussos P (2021). "dream: Powerful differential expression analysis for repeated measures designs." *Bioinformatics*, **37**(2), 192–201.

Examples

```
# load library
# library(variancePartition)

# Intialize parallel backend with 4 cores
library(BiocParallel)
register(SnowParam(4))

# load simulated data:
# geneExpr: matrix of gene expression values
# info: information/metadata about each sample
data(varPartData)

form <- ~ Batch + (1|Individual) + (1|Tissue)

# Fit linear mixed model for each gene
# run on just 10 genes for time
fit = dream( geneExpr[1:10,], form, info)

# view top genes
topTable( fit )

# get contrast matrix testing if the coefficient for Batch3 is
# different from coefficient for Batch2
# The variable of interest must be a fixed effect
L = makeContrastsDream(form, info, contrasts=c("Batch3 - Batch2"))

# plot contrasts
plotContrasts( L )

# Fit linear mixed model for each gene
# run on just 10 genes for time
# Note that that dream() is not compatible with eBayes()
fit2 = dream( geneExpr[1:10,], form, info, L)

# view top genes
topTable( fit2, coef="Batch3 - Batch2" )

# Parallel processing using multiple cores with reduced memory usage
param = SnowParam(4, "SOCK", progressBar=TRUE)
fit3 = dream( geneExpr[1:10,], form, info, L, BPPARAM = param)

# Fit fixed effect model for each gene
# Use lmFit in the backend
# Need to run eBayes afterward
form <- ~ Batch
fit4 = dream( geneExpr[1:10,], form, info, L)
fit4 = eBayes( fit4 )

# view top genes
topTable( fit4, coef="Batch3 - Batch2" )
```

```
# Compute residuals using dream
residuals(fit4)
```

dscchisq	<i>Scaled chi-square</i>
----------	--------------------------

Description

Scaled chi-square density using a gamma distribution

Usage

```
dscchisq(x, a, b)
```

Arguments

x	vector of quantiles.
a	scale
b	degrees of freedom

eBayes,MArrayLM2-method	<i>eBayes for MArrayLM2</i>
-------------------------	-----------------------------

Description

eBayes for result of linear mixed model for with dream() using residual degrees of freedom approximated with rdf.merMod()

Usage

```
## S4 method for signature 'MArrayLM2'
eBayes(
  fit,
  proportion = 0.01,
  stdev.coef.lim = c(0.1, 4),
  trend = FALSE,
  robust = FALSE,
  winsor.tail.p = c(0.05, 0.1)
)
```

Arguments

fit	fit
proportion	proportion
stdev.coef.lim	stdev.coef.lim
trend	trend
robust	robust
winsor.tail.p	winsor.tail.p

Value

results of eBayes using approximated residual degrees of freedom

See Also

dream.rdf.merMod

ESS

Effective sample size

Description

Compute effective sample size based on correlation structure in linear mixed model

Usage

```
ESS(fit, method = "full")

## S4 method for signature 'lmerMod'
ESS(fit, method = "full")
```

Arguments

fit	model fit from lmer()
method	"full" uses the full correlation structure of the model. The "approximate" method makes the simplifying assumption that the study has a mean of m samples in each of k groups, and computes m based on the study design. When the study design is evenly balanced (i.e. the assumption is met), this gives the same results as the "full" method.

Details

Effective sample size calculations are based on:

Liu, G., and Liang, K. Y. (1997). Sample size calculations for studies with correlated observations. *Biometrics*, 53(3), 937-47.

"full" method: if

$$V_x = \text{var}(Y; x)$$

is the variance-covariance matrix of Y, the response, based on the covariate x, then the effective sample size corresponding to this covariate is

$$\sum_{i,j} (V_x^{-1})_{i,j}$$

. In R notation, this is: `sum(solve(V_x))`. In practice, this can be evaluated as `sum(w)`, where R

"approximate" method: Letting m be the mean number of samples per group,

$$k$$

be the number of groups, and

$$\rho$$

be the intraclass correlation, the effective sample size is

$$mk / (1 + \rho(m - 1))$$

Note that these values are equal when there are exactly m samples in each group. If m is only an average then this an approximation.

Value

effective sample size for each random effect in the model

Examples

```
library(lme4)
data(varPartData)

# Linear mixed model
fit <- lmer( geneExpr[1,] ~ (1|Individual) + (1|Tissue) + Age, info)

# Effective sample size
ESS( fit )
```

extractVarPart	<i>Extract variance statistics</i>
----------------	------------------------------------

Description

Extract variance statistics from list of models fit with `lm()` or `lmer()`

Usage

```
extractVarPart(modelList, showWarnings = TRUE, ...)
```

Arguments

<code>modelList</code>	list of <code>lmer()</code> model fits
<code>showWarnings</code>	show warnings about model fit (default TRUE)
<code>...</code>	other arguments

Value

data.frame of fraction of variance explained by each variable, after correcting for all others.

Examples

```
# library(variancePartition)

# Intialize parallel backend with 4 cores
library(BiocParallel)
register(SnowParam(4))

# load simulated data:
# geneExpr: matrix of gene expression values
# info: information/metadata about each sample
data(varPartData)

# Specify variables to consider
# Age is continuous so we model it as a fixed effect
# Individual and Tissue are both categorical, so we model them as random effects
form <- ~ Age + (1|Individual) + (1|Tissue)

# Step 1: fit linear mixed model on gene expression
# If categorical variables are specified, a linear mixed model is used
# If all variables are modeled as continuous, a linear model is used
# each entry in results is a regression model fit on a single gene
# Step 2: extract variance fractions from each model fit
# for each gene, returns fraction of variation attributable to each variable
# Interpretation: the variance explained by each variable
# after correction for all other variables
varPart <- fitExtractVarPartModel( geneExpr, form, info )
```

```

# violin plot of contribution of each variable to total variance
plotVarPart( sortCols( varPart ) )

# Advanced:
# Fit model and extract variance in two separate steps
# Step 1: fit model for each gene, store model fit for each gene in a list
results <- fitVarPartModel( geneExpr, form, info )

# Step 2: extract variance fractions
varPart <- extractVarPart( results )

```

```
fitExtractVarPartModel
```

Fit linear (mixed) model, report variance fractions

Description

Fit linear (mixed) model to estimate contribution of multiple sources of variation while simultaneously correcting for all other variables. Report fraction of variance attributable to each variable

Usage

```

fitExtractVarPartModel(
  exprObj,
  formula,
  data,
  REML = FALSE,
  useWeights = TRUE,
  weightsMatrix = NULL,
  showWarnings = TRUE,
  control = lme4::lmerControl(calc.derivs = FALSE, check.rankX = "stop.deficient"),
  quiet = FALSE,
  BPPARAM = bpparam(),
  ...
)

## S4 method for signature 'matrix'
fitExtractVarPartModel(
  exprObj,
  formula,
  data,
  REML = FALSE,
  useWeights = TRUE,
  weightsMatrix = NULL,
  showWarnings = TRUE,
  control = lme4::lmerControl(calc.derivs = FALSE, check.rankX = "stop.deficient"),

```

```
    quiet = FALSE,
    BPPARAM = bpparam(),
    ...
)

## S4 method for signature 'data.frame'
fitExtractVarPartModel(
  exprObj,
  formula,
  data,
  REML = FALSE,
  useWeights = TRUE,
  weightsMatrix = NULL,
  showWarnings = TRUE,
  control = lme4::lmerControl(calc.derivs = FALSE, check.rankX = "stop.deficient"),
  quiet = FALSE,
  BPPARAM = bpparam(),
  ...
)

## S4 method for signature 'EList'
fitExtractVarPartModel(
  exprObj,
  formula,
  data,
  REML = FALSE,
  useWeights = TRUE,
  weightsMatrix = NULL,
  showWarnings = TRUE,
  control = lme4::lmerControl(calc.derivs = FALSE, check.rankX = "stop.deficient"),
  quiet = FALSE,
  BPPARAM = bpparam(),
  ...
)

## S4 method for signature 'ExpressionSet'
fitExtractVarPartModel(
  exprObj,
  formula,
  data,
  REML = FALSE,
  useWeights = TRUE,
  weightsMatrix = NULL,
  showWarnings = TRUE,
  control = lme4::lmerControl(calc.derivs = FALSE, check.rankX = "stop.deficient"),
  quiet = FALSE,
  BPPARAM = bpparam(),
  ...
)
```

```

)

## S4 method for signature 'sparseMatrix'
fitExtractVarPartModel(
  exprObj,
  formula,
  data,
  REML = FALSE,
  useWeights = TRUE,
  weightsMatrix = NULL,
  showWarnings = TRUE,
  control = lme4::lmerControl(calc.derivs = FALSE, check.rankX = "stop.deficient"),
  quiet = FALSE,
  BPPARAM = bpparam(),
  ...
)

```

Arguments

<code>exprObj</code>	matrix of expression data (g genes x n samples), or <code>ExpressionSet</code> , or <code>EList</code> returned by <code>voom()</code> from the <code>limma</code> package
<code>formula</code>	specifies variables for the linear (mixed) model. Must only specify covariates, since the rows of <code>exprObj</code> are automatically used as a response. e.g.: <code>~ a + b + (1 c)</code>
<code>data</code>	<code>data.frame</code> with columns corresponding to <code>formula</code>
<code>REML</code>	use restricted maximum likelihood to fit linear mixed model. default is <code>FALSE</code> . See Details.
<code>useWeights</code>	if <code>TRUE</code> , analysis uses heteroskedastic error estimates from <code>voom()</code> . Value is ignored unless <code>exprObj</code> is an <code>EList()</code> from <code>voom()</code> or <code>weightsMatrix</code> is specified
<code>weightsMatrix</code>	matrix the same dimension as <code>exprObj</code> with observation-level weights from <code>voom()</code> . Used only if <code>useWeights</code> is <code>TRUE</code>
<code>showWarnings</code>	show warnings about model fit (default <code>TRUE</code>)
<code>control</code>	control settings for <code>lmer()</code>
<code>quiet</code>	suppress message, default <code>FALSE</code>
<code>BPPARAM</code>	parameters for parallel evaluation
<code>...</code>	Additional arguments for <code>lmer()</code> or <code>lm()</code>

Details

A linear (mixed) model is fit for each gene in `exprObj`, using `formula` to specify variables in the regression. If categorical variables are modeled as random effects (as is recommended), then a linear mixed model is used. For example if `formula` is `~ a + b + (1|c)`, then the model is

```
fit <- lmer( exprObj[j,] ~ a + b + (1|c), data=data)
```

If there are no random effects, so `formula` is `~ a + b + c`, a 'standard' linear model is used:

```
fit <-lm( exprObj[j,] ~ a + b + c,data=data)
```

In both cases, useWeights=TRUE causes weightsMatrix[j,] to be included as weights in the regression model.

Note: Fitting the model for 20,000 genes can be computationally intensive. To accelerate computation, models can be fit in parallel using BiocParallel to run in parallel. Parallel processing must be enabled before calling this function. See below.

The regression model is fit for each gene separately. Samples with missing values in either gene expression or metadata are omitted by the underlying call to lm/lmer.

REML=FALSE uses maximum likelihood to estimate variance fractions. This approach produced unbiased estimates, while REML=TRUE can show substantial bias. See Vignette "3) Theory and practice of random effects and REML"

Value

list() of where each entry is a model fit produced by lmer() or lm()

Examples

```
# load library
# library(variancePartition)

# Intialize parallel backend with 4 cores
library(BiocParallel)
register(SnowParam(4))

# load simulated data:
# geneExpr: matrix of gene expression values
# info: information/metadata about each sample
data(varPartData)

# Specify variables to consider
# Age is continuous so we model it as a fixed effect
# Individual and Tissue are both categorical, so we model them as random effects
form <- ~ Age + (1|Individual) + (1|Tissue)

# Step 1: fit linear mixed model on gene expression
# If categorical variables are specified, a linear mixed model is used
# If all variables are modeled as continuous, a linear model is used
# each entry in results is a regression model fit on a single gene
# Step 2: extract variance fractions from each model fit
# for each gene, returns fraction of variation attributable to each variable
# Interpretation: the variance explained by each variable
# after correction for all other variables
varPart <- fitExtractVarPartModel( geneExpr, form, info )

# violin plot of contribution of each variable to total variance
plotVarPart( sortCols( varPart ) )

# Note: fitExtractVarPartModel also accepts ExpressionSet
data(sample.ExpressionSet, package="Biobase")
```

```
# ExpressionSet example
form <- ~ (1|sex) + (1|type) + score
info2 <- Biobase::pData(sample.ExpressionSet)
varPart2 <- fitExtractVarPartModel( sample.ExpressionSet, form, info2 )
```

fitVarPartModel	<i>Fit linear (mixed) model</i>
-----------------	---------------------------------

Description

Fit linear (mixed) model to estimate contribution of multiple sources of variation while simultaneously correcting for all other variables.

Usage

```
fitVarPartModel(
  exprObj,
  formula,
  data,
  REML = FALSE,
  useWeights = TRUE,
  weightsMatrix = NULL,
  showWarnings = TRUE,
  fxn = identity,
  control = lme4::lmerControl(calc.derivs = FALSE, check.rankX = "stop.deficient"),
  quiet = FALSE,
  BPPARAM = bpparam(),
  ...
)
```

S4 method for signature 'matrix'

```
fitVarPartModel(
  exprObj,
  formula,
  data,
  REML = FALSE,
  useWeights = TRUE,
  weightsMatrix = NULL,
  showWarnings = TRUE,
  fxn = identity,
  control = lme4::lmerControl(calc.derivs = FALSE, check.rankX = "stop.deficient"),
  quiet = FALSE,
  BPPARAM = bpparam(),
  ...
)
```

```
)

## S4 method for signature 'data.frame'
fitVarPartModel(
  exprObj,
  formula,
  data,
  REML = FALSE,
  useWeights = TRUE,
  weightsMatrix = NULL,
  showWarnings = TRUE,
  fxn = identity,
  control = lme4::lmerControl(calc.derivs = FALSE, check.rankX = "stop.deficient"),
  quiet = FALSE,
  BPPARAM = bpparam(),
  ...
)

## S4 method for signature 'EList'
fitVarPartModel(
  exprObj,
  formula,
  data,
  REML = FALSE,
  useWeights = TRUE,
  weightsMatrix = NULL,
  showWarnings = TRUE,
  fxn = identity,
  control = lme4::lmerControl(calc.derivs = FALSE, check.rankX = "stop.deficient"),
  quiet = FALSE,
  BPPARAM = bpparam(),
  ...
)

## S4 method for signature 'ExpressionSet'
fitVarPartModel(
  exprObj,
  formula,
  data,
  REML = FALSE,
  useWeights = TRUE,
  weightsMatrix = NULL,
  showWarnings = TRUE,
  fxn = identity,
  control = lme4::lmerControl(calc.derivs = FALSE, check.rankX = "stop.deficient"),
  quiet = FALSE,
  BPPARAM = bpparam(),
  ...
)
```

```

)

## S4 method for signature 'sparseMatrix'
fitVarPartModel(
  exprObj,
  formula,
  data,
  REML = FALSE,
  useWeights = TRUE,
  weightsMatrix = NULL,
  showWarnings = TRUE,
  fxn = identity,
  control = lme4::lmerControl(calc.derivs = FALSE, check.rankX = "stop.deficient"),
  quiet = FALSE,
  BPPARAM = bpparam(),
  ...
)

```

Arguments

exprObj	matrix of expression data (g genes x n samples), or ExpressionSet, or EList returned by voom() from the limma package
formula	specifies variables for the linear (mixed) model. Must only specify covariates, since the rows of exprObj are automatically used as a response. e.g.: ~ a + b + (1 c)
data	data.frame with columns corresponding to formula
REML	use restricted maximum likelihood to fit linear mixed model. default is FALSE. See Details.
useWeights	if TRUE, analysis uses heteroskedastic error estimates from voom(). Value is ignored unless exprObj is an EList() from voom() or weightsMatrix is specified
weightsMatrix	matrix the same dimension as exprObj with observation-level weights from voom(). Used only if useWeights is TRUE
showWarnings	show warnings about model fit (default TRUE)
fxn	apply function to model fit for each gene. Defaults to identity function so it returns the model fit itself
control	control settings for lmer()
quiet	suppress message, default FALSE
BPPARAM	parameters for parallel evaluation
...	Additional arguments for lmer() or lm()

Details

A linear (mixed) model is fit for each gene in exprObj, using formula to specify variables in the regression. If categorical variables are modeled as random effects (as is recommended), then a linear mixed model is used. For example if formula is ~ a + b + (1|c), then the model is


```
fit <- lmer( exprObj[j,] ~ a + b + (1|c), data=data)
```

If there are no random effects, so formula is $\sim a + b + c$, a 'standard' linear model is used:

```
fit <- lm( exprObj[j,] ~ a + b + c, data=data)
```

In both cases, `useWeights=TRUE` causes `weightsMatrix[j,]` to be included as weights in the regression model.

Note: Fitting the model for 20,000 genes can be computationally intensive. To accelerate computation, models can be fit in parallel using `BiocParallel` to run in parallel. Parallel processing must be enabled before calling this function. See below.

The regression model is fit for each gene separately. Samples with missing values in either gene expression or metadata are omitted by the underlying call to `lm/lmer`.

Since this function returns a list of each model fit, using this function is slower and uses more memory than `fitExtractVarPartModel()`.

`REML=FALSE` uses maximum likelihood to estimate variance fractions. This approach produced unbiased estimates, while `REML=TRUE` can show substantial bias. See Vignette "3) Theory and practice of random effects and REML"

Value

`list()` of where each entry is a model fit produced by `lmer()` or `lm()`

Examples

```
# load library
# library(variancePartition)

# Intialize parallel backend with 4 cores
library(BiocParallel)
register(SnowParam(4))

# load simulated data:
# geneExpr: matrix of gene expression values
# info: information/metadata about each sample
data(varPartData)

# Specify variables to consider
# Age is continuous so we model it as a fixed effect
# Individual and Tissue are both categorical, so we model them as random effects
form <- ~ Age + (1|Individual) + (1|Tissue)

# Step 1: fit linear mixed model on gene expression
# If categorical variables are specified, a linear mixed model is used
# If all variables are modeled as continuous, a linear model is used
# each entry in results is a regression model fit on a single gene
# Step 2: extract variance fractions from each model fit
# for each gene, returns fraction of variation attributable to each variable
# Interpretation: the variance explained by each variable
# after correction for all other variables
varPart <- fitExtractVarPartModel( geneExpr, form, info )
```

```

# violin plot of contribution of each variable to total variance
# also sort columns
plotVarPart( sortCols( varPart ) )

# Advanced:
# Fit model and extract variance in two separate steps
# Step 1: fit model for each gene, store model fit for each gene in a list
results <- fitVarPartModel( geneExpr, form, info )

# Step 2: extract variance fractions
varPart <- extractVarPart( results )

# Note: fitVarPartModel also accepts ExpressionSet
data(sample.ExpressionSet, package="Biobase")

# ExpressionSet example
form <- ~ (1|sex) + (1|type) + score
info2 <- Biobase::pData(sample.ExpressionSet)
results2 <- fitVarPartModel( sample.ExpressionSet, form, info2 )

```

getContrast

Extract contrast matrix for linear mixed model

Description

Extract contrast matrix, L, testing a single variable. Contrasts involving more than one variable can be constructed by modifying L directly

Usage

```
getContrast(exprObj, formula, data, coefficient)
```

Arguments

exprObj	matrix of expression data (g genes x n samples), or ExpressionSet, or EList returned by voom() from the limma package
formula	specifies variables for the linear (mixed) model. Must only specify covariates, since the rows of exprObj are automatically used as a response. e.g.: ~ a + b + (1 c) Formulas with only fixed effects also work
data	data.frame with columns corresponding to formula
coefficient	the coefficient to use in the hypothesis test

Value

Contrast matrix testing one variable

Examples

```

# load simulated data:
# geneExpr: matrix of gene expression values
# info: information/metadata about each sample
data(varPartData)

# get contrast matrix testing if the coefficient for Batch2 is zero
# The variable of interest must be a fixed effect
form <- ~ Batch + (1|Individual) + (1|Tissue)
L = getContrast( geneExpr, form, info, "Batch3")

# get contrast matrix testing if Batch3 - Batch2 = 0
form <- ~ Batch + (1|Individual) + (1|Tissue)
L = getContrast( geneExpr, form, info, c("Batch3", "Batch2"))

# To test against Batch1 use the formula:
# ~ 0 + Batch + (1|Individual) + (1|Tissue)
# to estimate Batch1 directly instead of using it as the baseline

```

get_prediction	<i>Compute predicted value of formula for linear (mixed) model</i>
----------------	--

Description

Compute predicted value of formula for linear (mixed) model for with `lm` or `lmer`

Usage

```

get_prediction(fit, formula)

## S4 method for signature 'lmerMod'
get_prediction(fit, formula)

## S4 method for signature 'lm'
get_prediction(fit, formula)

```

Arguments

fit	model fit with <code>lm</code> or <code>lmer</code>
formula	formula of fixed and random effects to predict

Details

Similar motivation as `lme4:::predict.merMod()`, but that function cannot use just a subset of the fixed effects: it either uses none or all. Note that the intercept is included in the formula by default. To exclude it from the prediction use `~ 0 + ...` syntax

Value

Predicted values from formula using parameter estimates from fit linear (mixed) model

Examples

```
library(lme4)

# Linear model
fit <- lm(Reaction ~ Days, sleepstudy)

# prediction of intercept
get_prediction( fit, ~ 1)

# prediction of Days without intercept
get_prediction( fit, ~ 0 + Days)

# Linear mixed model

# fit model
fm1 <- lmer(Reaction ~ Days + (Days | Subject), sleepstudy)

# predict Days, but exclude intercept
get_prediction( fm1, ~ 0 + Days)

# predict Days and (Days | Subject) random effect, but exclude intercept
get_prediction( fm1, ~ 0 + Days + (Days | Subject))
```

ggColorHue

Default colors for ggplot

Description

Return an array of n colors the same as the default used by ggplot2

Usage

```
ggColorHue(n)
```

Arguments

n number of colors

Value

array of colors of length n

Examples

```
ggColorHue(4)
```

isRunnableFormula	<i>Test if formula is full rank on this dataset</i>
-------------------	---

Description

Test if formula is full rank on this dataset

Usage

```
isRunnableFormula(exprObj, formula, data)
```

Arguments

exprObj	expression object
formula	formula
data	data

makeContrastsDream	<i>Construct Matrix of Custom Contrasts</i>
--------------------	---

Description

Construct the contrast matrix corresponding to specified contrasts of a set of parameters.

Usage

```
makeContrastsDream(
  formula,
  data,
  ...,
  contrasts = NULL,
  suppressWarnings = FALSE
)
```

Arguments

formula	specifies variables for the linear (mixed) model. Must only specify covariates, since the rows of exprObj are automatically used as a response. e.g.: $\sim a + b + (1 c)$ Formulas with only fixed effects also work
data	data.frame with columns corresponding to formula
...	expressions, or character strings which can be parsed to expressions, specifying contrasts
contrasts	character vector specifying contrasts
suppressWarnings	(default FALSE). suppress warnings for univariate contrasts

Details

This function expresses contrasts between a set of parameters as a numeric matrix. The parameters are usually the coefficients from a linear (mixed) model fit, so the matrix specifies which comparisons between the coefficients are to be extracted from the fit. The output from this function is usually used as input to `dream()`.

This function is inspired by `limma::makeContrasts()` but is designed to be compatible with linear mixed models for `dream()`

Names in `...` and `contrasts` will be used as column names in the returned value.

Value

matrix of linear contrasts between regression coefficients

Examples

```
# load library
# library(variancePartition)

# Intialize parallel backend with 4 cores
library(BiocParallel)
register(SnowParam(4))

# load simulated data:
# geneExpr: matrix of gene expression values
# info: information/metadata about each sample
data(varPartData)

form <- ~ 0 + Batch + (1|Individual) + (1|Tissue)

# Define contrasts
L = makeContrastsDream( form, info, contrasts = c(Batch1_vs_2 = "Batch1 - Batch2", Batch3_vs_4 = "Batch3 - Batch4"),

# show contrasts matrix
L

# Plot to visualize contrasts matrix
plotContrasts(L)

# Fit linear mixed model for each gene
# run on just 10 genes for time
fit = dream( geneExpr[1:10,], form, info, L=L)

# examine contrasts after fitting
head(coef(fit))

# show results from first contrast
topTable(fit, coef="Batch1_vs_2")

# show results from second contrast
topTable(fit, coef="Batch3_vs_4")
```

```
# show results from third contrast
topTable(fit, coef="Batch1_vs_34")
```

MArrayLM2-class	<i>Class MArrayLM2</i>
-----------------	------------------------

Description

Class MArrayLM2

plotCompareP	<i>Compare p-values from two analyses</i>
--------------	---

Description

Plot $-\log_{10}$ p-values from two analyses and color based on donor component from variancePartition analysis

Usage

```
plotCompareP(
  p1,
  p2,
  vpDonor,
  dupcorvalue,
  fraction = 0.2,
  xlabel = bquote(duplicateCorrelation ~ (-log[10] ~ p)),
  ylabel = bquote(dream ~ (-log[10] ~ p))
)
```

Arguments

p1	p-value from first analysis
p2	p-value from second analysis
vpDonor	donor component for each gene from variancePartition analysis
dupcorvalue	scalar donor component from duplicateCorrelation
fraction	fraction of highest/lowest values to use for best fit lines
xlabel	for x-axis
ylabel	label for y-axis

Value

ggplot2 plot

Examples

```
# load library
# library(variancePartition)

# Intialize parallel backend with 4 cores
library(BiocParallel)
register(SnowParam(4))

# load simulated data:
# geneExpr: matrix of gene expression values
# info: information/metadata about each sample
data(varPartData)

# Perform very simple analysis for demonstration

# Analysis 1
form <- ~ Batch
fit = dream( geneExpr, form, info)
fit = eBayes( fit )
res = topTable( fit, number=Inf, coef="Batch3" )

# Analysis 2
form <- ~ Batch + (1|Tissue)
fit2 = dream( geneExpr, form, info)
res2 = topTable( fit2, number=Inf, coef="Batch3" )

# Compare p-values
plotCompareP( res$P.Value, res2$P.Value, runif(nrow(res)), .3 )
```

plotContrasts

Plot representation of contrast matrix

Description

Plot contrast matrix to clarify interpretation of hypothesis tests with linear contrasts

Usage

```
plotContrasts(L)
```

Arguments

L contrast matrix

Value

ggplot2 object

Examples

```

# load library
# library(variancePartition)

# load simulated data:
# geneExpr: matrix of gene expression values
# info: information/metadata about each sample
data(varPartData)

# 1) get contrast matrix testing if the coefficient for Batch2 is zero
# 2) get contrast matrix testing if the coefficient for Batch2 is different from Batch3
form <- ~ Batch + (1|Individual) + (1|Tissue)
L = makeContrastsDream(form, info, contrasts=c("Batch2", Batch_3_vs_2 = "Batch3 - Batch2"))

# plot contrasts
plotContrasts( L )

```

plotCorrMatrix

plotCorrMatrix

Description

Plot correlation matrix

Usage

```

plotCorrMatrix(
  C,
  dendrogram = "both",
  sort = TRUE,
  margins = c(13, 13),
  key.xlab = "correlation",
  ...
)

```

Arguments

C	correlation matrix: R or R ² matrix
dendrogram	character string indicating whether to draw 'both' or none'
sort	sort rows and columns based on clustering
margins	spacing of plot
key.xlab	label of color gradient
...	additional arguments to heatmap.2

Details

Plots image of correlation matrix using customized call to heatmap.2

Value

Image of correlation matrix

Examples

```
# simulate simple matrix of 10 variables
mat = matrix(rnorm(1000), ncol=10)

# compute correlation matrix
C = cor(mat)

# plot correlations
plotCorrMatrix( C )

# plot squared correlations
plotCorrMatrix( C^2, dendrogram="none" )
```

`plotCorrStructure` *plotCorrStructure*

Description

Plot correlation structure of a gene based on random effects

Usage

```
plotCorrStructure(
  fit,
  varNames = names(coef(fit)),
  reorder = TRUE,
  pal = colorRampPalette(c("white", "red", "darkred")),
  hclust.method = "complete"
)
```

Arguments

<code>fit</code>	linear mixed model fit of a gene produced by <code>lmer()</code> or <code>fitVarPartModel()</code>
<code>varNames</code>	variables in the metadata for which the correlation structure should be shown. Variables must be random effects
<code>reorder</code>	how to reorder the rows/columns of the correlation matrix. <code>reorder=FALSE</code> gives no reorder. <code>reorder=TRUE</code> reorders based on <code>hclust</code> . <code>reorder</code> can also be an array of indices to reorder the samples manually
<code>pal</code>	color palette
<code>hclust.method</code>	clustering methods for <code>hclust</code>

Value

Image of correlation structure between each pair of experiments for a single gene

Examples

```
# load library
# library(variancePartition)

# Intialize parallel backend with 4 cores
library(BiocParallel)
register(SnowParam(4))

# load simulated data:
data(varPartData)

# specify formula
form <- ~ Age + (1|Individual) + (1|Tissue)

# fit and return linear mixed models for each gene
fitList <- fitVarPartModel( geneExpr[1:10,], form, info )

# Focus on the first gene
fit = fitList[[1]]

# plot correlation sturcture based on Individual, reordering samples with hclust
plotCorrStructure( fit, "Individual" )

# don't reorder
plotCorrStructure( fit, "Individual", reorder=FALSE )

# plot correlation sturcture based on Tissue, reordering samples with hclust
plotCorrStructure( fit, "Tissue" )

# don't reorder
plotCorrStructure( fit, "Tissue", FALSE )

# plot correlation structure based on all random effects
# reorder manually by Tissue and Individual
idx = order(info$Tissue, info$Individual)
plotCorrStructure( fit, reorder=idx )

# plot correlation structure based on all random effects
# reorder manually by Individual, then Tissue
idx = order(info$Individual, info$Tissue)
plotCorrStructure( fit, reorder=idx )
```

Description

Bar plot of variance fractions for a subset of genes

Usage

```
plotPercentBars(  
  varPart,  
  col = c(ggColorHue(ncol(varPart) - 1), "grey85"),  
  width = NULL  
)
```

Arguments

varPart	object returned by extractVarPart() or fitExtractVarPartModel()
col	color of bars for each variable
width	specify width of bars

Value

Returns ggplot2 barplot

Examples

```
# library(variancePartition)  
  
# Intialize parallel backend with 4 cores  
library(BiocParallel)  
register(SnowParam(4))  
  
# load simulated data:  
# geneExpr: matrix of gene expression values  
# info: information/metadata about each sample  
data(varPartData)  
  
# Specify variables to consider  
form <- ~ Age + (1|Individual) + (1|Tissue)  
  
# Fit model  
varPart <- fitExtractVarPartModel( geneExpr, form, info )  
  
# Bar plot for a subset of genes showing variance fractions  
plotPercentBars( varPart[1:5,] )  
  
# Move the legend to the top  
plotPercentBars( varPart[1:5,] ) + theme(legend.position="top")
```

plotStratify	<i>plotStratify</i>
--------------	---------------------

Description

Plot gene expression stratified by another variable

Usage

```
plotStratify(
  formula,
  data,
  xlab,
  ylab,
  main,
  sortBy,
  colorBy,
  sort = TRUE,
  text = NULL,
  text.y = 1,
  text.size = 5,
  pts.cex = 1,
  ylim = NULL,
  legend = TRUE,
  x.labels = FALSE
)
```

Arguments

formula	specify variables shown in the x- and y-axes. Y-axis should be continuous variable, x-axis should be discrete.
data	data.frame storing continuous and discrete variables specified in formula
xlab	label x-axis. Defaults to value of xval
ylab	label y-axis. Defaults to value of yval
main	main label
sortBy	name of column in geneExpr to sort samples by. Defaults to xval
colorBy	name of column in geneExpr to color box plots. Defaults to xval
sort	if TRUE, sort boxplots by median value, else use default ordering
text	plot text on the top left of the plot
text.y	indicate position of the text on the y-axis as a fraction of the y-axis range
text.size	size of text
pts.cex	size of points
ylim	specify range of y-axis
legend	show legend
x.labels	show x axis labels

Value

ggplot2 object

Examples

```
# Note: This is a newer, more convient interface to plotStratifyBy()

# load library
# library(variancePartition)

# load simulated data:
data(varPartData)

# Create data.frame with expression and Tissue information for each sample
GE = data.frame( Expression = geneExpr[1,], Tissue = info$Tissue)

# Plot expression stratified by Tissue
plotStratify( Expression ~ Tissue, GE )

# Omit legend and color boxes grey
plotStratify( Expression ~ Tissue, GE, colorBy = NULL)

# Specify colors
col = c( B="green", A="red", C="yellow")
plotStratify( Expression ~ Tissue, GE, colorBy=col, sort=FALSE)
```

plotStratifyBy

plotStratifyBy

Description

Plot gene expression stratified by another variable

Usage

```
plotStratifyBy(
  geneExpr,
  xval,
  yval,
  xlab = xval,
  ylab = yval,
  main = NULL,
  sortBy = xval,
  colorBy = xval,
  sort = TRUE,
  text = NULL,
  text.y = 1,
```

```

    text.size = 5,
    pts.cex = 1,
    ylim = NULL,
    legend = TRUE,
    x.labels = FALSE
  )

```

Arguments

geneExpr	data.frame of gene expression values and another variable for each sample. If there are multiple columns, the user can specify which one to use
xval	name of column in geneExpr to be used along x-axis to stratify gene expression
yval	name of column in geneExpr indicating gene expression
xlab	label x-axis. Defaults to value of xval
ylab	label y-axis. Defaults to value of yval
main	main label
sortBy	name of column in geneExpr to sort samples by. Defaults to xval
colorBy	name of column in geneExpr to color box plots. Defaults to xval
sort	if TRUE, sort boxplots by median value, else use default ordering
text	plot text on the top left of the plot
text.y	indicate position of the text on the y-axis as a fraction of the y-axis range
text.size	size of text
pts.cex	size of points
ylim	specify range of y-axis
legend	show legend
x.labels	show x axis labels

Value

ggplot2 object

Examples

```

# load library
# library(variancePartition)

# load simulated data:
data(varPartData)

# Create data.frame with expression and Tissue information for each sample
GE = data.frame( Expression = geneExpr[1,], Tissue = info$Tissue)

# Plot expression stratified by Tissue
plotStratifyBy( GE, "Tissue", "Expression")

# Omit legend and color boxes grey

```

```

plotStratifyBy( GE, "Tissue", "Expression", colorBy = NULL)

# Specify colors
col = c( B="green", A="red", C="yellow")
plotStratifyBy( GE, "Tissue", "Expression", colorBy=col, sort=FALSE)

```

plotVarianceEstimates *Plot Variance Estimates*

Description

Plot Variance Estimates

Usage

```

plotVarianceEstimates(
  fit,
  fitEB,
  var_true = NULL,
  xmax = quantile(fit$sigma^2, 0.999)
)

```

Arguments

fit	model fit from dream()
fitEB	model fit from eBayes()
var_true	array of true variance values from simulation (optional)
xmax	maximum value on the x-axis

plotVarPart *Violin plot of variance fractions*

Description

Violin plot of variance fraction for each gene and each variable

Usage

```
plotVarPart(  
  obj,  
  col = c(ggColorHue(ncol(obj) - 1), "grey85"),  
  label.angle = 20,  
  main = "",  
  ylab = "",  
  convertToPercent = TRUE,  
  ...  
)  
  
## S4 method for signature 'matrix'  
plotVarPart(  
  obj,  
  col = c(ggColorHue(ncol(obj) - 1), "grey85"),  
  label.angle = 20,  
  main = "",  
  ylab = "",  
  convertToPercent = TRUE,  
  ...  
)  
  
## S4 method for signature 'data.frame'  
plotVarPart(  
  obj,  
  col = c(ggColorHue(ncol(obj) - 1), "grey85"),  
  label.angle = 20,  
  main = "",  
  ylab = "",  
  convertToPercent = TRUE,  
  ...  
)  
  
## S4 method for signature 'varPartResults'  
plotVarPart(  
  obj,  
  col = c(ggColorHue(ncol(obj) - 1), "grey85"),  
  label.angle = 20,  
  main = "",  
  ylab = "",  
  convertToPercent = TRUE,  
  ...  
)
```

Arguments

obj	varParFrac object returned by fitExtractVarPart or extractVarPart
col	vector of colors

label.angle	angle of labels on x-axis
main	title of plot
ylab	text on y-axis
convertToPercent	
	multiply fractions by 100 to convert to percent values
...	additional arguments

Value

Makes violin plots of variance components model. This function uses the graphics interface from ggplot2. Warnings produced by this function usually ggplot2 warning that the window is too small.

Examples

```
# load library
# library(variancePartition)

# Intialize parallel backend with 4 cores
library(BiocParallel)
register(SnowParam(4))

# load simulated data:
# geneExpr: matrix of gene expression values
# info: information/metadata about each sample
data(varPartData)

# Specify variables to consider
# Age is continuous so we model it as a fixed effect
# Individual and Tissue are both categorical, so we model them as random effects
form <- ~ Age + (1|Individual) + (1|Tissue)

varPart <- fitExtractVarPartModel( geneExpr, form, info )

# violin plot of contribution of each variable to total variance
plotVarPart( sortCols( varPart ) )
```

rdf.merMod

Approximate residual degrees of freedom

Description

For a linear model with n samples and p covariates, $RSS/\sigma^2 \sim \chi_\nu^2$ where $\nu = n - p$ is the residual degrees of freedom. In the case of a linear mixed model, the distribution is no longer exactly a chi-square distribution, but can be approximated with a chi-square distribution.

Given the hat matrix, H , that maps between observed and fitted responses, the approximate residual degrees of freedom is $\nu = \text{tr}((I - H)^T(I - H))$. For a linear model, this simplifies to the well

known form $\nu = n - p$. In the more general case, such as a linear mixed model, the original form simplifies only to $n - 2tr(H) + tr(HH)$ and is an approximation rather than being exact. The third term here is quadratic time in the number of samples, n , and can be computationally expensive to evaluate for larger datasets. Here we develop a linear time algorithm that takes advantage of the fact that H is low rank.

H is computed as $A^T A + B^T B$ for $A=CL$ and $B=CR$ defined in the code. Since A and B are low rank, there is no need to compute H directly. Instead, the terms $tr(H)$ and $tr(HH)$ can be computed using the eigen decompositions of AA^T and BB^T which is linear time in the number of samples.

Usage

```
rdf.merMod(model, method = c("linear", "quadratic"))
```

Arguments

model	An object of class merMod
method	Use algorithm that is "linear" (default) or quadratic time in the number of samples

Details

Compute the approximate residual degrees of freedom from a linear mixed model.

See Also

`rdf_from_matrices`

Examples

```
library(lme4)

# Fit linear mixed model
fit <- lmer(Reaction ~ Days + (Days | Subject), sleepstudy)

# Evaluate the approximate residual degrees of freedom
rdf.merMod(fit)
```

`rdf_from_matrices` *Fast approximate residual degrees of freedom*

Description

Defining $H = A^T A + B^T B$ where A and B are low rank, compute $n - 2tr(H) + tr(HH)$ in $O(np^2)$ instead of $O(n^2p^2)$.

Usage

```
rdf_from_matrices(A, B)
```

Arguments

A a matrix or sparseMatrix
 B a matrix or sparseMatrix

See Also

rdf.merMod

reOnly *Adapted from lme4:::reOnly*

Description

Adapted from lme4:::reOnly

Usage

```
reOnly(f, response = FALSE)
```

Arguments

f formula
 response (FALSE) is there a response in the formula

residuals,MArrayLM-method
residuals for MArrayLM

Description

residuals for MArrayLM

Usage

```
## S4 method for signature 'MArrayLM'  

residuals(object, ...)
```

Arguments

object MArrayLM object from dream
 ... other arguments, currently ignored

Value

results of residuals

residuals,MArrayLM2-method
residuals for MArrayLM2

Description

residuals for MArrayLM2

Usage

```
## S4 method for signature 'MArrayLM2'  
residuals(object, ...)
```

Arguments

object	MArrayLM2 object from dream
...	other arguments, currently ignored

Value

results of residuals

residuals,VarParFitList-method
Residuals from model fit

Description

Extract residuals for each gene from model fit with fitVarPartModel()

Usage

```
## S4 method for signature 'VarParFitList'  
residuals(object, ...)
```

Arguments

object	object produced by fitVarPartModel()
...	other arguments.

Details

If model is fit with missing data, residuals returns NA for entries that were missing in the original data

Value

Residuals extracted from model fits stored in object

Examples

```
# load library
# library(variancePartition)

# Intialize parallel backend with 4 cores
library(BiocParallel)
register(SnowParam(4))

# load simulated data:
# geneExpr: matrix of gene expression values
# info: information/metadata about each sample
data(varPartData)

# Specify variables to consider
# Age is continuous so we model it as a fixed effect
# Individual and Tissue are both categorical, so we model them as random effects
form <- ~ Age + (1|Individual) + (1|Tissue)

# Fit model
modelFit <- fitVarPartModel( geneExpr, form, info )

# Extract residuals of model fit
res <- residuals( modelFit )
```

residuals.MArrayLM2 *Residuals for result of dream*

Description

Residuals for result of dream

Usage

```
residuals.MArrayLM2(object, ...)
```

Arguments

object	See ?stats::residuals
...	See ?stats::residuals

shrinkageMetric	<i>Shrinkage metric for eBayes</i>
-----------------	------------------------------------

Description

Evaluates the coefficient from the linear regression of $s2.post \sim \sigma Sq$. When there is no shrinkage, this value is 1. Values less than 1 indicate the amount of shrinkage.

Usage

```
shrinkageMetric(sigmaSq, s2.post)
```

Arguments

sigmaSq	maximum likelihood residual variance for every gene
s2.post	empirical Bayes posterior estimate of residual variance for every gene

Details

Shrinkage metric for eBayes quantifying the amount of shrinkage that is applied to shrink the maximum likelihood residual variance to the empirical Bayes posterior estimate

sortCols	<i>Sort variance partition statistics</i>
----------	---

Description

Sort columns returned by `extractVarPart()` or `fitExtractVarPartModel()`

Usage

```
sortCols(
  x,
  FUN = median,
  decreasing = TRUE,
  last = c("Residuals", "Measurement.error"),
  ...
)

## S4 method for signature 'matrix'
sortCols(
  x,
  FUN = median,
  decreasing = TRUE,
  last = c("Residuals", "Measurement.error"),
```

```

    ...
  )

  ## S4 method for signature 'data.frame'
  sortCols(
    x,
    FUN = median,
    decreasing = TRUE,
    last = c("Residuals", "Measurement.error"),
    ...
  )

  ## S4 method for signature 'varPartResults'
  sortCols(
    x,
    FUN = median,
    decreasing = TRUE,
    last = c("Residuals", "Measurement.error"),
    ...
  )

```

Arguments

x	object returned by <code>extractVarPart()</code> or <code>fitExtractVarPartModel()</code>
FUN	function giving summary statistic to sort by. Defaults to <code>median</code>
decreasing	logical. Should the sorting be increasing or decreasing?
last	columns to be placed on the right, regardless of values in these columns
...	other arguments to sort

Value

data.frame with columns sorted by mean value, with Residuals in last column

Examples

```

# library(variancePartition)

# Intialize parallel backend with 4 cores
library(BiocParallel)
register(SnowParam(4))

# load simulated data:
# geneExpr: matrix of gene expression values
# info: information/metadata about each sample
data(varPartData)

# Specify variables to consider
# Age is continuous so we model it as a fixed effect
# Individual and Tissue are both categorical, so we model them as random effects

```



```

form <- ~ Age + (1|Individual) + (1|Tissue)

# Step 1: fit linear mixed model on gene expression
# If categorical variables are specified, a linear mixed model is used
# If all variables are modeled as continuous, a linear model is used
# each entry in results is a regression model fit on a single gene
# Step 2: extract variance fractions from each model fit
# for each gene, returns fraction of variation attributable to each variable
# Interpretation: the variance explained by each variable
# after correction for all other variables
varPart <- fitExtractVarPartModel( geneExpr, form, info )

# violin plot of contribution of each variable to total variance
# sort columns by median value
plotVarPart( sortCols( varPart ) )

```

VarParCList-class *Class VarParCList*

Description

Class VarParCList

VarParFitList-class *Class VarParFitList*

Description

Class VarParFitList

varParFrac-class *Class varParFrac*

Description

Class varParFrac

varPartConfInf *Linear mixed model confidence intervals*

Description

Fit linear mixed model to estimate contribution of multiple sources of variation while simultaneously correcting for all other variables. Then perform parametric bootstrap sampling to get a 95% confidence intervals for each variable for each gene.

Usage

```
varPartConfInf(
  exprObj,
  formula,
  data,
  REML = FALSE,
  useWeights = TRUE,
  weightsMatrix = NULL,
  showWarnings = TRUE,
  colinearityCutoff = 0.999,
  control = lme4::lmerControl(calc.derivs = FALSE, check.rankX = "stop.deficient"),
  nsim = 1000,
  ...
)
```

Arguments

exprObj	matrix of expression data (g genes x n samples), or ExpressionSet, or EList returned by voom() from the limma package
formula	specifies variables for the linear (mixed) model. Must only specify covariates, since the rows of exprObj are automatically used as a response. e.g.: ~ a + b + (1 c)
data	data.frame with columns corresponding to formula
REML	use restricted maximum likelihood to fit linear mixed model. default is FALSE. Strongly discourage against changing this option, but here for compatibility.
useWeights	if TRUE, analysis uses heteroskedastic error estimates from voom(). Value is ignored unless exprObj is an EList from voom() or weightsMatrix is specified
weightsMatrix	matrix the same dimension as exprObj with observation-level weights from voom(). Used only if useWeights is TRUE
showWarnings	show warnings about model fit (default TRUE)
colinearityCutoff	cutoff used to determine if model is computationally singular
control	control settings for lmer()
nsim	number of bootstrap datasets
...	Additional arguments for lmer() or lm()

Details

A linear mixed model is fit for each gene, and `bootMer()` is used to generate parametric bootstrap confidence intervals. `use.u=TRUE` is used so that the

$$\hat{u}$$

values from the random effects are used as estimated and are not re-sampled. This gives confidence intervals as if additional data were generated from these same current samples. Conversely, `use.u=FALSE` assumes that this dataset is a sample from a larger population. Thus it simulates

$$\hat{u}$$

based on the estimated variance parameter. This approach gives confidence intervals as if additional data were collected from the larger population from which this dataset is sampled. Overall, `use.u=TRUE` gives smaller confidence intervals that are appropriate in this case.

Value

`list()` of where each entry is the result for a gene. Each entry is a matrix of the 95% confidence interval of the variance fraction for each variable

Examples

```
# load library
# library(variancePartition)

# Intialize parallel backend with 4 cores
library(BiocParallel)
register(SnowParam(4))

# load simulated data:
# geneExpr: matrix of gene expression values
# info: information/metadata about each sample
data(varPartData)

# Specify variables to consider
# Age is continuous so we model it as a fixed effect
# Individual and Tissue are both categorical, so we model them as random effects
form <- ~ Age + (1|Individual) + (1|Tissue)

# Compute bootstrap confidence intervals for each variable for each gene
resCI <- varPartConfInf( geneExpr[1:5,], form, info, nsim=100 )
```

varPartData	<i>Simulation dataset for examples</i>
-------------	--

Description

A simulated dataset of gene expression and metadata

A simulated dataset of gene counts

info about study design

Normalized expression data

Usage

```
data(varPartData)
```

```
data(varPartData)
```

```
data(varPartData)
```

```
data(varPartData)
```

Format

A dataset of 100 samples and 200 genes

A dataset of 100 samples and 200 genes

A dataset of 100 samples and 200 genes

A dataset of 100 samples and 200 genes

Details

- geneCounts gene expression in the form of RNA-seq counts
- geneExpr gene expression on a continuous scale
- info metadata about the study design
- geneCounts gene expression in the form of RNA-seq counts
- geneExpr gene expression on a continuous scale
- info metadata about the study design
- geneCounts gene expression in the form of RNA-seq counts
- geneExpr gene expression on a continuous scale
- info metadata about the study design
- geneCounts gene expression in the form of RNA-seq counts
- geneExpr gene expression on a continuous scale
- info metadata about the study design

varPartDEdata	<i>Simulation dataset for dream example</i>
---------------	---

Description

Gene counts from RNA-seq
metadata matrix of sample information

Usage

```
data(varPartDEdata)  
  
data(varPartDEdata)
```

Format

A dataset of 24 samples and 19,364 genes
A dataset of 24 samples and 19,364 genes

Details

- countMatrix gene expression in the form of RNA-seq counts
- metadata metadata about the study design
- countMatrix gene expression in the form of RNA-seq counts
- metadata metadata about the study design

varPartResults-class	<i>Class varPartResults</i>
----------------------	-----------------------------

Description

Class varPartResults

voomWithDreamWeights *Transform RNA-Seq Data Ready for Linear Mixed Modelling with dream()*

Description

Transform count data to log₂-counts per million (logCPM), estimate the mean-variance relationship and use this to compute appropriate observation-level weights. The data are then ready for linear mixed modelling with `dream()`. This method is the same as `limma::voom()`, except that it allows random effects in the formula

Usage

```
voomWithDreamWeights(
  counts,
  formula,
  data,
  lib.size = NULL,
  normalize.method = "none",
  span = 0.5,
  weights = NULL,
  plot = FALSE,
  save.plot = FALSE,
  quiet = FALSE,
  BPPARAM = bpparam(),
  ...
)
```

Arguments

<code>counts</code>	a numeric matrix containing raw counts, or an <code>ExpressionSet</code> containing raw counts, or a <code>DGEList</code> object. Counts must be non-negative and NAs are not permitted.
<code>formula</code>	specifies variables for the linear (mixed) model. Must only specify covariates, since the rows of <code>exprObj</code> are automatically used as a response. e.g.: <code>~ a + b + (1 c)</code> Formulas with only fixed effects also work, and <code>lmFit()</code> followed by <code>contrasts.fit()</code> are run.
<code>data</code>	<code>data.frame</code> with columns corresponding to formula
<code>lib.size</code>	numeric vector containing total library sizes for each sample. Defaults to the normalized (effective) library sizes in counts if <code>counts</code> is a <code>DGEList</code> or to the columnwise count totals if <code>counts</code> is a matrix.
<code>normalize.method</code>	the microarray-style normalization method to be applied to the logCPM values (if any). Choices are as for the <code>method</code> argument of <code>normalizeBetweenArrays</code> when the data is single-channel. Any normalization factors found in counts will still be used even if <code>normalize.method="none"</code> .

span	width of the lowess smoothing window as a proportion.
weights	Can be a numeric matrix of individual weights of same dimensions as the counts, or a numeric vector of sample weights with length equal to ncol(counts)
plot	logical, should a plot of the mean-variance trend be displayed?
save.plot	logical, should the coordinates and line of the plot be saved in the output?
quiet	suppress message, default FALSE
BPPARAM	parameters for parallel evaluation
...	other arguments are passed to lmer.

Details

Adapted from vomm() in limma v3.40.2

Value

An EList object just like the result of limma::voom()

See Also

limma::voom()

Examples

```
# library(variancePartition)
library(edgeR)
library(BiocParallel)

data(varPartDEdata)

# normalize RNA-seq counts
dge = DGEList(counts = countMatrix)
dge = calcNormFactors(dge)

# specify formula with random effect for Individual
form <- ~ Disease + (1|Individual)

# compute observation weights
vobj = voomWithDreamWeights( dge[1:20,], form, metadata)

# fit dream model
res = dream( vobj, form, metadata)

# extract results
topTable(res, coef="Disease1")
```

Index

- * **datasets**
 - varPartData, [52](#)
 - varPartDEdata, [53](#)
- as.data.frame.varPartResults, [3](#)
- as.matrix
 - (as.matrix, varPartResults-method), [4](#)
- as.matrix, varPartResults-method, [4](#)
- calcVarPart, [5](#)
- calcVarPart, glm-method (calcVarPart), [5](#)
- calcVarPart, glmer-method (calcVarPart), [5](#)
- calcVarPart, glmerMod-method (calcVarPart), [5](#)
- calcVarPart, lm-method (calcVarPart), [5](#)
- calcVarPart, lmerMod-method (calcVarPart), [5](#)
- calcVarPart, negbin-method (calcVarPart), [5](#)
- canCorPairs, [7](#)
- classifyTestsF, [8](#)
- classifyTestsF, MArrayLM2-method, [9](#)
- colinearityScore, [10](#)
- countMatrix (varPartDEdata), [53](#)

- dream, [11](#)
- dscchisq, [14](#)

- eBayes, MArrayLM2-method, [14](#)
- ESS, [15](#)
- ESS, lmerMod-method (ESS), [15](#)
- extractVarPart, [17](#)

- fitExtractVarPartModel, [18](#)
- fitExtractVarPartModel, data.frame-method (fitExtractVarPartModel), [18](#)
- fitExtractVarPartModel, EList-method (fitExtractVarPartModel), [18](#)
- fitExtractVarPartModel, ExpressionSet-method (fitExtractVarPartModel), [18](#)
- fitExtractVarPartModel, matrix-method (fitExtractVarPartModel), [18](#)
- fitExtractVarPartModel, sparseMatrix-method (fitExtractVarPartModel), [18](#)
- fitVarPartModel, [22](#)
- fitVarPartModel, data.frame-method (fitVarPartModel), [22](#)
- fitVarPartModel, EList-method (fitVarPartModel), [22](#)
- fitVarPartModel, ExpressionSet-method (fitVarPartModel), [22](#)
- fitVarPartModel, matrix-method (fitVarPartModel), [22](#)
- fitVarPartModel, sparseMatrix-method (fitVarPartModel), [22](#)

- geneCounts (varPartData), [52](#)
- geneExpr (varPartData), [52](#)
- get_prediction, [27](#)
- get_prediction, lm-method (get_prediction), [27](#)
- get_prediction, lmerMod-method (get_prediction), [27](#)
- getContrast, [26](#)
- ggColorHue, [28](#)

- info (varPartData), [52](#)
- isRunnableFormula, [29](#)

- makeContrastsDream, [29](#)
- MArrayLM2-class, [31](#)
- metadata (varPartDEdata), [53](#)

- plotCompareP, [31](#)
- plotContrasts, [32](#)
- plotCorrMatrix, [33](#)
- plotCorrStructure, [34](#)
- plotPercentBars, [35](#)

plotStratify, [37](#)
plotStratifyBy, [38](#)
plotVarianceEstimates, [40](#)
plotVarPart, [40](#)
plotVarPart, data.frame-method
 (plotVarPart), [40](#)
plotVarPart, matrix-method
 (plotVarPart), [40](#)
plotVarPart, varPartResults-method
 (plotVarPart), [40](#)

rdf.merMod, [42](#)
rdf_from_matrices, [43](#)
reOnly, [44](#)
residuals, MArrayLM-method, [44](#)
residuals, MArrayLM2-method, [45](#)
residuals, VarParFitList-method, [45](#)
residuals.MArrayLM2, [46](#)

shrinkageMetric, [47](#)
sortCols, [47](#)
sortCols, data.frame-method (sortCols),
 [47](#)
sortCols, matrix-method (sortCols), [47](#)
sortCols, varPartResults-method
 (sortCols), [47](#)

VarParCIList-class, [49](#)
VarParFitList-class, [49](#)
varParFrac-class, [49](#)
varPartConfInf, [50](#)
varPartData, [52](#)
varPartDEdata, [53](#)
varPartResults-class, [53](#)
vroomWithDreamWeights, [54](#)