

gaga

February 8, 2012

<code>buildPatterns</code>	<i>Build a matrix with all possible patterns given a number of groups where samples may belong to.</i>
----------------------------	--

Description

Creates a matrix indicating which groups are put together under each pattern. The number of possible patterns increases very fast as the number of groups increases. This function provides an easy way to compute all possible patterns. The output of this function is usually used for the `patterns` parameter of the `lmFit` function.

Usage

```
buildPatterns(groups)
```

Arguments

<code>groups</code>	Character containing the names of the groups at which samples may belong to. If the output of the function is going to be used in <code>fitGG</code> it must match the group levels specified in the <code>groups</code> parameter that will be given to <code>fitGG</code> .
---------------------	---

Examples

```
buildPatterns(groups=c('GroupControl', 'GroupA', 'GroupB'))
```

<code>checkfit</code>	<i>Check goodness-of-fit of GaGa and MiGaGa models</i>
-----------------------	--

Description

Produces plots to check fit of GaGa and MiGaGa model. Compares observed data with posterior predictive distribution of the model. Can also compare posterior distribution of parameters with method of moments estimates.

Usage

```
checkfit(gg.fit, x, groups, type='data', logexpr=FALSE, xlab, ylab, main, lty, l
```

Arguments

<code>gg.fit</code>	GaGa or MiGaGa fit (object of type <code>gagafit</code> , as returned by <code>fitGG</code>).
<code>x</code>	<code>ExpressionSet</code> , <code>exprSet</code> , data frame or matrix containing the gene expression measurements used to fit the model.
<code>groups</code>	If <code>x</code> is of type <code>ExpressionSet</code> or <code>exprSet</code> , <code>groups</code> should be the name of the column in <code>pData(x)</code> with the groups that one wishes to compare. If <code>x</code> is a matrix or a data frame, <code>groups</code> should be a vector indicating to which group each column in <code>x</code> corresponds to.
<code>type</code>	<code>data</code> checks marginal density of the data; <code>shape</code> checks shape parameter; <code>mean</code> checks mean parameter; <code>shapemean</code> checks the joint of shape and mean parameters
<code>logexpr</code>	If set to <code>TRUE</code> , the expression values are in log2 scale.
<code>xlab</code>	Passed on to <code>plot</code>
<code>ylab</code>	Passed on to <code>plot</code>
<code>main</code>	Passed on to <code>plot</code>
<code>lty</code>	Ignored.
<code>lwd</code>	Ignored.
<code>...</code>	Other arguments to be passed to <code>plot</code>

Details

The routine generates random draws from the posterior and posterior predictive distributions, fixing the hyper-parameters at their estimated value (posterior mean if model was fit with `method=='Bayes'` or maximum likelihood estimate if model was fit with `method=='EBayes'`).

Value

Produces a plot.

Note

Posterior and posterior predictive checks can lack sensitivity to detect model misfit, since they are susceptible to over-fitting. An alternative is to perform prior predictive checks by generating parameters and data with `simGG`.

Author(s)

David Rossell

References

Rossell D. GaGa: a simple and flexible hierarchical model for microarray data analysis. <http://rosselldavid.googlepages.com>.

See Also

`simGG` to simulate samples from the prior-predictive distribution, `simnewsamples` to generate parameters and observations from the posterior predictive, which is useful to check goodness-of-fit individually a desired gene.

classpred

*Predict the class that a new sample belongs to.***Description**

Computes the posterior probability that a new sample belongs to each group and classifies it into the group with highest probability.

Usage

```
classpred(gg.fit, xnew, x, groups, prgroups, ngene=100)
```

Arguments

gg.fit	GaGa or MiGaGa fit (object of type gagafit, as returned by fitGG).
xnew	Expression levels of the sample to be classified. Only the subset of the genes indicated by ngene is used.
x	ExpressionSet, exprSet, data frame or matrix containing the gene expression measurements used to fit the model.
groups	If x is of type ExpressionSet or exprSet, groups should be the name of the column in pData(x) with the groups that one wishes to compare. If x is a matrix or a data frame, groups should be a vector indicating to which group each column in x corresponds to.
prgroups	Vector specifying prior probabilities for each group. Defaults to equally probable groups.
ngene	Number of genes to use to build the classifier. Genes with smaller probability of being equally expressed are selected first.

Details

The classifier weights each gene according to the posterior probability that it is differentially expressed. Hence, adding genes that are unlikely to be differentially expressed does not affect the performance of the classifier, but it does increase the computational cost. All computations are performed by fixing the hyper-parameters to their estimated value (posterior mean if model was fit with method=='Bayes' or maximum likelihood estimate if model was fit with method=='EBayes').

Value

List with the following elements:

d	Numeric value indicating the group that the new sample is classified into, i.e. where the maximum in posgroups is.
posgroups	Vector giving the posterior probability that the xnew belongs to each of the groups.

Author(s)

David Rossell

References

Rossell D. GaGa: a simple and flexible hierarchical model for microarray data analysis. <http://rosselldavid.googlepages.com>.

See Also

[fitGG](#), [parest](#)

Examples

```
#Not run. Example from the help manual
#library(gaga)
#set.seed(10)
#n <- 100; m <- c(6,6)
#a0 <- 25.5; nu <- 0.109
#balpha <- 1.183; nualpha <- 1683
#probpatt <- c(.95,.05)
#xsim <- simGG(n,m,p.de=probpatt[2],a0,nu,balpha,nualpha)
#
#ggfit <- fitGG(xsim$x[,c(-6,-12)],groups,patterns=patterns,nclust=1)
#ggfit <- parest(ggfit,x=xsim$x[,c(-6,-12)],groups,burnin=100,alpha=.05)
#
#pred1 <- classpred(ggfit,xnew=xsim$x[,6],x=xsim$x[,c(-6,-12)],groups)
#pred2 <- classpred(ggfit,xnew=xsim$x[,12],x=xsim$x[,c(-6,-12)],groups)
#pred1
#pred2
```

dgamma

Approximate gamma shape distribution

Description

`dgamma` approximates density of a gamma shape distribution with a gamma density. `rcgamma` obtains random draws from the approximation. `mgamma` computes approximated mean, variance and normalization constant.

Usage

```
dgamma(x, a, b, c, d, r, s, newton = TRUE)
rcgamma(n, a, b, c, d, r, s, newton = TRUE)
mgamma(a, b, c, d, r, s, newton = TRUE)
```

Arguments

<code>x</code>	Vector indicating the values at which to evaluate the density.
<code>n</code>	Number of random draws to obtain.
<code>a, b, c, d, r, s</code>	Parameter values.
<code>newton</code>	Set to TRUE to try to locate the mode by taking a few Newton-Raphson steps.

Details

The density of a gamma shape distribution is given by $C(a, b, c, d, r, s) \frac{\Gamma(a+x+d)}{\Gamma(x)^a} \left(\frac{x}{r+s+x}\right)^{a+x+d} x^{b-d-1} \exp(-x*c)$ for $x \geq 0$, and 0 otherwise, where $C()$ is the normalization constant. The gamma approximation is $Ga(a/2+b-1/2, c+a \log(s/a))$. The approximate normalization constant is obtained by taking the ratio of the exact density and the approximation at the maximum, as described in Rossell (2007).

Value

`dcgamma` returns a vector with approximate density. `rcgamma` returns a vector with draws from the approximating gamma. `mcgamma` returns a list with components:

<code>m</code>	Approximate mean
<code>v</code>	Approximate variance
<code>normk</code>	Approximate normalization constant

Note

For general values of the parameters the gamma approximation may be poor. In such a case one could use this function to obtain draws from the proposal distribution in a Metropolis-Hastings step.

Author(s)

David Rossell

References

Rossell D. GaGa: a simple and flexible hierarchical model for microarray data analysis. <http://rosselldavid.googlepages.com>.

See Also

[dgamma](#), [rgamma](#)

findgenes

Find differentially expressed genes after GaGa fit.

Description

Obtains a list of differentially expressed genes using the posterior probabilities from a GaGa or MiGaGa fit. For `parametric==TRUE` the procedure controls the Bayesian FDR below `fdrmax`. For `parametric==FALSE` it controls the estimated frequentist FDR.

Usage

```
findgenes(gg.fit, x, groups, fdrmax=.05, parametric=TRUE, B=500)
```

Arguments

<code>gg.fit</code>	GaGa or MiGaGa fit (object of type <code>gagafit</code> , as returned by <code>fitGG</code>).
<code>x</code>	ExpressionSet, <code>exprSet</code> , data frame or matrix containing the gene expression measurements used to fit the model.
<code>groups</code>	If <code>x</code> is of type <code>ExpressionSet</code> or <code>exprSet</code> , <code>groups</code> should be the name of the column in <code>pData(x)</code> with the groups that one wishes to compare. If <code>x</code> is a matrix or a data frame, <code>groups</code> should be a vector indicating to which group each column in <code>x</code> corresponds to.
<code>fdrmax</code>	Upper bound on FDR..
<code>parametric</code>	Set to <code>TRUE</code> to use the Bayes rule. Set to <code>FALSE</code> to estimate the frequentist FDR non-parametrically.
<code>B</code>	Number of bootstrap samples to estimate FDR non-parametrically (ignored if <code>parametric==TRUE</code>)

Details

The Bayes rule to minimize expected FNR subject to $FDR \leq fdrmax$ declares differentially expressed all genes with posterior probability of being equally expressed below a certain threshold. The value of the threshold is computed exactly for `parametric==TRUE`, FDR being defined in a Bayesian sense. For `parametric==FALSE` the FDR is defined in a frequentist sense.

Value

List with components:

<code>truePos</code>	Expected number of true positives.
<code>d</code>	Vector indicating the pattern that each gene is assigned to.
<code>fdr</code>	Frequentist estimated FDR that is closest to <code>fdrmax</code> .
<code>fdrpar</code>	Bayesian FDR. If <code>parametric==TRUE</code> , this is equal to <code>fdrmax</code> . If <code>parametric==FALSE</code> , it's the Bayesian FDR needed to achieve frequentist estimated $FDR=fdrmax$.
<code>fdrest</code>	Data frame with estimated frequentist FDR for each target Bayesian FDR
<code>fnr</code>	Bayesian FNR
<code>power</code>	Bayesian power as estimated by expected number of true positives divided by the expected number of differentially expressed genes
<code>threshold</code>	Optimal threshold for posterior probability of equal expression (genes with probability $< threshold$ are declared DE)

Author(s)

David Rossell

References

Rossell D. GaGa: a simple and flexible hierarchical model for microarray data analysis. <http://rosselldavid.googlepages.com>.

See Also

[fitGG](#), [parest](#)

Examples

```

#Not run. Example from the help manual
#library(gaga)
#set.seed(10)
#n <- 100; m <- c(6,6)
#a0 <- 25.5; nu <- 0.109
#balpha <- 1.183; nualpha <- 1683
#probpatt <- c(.95,.05)
#xsim <- simGG(n,m,p.de=probpatt[2],a0,nu,balphi,nualpha)
#
#ggfit <- fitGG(xsim$x[,c(-6,-12)],groups,patterns=patterns,nclust=1)
#ggfit <- parest(ggfit,x=xsim$x[,c(-6,-12)],groups,burnin=100,alpha=.05)
#
#d <- findgenes(ggfit,xsim$x[,c(-6,-12)],groups,fdrmax=.05,parametric=TRUE)
#dtrue <- (xsim$l[,1]!=xsim$l[,2])
#table(d$d,dtrue)

```

fitGG

*Fit GaGa hierarchical model***Description**

Fits GaGa or MiGaGa hierarchical models, either via a fully Bayesian approach or via maximum likelihood.

Usage

```
fitGG(x, groups, patterns, equalcv = TRUE, nclust = 1, method = "quickEM", B, pr
```

Arguments

<code>x</code>	ExpressionSet, exprSet, data frame or matrix containing the gene expression measurements used to fit the model.
<code>groups</code>	If <code>x</code> is of type ExpressionSet or exprSet, <code>groups</code> should be the name of the column in <code>pData(x)</code> with the groups that one wishes to compare. If <code>x</code> is a matrix or a data frame, <code>groups</code> should be a vector indicating to which group each column in <code>x</code> corresponds to.
<code>patterns</code>	Matrix indicating which groups are put together under each pattern, i.e. the hypotheses to consider for each gene. <code>colnames(patterns)</code> must match the group levels specified in <code>groups</code> . Defaults to two hypotheses: null hypothesis of all groups being equal and full alternative of all groups being different. The function <code>buildPatterns</code> can be used to construct a matrix with all possible patterns.
<code>equalcv</code>	<code>equalcv==TRUE</code> fits model assuming constant CV across groups. <code>equalcv==FALSE</code> compares cv as well as mean expression levels between groups
<code>nclust</code>	Number of clusters in the MiGaGa model. <code>nclust</code> corresponds to the GaGa model.
<code>method</code>	<code>method=='MH'</code> fits a fully Bayesian model via Metropolis-Hastings posterior sampling. <code>method=='Gibbs'</code> does the same using Gibbs sampling. <code>method=='SA'</code> uses Simulated Annealing to find the posterior mode. <code>method=='EM'</code>

	finds maximum-likelihood estimates via the expectation-maximization algorithm, but this is currently only implemented for <code>nclust > 1</code> . <code>method == 'quickEM'</code> is a quicker implementation that only performs 2 optimization steps (see details).
<code>B</code>	Number of iterations. For <code>method == 'MH'</code> and <code>method == 'Gibbs'</code> , <code>B</code> is the number of MCMC iterations (defaults to 1000). For <code>method == 'SA'</code> , <code>B</code> is the number of iterations in the Simulated Annealing scheme (defaults to 200). For <code>method == 'EM'</code> , <code>B</code> is the maximum number of iterations (defaults to 20).
<code>priorpar</code>	List with prior parameter values. It must have components <code>a.alpha0</code> , <code>b.alpha0</code> , <code>a.nu</code> , <code>b.nu</code> , <code>a</code> and <code>p.probpat</code> . If missing they are set to non-informative values that are usually reasonable for RMA and GCRMA normalized data.
<code>parini</code>	list with components <code>a0</code> , <code>nu</code> , <code>balpha</code> , <code>nualpha</code> , <code>probclus</code> and <code>probpat</code> indicating the starting values for the hyper-parameters. If not specified, a method of moments estimate is used.
<code>trace</code>	For <code>trace == TRUE</code> the progress of the model fitting routine is printed.

Details

An approximation is used to sample faster from the posterior distribution of the gamma shape parameters and to compute the normalization constants (needed to evaluate the likelihood). These approximations are implemented in `rcgamma` and `mccgamma`.

The cooling scheme in `method == 'SA'` uses a temperature equal to $1/\log(1+i)$, where `i` is the iteration number.

The EM implementation in `method == 'quickEM'` is a quick EM algorithm that usually delivers hyper-parameter estimates very similar to those obtained via the slower `method == 'EM'`. Additionally, the GaGa model inference has been seen to be robust to moderate changes in the hyper-parameter estimates in most datasets.

Value

An object of class `gagafit`, with components

<code>parest</code>	Hyper-parameter estimates. Only returned if <code>method == 'EBayes'</code> , for <code>method == 'Bayes'</code> one must call the function <code>parest</code> after <code>fitGG</code>
<code>mcmc</code>	Object of class <code>mcmc</code> with posterior draws for hyper-parameters. Only returned if <code>method == 'Bayes'</code> .
<code>lhood</code>	For <code>method == 'Bayes'</code> it is the log-likelihood evaluated at each MCMC iteration. For <code>method == 'EBayes'</code> it is the log-likelihood evaluated at the maximum.
<code>nclust</code>	Same as input argument.
<code>patterns</code>	Same as input argument, converted to object of class <code>gagahyp</code> .

Author(s)

David Rossell

References

Rossell D. GaGa: a simple and flexible hierarchical model for microarray data analysis. <http://rosselldavid.googlepages.com>.

See Also

[parest](#) to estimate hyper-parameters and compute posterior probabilities after a GaGa or MiGaGa fit. [findgenes](#) to find differentially expressed genes. [classpred](#) to predict the group that a new sample belongs to.

Examples

```
library(gaga)
set.seed(10)
n <- 100; m <- c(6,6)
a0 <- 25.5; nu <- 0.109
balpha <- 1.183; nualpha <- 1683
probpatt <- c(.95, .05)
xsim <- simGG(n,m,p.de=probpatt[2],a0,nu,balphi,nualpha,equalcv=TRUE)
x <- exprs(xsim)

#Frequentist fit: EM algorithm to obtain MLE
groups <- pData(xsim)$group[c(-6,-12)]
patterns <- matrix(c(0,0,0,1),2,2)
colnames(patterns) <- c('group 1','group 2')
gg1 <- fitGG(x[,c(-6,-12)],groups,patterns=patterns,method='EM',trace=FALSE)
gg1 <- parest(gg1,x=x[,c(-6,-12)],groups)
gg1
```

forwsimDiffExpr *Forward simulation for differential expression.*

Description

Forward simulation allows to evaluate the expected utility for sequential designs. Here the utility is the expected number of true discoveries minus a sampling cost. The routine simulates future data either from the prior predictive or using a set of pilot data and a GaGa model fit. At each future time point, it computes a summary statistic that will be used to determine when to stop the experiment.

Usage

```
forwsimDiffExpr(gg.fit, x, groups, ngenes, maxBatch, batchSize, fdrmax = 0.05, g
Bsummary = 100, trace = TRUE, randomSeed)
```

Arguments

<code>gg.fit</code>	GaGa or MiGaGa fit (object of type <code>gagafit</code> , as returned by <code>fitGG</code>).
<code>x</code>	<code>ExpressionSet</code> , <code>exprSet</code> , data frame or matrix containing the gene expression measurements used to fit the model.
<code>groups</code>	If <code>x</code> is of type <code>ExpressionSet</code> or <code>exprSet</code> , <code>groups</code> should be the name of the column in <code>pData(x)</code> with the groups that one wishes to compare. If <code>x</code> is a matrix or a data frame, <code>groups</code> should be a vector indicating to which group each column in <code>x</code> corresponds to.
<code>ngenes</code>	Number of genes to simulate data for. If <code>x</code> is specified this argument is set to <code>nrow(x)</code> and data is simulated from the posterior predictive conditional on <code>x</code> . If <code>x</code> not specified simulation is from the prior predictive.

<code>maxBatch</code>	Maximum number of batches, i.e. the routine simulates <code>batchSize*maxBatch</code> samples per group.
<code>batchSize</code>	Batch size, i.e. number of observations per group to simulate at each time point. Defaults to <code>ncol(x)/length(unique(groups))</code> .
<code>fdrmax</code>	Upper bound on FDR.
<code>genelimit</code>	Only the <code>genelimit</code> genes with the lowest probability of being equally expressed across all groups will be simulated. Setting this limit can significantly increase the computational speed.
<code>v0thre</code>	Only genes with posterior probability of being equally expressed $< v0thre$ will be simulated. Setting this limit can significantly increase the computational speed.
<code>B</code>	Number of forward simulations.
<code>Bsummary</code>	Number of simulations for estimating the summary statistic.
<code>trace</code>	For <code>trace==TRUE</code> iteration progress is displayed.
<code>randomSeed</code>	Integer value used to set random number generator seed. Defaults to <code>as.numeric(Sys.time())</code> modulus 10^6 .

Details

To improve computational speed hyper-parameters are not re-estimated as new data is simulated.

Value

A `data.frame` with the following columns:

<code>simid</code>	Simulation number.
<code>j</code>	Time (sample size).
<code>u</code>	Expected number of true positives if we were to stop experimentation at this time.
<code>fdr</code>	Expected FDR if we were to stop experimentation at this time.
<code>fnr</code>	Expected FNR if we were to stop experimentation at this time.
<code>power</code>	Expected power (as estimated by $E(TP)/E(\text{positives})$) if we were to stop experimentation at this time.
<code>summary</code>	Summary statistic: increase in expected true positives if we were to obtain one more data batch.

Author(s)

David Rossell.

References

- Rossell D., Mueller P. Sequential sample sizes for high-throughput hypothesis testing experiments. <http://sites.google.com/site/rosselldavid/home>.
- Rossell D. GaGa: a simple and flexible hierarchical model for microarray data analysis. *Annals of Applied Statistics*, 2009, 3, 1035-1051.

See Also

[fitGG](#) for fitting a GaGa model, [seqBoundariesGrid](#) for finding the optimal design based on the forwards simulation output. [powfindgenes](#) for fixed sample size calculations.

Examples

```
#Simulate data and fit GaGa model
set.seed(1)
x <- simGG(n=20,m=2,p.de=.5,a0=3,nu=.5,balpa=.5,nualpha=25)
gg1 <- fitGG(x,groups=1:2,method='EM')
gg1 <- parest(gg1,x=x,groups=1:2)

#Run forward simulation
fs1 <- forwsimDiffExpr(gg1, x=x, groups=1:2,
maxBatch=2,batchSize=1,fdrmax=0.05, B=100, Bsummary=100, randomSeed=1)

#Expected number of true positives for each sample size
tapply(fs1$u,fs1$time,'mean')

#Expected utility for each sample size
samplingCost <- 0.01
tapply(fs1$u,fs1$time,'mean') - samplingCost*(0:2)

#Optimal sequential design
b0seq <- seq(0,20,length=200); b1seq <- seq(0,40,length=200)
bopt <- seqBoundariesGrid(b0=b0seq,b1=b1seq,forwsim=fs1,samplingCost=samplingCost,powmin=0)
bopt <- bopt$opt

plot(fs1$time,fs1$u,xlab='Additional batches',ylab='E(newly discovered DE genes)')
abline(bopt['b0'],bopt['b1'])
text(.2,bopt['b0'],'Continue',pos=3)
text(.2,bopt['b0'],'Stop',pos=1)
```

geneclus

Cluster genes into expression patterns.

Description

Performs supervised gene clustering. Clusters genes into the expression pattern with highest posterior probability, according to a GaGa or MiGaGa fit.

Usage

```
geneclus(gg.fit, method='posprob')
```

Arguments

<code>gg.fit</code>	GaGa or MiGaGa fit (object of type <code>gagafit</code> , as returned by <code>fitGG</code>).
<code>method</code>	For <code>method==1</code> samples are assigned to pattern with highest posterior probability, and for <code>method==2</code> to the pattern with highest likelihood (e.g. assuming equal a priori prob for all patterns)

Details

Each gene is assigned to the pattern with highest posterior probability. This is similar to routine `findgenes`, which also assigns genes to the pattern with highest posterior probability, although `findgenes` applies an FDR-based correction i.e. tends to assign more genes to the null pattern of no differential expression.

Value

List with components:

`d` Vector indicating the pattern that each gene is assigned to.
`posprob` Vector with posterior probabilities of the assigned patterns.

Author(s)

David Rossell

References

Rossell D. GaGa: a simple and flexible hierarchical model for microarray data analysis. <http://rosselldavid.googlepages.com>.

See Also

`fitGG`, `parest`

Examples

```
#Not run. Example from the help manual
#library(gaga)
#set.seed(10)
#n <- 100; m <- c(6,6)
#a0 <- 25.5; nu <- 0.109
#balpha <- 1.183; nualpha <- 1683
#probpat <- c(.95,.05)
#xsim <- simGG(n,m,p.de=probpat[2],a0,nu,balphi,nualpha)
#
#ggfit <- fitGG(xsim$x[,c(-6,-12)],groups,patterns=patterns,nclust=1)
#ggfit <- parest(ggfit,x=xsim$x[,c(-6,-12)],groups,burnin=100,alpha=.05)
#
#dclus <- geneclus(ggfit) #not use FDR correction
#dfdr <- findgenes(ggfit,xsim$x[,c(-6,-12)],groups,fdrmax=.05,parametric=TRUE) #use FDR c
#table(dfdr$d,dclus$d) #compare results
```

getpar

Extract hyper-parameter estimates from a gagafit object

Description

Extracts the hyper-parameter estimates from a `gagafit` object and puts them in a list.

Usage

```
getpar(gg.fit)
```

Arguments

`gg.fit` Object of class `gagafit`, as returned by `parest`.

Details

This routine simply evaluates the component `gg.fit$parest` from a `gagafit` object, which causes an error if this component is not available. This routine is used internally by a number of other routines.

Value

A list with components:

<code>a0</code>	Estimated value of hyper-parameter <code>a0</code>
<code>nu</code>	Estimated value of hyper-parameter <code>nu</code>
<code>balpha</code>	Estimated value of hyper-parameter <code>balpha</code>
<code>nualpha</code>	Estimated value of hyper-parameter <code>nualpha</code>
<code>probclus</code>	Estimated cluster probabilities
<code>probpatt</code>	Estimated prior probability of each expression pattern

Author(s)

David Rossell

References

Rossell D. GaGa: a simple and flexible hierarchical model for microarray data analysis. <http://rosselldavid.googlepages.com>.

See Also

[fitGG](#), [parest](#)

`parest`

Parameter estimates and posterior probabilities of differential expression for GaGa and MiGaGa model

Description

Obtains parameter estimates and posterior probabilities of differential expression after a GaGa or MiGaGa model has been fit with the function `fitGG`.

Usage

```
parest(gg.fit, x, groups, burnin, alpha=.05)
```

Arguments

<code>gg.fit</code>	GaGa or MiGaGa fit (object of type <code>gagafit</code> , as returned by <code>fitGG</code>).
<code>x</code>	ExpressionSet, <code>exprSet</code> , data frame or matrix containing the gene expression measurements used to fit the model.
<code>groups</code>	If <code>x</code> is of type ExpressionSet or <code>exprSet</code> , <code>groups</code> should be the name of the column in <code>pData(x)</code> with the groups that one wishes to compare. If <code>x</code> is a matrix or a data frame, <code>groups</code> should be a vector indicating to which group each column in <code>x</code> corresponds to.
<code>burnin</code>	Number of MCMC samples to discard. Ignored if <code>gg.fit</code> was fit with the option <code>method=='EBayes'</code> .
<code>alpha</code>	If <code>gg.fit</code> was fit with the option <code>method=='Bayes'</code> , <code>parest</code> also computes $1-\alpha$ posterior credibility intervals.

Details

If `gg.fit` was fit via MCMC posterior sampling (option `method=='Bayes'`), `parest` discards the first `burnin` iterations and uses the rest to obtain point estimates and credibility intervals for the hyper-parameters. To compute posterior probabilities of differential expression the hyper-parameters are fixed to their estimated value, i.e. not averaged over MCMC iterations.

Value

An object of class `gagafit`, with components:

<code>parest</code>	Hyper-parameter estimates.
<code>mcmc</code>	Object of class <code>mcmc</code> with posterior draws for hyper-parameters. Only returned if <code>method=='Bayes'</code> .
<code>lhood</code>	For <code>method=='Bayes'</code> it is the posterior mean of the log-likelihood. For <code>method=='EBayes'</code> it is the log-likelihood evaluated at the maximum.
<code>nclust</code>	Number of clusters.
<code>patterns</code>	Object of class <code>gagahyp</code> indicating which hypotheses (expression patterns) were tested.
<code>pp</code>	Matrix with posterior probabilities of differential expression for each gene. Genes are in rows and expression patterns are in columns (e.g. for 2 hypotheses, 1st column is the probability of the null hypothesis and 2nd column for the alternative).

Author(s)

David Rossell

References

Rossell D. GaGa: a simple and flexible hierarchical model for microarray data analysis. <http://rosselldavid.googlepages.com>.

See Also

`fitGG` to fit a GaGa or MiGaGa model, `findgenes` to find differentially expressed genes and `posmeansGG` to obtain posterior expected expression values. `classpred` performs class prediction.

Examples

```
#Not run
#library(EBarrays); data(gould)
#x <- log(exprs(gould)[,-1]) #exclude 1st array
#groups <- pData(gould)[-1,1]
#patterns <- rbind(rep(0,3),c(0,0,1),c(0,1,1),0:2) #4 hypothesis
#gg <- fitGG(x,groups,patterns,method='EBayes')
#gg
#gg <- parest(gg,x,groups)
#gg
```

posmeansGG	<i>Gene-specific posterior means</i>
------------	--------------------------------------

Description

Computes posterior means for the gene expression levels using a GaGa or MiGaGa model.

Usage

```
posmeansGG(gg.fit, x, groups, sel, underpattern)
```

Arguments

gg.fit	GaGa or MiGaGa fit (object of type gagafit, as returned by fitGG).
x	ExpressionSet, exprSet, data frame or matrix containing the gene expression measurements used to fit the model.
groups	If x is of type ExpressionSet or exprSet, groups should be the name of the column in pData(x) with the groups that one wishes to compare. If x is a matrix or a data frame, groups should be a vector indicating to which group each column in x corresponds to.
sel	Numeric vector with the indexes of the genes we want to draw new samples for (defaults to all genes). If a logical vector is indicated, it is converted to (1:nrow(x))[sel].
underpattern	Expression pattern assumed to be true (defaults to last pattern in gg.fit\$patterns). Posterior means are computed under this pattern. For example, if only the null pattern that all groups are equal and the full alternative that all groups are different are considered, underpattern=1 returns the posterior means under the assumption that groups are different from each other (underpattern=0 returns the same mean for all groups).

Details

The posterior distribution of the mean parameters actually depends on the gene-specific shape parameter(s), which is unknown. To speed up computations, a gamma approximation to the shape parameter posterior is used (see `rcgamma` for details) and the shape parameter is fixed to its mode a posteriori.

Value

Matrix with mean expression values a posteriori, for each selected gene and each group. Genes are in rows and groups in columns.

Author(s)

David Rossell

References

Rossell D. GaGa: a simple and flexible hierarchical model for microarray data analysis. <http://rosselldavid.googlepages.com>.

See Also

[fitGG](#) for fitting GaGa and MiGaGa models, [parest](#) for computing posterior probabilities of each expression pattern.

powclasspred

Expected probability that a future sample is correctly classified.

Description

Estimates posterior expected probability that a future sample is correctly classified when performing class prediction. The estimate is obtained via Monte Carlo simulation from the posterior predictive.

Usage

```
powclasspred(gg.fit, x, groups, prgroups, v0thre=1, ngene=100, B=100)
```

Arguments

<code>gg.fit</code>	GaGa or MiGaGa fit (object of type <code>gagafit</code> , as returned by <code>fitGG</code>).
<code>x</code>	<code>ExpressionSet</code> , <code>exprSet</code> , data frame or matrix containing the gene expression measurements used to fit the model.
<code>groups</code>	If <code>x</code> is of type <code>ExpressionSet</code> or <code>exprSet</code> , <code>groups</code> should be the name of the column in <code>pData(x)</code> with the groups that one wishes to compare. If <code>x</code> is a matrix or a data frame, <code>groups</code> should be a vector indicating to which group each column in <code>x</code> corresponds to.
<code>prgroups</code>	Vector specifying prior probabilities for each group. Defaults to equally probable groups.
<code>v0thre</code>	Only genes with posterior probability of being equally expressed below <code>v0thre</code> are used.
<code>ngene</code>	Number of genes to use to build the classifier. Genes with smaller probability of being equally expressed are selected first.
<code>B</code>	Number of Monte Carlo samples to be used.

Details

The routine simulates future samples (microarrays) from the posterior predictive distribution of a given group (e.g. control/cancer). Then it computes the posterior probability that the new sample belongs to each of the groups and classifies the sample into the group with highest probability. This process is repeated B times, and the proportion of correctly classified samples is reported for each group. The standard error is obtained via the usual normal approximation (i.e. SD/B). The overall probability of correct classification is also provided (i.e. for all groups together), but using a more efficient variant of the algorithm. Instead of reporting the observed proportion of correctly classified samples, it reports the expected proportion of correctly classified samples (i.e. the average posterior probability of the class that the sample is assigned to).

Value

List with components:

<code>ccall</code>	Estimated expected probability of correctly classifying a future sample.
<code>seccall</code>	Estimated standard error of <code>ccall</code> .
<code>ccgroup</code>	Vector with the estimated probability of correctly classifying a sample from each group.
<code>seggroup</code>	Estimated standard error of <code>ccgroup</code> .

Author(s)

David Rossell

References

Rossell D. GaGa: a simple and flexible hierarchical model for microarray data analysis. <http://rosselldavid.googlepages.com>.

See Also

`classpred`, `fitGG`, `parest`. See `powfindgenes` for differential expression power calculations.

`powfindgenes`

Power computations for differential expression

Description

`powfindgenes` evaluates the posterior expected number of true positives (e.g. true gene discoveries) if one were to obtain an additional batch of data. It uses a gaga model fit based on a set of pilot data.

Usage

```
powfindgenes(gg.fit, x, groups, batchSize = 1, fdmax = 0.05, genelimit, v0thre
```

Arguments

<code>gg.fit</code>	GaGa or MiGaGa fit (object of type <code>gagafit</code> , as returned by <code>fitGG</code>).
<code>x</code>	ExpressionSet, <code>exprSet</code> , data frame or matrix containing the gene expression measurements used to fit the model.
<code>groups</code>	If <code>x</code> is of type <code>ExpressionSet</code> or <code>exprSet</code> , <code>groups</code> should be the name of the column in <code>pData(x)</code> with the groups that one wishes to compare. If <code>x</code> is a matrix or a data frame, <code>groups</code> should be a vector indicating to which group each column in <code>x</code> corresponds to.
<code>batchSize</code>	Number of additional samples to obtain per group.
<code>fdrmax</code>	Upper bound on FDR..
<code>genelimit</code>	Only the <code>genelimit</code> genes with the lowest probability of being equally expressed across all groups will be simulated. Setting this limit can significantly increase the computational speed.
<code>v0thre</code>	Only genes with posterior probability of being equally expressed $< v0thre$ will be simulated. Setting this limit can significantly increase the computational speed.
<code>B</code>	Number of simulations from the GaGa predictive distribution to be used to estimate the posterior expected number of true positives.

Details

The routine simulates data from the posterior predictive distribution of a GaGa model. That is, first it simulates parameter values (differential expression status, mean expression levels etc.) from the posterior distribution. Then it simulates data using Gamma distributions and the parameter values drawn from the posterior. Finally the simulated data is used to determine the differential status of each gene, controlling the Bayesian FDR at the `fdrmax` level, as implemented in `findgenes`. As the differential expression status is known for each gene, one can evaluate the number of true discoveries in the reported gene list.

Value

<code>m</code>	Posterior expected number of true positives (as estimated by the sample mean of <code>B</code> simulations)
<code>s</code>	Standard error of the estimate i.e. $SD \text{ of the simulations} / \sqrt{B}$

Author(s)

David Rossell

References

Rossell D. GaGa: a simple and flexible hierarchical model for microarray data analysis. <http://rosselldavid.googlepages.com>.

See Also

`findgenes`, `fitGG`, `parest`. See `powclasspred` for power calculations for sample classification.

Examples

```
#Simulate data and fit GaGa model
set.seed(1)
x <- simGG(n=20,m=2,p.de=.5,a0=3,nu=.5,balpha=.5,nalpha=25)
gg1 <- fitGG(x,groups=1:2,method='EM')
gg1 <- parest(gg1,x=x,groups=1:2)

#Expected nb of TP for 1 more sample per group
powfindgenes(gg1,x=x,groups=1:2,batchSize=1,fdrmax=.05)$m

#Expected nb of TP for 10 more samples per group
powfindgenes(gg1,x=x,groups=1:2,batchSize=10,fdrmax=.05)$m
```

print.gagaclus *Print an object of class gagaclus*

Description

Prints an object of class `gagaclus`, which contains the result of clustering genes into expression patterns.

Usage

```
## S3 method for class 'gagaclus'
print(x, ...)
```

Arguments

`x` Object of type `gagaclus`.
`...` Other arguments to be passed on to the generic print function.

Value

Displays the expression patterns and the number of genes classified into each of them.

Author(s)

David Rossell

References

Rossell D. GaGa: a simple and flexible hierarchical model for microarray data analysis. <http://rosselldavid.googlepages.com>.

See Also

[fitGG](#), [geneclus](#)

print.gagafit	<i>Print an object of class gagafit</i>
---------------	---

Description

Prints an object of class `gagafit`, as returned by `fitGG` or `parest`. Provides general information and hyper-parameter estimates, if available.

Usage

```
## S3 method for class 'gagafit'  
print(x, ...)
```

Arguments

<code>x</code>	Object of type <code>gagafit</code> , as returned by <code>fitGG</code> or <code>parest</code> .
<code>...</code>	Other arguments to be passed on to the generic <code>print</code> function.

Details

`fitGG` does not create a complete `gagafit` object. The complete object is returned by `parest`, which computes the posterior probabilities of differential expression and obtain hyper-parameter estimates (these are only provided by `fitGG` when the option `method='EBayes'` is used).

Value

Prints number of genes, hypotheses, details about the model fitting and hyper-parameter estimates (when available).

Author(s)

David Rossell

References

Rossell D. GaGa: a simple and flexible hierarchical model for microarray data analysis. <http://rosselldavid.googlepages.com>.

See Also

[fitGG](#), [parest](#)

```
print.gagahyp
```

Print an object of class gagahyp

Description

Prints an object of class `gagahyp`, which contains information on the hypotheses (expression patterns) from a GaGa or MiGaGa model.

Usage

```
## S3 method for class 'gagahyp'
print(x, probpat=NA, ...)
```

Arguments

<code>x</code>	Object of type <code>gagahyp</code> .
<code>probp</code>	Vector with either estimated probabilities of each hypothesis, or with number of genes classified into each expression pattern.
<code>...</code>	Other arguments to be passed on to the generic print function.

Value

Prints hypotheses. When available, also displays estimated proportion of genes following each expression pattern or the number of genes classified into each expression pattern.

Author(s)

David Rossell

References

Rossell D. GaGa: a simple and flexible hierarchical model for microarray data analysis. <http://rosselldavid.googlepages.com>.

See Also

[fitGG](#), [geneclus](#)

```
seqBoundariesGrid
```

Evaluate expected utility for parametric sequential stopping boundaries.

Description

Estimate the expected utility for sequential boundaries parameterized by (b_0, b_1) . Expected utility is estimated on a grid of (b_0, b_1) values based on a forward simulation output such as that generated by the function `forwsimDiffExpr`.

Usage

```
seqBoundariesGrid(b0, b1, forwsim, samplingCost, powmin = 0, f = "linear", ineq
```

Arguments

<code>b0</code>	Vector with <code>b0</code> values. Expected utility is evaluated for a grid defined by all combinations of (<code>b0</code> , <code>b1</code>) values.
<code>b1</code>	Vector with <code>b1</code> values.
<code>forwsim</code>	<code>data.frame</code> with forward simulation output, such as that returned by the function <code>forwsimDiffExpr</code> . It must have columns named <code>simid</code> , <code>time</code> , <code>u</code> , <code>fdr</code> , <code>fnr</code> , <code>power</code> and <code>summary</code> . See <code>forwsimDiffExpr</code> for details on the meaning of each column.
<code>samplingCost</code>	Cost of obtaining one more data batch, in terms of the number of new truly differentially expressed discoveries that would make it worthwhile to obtain one more data batch.
<code>powmin</code>	Constraint on power. Optimization chooses the optimal <code>b0</code> , <code>b1</code> satisfying <code>power >= powermin</code> (if such <code>b0</code> , <code>b1</code> exists).
<code>f</code>	Parametric form for the stopping boundary. Currently only 'linear' and 'invsqrt' are implemented. For 'linear', the boundary is $b_0 + b_1 * \text{time}$. For 'invsqrt', the boundary is $b_0 + b_1 / \sqrt{\text{time}}$, where <code>time</code> is the sample size measured as number of batches.
<code>ineq</code>	For <code>ineq == 'less'</code> the trial stops when <code>summary</code> is below the stopping boundary. This is appropriate whenever <code>summary</code> measures the potential benefit of obtaining one more data batch. For <code>ineq == 'greater'</code> the trial stops when <code>summary</code> is above the stopping boundary. This is appropriate whenever <code>summary</code> measures the potential costs of obtaining one more data batch.

Details

Intuitively, the goal is to stop collecting new data when the expected benefit of obtaining one more data batch is small, i.e. below a certain boundary. We consider two simple parametric forms for such a boundary (linear and inverse square root), which allows to easily evaluate the expected utility for each boundary within a grid of parameter values. The optimal boundary is defined by the parameter values achieving the largest expected utility, restricted to parameter values with an estimated power greater or equal than `powmin`. Here power is defined as the expected number of true discoveries divided by the expected number of differentially expressed entities.

The routine evaluates the expected utility, as well as expected FDR, FNR, power and sample size for each specified boundary, and also reports the optimal boundary.

Value

A list with two components:

<code>opt</code>	Vector with optimal stopping boundary (<code>b</code>), estimated expected utility (<code>u</code>), false discovery rate (<code>fdr</code>), false negative rate (<code>fnr</code>), power (<code>power</code>) and the expected sample size measured as the number of batches (<code>time</code>).
<code>grid</code>	<code>data.frame</code> with all evaluated boundaries (columns <code>b0</code> and <code>b1</code>) and their respective estimated expected utility, false discovery rate, false negative rate, power and expected sample size (measured as the number of batches).

Author(s)

David Rossell.

References

Rossell D., Mueller P. Sequential sample sizes for high-throughput hypothesis testing experiments. <http://sites.google.com/site/rosselldavid/home>.

Rossell D. GaGa: a simple and flexible hierarchical model for microarray data analysis. *Annals of Applied Statistics*, 2009, 3, 1035-1051.

See Also

[forwsimDiffExpr](#)

 simGG

Prior predictive simulation

Description

Simulates parameters and data from the prior-predictive of GaGa or MiGaGa model with several groups, fixing the hyper-parameters.

Usage

```
simGG(n, m, p.de=.1, a0, nu, balpha, nualpha, equalcv = TRUE, probclus
= 1, a = NA, l = NA, useal = FALSE)
```

Arguments

n	Number of genes.
m	Vector indicating number of observations to be simulated for each group.
p.de	Probability that a gene is differentially expressed.
a0, nu	Mean expression for each gene is generated from $1/\text{rgamma}(a0, a0/\text{nu})$ if <code>probclus</code> is of length 1, and from a mixture if <code>length(probclus) > 1</code> .
balpha, nualpha	Shape parameter for each gene is generated from $\text{rgamma}(\text{balpha}, \text{balpha}/\text{nualpha})$.
equalcv	If <code>equalcv==TRUE</code> the shape parameter is simulated to be constant across groups.
probclus	Vector with the probability of each component in the mixture. Set to 1 for the GaGa model.
a, l	Optionally, if <code>useal==TRUE</code> the parameter values are not generated, only the data is generated. <code>a</code> is a matrix with the shape parameters of each gene and group and <code>l</code> is a matrix with the mean expressions.
useal	For <code>useal==TRUE</code> the parameter values specified in <code>a</code> and <code>l</code> are used, instead of being generated.

Details

The shape parameters are actually drawn from a gamma approximation to their posterior distribution. The function `rcgamma` implements this approximation.

Value

Object of class 'ExpressionSet'. Expression values can be accessed via `exprs(object)` and the parameter values used to generate the expression values can be accessed via `fData(object)`.

Note

Currently, the routine only implements prior predictive simulation for the 2 hypothesis case.

Author(s)

David Rossell

References

Rossell D. GaGa: a simple and flexible hierarchical model for microarray data analysis. <http://rosselldavid.googlepages.com>.

See Also

`simnewsamples` to simulate from the posterior predictive, `checkfit` for graphical posterior predictive checks.

Examples

```
#Not run. Example from the help manual
#library(gaga)
#set.seed(10)
#n <- 100; m <- c(6,6)
#a0 <- 25.5; nu <- 0.109
#balpha <- 1.183; nualpha <- 1683
#probpatt <- c(.95,.05)
#xsim <- simGG(n,m,p.de=probpatt[2],a0,nu,balpha,nualpha)
#
#plot(density(xsim$x),main='')
#plot(xsim$l,xsim$a,ylab='Shape',xlab='Mean')
```

simnewsamples

Posterior predictive simulation

Description

Simulates parameters and data from the posterior and posterior predictive distributions, respectively, of a GaGa or MiGaGa model.

Usage

```
simnewsamples(gg.fit, groupsnew, sel, x, groups)
```

Arguments

<code>gg.fit</code>	GaGa or MiGaGa fit (object of type <code>gagafit</code> , as returned by <code>fitGG</code>).
<code>groupsnew</code>	Vector indicating the group that each new sample should belong to. <code>length(groupsnew)</code> is the number of new samples that will be generated.
<code>sel</code>	Numeric vector with the indexes of the genes we want to draw new samples for (defaults to all genes). If a logical vector is indicated, it is converted to <code>(1:nrow(x))[sel]</code> .
<code>x</code>	ExpressionSet, <code>exprSet</code> , data frame or matrix containing the gene expression measurements used to fit the model.
<code>groups</code>	If <code>x</code> is of type ExpressionSet or <code>exprSet</code> , <code>groups</code> should be the name of the column in <code>pData(x)</code> with the groups that one wishes to compare. If <code>x</code> is a matrix or a data frame, <code>groups</code> should be a vector indicating to which group each column in <code>x</code> corresponds to.

Details

The shape parameters are actually drawn from a gamma approximation to their posterior distribution. The function `rcgamma` implements this approximation.

Value

Object of class 'ExpressionSet'. Expression values can be accessed via `exprs(object)` and the parameter values used to generate the expression values can be accessed via `fData(object)`.

Author(s)

David Rossell

References

Rossell D. GaGa: a simple and flexible hierarchical model for microarray data analysis. <http://rosselldavid.googlepages.com>.

See Also

`checkfit` for posterior predictive plot, `simGG` for prior predictive simulation.

Index

*Topic **design**

forwsimDiffExpr, 9
seqBoundariesGrid, 21

*Topic **distribution**

checkfit, 1
dgamma, 4
posmeansGG, 15
simGG, 23
simnewsamples, 24

*Topic **htest**

classpred, 3
findgenes, 5
forwsimDiffExpr, 9
geneclus, 11
powclasspred, 16
powfindgenes, 17
seqBoundariesGrid, 21

*Topic **logic**

buildPatterns, 1

*Topic **models**

checkfit, 1
classpred, 3
findgenes, 5
fitGG, 7
geneclus, 11
getpar, 12
parest, 13
posmeansGG, 15
powclasspred, 16
powfindgenes, 17
simGG, 23
simnewsamples, 24

*Topic **print**

print.gagaclus, 19
print.gagafit, 20
print.gagahyp, 21

buildPatterns, 1

checkfit, 1, 24, 25
classpred, 3, 9, 14, 17

dgamma, 4
dgamma, 5

findgenes, 5, 9, 14, 18
fitGG, 4, 6, 7, 11–14, 16–21
forwsimDiffExpr, 9, 23

geneclus, 11, 19, 21
getpar, 12

mcgamma (*dgamma*), 4

parest, 4, 6, 9, 12, 13, 13, 16–18, 20
posmeansGG, 14, 15
powclasspred, 16, 18
powfindgenes, 11, 17, 17
print.gagaclus, 19
print.gagafit, 20
print.gagahyp, 21

rcgamma (*dgamma*), 4
rgamma, 5

seqBoundariesGrid, 11, 21
simGG, 2, 23, 25
simnewsamples, 2, 24, 24