

crlmm

October 25, 2011

AssayData-methods *Methods for class "AssayData" in crlmm*

Description

The `batchStatistics` slot in a `CNSet` object is an instance of the `AssayData` slot. In general, the accessors for `AssayData` are called indirectly by the corresponding method for the `CNSet` class and not called directly by the user.

Methods

Ns signature(object="AssayData"): ...
corr signature(object="AssayData"): ...
mads signature(x="AssayData"): ...
medians signature(object="AssayData"): ...
tau2 signature(object="AssayData"): ...

See Also

[CNSet-class](#), [Ns](#), [tau2](#), [corr](#), [mads](#), [medians](#)

CNSet-methods *crlmm methods for class "CNSet"*

Description

`CNSet` is a container defined in the `oligoClasses` package for storing normalized intensities for genotyping platforms, genotype calls, and parameters estimated for copy number. Accessors for data that an object of this class contains are largely defined in the package `oligoClasses`. `CNSet` methods that involve more complex calculations that are specific to the `crlmm` package, such as computing allele-specific copy number, are included in `crlmm` and described here.

Methods

CA signature(object="CNSet"): ...
CB signature(object="CNSet"): ...
lines signature(x="CNSet"): ...
totalCopynumber signature(object="CNSet"): ...
rawCopynumber signature(object="CNSet"): ...
nuA signature(object="CNSet"): ...
nuB signature(object="CNSet"): ...
phiA signature(object="CNSet"): ...
phiB signature(object="CNSet"): ...
Ns signature(object="CNSet"): ...
corr signature(object="CNSet"): ...
mads signature(x="CNSet"): ...
medians signature(object="CNSet"): ...
tau2 signature(object="CNSet"): ...

See Also

[CNSet-class](#), [CA](#), [CB](#), [totalCopynumber](#), [rawCopynumber](#)

batchStatisticAccessors

Accessors for batch-specific summary statistics.

Description

The summary statistics stored here are used by the tools for copy number estimation.

Usage

```
corr(object, ...)
tau2(object, ...)
mads(object, ...)
medians(object, ...)
Ns(object, ...)
```

Arguments

object An object of class CNSet.
... An additional argument named 'i' can be passed to subset the markers and an argument 'j' can be passed to subset the batches. Other arguments are ignored.

Value

An array with dimension $R \times A \times G \times C$, or $R \times G \times C$.

R: number of markers **A**: number of alleles (2) **G**: number of biallelic genotypes (3) **C**: number of batches

`Ns` returns an array of genotype frequencies stratified by batch. Dimension $R \times G \times C$.

`corr` returns an array of within-genotype correlations (log2-scale) stratified by batch. Dimension $R \times G \times C$.

`medians` returns an array of the within-genotype medians (intensity-scale) stratified by batch and allele. Dimension $R \times A \times G \times C$.

`mads` returns an array of the within-genotype median absolute deviations (intensity-scale) stratified by batch and allele. Dimension is the same as for `medians`.

`tau2` returns an array of the squared within-genotype median absolute deviation on the log-scale. Only the `mads` for AA and BB genotypes are stored. Dimension is $R \times A \times G \times C$, where **G** is AA or BB. Note that the mad for allele A/B for subjects with genotype BB/AA is a robust estimate of the background variance, whereas the the mad for allele A/B for subjects with genotype AA/BB is a robust estimate of the variance for copy number greater than 0 (we assume that on the log-scale the variance is roughly constant for CA, CB > 0).

See Also

[batchStatistics](#)

Examples

```
data(sample.CNSet)
## All NAs. Need to replace sample.CNSet with a HapMap example
Ns(cnSet, i=1:5, j=1:2)
corr(cnSet, i=1:5, j=1:2)
medians(cnSet, i=1:5, j=1:2)
mads(cnSet, i=1:5, j=1:2)
tau2(cnSet, i=1:5, j=1:2)
```

celDates

Extract dates from the cel file header

Description

Extract dates from the cel file header.

Usage

```
celDates(celfiles)
```

Arguments

`celfiles` CEL file names. Must specify the complete path.

Value

date-time class POSIXt

Author(s)

R. Scharpf

See Also[read.celfile.header](#), [POSIXt](#)

`constructInf`*Instantiate an object of class CNSet for the Infinium platforms.*

Description

Instantiates an object of class CNSet for the Infinium platforms. Elements of `assayData` and `batchStatistics` will be `ff` objects. See details.

Usage

```
constructInf(sampleSheet = NULL, arrayNames = NULL, path = ".", arrayInfoColName
```

Arguments

- `sampleSheet` `data.frame` containing Illumina sample sheet information (for required columns, refer to BeadStudio Genotyping guide - Appendix A).
- `arrayNames` character vector containing names of arrays to be read in. If `NULL`, all arrays that can be found in the specified working directory will be read in.
- `path` character string specifying the location of files to be read by the function
- `arrayInfoColNames` (used when `sampleSheet` is specified) list containing elements 'barcode' which indicates column names in the `sampleSheet` which contains the arrayNumber/barcode number and 'position' which indicates the strip number. In older style sample sheets, this information is combined (usually in a column named 'SentryPosition') and this should be specified as `list(barcode=NULL, position="SentryPosition")`
- `highDensity` logical (used when `sampleSheet` is specified). If `TRUE`, array extensions '_A', '_B' in `sampleSheet` are replaced with 'R01C01', 'R01C02' etc.
- `sep` character string specifying separator used in `.idat` file names.
- `fileExt` list containing elements 'Green' and 'Red' which specify the `.idat` file extension for the Cy3 and Cy5 channels.
- `cdfName` annotation package (see also `validCdfNames`)
- `verbose` 'logical.' Whether to print descriptive messages during processing.
- `batch` batch variable. See details.
- `saveDate` 'logical'. Should the dates from each `.idat` be saved with sample information?

Details

This function initializes a container for storing the normalized intensities for the A and B alleles at polymorphic loci and the normalized intensities for the 'A' allele at nonpolymorphic loci. CRLMM genotype calls and confidence scores are also stored in the `assayData`. This function does not do any preprocessing or genotyping – it only creates an object of the appropriate size. The initialized values will all be 'NA'.

The `ff` package provides infrastructure for accessing and writing data to disk instead of keeping data in memory. Each element of the `assayData` and `batchStatistics` slot are `ff` objects. `ff` objects in the R workspace contain pointers to several files with the '.ff' extension on disk. The location of where the data is stored on disk can be specified by use of the `ldPath` function. Users should not move or rename this directory. If only output files are stored in `ldPath`, one can either remove the entire directory prior to rerunning the analysis or all of the '.ff' files. Otherwise, one would accumulate a large number of '.ff' files on disk that are no longer in use.

We have adopted the `ff` package in order to reduce `crlmm`'s memory footprint. The memory usage can be fine-tuned by the utilities `ocSamples` and `ocProbesets` provided in the `oligoClasses` package. In most instances, the user-level interface will be no different than accessing data from ordinary matrices in R. However, the differences in the underlying representation can become more noticeable for very large datasets in which the I/O for accessing data from the disk can be substantial.

Value

A `CNSet` object

Author(s)

R. Scharpf

See Also

[ldPath](#), [ocSamples](#), [ocProbesets](#), [CNSet-class](#), [preprocessInf](#), [genotypeInf](#)

Examples

```
## See the illumina_copynumber.Rnw vignette in inst/scripts of the
## source package for an example
```

copynumberAccessors

Accessors for allele-specific or total copy number

Description

These methods can be applied after an object of class `CNSet` has been generated by the `crlmmCopynumber` function.

Usage

```
CA(object, ...)
CB(object, ...)
nuA(object)
nuB(object)
phiA(object)
phiB(object)
totalCopynumber(object, ...)
rawCopynumber(object, ...)
```

Arguments

<code>object</code>	An object of class <code>CNSet</code> .
<code>...</code>	An additional argument named 'i' can be passed to subset the markers and an argument 'j' can be passed to subset the samples. Other arguments are ignored.

Details

At polymorphic markers, `nuA` and `nuB` provide the intercept coefficient (the estimated background intensity) for the A and B alleles, respectively. `phiA` and `phiB` provide the slope coefficients for the A and B alleles, respectively.

At nonpolymorphic markers, `nuB` and `phiB` are 'NA'.

These functions can be used to translate the normalized intensities to the copy number scale. Plotting the copy number estimates as a function of physical position can be used to guide downstream algorithms that smooth, as well as to assess possible mosaicism.

Value

`nu[A/B]` and `phi[A/B]` return matrices of the intercept and slope coefficients, respectively.

`CA` and `CB` return matrices of allele-specific copy number.

`totalCopynumber` (or `rawCopynumber`) returns a matrix of `CA+CB`.

See Also

[crlmmCopynumber](#), [CNSet-class](#)

Examples

```
## Version 1.6* of crlmm used CNSetLM objects.
data(sample.CNSet)
all(isCurrent(cnSet)) ## is the cnSet object current?

## -----
## calculating allele-specific copy number
## -----
## copy number for allele A, first 5 markers, first 2 samples
(ca <- CA(cnSet, i=1:5, j=1:2))
## copy number for allele B, first 5 markers, first 2 samples
(cb <- CB(cnSet, i=1:5, j=1:2))
## total copy number for first 5 markers, first 2 samples
(cn1 <- ca+cb)

## total copy number at first 5 nonpolymorphic loci
```

```

index <- which(!isSnp(cnSet))[1:5]
cn2 <- CA(cnSet, i=index, j=1:2)
## note, cb is NA at nonpolymorphic loci
(cb <- CB(cnSet, i=index, j=1:2))
## note, ca+cb will give NAs at nonpolymorphic loci
CA(cnSet, i=index, j=1:2) + cb
## A shortcut for total copy number
cn3 <- totalCopynumber(cnSet, i=1:5, j=1:2)
all.equal(cn3, cn1)
cn4 <- totalCopynumber(cnSet, i=index, j=1:2)
all.equal(cn4, cn2)

## markers 1-5, all samples
cn5 <- totalCopynumber(cnSet, i=1:5)
## all markers, samples 1-5
cn6 <- totalCopynumber(cnSet, j=1:5)

## NOTE: subsetting the object before extracting copy number
##       can be very inefficient when the data set is very large,
##       particularly if using ff objects. IN particular, subsetting
##       the CNSet object will subset all of the assay data elements
##       and all of the elements in the LinearModelParameter slot
## Not run:
##       do not do the following
cn <- CA(cnSet[1:5, ], "A")

## End(Not run)

```

crlmm-package

Genotype Calling via CRLMM Algorithm

Description

Faster implementation of CRLMM specific to SNP 5.0 and 6.0 arrays.

Details

Index:

crlmm-package	New implementation of the CRLMM Algorithm.
crlmm	Genotype SNP 5.0 or 6.0 samples.
calls	Accessor for genotype calls.
confs	Accessor for confidences.

The 'crlmm' package reimplements the CRLMM algorithm present in the 'oligo' package. This implementation primes for efficient genotyping of samples on SNP 5.0 and SNP 6.0 Affymetrix arrays.

To use this package, the user must have additional data packages: 'genomewidesnp5Crlmm' - SNP 5.0 arrays 'genomewidesnp6Crlmm' - SNP 6.0 arrays

Author(s)

Rafael A Irizarry Maintainer: Benilton S Carvalho <carvalho@bclab.org>

References

Carvalho BS, Louis TA, Irizarry RA. Quantifying uncertainty in genotype calls. *Bioinformatics*. 2010 Jan 15;26(2):242-9. Epub 2009 Nov 11.

crlmm

Genotype oligonucleotide arrays with CRLMM

Description

This is a faster and more efficient implementation of the CRLMM algorithm, especially designed for Affymetrix SNP 5 and 6 arrays (to be soon extended to other platforms).

Usage

```
crlmm(filenamees, row.names=TRUE, col.names=TRUE,
       probs=c(1/3, 1/3, 1/3), DF=6, SNRMin=5,
       gender=NULL, save.it=FALSE, load.it=FALSE,
       intensityFile, mixtureSampleSize=10^5,
       eps=0.1, verbose=TRUE, cdfName, sns, recallMin=10,
       recallRegMin=1000, returnParams=FALSE, badSNP=0.7)
crlmm2(filenamees, row.names=TRUE, col.names=TRUE,
        probs=c(1/3, 1/3, 1/3), DF=6, SNRMin=5,
        gender=NULL, save.it=FALSE, load.it=FALSE,
        intensityFile, mixtureSampleSize=10^5,
        eps=0.1, verbose=TRUE, cdfName, sns, recallMin=10,
        recallRegMin=1000, returnParams=FALSE, badSNP=0.7)
```

Arguments

filenamees	'character' vector with CEL files to be genotyped.
row.names	'logical'. Use rownames - SNP names?
col.names	'logical'. Use colnames - Sample names?
probs	'numeric' vector with priors for AA, AB and BB.
DF	'integer' with number of degrees of freedom to use with t-distribution.
SNRMin	'numeric' scalar defining the minimum SNR used to filter out samples.
gender	'integer' vector, with same length as 'filenamees', defining sex. (1 - male; 2 - female)
save.it	'logical'. Save preprocessed data?
load.it	'logical'. Load preprocessed data to speed up analysis?
intensityFile	'character' with filename to be saved/loaded - preprocessed data.
mixtureSampleSize	Number of SNP's to be used with the mixture model.
eps	Minimum change for mixture model.
verbose	'logical'.
cdfName	'character' defining the CDF name to use ('GenomeWideSnps5', 'GenomeWideSnps6')

sns 'character' vector with sample names to be used.
 recallMin Minimum number of samples for recalibration.
 recallRegMin Minimum number of SNP's for regression.
 returnParams 'logical'. Return recalibrated parameters.
 badSNP 'numeric'. Threshold to flag as bad SNP (affects batchQC)

Details

'crlmm2' allows one to genotype very large datasets (via ff package) and also permits the use of clusters or multiple cores (via snow package) to speed up genotyping.

As noted above, the call probabilities are stored using an integer representation to reduce file size using the transformation $\text{round}(-1000 \cdot \log_2(1-p))$, where p is the probability. The function `i2P` can be used to convert the integers back to the scale of probabilities.

Value

A SnpSet object.

calls Genotype calls (1 - AA, 2 - AB, 3 - BB)
 confs Confidence scores $\text{round}(-1000 \cdot \log_2(1-p))$
 SNPQC SNP Quality Scores
 batchQC Batch Quality Score
 params Recalibrated parameters

References

Carvalho B, Bengtsson H, Speed TP, Irizarry RA. Exploration, normalization, and genotype calls of high-density oligonucleotide SNP array data. *Biostatistics*. 2007 Apr;8(2):485-99. Epub 2006 Dec 22. PMID: 17189563.

Carvalho BS, Louis TA, Irizarry RA. Quantifying uncertainty in genotype calls. *Bioinformatics*. 2010 Jan 15;26(2):242-9.

See Also

[i2p](#), [snpCall](#), [snpCallProbability](#)

Examples

```
## this can be slow
if (require(genomewidesnp6Crlmm) & require(hapmapsnp6)){
  path <- system.file("celFiles", package="hapmapsnp6")

  ## the filenames with full path...
  ## very useful when genotyping samples not in the working directory
  cels <- list.celfiles(path, full.names=TRUE)
  (crlmmOutput <- crlmm(cels))

  ## If gender is known, one should check that the assigned gender is
  ## correct, or pass the integer coding of gender as an argument to the
  ## crlmm function as done below
  gender <- c("female", "female", "male")
  gender[gender == "female"] <- 2
```

```

gender[gender == "male"] <- 1
## Not run: (crlmmOutput <- crlmm(cels, gender=gender))
}

## Not run:
## HPC Example
library(ff)
library(snow)
library(crlmm)
## genotype 50K SNPs at a time
ocProbesets(50000)
## setup cluster - 8 cores on the machine
setCluster(8, "SOCK")

path <- system.file("celFiles", package="hapmapsnp6")
cels <- list.celfiles(path, full.names=TRUE)
crlmmOutput <- crlmm2(cels)

## End(Not run)

```

crlmmCopynumber *Locus- and allele-specific estimation of copy number*

Description

Locus- and allele-specific estimation of copy number.

Usage

```

crlmmCopynumber(object, MIN.SAMPLES=10, SNRMin = 5, MIN.OBS = 1,
  DF.PRIOR = 50, bias.adj = FALSE,
  prior.prob = rep(1/4, 4), seed = 1, verbose = TRUE,
  GT.CONF.THR = 0.80, MIN.NU = 2^3, MIN.PHI = 2^3,
  THR.NU.PHI = TRUE, type=c("SNP", "NP", "X.SNP", "X.NP"))

```

Arguments

object	object of class CNSet.
MIN.SAMPLES	'Integer'. The minimum number of samples in a batch. Batches with fewer than MIN.SAMPLES are skipped. Therefore, samples in batches with fewer than MIN.SAMPLES have NA's for the allele-specific copy number and NA's for the linear model parameters.
SNRMin	Samples with low signal to noise ratios are excluded.
MIN.OBS	For a SNP with with fewer than MIN.OBS of a genotype in a given batch, the within-genotype median is imputed. The imputation is based on a regression using SNPs for which all three biallelic genotypes are observed. For example, assume at a given SNP genotypes AA and AB were observed and BB is an unobserved genotype. For SNPs in which all 3 genotypes were observed, we fit the model $E(\text{mean_BB}) = \beta_0 + \beta_1 * \text{mean_AA} + \beta_2 * \text{mean_AB}$, obtaining estimates; of β_0 , β_1 , and β_2 . The imputed mean at the SNP with unobserved BB is then $\hat{\beta}_0 + \hat{\beta}_1 * \text{mean_AA} + \hat{\beta}_2 * \text{mean_AB}$.

<code>DF.PRIOR</code>	The 2 x 2 covariance matrix of the background and signal variances is estimated from the data at each locus. This matrix is then smoothed towards a common matrix estimated from all of the loci. <code>DF.PRIOR</code> controls the amount of smoothing towards the common matrix, with higher values corresponding to greater smoothing. Currently, <code>DF.PRIOR</code> is not estimated from the data. Future versions may estimate <code>DF.PRIOR</code> empirically.
<code>bias.adj</code>	<code>bias.adj</code> is currently ignored (as well as the <code>prior.prob</code> argument). We plan to add this feature back to the <code>crlmm</code> package in the near future. This feature, when <code>TRUE</code> , updated initial estimates from the linear model after excluding samples with a low posterior probability of normal copy number. Excluding samples that have a low posterior probability can be helpful at loci in which a substantial fraction of the samples have a copy number alteration. For additional information, see Scharpf et al., 2010.
<code>prior.prob</code>	This argument is currently ignored. A numerical vector providing prior probabilities for copy number states corresponding to homozygous deletion, hemizygous deletion, normal copy number, and amplification, respectively.
<code>seed</code>	Seed for random number generation.
<code>verbose</code>	Logical.
<code>GT.CONF.THR</code>	Confidence threshold for genotype calls (0, 1). Calls with confidence scores below this threshold are not used to estimate the within-genotype medians. See Carvalho et al., 2007 for information regarding confidence scores of biallelic genotypes.
<code>MIN.NU</code>	numeric. Minimum value for background intensity. Ignored if <code>THR.NU.PHI</code> is <code>FALSE</code> .
<code>MIN.PHI</code>	numeric. Minimum value for slope. Ignored if <code>THR.NU.PHI</code> is <code>FALSE</code> .
<code>THR.NU.PHI</code>	If <code>THR.NU.PHI</code> is <code>FALSE</code> , <code>MIN.NU</code> and <code>MIN.PHI</code> are ignored. When <code>TRUE</code> , background (<code>nu</code>) and slope (<code>phi</code>) coefficients below <code>MIN.NU</code> and <code>MIN.PHI</code> are set to <code>MIN.NU</code> and <code>MIN.PHI</code> , respectively.
<code>type</code>	Character string vector that must be one or more of "SNP", "NP", "X.SNP", or "X.NP". <code>Type</code> refers to a set of markers. See details below

Details

We suggest a minimum of 10 samples per batch for using `crlmmCopynumber`. 50 or more samples per batch is preferred and will improve the estimates.

The functions `crlmmCopynumberLD` and `crlmmCopynumber2` have been deprecated.

The argument `type` can be used to specify a subset of markers for which the copy number estimation algorithm is run. One or more of the following possible entries are valid: 'SNP', 'NP', 'X.SNP', and 'X.NP'.

'SNP' refers to autosomal SNPs.

'NP' refers to autosomal nonpolymorphic markers.

'X.SNP' refers to SNPs on chromosome X.

'X.NP' refers to autosomes on chromosome X.

However, users must run 'SNP' prior to running 'NP' and 'X.NP', or specify `type = c('SNP', 'X.NP')`.

Value

The value returned by the `crlmmCopynumber` function depends on whether the data is stored in RAM or whether the data is stored on disk using the R package `ff` for reading / writing. If uncertain, the first line of the `show` method defined for `CNSet` objects prints whether the `assayData` elements are derived from the `ff` package in the first line. Specifically,

- if the elements of the `batchStatistics` slot in the `CNSet` object have the class "ff_matrix" or "ffdf", then the `crlmmCopynumber` function updates the data stored on disk and returns the value `TRUE`.

- if the elements of the `batchStatistics` slot in the `CNSet` object have the class 'matrix', then the `crlmmCopynumber` function returns an object of class `CNSet` with the elements of `batchStatistics` updated.

Author(s)

R. Scharpf

References

Carvalho B, Bengtsson H, Speed TP, Irizarry RA. Exploration, normalization, and genotype calls of high-density oligonucleotide SNP array data. *Biostatistics*. 2007 Apr;8(2):485-99. Epub 2006 Dec 22. PMID: 17189563.

Carvalho BS, Louis TA, Irizarry RA. Quantifying uncertainty in genotype calls. *Bioinformatics*. 2010 Jan 15;26(2):242-9.

Scharpf RB, Ruczinski I, Carvalho B, Doan B, Chakravarti A, and Irizarry RA, *Biostatistics*. *Biostatistics*, Epub July 2010.

crlmmIllumina

Genotype Illumina Infinium II BeadChip data with CRLMM

Description

Implementation of the CRLMM algorithm for data from Illumina's Infinium II BeadChips.

Usage

```
crlmmIllumina(RG, XY, stripNorm=TRUE,
              useTarget=TRUE, row.names=TRUE, col.names=TRUE,
              probs=c(1/3, 1/3, 1/3), DF=6, SNRMin=5,
              gender=NULL, seed=1, mixtureSampleSize=10^5,
              eps=0.1, verbose=TRUE, cdfName, sns, recallMin=10,
              recallRegMin=1000, returnParams=FALSE, badSNP=0.7)
```

Arguments

RG	NChannelSet containing R and G bead intensities
XY	NChannelSet containing X and Y bead intensities
stripNorm	'logical'. Should the data be strip-level normalized?

useTarget	'logical' (only used when stripNorm=TRUE). Should the reference HapMap intensities be used in strip-level normalization?
row.names	'logical'. Use rownames - SNP names?
col.names	'logical'. Use colnames - Sample names?
probs	'numeric' vector with priors for AA, AB and BB.
DF	'integer' with number of degrees of freedom to use with t-distribution.
SNRMin	'numeric' scalar defining the minimum SNR used to filter out samples.
gender	'integer' vector, with same length as 'filenames', defining sex. (1 - male; 2 - female)
seed	'integer' scalar for random number generator (used to sample mixtureSampleSize SNPs for mixture model).
mixtureSampleSize	'integer'. The number of SNP's to be used when fitting the mixture model.
eps	Minimum change for mixture model.
verbose	'logical'.
cdfName	'character' defining the chip annotation (manifest) to use ('human370v1c', 'human550v3b', 'human650v3a', 'human1mv1c', 'human370quadv3c', 'human610quadv1b', 'human660quadv1a', 'human1mduov3b', 'humanomni1quadv1b', 'humanomniexpress12v1b')
sns	'character' vector with sample names to be used.
recallMin	'integer'. Minimum number of samples for recalibration.
recallRegMin	'integer'. Minimum number of SNP's for regression.
returnParams	'logical'. Return recalibrated parameters.
badSNP	'numeric'. Threshold to flag as bad SNP (affects batchQC)

Details

Note: The user should specify either the RG or XY intensities, not both.

Value

A SnpSet object which contains

calls Genotype calls (1 - AA, 2 - AB, 3 - BB)
 callProbability confidence scores 'round(-1000*log2(1-p))'

in the assayData slot and

SNPQC SNP Quality Scores
 batchQC Batch Quality Scores

along with center and scale parameters when returnParams=TRUE in the featureData slot.

Author(s)

Matt Ritchie

References

- Ritchie ME, Carvalho BS, Hetrick KN, Tavar'e S, Irizarry RA. R/Bioconductor software for Illumina's Infinium whole-genome genotyping BeadChips. *Bioinformatics*. 2009 Oct 1;25(19):2621-3.
- Carvalho B, Bengtsson H, Speed TP, Irizarry RA. Exploration, normalization, and genotype calls of high-density oligonucleotide SNP array data. *Biostatistics*. 2007 Apr;8(2):485-99. Epub 2006 Dec 22. PMID: 17189563.
- Carvalho BS, Louis TA, Irizarry RA. Quantifying uncertainty in genotype calls. *Bioinformatics*. 2010 Jan 15;26(2):242-9.

Examples

```
## crlmmOut = crlmmIllumina(RG)
```

crlmmIlluminaV2 *Read and Genotype Illumina Infinium II BeadChip data with CRLMM*

Description

Implementation of the CRLMM algorithm for data from Illumina's Infinium II BeadChips.

Usage

```
crlmmIlluminaV2(sampleSheet=NULL, arrayNames=NULL, ids=NULL, path=".",
  arrayInfoColNames=list(barcode="SentrixBarcode_A", position="SentrixPosition"),
  highDensity=FALSE, sep="_", fileExt=list(green="Grn.idat", red="Red.idat"),
  saveDate=FALSE, stripNorm=TRUE, useTarget=TRUE,
  row.names=TRUE, col.names=TRUE, probs=c(1/3, 1/3, 1/3),
  DF=6, SNRMin=5, gender=NULL, seed=1, mixtureSampleSize=10^5,
  eps=0.1, verbose=TRUE, cdfName, sns, recallMin=10,
  recallRegMin=1000, returnParams=FALSE, badSNP=.7)
```

Arguments

- | | |
|--------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <code>sampleSheet</code> | <code>data.frame</code> containing Illumina sample sheet information (for required columns, refer to <i>BeadStudio Genotyping guide - Appendix A</i>). |
| <code>arrayNames</code> | character vector containing names of arrays to be read in. If <code>NULL</code> , all arrays that can be found in the specified working directory will be read in. |
| <code>ids</code> | vector containing ids of probes to be read in. If <code>NULL</code> all probes found on the first array are read in. |
| <code>path</code> | character string specifying the location of files to be read by the function |
| <code>arrayInfoColNames</code> | (used when <code>sampleSheet</code> is specified) list containing elements 'barcode' which indicates column names in the <code>sampleSheet</code> which contains the arrayNumber/barcode number and 'position' which indicates the strip number. In older style sample sheets, this information is combined (usually in a column named 'SentrixPosition') and this should be specified as <code>list(barcode=NULL, position="SentrixPosition")</code> |

highDensity	logical (used when sampleSheet is specified). If TRUE, array extensions '_A', '_B' in sampleSheet are replaced with 'R01C01', 'R01C02' etc.
sep	character string specifying separator used in .idat file names.
fileExt	list containing elements 'Green' and 'Red' which specify the .idat file extension for the Cy3 and Cy5 channels.
saveDate	'logical'. Should the dates from each .idat be saved with sample information?
stripNorm	'logical'. Should the data be strip-level normalized?
useTarget	'logical' (only used when stripNorm=TRUE). Should the reference HapMap intensities be used in strip-level normalization?
row.names	'logical'. Use rownames - SNP names?
col.names	'logical'. Use colnames - Sample names?
probs	'numeric' vector with priors for AA, AB and BB.
DF	'integer' with number of degrees of freedom to use with t-distribution.
SNRMin	'numeric' scalar defining the minimum SNR used to filter out samples.
gender	'integer' vector, with same length as 'filenames', defining sex. (1 - male; 2 - female)
seed	'integer' scalar for random number generator (used to sample mixtureSampleSize SNPs for mixture model.
mixtureSampleSize	'integer'. The number of SNP's to be used when fitting the mixture model.
eps	Minimum change for mixture model.
verbose	'logical'.
cdfName	'character' defining the chip annotation (manifest) to use ('human370v1c', 'human550v3b', 'human650v3a', 'human1mv1c', 'human370quadv3c', 'human610quadv1b', 'human660quadv1a', 'human1mduov3b', 'humanomni1quadv1b', 'humanomniexpress12v1b')
sns	'character' vector with sample names to be used.
recallMin	'integer'. Minimum number of samples for recalibration.
recallRegMin	'integer'. Minimum number of SNP's for regression.
returnParams	'logical'. Return recalibrated parameters.
badSNP	'numeric'. Threshold to flag as bad SNP (affects batchQC)

Details

This function combines the reading of data from idat files using `readIdatFiles` and genotyping to reduce memory usage.

Value

A `SnpSet` object which contains

`calls` Genotype calls (1 - AA, 2 - AB, 3 - BB)
`callProbability` confidence scores `'round(-1000*log2(1-p))'`

in the `assayData` slot and

`SNPQC` SNP Quality Scores
`batchQC` Batch Quality Scores

along with center and scale parameters when `returnParams=TRUE` in the `featureData` slot.

Author(s)

Matt Ritchie

References

Ritchie ME, Carvalho BS, Hetrick KN, Tavar'e S, Irizarry RA. R/Bioconductor software for Illumina's Infinium whole-genome genotyping BeadChips. *Bioinformatics*. 2009 Oct 1;25(19):2621-3.

Carvalho B, Bengtsson H, Speed TP, Irizarry RA. Exploration, normalization, and genotype calls of high-density oligonucleotide SNP array data. *Biostatistics*. 2007 Apr;8(2):485-99. Epub 2006 Dec 22. PMID: 17189563.

Carvalho BS, Louis TA, Irizarry RA. Quantifying uncertainty in genotype calls. *Bioinformatics*. 2010 Jan 15;26(2):242-9.

See Also

[crlmmIllumina](#)

Examples

```
## crlmmOut = crlmmIlluminaV2(samples, path=path, arrayInfoColNames=list(barcode="Chip", position="A",
##                               saveDate=TRUE, cdfName="human370v1c", returnParams=TRUE))
```

genotype.Illumina *Preprocessing and genotyping of Illumina Infinium II arrays.*

Description

Preprocessing and genotyping of Illumina Infinium II arrays.

Usage

```
genotype.Illumina(sampleSheet=NULL, arrayNames=NULL, ids=NULL, path=".",
  arrayInfoColNames=list(barcode="SentrixBarcode_A", position="SentrixPosition",
  highDensity=FALSE, sep="_", fileExt=list(green="Grn.idat", red="Red.idat"),
  cdfName, copynumber=TRUE, batch, saveDate=TRUE, stripNorm=TRUE, useTarget=FALSE,
  mixtureSampleSize=10^5, fitMixture=TRUE, eps = 0.1, verbose = TRUE, seed =
  sns, probs = rep(1/3, 3), DF = 6, SNRMin = 5, recallMin = 10, recallRegMin = 10,
  gender = NULL, returnParams = TRUE, badSNP = 0.7)
```

Arguments

sampleSheet	data.frame containing Illumina sample sheet information (for required columns, refer to BeadStudio Genotyping guide - Appendix A).
arrayNames	character vector containing names of arrays to be read in. If NULL, all arrays that can be found in the specified working directory will be read in.
ids	vector containing ids of probes to be read in. If NULL all probes found on the first array are read in.
path	character string specifying the location of files to be read by the function

arrayInfoColNames	(used when <code>sampleSheet</code> is specified) list containing elements 'barcode' which indicates column names in the <code>sampleSheet</code> which contains the arrayNumber/barcode number and 'position' which indicates the strip number. In older style sample sheets, this information is combined (usually in a column named 'SentryPosition') and this should be specified as <code>list (barcode=NULL, position="SentryPosition")</code>
highDensity	logical (used when <code>sampleSheet</code> is specified). If TRUE, array extensions '_A', '_B' in <code>sampleSheet</code> are replaced with 'R01C01', 'R01C02' etc.
sep	character string specifying separator used in .dat file names.
fileExt	list containing elements 'Green' and 'Red' which specify the .dat file extension for the Cy3 and Cy5 channels.
cdfName	annotation package (see also <code>validCdfNames</code>)
copynumber	'logical.' Whether to store copy number intensities with SNP output.
batch	batch variable. See details.
saveDate	'logical'. Should the dates from each .dat be saved with sample information?
stripNorm	'logical'. Should the data be strip-level normalized?
useTarget	'logical' (only used when <code>stripNorm=TRUE</code>). Should the reference HapMap intensities be used in strip-level normalization?
mixtureSampleSize	Sample size to be use when fitting the mixture model.
fitMixture	'logical.' Whether to fit per-array mixture model.
eps	Stop criteria.
verbose	'logical.' Whether to print descriptive messages during processing.
seed	Seed to be used when sampling. Useful for reproducibility
sns	The sample identifiers. If missing, the default sample names are <code>basename (filenames)</code>
probs	'numeric' vector with priors for AA, AB and BB.
DF	'integer' with number of degrees of freedom to use with t-distribution.
SNRMin	'numeric' scalar defining the minimum SNR used to filter out samples.
recallMin	Minimum number of samples for recalibration.
recallRegMin	Minimum number of SNP's for regression.
gender	integer vector (male = 1, female =2) or missing, with same length as filenames. If missing, the gender is predicted.
returnParams	'logical'. Return recalibrated parameters from <code>crlmm</code> .
badSNP	'numeric'. Threshold to flag as bad SNP (affects <code>batchQC</code>)

Details

For large datasets it is important to utilize the large data support by installing and loading the `ff` package before calling the `genotype` function. In previous versions of the `crlmm` package, we used different functions for genotyping depending on whether the `ff` package is loaded, namely `genotype` and `genotype2`. The `genotype` function now handles both instances.

`genotype.Illumina` is a wrapper of the `crlmm` function for genotyping. Differences include (1) that the copy number probes (if present) are also quantile-normalized and (2) the class of object returned by this function, `CNSet`, is needed for subsequent copy number estimation. Note that the `batch` variable that must be passed to this function has no effect on the normalization or genotyping steps. Rather, `batch` is required in order to initialize a `CNSet` container with the appropriate dimensions.

Value

A `SnpSuperSet` instance.

Note

For large datasets, load the 'ff' package prior to genotyping – this will greatly reduce the RAM required for big jobs. See `ldPath` and `ocSamples`.

Author(s)

Matt Ritchie

References

Ritchie ME, Carvalho BS, Hetrick KN, Tavar'e S, Irizarry RA. R/Bioconductor software for Illumina's Infinium whole-genome genotyping BeadChips. *Bioinformatics*. 2009 Oct 1;25(19):2621-3.

Carvalho B, Bengtsson H, Speed TP, Irizarry RA. Exploration, normalization, and genotype calls of high-density oligonucleotide SNP array data. *Biostatistics*. 2007 Apr;8(2):485-99. Epub 2006 Dec 22. PMID: 17189563.

Carvalho BS, Louis TA, Irizarry RA. Quantifying uncertainty in genotype calls. *Bioinformatics*. 2010 Jan 15;26(2):242-9.

See Also

[crlmmIlluminaV2](#), [ocSamples](#), [ldOpts](#)

Examples

```
##
```

genotype

Preprocessing and genotyping of Affymetrix arrays.

Description

Preprocessing and genotyping of Affymetrix arrays.

Usage

```
genotype(filename, cdfName, batch, mixtureSampleSize = 10^5, eps = 0.1,
          verbose = TRUE, seed = 1, sns, probs = rep(1/3, 3),
          DF = 6, SNRMin = 5, recallMin = 10, recallRegMin = 1000,
          gender = NULL, returnParams = TRUE, badSNP = 0.7)
```

Arguments

<code>filenames</code>	complete path to CEL files
<code>cdfName</code>	annotation package (see also <code>validCdfNames</code>)
<code>batch</code>	batch variable. See details.
<code>mixtureSampleSize</code>	Sample size to be use when fitting the mixture model.
<code>eps</code>	Stop criteria.
<code>verbose</code>	Logical. Whether to print descriptive messages during processing.
<code>seed</code>	Seed to be used when sampling. Useful for reproducibility
<code>sns</code>	The sample identifiers. If missing, the default sample names are <code>basename (filenames)</code>
<code>probs</code>	'numeric' vector with priors for AA, AB and BB.
<code>DF</code>	'integer' with number of degrees of freedom to use with t-distribution.
<code>SNRMin</code>	'numeric' scalar defining the minimum SNR used to filter out samples.
<code>recallMin</code>	Minimum number of samples for recalibration.
<code>recallRegMin</code>	Minimum number of SNP's for regression.
<code>gender</code>	integer vector (male = 1, female =2) or missing, with same length as <code>filenames</code> . If missing, the gender is predicted.
<code>returnParams</code>	'logical'. Return recalibrated parameters from <code>crlmm</code> .
<code>badSNP</code>	'numeric'. Threshold to flag as bad SNP (affects <code>batchQC</code>)

Details

For large datasets it is important to utilize the large data support by installing and loading the `ff` package before calling the `genotype` function. In previous versions of the `crlmm` package, we used different functions for genotyping depending on whether the `ff` package is loaded, namely `genotype` and `genotype2`. The `genotype` function now handles both instances.

`genotype` is essentially a wrapper of the `crlmm` function for genotyping. Differences include (1) that the copy number probes (if present) are also quantile-normalized and (2) the class of object returned by this function, `CNSet`, is needed for subsequent copy number estimation. Note that the `batch` variable that must be passed to this function has no effect on the normalization or genotyping steps. Rather, `batch` is required in order to initialize a `CNSet` container with the appropriate dimensions.

Value

A `SnpSuperSet` instance.

Note

For large datasets, load the 'ff' package prior to genotyping – this will greatly reduce the RAM required for big jobs. See `ldPath` and `ocSamples`.

Author(s)

R. Scharpf

References

Carvalho B, Bengtsson H, Speed TP, Irizarry RA. Exploration, normalization, and genotype calls of high-density oligonucleotide SNP array data. *Biostatistics*. 2007 Apr;8(2):485-99. Epub 2006 Dec 22. PMID: 17189563.

Carvalho BS, Louis TA, Irizarry RA. Quantifying uncertainty in genotype calls. *Bioinformatics*. 2010 Jan 15;26(2):242-9.

See Also

[snprma](#), [crlmm](#), [ocSamples](#), [ldOpts](#), [batch](#), [crlmmCopynumber](#)

Examples

```
if (require(ff) & require(genomewidesnp6Crlmm) & require(hapmapsnp6)){

  path <- system.file("celFiles", package="hapmapsnp6")
  ## the filenames with full path...
  ## very useful when genotyping samples not in the working directory
  cels <- list.celfiles(path, full.names=TRUE)

  ## Note: one would need at least 10 CEL files for copy number estimation
  ## To use less RAM, specify a smaller argument to ocProbesets
  ocProbesets(50e3)
  batch <- as.factor(rep("A", length(cels)))
  (cnSet <- genotype(cels, cdfName="genomewidesnp6", batch=batch))

  ## when gender is not specified (as in the above example), crlmm tries
  ## to predict the gender from SNPs on chromosome X
  cnSet$gender

  ## If gender is known, one should check that the assigned gender is
  ## correct. Alternatively, one can pass gender as an argument to the
  ## genotype function.
  gender <- c("female", "female", "male")
  gender[gender == "female"] <- 2
  gender[gender == "male"] <- 1
  ## Not run:
  cnSet2 <- (cnSet <- genotype(cels, cdfName="genomewidesnp6", batch=batch, gender=as.int

  ## End(Not run)
  dim(cnSet)
  table(isSnp(cnSet))
}
```

genotypeInf

Genotyping of Illumina Infinium II arrays.

Description

Genotyping of Illumina Infinium II arrays. This function provides CRLMM genotypes and confidence scores for the the polymorphic markers and is a required step prior to copy number estimation.

Usage

```
genotypeInf(cnSet, mixtureParams, probs = rep(1/3, 3), SNRMin = 5, recallMin = 1
```

Arguments

cnSet	An object of class CNSet
mixtureParams	data.frame containing mixture model parameters needed for genotyping. The mixture model parameters are estimated from the preprocessInf function.
probs	'numeric' vector with priors for AA, AB and BB.
SNRMin	'numeric' scalar defining the minimum SNR used to filter out samples.
recallMin	Minimum number of samples for recalibration.
recallRegMin	Minimum number of SNP's for regression.
verbose	'logical.' Whether to print descriptive messages during processing.
returnParams	'logical'. Return recalibrated parameters from crlmm.
badSNP	'numeric'. Threshold to flag as bad SNP (affects batchQC)
gender	integer vector (male = 1, female =2) or missing, with same length as filenames. If missing, the gender is predicted.
DF	'integer' with number of degrees of freedom to use with t-distribution.

Details

The CRLMM genotype calls and confidence scores are written to file using *ff* protocols for I/O. For the most part, the calls and confidence scores can be accessed as though the data is in memory through the methods `snpCall` and `snpCallProbability`, respectively.

The genotype calls are stored using an integer representation: 1 - AA, 2 - AB, 3 - BB. Similarly, the call probabilities are stored using an integer representation to reduce file size using the transformation `'round(-1000*log2(1-p))'`, where *p* is the probability. The function `i2P` can be used to convert the integers back to the scale of probabilities.

Value

Logical. If the genotyping is completed, the value 'TRUE' is returned. Note that `assayData` elements 'call' and 'callProbability' are updated on disk. Therefore, the genotypes and confidence scores can be retrieved using accessors for the CNSet class.

Author(s)

R. Scharpf

See Also

[crlmm](#), [snpCall](#), [snpCallProbability](#)

Examples

```
## See the 'illumina_copynumber' vignette in inst/scripts of
## the source package
```

```
preprocessInf      Preprocessing of Illumina Infinium II arrays.
```

Description

This function normalizes the intensities for the 'A' and 'B' alleles for a `CNSet` object and estimates mixture parameters used for subsequent genotyping. See details for how the normalized intensities are written to file. This step is required for subsequent genotyping and copy number estimation.

Usage

```
preprocessInf(cnSet, sampleSheet=NULL, arrayNames = NULL, ids = NULL, path = ".")
```

Arguments

<code>cnSet</code>	object of class <code>CNSet</code>
<code>sampleSheet</code>	<code>data.frame</code> containing Illumina sample sheet information (for required columns, refer to BeadStudio Genotyping guide - Appendix A).
<code>arrayNames</code>	character vector containing names of arrays to be read in. If <code>NULL</code> , all arrays that can be found in the specified working directory will be read in.
<code>ids</code>	vector containing ids of probes to be read in. If <code>NULL</code> all probes found on the first array are read in.
<code>path</code>	character string specifying the location of files to be read by the function
<code>arrayInfoColNames</code>	(used when <code>sampleSheet</code> is specified) list containing elements 'barcode' which indicates column names in the <code>sampleSheet</code> which contains the arrayNumber/barcode number and 'position' which indicates the strip number. In older style sample sheets, this information is combined (usually in a column named 'SentryPosition') and this should be specified as <code>list(barcode=NULL, position="SentryPosition")</code>
<code>highDensity</code>	logical (used when <code>sampleSheet</code> is specified). If <code>TRUE</code> , array extensions '_A', '_B' in <code>sampleSheet</code> are replaced with 'R01C01', 'R01C02' etc.
<code>sep</code>	character string specifying separator used in <code>.idat</code> file names.
<code>fileExt</code>	list containing elements 'Green' and 'Red' which specify the <code>.idat</code> file extension for the Cy3 and Cy5 channels.
<code>saveDate</code>	'logical'. Should the dates from each <code>.idat</code> be saved with sample information?
<code>stripNorm</code>	'logical'. Should the data be strip-level normalized?
<code>useTarget</code>	'logical' (only used when <code>stripNorm=TRUE</code>). Should the reference HapMap intensities be used in strip-level normalization?
<code>mixtureSampleSize</code>	Sample size to be use when fitting the mixture model.
<code>fitMixture</code>	'logical.' Whether to fit per-array mixture model.
<code>eps</code>	Stop criteria.
<code>verbose</code>	'logical.' Whether to print descriptive messages during processing.
<code>seed</code>	Seed to be used when sampling. Useful for reproducibility

Details

The normalized intensities are written to disk using package `ff` protocols for writing/reading to disk. Note that the object `CNSet` containing the `ff` objects in the `assayData` slot will be updated after applying this function.

Value

A `ff_matrix` object containing parameters for fitting the mixture model. Note that while the `CNSet` object is not returned by this function, the object will be updated as the normalized intensities are written to disk. In particular, after applying this function the normalized intensities in the `alleleA` and `alleleB` elements of `assayData` are now available.

Author(s)

R. Scharpf

See Also

[CNSet-class](#), [A](#), [B](#), [constructInf](#), [genotypeInf](#)

Examples

```
## See the 'illumina_copynumber' vignette in inst/scripts of
## the source package
```

readIdatFiles	<i>Reads Idat Files from Infinium II Illumina BeadChips</i>
---------------	-------------------------------------------------------------

Description

Reads intensity information for each bead type from `.idat` files of Infinium II genotyping BeadChips

Usage

```
readIdatFiles(sampleSheet=NULL, arrayNames=NULL, ids=NULL, path="",
              arrayInfoColNames=list(barcode="SentrixBarcode_A",
                                     position="SentrixPosition_A"),
              highDensity=FALSE, sep="_",
              fileExt=list(green="Grn.idat", red="Red.idat"),
              saveDate=FALSE, verbose=FALSE)
```

Arguments

<code>sampleSheet</code>	data.frame containing Illumina sample sheet information (for required columns, refer to <i>BeadStudio Genotyping guide - Appendix A</i>).
<code>arrayNames</code>	character vector containing names of arrays to be read in. If <code>NULL</code> , all arrays that can be found in the specified working directory will be read in.
<code>ids</code>	vector containing ids of probes to be read in. If <code>NULL</code> all probes found on the first array are read in.
<code>path</code>	character string specifying the location of files to be read by the function

<code>arrayInfoColNames</code>	(used when <code>sampleSheet</code> is specified) list containing elements 'barcode' which indicates column names in the <code>sampleSheet</code> which contains the arrayNumber/barcode number and 'position' which indicates the strip number. In older style sample sheets, this information is combined (usually in a column named 'SentryPosition') and this should be specified as <code>list (barcode=NULL, position="SentryPosition")</code>
<code>highDensity</code>	logical (used when <code>sampleSheet</code> is specified). If TRUE, array extensions '_A', '_B' in <code>sampleSheet</code> are replaced with 'R01C01', 'R01C02' etc.
<code>sep</code>	character string specifying separator used in .idat file names.
<code>fileExt</code>	list containing elements 'Green' and 'Red' which specify the .idat file extension for the Cy3 and Cy5 channels.
<code>saveDate</code>	logical. Should the dates from each .idat be saved with sample information?
<code>verbose</code>	logical. Should processing information be displayed as data is read in?

Details

The summarised Cy3 (G) and Cy5 (R) intensities (on the original scale) are read in from the .idat files.

Where available, a `sampleSheet` data.frame, in the same format as used by BeadStudio (columns 'Sample_ID', 'SentryBarcode_A' and 'SentryPosition_A' are required) which keeps track of sample information can be specified.

Thanks to Keith Baggerly who provided the code to read in the binary .idat files.

Value

NChannelSet with intensity data (R, G), and indicator for SNPs with 0 beads (`zero`) for each bead type.

Author(s)

Matt Ritchie

References

Ritchie ME, Carvalho BS, Hetrick KN, Tavar'e S, Irizarry RA. R/Bioconductor software for Illumina's Infinium whole-genome genotyping BeadChips. *Bioinformatics*. 2009 Oct 1;25(19):2621-3.

Examples

```
#RG = readIdatFiles()
```

sample.CNSet	<i>Object of class 'CNSet'</i>
--------------	--------------------------------

Description

The data for the first 16 polymorphic markers in the HapMap analysis.

Usage

```
data(sample.CNSet)
```

Format

The data illustrates the `CNSet`-class, with `assayData` containing the quantile-normalized intensities for the A and B alleles, `genotype` calls and `confidence` scores. New slots that specific to copy number estimation are `batch` and `batchStatistics`.

Details

This object was created from the `copynumber` vignette in `inst/scripts`.

Examples

```
data(sample.CNSet)
## -----
## accessors for the feature-level info
## -----
chromosome(cnSet)[1:5]
position(cnSet)[1:5]
isSnp(cnSet)[1:5]
table(isSnp(cnSet))
## -----
## sample-level statistics computed by crlmm
## -----
varLabels(cnSet)
## accessors for sample-level statistics
## The signal to noise ratio (SNR)
cnSet$SNR[1:5]
## the skew
cnSet$SKW[1:5]
## the gender (gender is imputed unless specified in the call to crlmm)
table(cnSet$gender) ## 1=male, 2=female
## -----
## batchStatistics
## ----- estimate of
## intercept from linear model
dim(nu(cnSet, "A"))
## background for the A allele in the 2 batches for the
## first 5 markers
nu(cnSet, "A")[1:5, ]
## background for the B allele in the 2 batches for the
## first 5 markers
nu(cnSet, "B")[1:5, ]
## the slope
```

```

phi(cnSet, "A")[1:5, ]
## correlation within genotype cluster AA
##corr(cnSet, "AA")[1:5, ]
#### correlation within genotype cluster AB
##corr(cnSet, "AB")[1:5, ]
#### correlation within genotype cluster BB
##corr(cnSet, "BB")[1:5, ]
## -----

## -----
## calculating allele-specific copy number
## -----
## copy number for allele A, first 5 markers, first 2 samples
(ca <- CA(cnSet, i=1:5, j=1:2))
## copy number for allele B, first 5 markers, first 2 samples
(cb <- CB(cnSet, i=1:5, j=1:2))
## total copy number for first 5 markers, first 2 samples
(cn1 <- ca+cb)

## total copy number at first 5 nonpolymorphic loci
index <- which(!isSnp(cnSet))[1:5]
cn2 <- CA(cnSet, i=index, j=1:2)
## note, cb is NA at nonpolymorphic loci
(cb <- CB(cnSet, i=index, j=1:2))
## note, ca+cb will give NAs at nonpolymorphic loci
CA(cnSet, i=index, j=1:2) + cb
## A shortcut for total copy number
cn3 <- totalCopynumber(cnSet, i=1:5, j=1:2)
all.equal(cn3, cn1)
cn4 <- totalCopynumber(cnSet, i=index, j=1:2)
all.equal(cn4, cn2)

## markers 1-5, all samples
cn5 <- totalCopynumber(cnSet, i=1:5)
## all markers, samples 1-5
cn6 <- totalCopynumber(cnSet, j=1:5)

```

snprma

Preprocessing tool for SNP arrays.

Description

SNPRMA will preprocess SNP chips. The preprocessing consists of quantile normalization to a known target distribution and summarization to the SNP-Allele level.

Usage

```

snprma(filenamees, mixtureSampleSize = 10^5, fitMixture = FALSE, eps = 0.1, verbo
snprma2(filenamees, mixtureSampleSize = 10^5, fitMixture = FALSE, eps = 0.1, verb

```

Arguments

filenamees 'character' vector with file names.

mixtureSampleSize	Sample size to be use when fitting the mixture model.
fitMixture	'logical'. Fit the mixture model?
eps	Stop criteria.
verbose	'logical'.
seed	Seed to be used when sampling.
cdfName	cdfName: 'GenomeWideSnp_5', 'GenomeWideSnp_6'
sns	Sample names.

Details

'snprma2' allows one to genotype very large datasets (via ff package) and also permits the use of clusters or multiple cores (via snow package) to speed up preprocessing.

Value

A	Summarized intensities for Allele A
B	Summarized intensities for Allele B
sns	Sample names
gns	SNP names
SNR	Signal-to-noise ratio
SKW	Skewness
mixtureParams	Parameters from mixture model
cdfName	Name of the CDF

Examples

```
if (require(genomewidesnp6Crlmm) & require(hapmapsnp6) & require(oligoClasses)){
  path <- system.file("celFiles", package="hapmapsnp6")

  ## the filenames with full path...
  ## very useful when genotyping samples not in the working directory
  cels <- list.celfiles(path, full.names=TRUE)
  snprmaOutput <- snprma(cels)
  snprmaOutput[["A"]][1:10,]
  snprmaOutput[["B"]][1:10,]
}
## Not run:
## HPC Example
library(ff)
library(snow)
library(crlmm)
## genotype 50K SNPs at a time
ocProbesets(50000)
## setup cluster - 8 cores on the machine
setCluster(8, "SOCK")

path <- system.file("celFiles", package="hapmapsnp6")
cels <- list.celfiles(path, full.names=TRUE)
snprmaOutput <- snprma2(cels)
```

```
## End(Not run)
```

Index

*Topic **IO**

readIdatFiles, 23

*Topic **classif**

crlmm, 8

crlmmIllumina, 12

crlmmIlluminaV2, 14

genotype, 18

genotype.Illumina, 16

genotypeInf, 20

snprma, 26

*Topic **datasets**

sample.CNSet, 25

*Topic **manip**

AssayData-methods, 1

batchStatisticAccessors, 2

celDates, 3

constructInf, 4

copynumberAccessors, 5

crlmmCopynumber, 10

preprocessInf, 22

snprma, 26

*Topic **methods**

CNSet-methods, 1

*Topic **package**

crlmm-package, 7

A, 23

AssayData-methods, 1

B, 23

batch, 20

batchStatisticAccessors, 2

batchStatistics, 3

CA, 2

CA (copynumberAccessors), 5

CA, CNSet-method (CNSet-methods), 1

CB, 2

CB (copynumberAccessors), 5

CB, CNSet-method (CNSet-methods), 1

celDates, 3

cnSet (sample.CNSet), 25

CNSet-class, 1, 2, 5, 6, 23

CNSet-methods, 1

constructInf, 4, 23

copynumberAccessors, 5

corr, 1

corr (batchStatisticAccessors), 2

corr, AssayData-method
(AssayData-methods), 1

corr, CNSet-method
(CNSet-methods), 1

crlmm, 8, 20, 21

crlmm-package, 7

crlmm2 (crlmm), 8

crlmmCopynumber, 6, 10, 20

crlmmCopynumber2
(crlmmCopynumber), 10

crlmmCopynumberLD
(crlmmCopynumber), 10

crlmmIllumina, 12, 16

crlmmIlluminaV2, 14, 18

genotype, 18

genotype.Illumina, 16

genotype2 (genotype), 18

genotypeInf, 5, 20, 23

genotypeLD (genotype), 18

i2p, 9

ldOpts, 18, 20

ldPath, 5

lines, CNSet-method
(CNSet-methods), 1

mads, 1

mads (batchStatisticAccessors), 2

mads, AssayData-method
(AssayData-methods), 1

mads, CNSet-method
(CNSet-methods), 1

medians, 1

medians
(batchStatisticAccessors),
2

medians, AssayData-method
(AssayData-methods), 1

medians, C`NSet`-method
 (*C`NSet`-methods*), 1

Ns, 1
 Ns (*batchStatisticAccessors*), 2
 Ns, AssayData-method
 (*AssayData-methods*), 1
 Ns, C`NSet`-method (*C`NSet`-methods*), 1
 nuA (*copynumberAccessors*), 5
 nuA, C`NSet`-method (*C`NSet`-methods*),
 1
 nuB (*copynumberAccessors*), 5
 nuB, C`NSet`-method (*C`NSet`-methods*),
 1

ocProbesets, 5
 ocSamples, 5, 18, 20

phiA (*copynumberAccessors*), 5
 phiA, C`NSet`-method
 (*C`NSet`-methods*), 1
 phiB (*copynumberAccessors*), 5
 phiB, C`NSet`-method
 (*C`NSet`-methods*), 1
 POSIXt, 4
 preprocessInf, 5, 22

rawCopynumber, 2
 rawCopynumber
 (*copynumberAccessors*), 5
 rawCopynumber, C`NSet`-method
 (*C`NSet`-methods*), 1
 read.celfile.header, 4
 readIdatFiles, 23
 readIdatFiles2 (*readIdatFiles*), 23

sample.C`NSet`, 25
 snpCall, 9, 21
 snpCallProbability, 9, 21
 snprma, 20, 26
 snprma2 (*snprma*), 26

tau2, 1
 tau2 (*batchStatisticAccessors*), 2
 tau2, AssayData-method
 (*AssayData-methods*), 1
 tau2, C`NSet`-method
 (*C`NSet`-methods*), 1
 totalCopynumber, 2
 totalCopynumber
 (*copynumberAccessors*), 5
 totalCopynumber, C`NSet`-method
 (*C`NSet`-methods*), 1