

Package ‘Clonality’

December 5, 2021

Type Package

Title Clonality testing

Version 1.42.0

Date 2019-12-09

Author Irina Ostrovnaya

Maintainer Irina Ostrovnaya <ostrovni@mskcc.org>

Depends R (>= 2.12.2), DNACopy

Imports grDevices, graphics, stats, utils

Suggests gdata

Description Statistical tests for clonality versus independence of tumors from the same patient based on their LOH or genomewide copy number profiles

License GPL-3

LazyLoad yes

biocViews CopyNumber, Classification, aCGH, Mutations, Diagnosis, metastasis

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

git_branch RELEASE_3_14

git_last_commit 3792eeb

git_last_commit_date 2021-10-26

Date/Publication 2021-12-05

R topics documented:

Clonality-package	2
ave.adj.probes	3
chromosomePlots	4
clonality.analysis	5
clonEM	10
create.mutation.matrix	11

ECMtesting	12
freqdata	13
genomewidePlots	14
get.mutation.frequencies	15
grid.lik	17
histogramPlot	18
lcis	19
LOHclonality	19
LRtesting3or4tumors	21
model.lik	22
mutation.proba	23
mutation.rem	24
print.mutation.proba	26
print.mutation.rem	27
SNVtest	27
SNVtest2	29
splitChromosomes	31
xidens	32

Index	33
--------------	-----------

Clonality-package	<i>Clonality testing</i>
-------------------	--------------------------

Description

Statistical tests for clonality versus independence of tumors from the same patient based on their LOH, copy number or mutational profiles.

Details

Package:	Clonality
Type:	Package
Version:	0.99.3
Date:	2014-9-07
License:	GPL-3
LazyLoad:	yes

Author(s)

Irina Ostrovnyaya <ostrovni@mskcc.org>, Audrey Mauguen <mauguena@mskcc.org>

ave.adj.probes *Averaging of adjacent probes in copy number arrays*

Description

For each sample the log-ratios at each consecutive K number of probes are averaged.

Usage

```
ave.adj.probes(data, K)
```

Arguments

data	Copy Number Array object (output of function CNA() from the package DNA-copy). First column contains chromosomes, second column contains genomic locations. Each remaining column contains log-ratios from a particular tumor or sample.
K	Number of markers to be averaged. Should be selected so that the final resolution of the averaged data would be 5,000-10,000 markers.

Details

Averages log-ratios in every K consecutive markers. The purpose of this step is to reduce the noise in the data, eliminate possible very small germline copy number variations, and get rid of a possible wave effect.

Value

Returns CNA object of reduced resolution

Examples

```
# Same example as in clonality.analysis()

set.seed(100)
chrom<-rep(c(1:22),each=100)
maploc<- runif(2200)* 200000
chromarm<-splitChromosomes(chrom,maploc)

#Simulate the dataset with 10 pairs of tumors with 22 chromosomes, 100 markers each
#Simulated log-ratios are equal to signal + noise
#Signal: each chromosome has 50% chance to be normal, 30% to be whole-arm loss/gain, and 20% to be partial arm loss/g
#There are no chromosomes with recurrent losses/gains
#Noise: drawn from normal distribution with mean 0, standard deviation 0.25
#First 9 patients have independent tumors, last patient has two tumors with identical signal, independent noise

set.seed(100)
```

```

chrom<-paste("chr",rep(c(1:22),each=100),"p",sep="")
chrom[nchar(chrom)==5]<-paste("chr0",substr(chrom[nchar(chrom)==5],4,5),sep="")
maploc<- rep(c(1:100),22)
data<-NULL
for (pt in 1:9) #first 9 patients have independent tumors
{
tumor1<-tumor2<- NULL
mean1<- rnorm(22)
mean2<- rnorm(22)
for (chr in 1:22)
{
r<-runif(2)
if (r[1]<=0.5) tumor1<-c(tumor1,rep(0,100))
else if (r[1]>0.7) tumor1<-c(tumor1,rep(mean1[chr],100))
else { i<-sort(sample(1:100,2))
tumor1<-c(tumor1,mean1[chr]*c(rep(0, i[1]),rep(1, i[2]-i[1]), rep(0, 100-i[2])))
}
if (r[2]<=0.5) tumor2<-c(tumor2,rep(0,100))
else if (r[2]>0.7) tumor2<-c(tumor2,rep(mean2[chr],100))
else { i<-sort(sample(1:100,2))
tumor2<-c(tumor2,mean2[chr]*c(rep(0, i[1]),rep(1, i[2]-i[1]), rep(0, 100-i[2])))
}
}
data<-cbind(data,tumor1,tumor2)
}

#last patient has identical profiles
tumor1<- NULL
mean1<- rnorm(22)
for (chr in 1:22)
{
r<-runif(1)
if (r<=0.4) tumor1<-c(tumor1,rep(0,100))
else if (r>0.6) tumor1<-c(tumor1,rep(mean1[chr],100))
else { i<-sort(sample(1:100,2))
tumor1<-c(tumor1,mean1[chr]*c(rep(0, i[1]),rep(1, i[2]-i[1]), rep(0, 100-i[2])))
}
}
data<-cbind(data,tumor1,tumor1)

data<-data+matrix(rnorm( 44000,mean=0,sd=0.4) ,nrow=2200,ncol=20)
dataCNA<-CNA(data,chrom=chrom,maploc=maploc,sampleid=paste("pt",rep(1:10,each=2),rep(1:2,10)))
dim(dataCNA)
dataCNA2<-ave.adj.probes(dataCNA, 2)
dim(dataCNA2)

```

Description

The function produces a sequence of plots for each chromosome with one-step segmented data of all samples of a particular patient.

Usage

```
chromosomePlots(data.seg1, ptlist, ptname, nmad)
```

Arguments

<code>data.seg1</code>	Output of one-step segmentation - output <code>OneStepSeg</code> of <code>clonality.analysis()</code> .
<code>ptlist</code>	Vector of the patient IDs in the order of the samples appearing in the data. For example, if the first three tumors belong to patient A, and the following two belong to patient B, then <code>ptlist=c('ptA', 'ptA', 'ptA', 'ptB', 'ptB')</code> .
<code>ptname</code>	Name of the patient from <code>ptlist</code> for which the data should be plotted
<code>nmad</code>	Number of MADs (median absolute deviations) that is used for Gain/Loss calls. Used to mark the Gain/Loss threshold on the plots.

Details

The function produces a sequence of plots for each chromosome with one-step segmented data of all samples of a particular patient. The dotted horizontal lines denote the gain and loss thresholds.

Examples

```
# See example as in clonality.analysis()
```

`clonality.analysis` *Clonality testing using copy number data*

Description

Function to test clonality of two tumors from the same patient based on their genomewide copy number profiles. This function calculates likelihood ratios and the reference distribution under the hypothesis of independence.

Usage

```
clonality.analysis(data, ptlist, pfreq = NULL, refdata = NULL, nmad = 1.25, reference = TRUE, allpairs =
```

Arguments

data	Copy Number Array object (output of function CNA() from package DNACopy). First column contains chromosomes, second column contains genomic locations. Each remaining column contains log-ratios from a particular tumor or sample. Chromosomes X and Y should be removed prior to analysis, and chromosomes should be split into p and q arms to improve the power (use function splitChromosomes()).
ptlist	Vector of the patient IDs in the order of the samples appearing in the data. For example, if the first three tumors (columns 3, 4, 5 of data) belong to patient A, and the following two (columns 6, 7 of data) belong to patient B, then ptlist=c('ptA', 'ptA', 'ptA', 'ptB', 'ptB'). Note that while sample names in data should be unique the ptlist should have repeated labels.
pfreq	Marginal frequencies of Gains, Losses and Normals for all the chromosomes. If it is not known, pfreq should be set to NULL and frequencies will be estimated from all the samples in the dataset. If frequencies are known, pfreq should be a data frame with 4 columns: 1) chromosome arm in the format 'chr01p', probability of 2) gain, 3) loss and 4) normal.
refdata	If available, additional cohort of patients with the same disease that should be used to estimate the marginal gain/loss frequencies. If NULL, the original set of tumors is used, otherwise, refdata should be a CNA object. It will be segmented with 1 step CBS and each chromosome will be classified as gain/loss as described in the manuscript, leading to frequency estimates. No averaging or chromosome splitting is done for this dataset, so users should make sure refdata has chromosomes in the format 'chr01p' and that its resolution is similar to the one of the original data.
nmad	Number of MADs (median absolute deviations) that is used for Gain/Loss calls. For each array MAD of its residuals (that is, data minus segmentation means) is calculated. Residuals represent the array's noise level. Any segment of this array that has a mean at least nmad MADs above or below array's median is called a gain or a loss. We use value of 1.25, while values in the range of 0.5 to 2 can also be admissible depending on the resolution and presence of artifacts.
reference	If TRUE the reference distribution of likelihood ratios is created under hypothesis of independence by pairing (independent) tumors from different patients.
allpairs	If TRUE all possible pairs of tumors from different patients will be used for reference distribution. If two tumors in a pair are not exchangeable, for example primary tumor vs recurrence, or pre-cancerous lesion vs tumor, then allpairs should be set to FALSE and the 'first' tumor should always come earlier in the data before the 'second' tumor for all the patients. Then 'first' tumors of patients will only be paired with 'second' tumors of other patients for the reference distribution.
segmethod	The segmentation algorithm to be used. The default is "oneseg" which uses the built in function of the same name based on the CBS algorithm. An alternative segmentation algorithm can be used. A function should be created and the name passed as described in the vignette.
separ	The parameters necessary for the segmentation algorithm as a list. For "oneseg" you can specify alpha (default = 0.01) and nperm (default = 2000) necessary for

the CBS algorithm.

Details

The function implements the statistical procedure designed to distinguish whether the two tumors from the same patient are clonal (have the same progenitor cancer cell) or independent (developed from normal cells independently). At first data are segmented with one step CBS (Olshen, A. B., Venkatraman, E. S., Lucito, R., Wigler, M. (2004). Circular binary segmentation for the analysis of array-based DNA copy number data. *Biostatistics* 5: 557-572) that picks at most one copy number change per chromosome arm. Then each chromosome arm is classified as Gain/Loss/Normal based on a middle segment if there are 3 segments, or based on the most outstanding segment if there are 2 segments. The multinomial likelihood ratio comparing these classifications is computed (LR1). For each concordant partial arm gain or loss we also calculate likelihood ratio that this change is exactly the same in both tumors. These likelihood ratios are multiplied by LR1 to obtain our final statistic, LR2. If LR2 is much greater than 1, that indicates clonality. If LR2 is much smaller than 1, it indicates independence. The reference distribution of LR2 under the hypothesis of independence is obtained by pairing up tumors from different patients, which are independent by default.

Since only one gain/loss is admissible per chromosome arm it is highly recommended to apply this methodology to arrays with at most 10,000-15,000 markers. We suggest averaging blocks of consecutive probes for arrays with larger resolution, see function `ave.adj.probes`.

Value

If the reference is TRUE, function returns the list with 4 elements: LR, OneStepSeg, ChromClass, refLR.

LR - matrix with the within patient comparisons. Each row corresponds to a pair of samples being compared. Columns are: Sample1 - name of sample 1; Sample2 - name of sample 2; LR1 - likelihood ratio without comparisons of specific concordant gains/losses; LR2 - final likelihood ratio with individual comparisons; GGorLL - number of chromosome arms that are classified as Gains in both tumors or Losses in both tumors; NN - number of chromosome arms that are classified as Normal in both tumors; GL - number of chromosome arms that are classified as Gain in one tumors and Loss in another; GNorLN - number of chromosome arms that are classified as Gain(Loss) in one tumors and Normal in another; IndividualComparisons - list of chromosome arms that had comparisons of specific concordant gains/losses in both tumors and the corresponding likelihood ratio for them being exactly the same. p-value - quantile of the reference distribution under the null hypothesis (`refLR$LR2`) that the value of LR2 match.

OneStepSeg - is the output of one step segmentation of the data. It has the same structure as the output of 'segment' from `DNAcopy`, but only one most prominent change per arm is allowed.

ChromClass - is the matrix of chromosome classifications based on the one step segmentation. Rows correspond to chromosome arms, columns correspond to samples. Chromosome arms are classified by the middle segment if there are 3 segments, and by the most outstanding segment if there are 2 segments.

refLR - matrix with the between patient comparisons (reference distribution under the hypothesis of independence). Has the same structure as LR but the pairs of tumors are selected from different patients.

Note that calculating the reference distribution might take a long time.

If the reference is FALSE, there is no p-value column in LR and no refLR output.

Author(s)

Irina Ostrovnaya <ostrovni@mskcc.org>

References

Ostrovnya, I., Olshen, A. B., Seshan, V.E., Orlow, I., Albertson, D. G. and Begg, C. B. (2010), A metastasis or a second independent cancer? Evaluating the clonal origin of tumors using array copy number data. *Statistics in Medicine*, 29: 1608-1621

Ostrovnya, I. and Begg, C. Testing Clonal Relatedness of Tumors Using Array Comparative Genomic Hybridization: A Statistical Challenge *Clin Cancer Res* March 1, 2010 16:1358-1367

Venkatraman, E. S. and Olshen, A. B. (2007). A faster circular binary segmentation algorithm for the analysis of array CGH data. *Bioinformatics*, 23:657-63.

Olshen, A. B., Venkatraman, E. S., Lucito, R., Wigler, M. (2004). Circular binary segmentation for the analysis of array-based DNA copy number data. *Biostatistics* 5: 557-572.

Examples

```
#Analysis of simulated data
```

```
#Simulate the dataset with 10 pairs of tumors with 22 chromosomes, 100 markers each
#Simulated log-ratios are equal to signal + noise
#Signal: each chromosome has 50% chance to be normal, 30% to be whole-arm loss/gain, and 20% to be partial arm loss/g
#There are no chromosomes with recurrent losses/gains
#Noise: drawn from normal distribution with mean 0, standard deviation 0.4
#First 9 patients have independent tumors, last patient has two tumors with identical signal, independent noise
```

```
set.seed(100)
chrom<-paste("chr",rep(c(1:22),each=100),"p",sep="")
chrom[nchar(chrom)==5]<-paste("chr0",substr(chrom[nchar(chrom)==5],4,5),sep="")
maploc<- rep(c(1:100),22)
data<-NULL
for (pt in 1:9) #first 9 patients have independent tumors
{
  tumor1<-tumor2<- NULL
  mean1<- rnorm(22)
  mean2<- rnorm(22)
  for (chr in 1:22)
  {
    r<-runif(2)
    if (r[1]<=0.5) tumor1<-c(tumor1,rep(0,100))
    else if (r[1]>0.7) tumor1<-c(tumor1,rep(mean1[chr],100))
    else { i<-sort(sample(1:100,2))
          tumor1<-c(tumor1,mean1[chr]*c(rep(0, i[1]),rep(1, i[2]-i[1]), rep(0, 100-i[2])))
        }
    if (r[2]<=0.5) tumor2<-c(tumor2,rep(0,100))
    else if (r[2]>0.7) tumor2<-c(tumor2,rep(mean2[chr],100))
    else { i<-sort(sample(1:100,2))
          tumor2<-c(tumor2,mean2[chr]*c(rep(0, i[1]),rep(1, i[2]-i[1]), rep(0, 100-i[2])))
        }
  }
}
```



```

    }
  }
  data<-cbind(data,tumor1,tumor2)
}

#last patient has identical profiles
tumor1<- NULL
mean1<- rnorm(22)
for (chr in 1:22)
{
  r<-runif(1)
  if (r<=0.4) tumor1<-c(tumor1,rep(0,100))
  else if (r>0.6) tumor1<-c(tumor1,rep(mean1[chr],100))
  else { i<-sort(sample(1:100,2))
        tumor1<-c(tumor1,mean1[chr]*c(rep(0, i[1]),rep(1, i[2]-i[1]), rep(0, 100-i[2])))
      }
}

data<-cbind(data,tumor1,tumor1)

data<-data+matrix(rnorm( 44000,mean=0,sd=0.4) ,nrow=2200,ncol=20)
dataCNA<-CNA(data,chrom=chrom,maploc=maploc,sampleid=paste("pt",rep(1:10,each=2),rep(1:2,10)))
ptlist<- paste("pt",rep(1:10,each=2),sep=".")
samnms<-paste("pt",rep(1:10,each=2),rep(1:2,10),sep=".")
results<-clonality.analysis(dataCNA, ptlist, pfreq = NULL, refdata = NULL, nmad = 1,
  reference = TRUE, allpairs = TRUE)

#genomewide plots of pairs of tumors from the same patient
pdf("genomewideplots.pdf",height=7,width=11)
for (i in unique(ptlist))
{
  w<-which(ptlist==i)
  ns<- length(w)
  if (ns>1)
  {
    for (p1 in c(1:(ns-1)))
    for (p2 in c((p1+1):ns))
    genomewidePlots(results$OneStepSeg, results$ChromClass,ptlist , ptpair=samnms[c(w[p1],w[p2])],results$LR, plot.
  }
}
dev.off()

pdf("hist.pdf",height=7,width=11)
histogramPlot(results$LR[,4], results$refLR[,4])
dev.off()

for (i in unique(ptlist))
{

```

```
pdf(paste("Patient", i, ".pdf", sep=""), height=7, width=11)
chromosomePlots(results$OneStepSeg, ptlist, ptname=i, nmad=1.25)
dev.off()
}
```

clonEM

Auxiliary EM algorithm Function.

Description

This function optimizes the likelihood of the model using an EM algorithm.

Usage

```
clonEM(mutmat, init.para, xigrid, conv.crit, niter)
```

Arguments

mutmat	Matrix containing the data, with all mutations in rows and the tumor pairs in columns. The data are coded as 0=mutation not observed, 1=shared mutation (observed in both tumors), 2=private mutation (observed in one tumor only). The first column contains the probabilities of occurrence for each mutation.
init.para	Initial values of the parameters for the optimization. The order of the parameters is c(mu, sigma, pi), where mu and sigma are the mean and variance of the lognormal distribution of the random-effect xi, and pi is the proportion of clonal cases.
xigrid	Grid of the values of xi used to compute the integration; it corresponds to the domain of definition of xi.
conv.crit	Criteria used to defined convergence (on all three parameters).
niter	Maximum number of iterations used for the EM algorithm.

Value

Return the value of the parameters maximizing the likelihood, the number of iterations used, and the convergence status.

Author(s)

Audrey Mauguen <mauguena@mskcc.org> and Venkatraman E. Seshan.

`create.mutation.matrix`*Formatting matrix of mutations*

Description

This functions reformats matrix of mutations for subsequent analysis

Usage

```
create.mutation.matrix(maf,rem=FALSE)
```

Arguments

<code>maf</code>	<code>maf</code> is a mutations file in generic MAF (mutation annotation format) style: one row per mutation per sample. See this webpage for detailed description of the format: https://docs.gdc.cancer.gov/Data/File_Formats/MAF_Format/ . Object 'maf' should have the following columns: PatientID, Tumor_Sample_Barcode (sample ID), Chromosome, Start_Position, Reference_Allele, Tumor_Seq_Allele2 (reference allele. For compatibility with mutational frequencies that can be obtained using function "get.mutation.frequencies" chromosome should be a number 1-22 or X or Y; Start_Position is genomic location in GRCh37 build; Reference_Allele is a reference allele and Tumor_Seq_Allele2 is Alternative allele. Mutation IDs are created, e.g. '10 100003849 G A' is the mutation at chromosome 10, genomic location 100003849, where reference allele G is substituted with A, or '10 100011448 - CCGCTGCAAT' is the insertion of 'CCGCTGCAAT' at chromosome 10, location 100011448. The ref and alt alleles follow standard TCGA maf file notations.
<code>rem</code>	if TRUE, the mutationan matrix for random effect function will be prepared.

Details

if `rem=FALSE` (default), binary mutational matrix will be created, where each row is a mutation, each column is a sample, and values are binary taking values 0 if there is no mutation, 1 if there is a mutation in this sample. If `rem=TRUE` matrix with each possible pair of samples from the same patient will be created. Each row represents mutation, and each column - pair of samples, where value 0 denotes that mutation is not observed, 1 if shared mutation is observed in both tumors, and 2 if it is a private mutation observed in only one tumor.

Value

Data frame with matrix of mutations

Author(s)

Irina Ostrovnaya <ostrovni@mskcc.org>

References

Ostrovnyaya, Irina, Venkatraman E. Seshan, and Colin B. Begg. 2015. USING SOMATIC MUTATION DATA TO TEST TUMORS FOR CLONAL RELATEDNESS. *The Annals of Applied Statistics* 9 (3): 1533-48.

Examples

```
data(lcis)
#we want to analyze pair TK53IDC2.TK53LCIS2 that has only 1 shared mutation

mut.matrix<-create.mutation.matrix(lcis )
table(mut.matrix$TK53IDC2,mut.matrix$TK53LCIS2)

freq<-get.mutation.frequencies(rownames(mut.matrix),"BRCA")
SNVtest(mut.matrix$TK53IDC2,mut.matrix$TK53LCIS2,freq)
```

ECMtesting	<i>Clonality testing of ≥ 3 tumors using Extended Concordant Mutations (ECM) test based on LOH (Loss of Heterozygosity) profiles</i>
------------	--

Description

Function to test clonality of three and more tumors from the same patient based on their LOH profiles. This function implements Extended Concordant Mutations for all possible subsets of tumors from the same patient and minP multiplicity adjustment using simulated tumors.

Usage

```
ECMtesting(LOHtable,ptlist,noloh,loh1,loh2,Nsim=100)
```

Arguments

LOHtable	Matrix of LOH calls. Each row corresponds to a marker. First column contains the names of the markers. Each other column represents a sample and contains LOH calls.
ptlist	Vector of the patient IDs in the order the samples appear in the data. For example, if the first three tumors (columns 2, 3, 4 of data) belong to patient A, and the following two (columns 5, 6 of data) belong to patient B, then ptlist=c('ptA', 'ptA', 'ptA', 'ptB', 'ptB').
noloh	The string or a number that denotes absence of LOH.
loh1	The string or a number that denotes presence of LOH.
loh2	The string or a number that denotes presence of LOH that is discordant from loh1.
Nsim	Number of simulations used to calculate minP adjusted p-values

Details

Extended Concordant Mutations test is done for every subset of tumors. It uses number of concordant mutations in all tumors of the subset as a test statistic, and its reference distribution is calculated assuming fixed counts of LOH per tumor and equal probability of maternal and paternal alleles being affected. Note that ECM test for 2 tumors and original CM test will give slightly different p-values since continuity correction is done in ECM test.

Value

The function returns a list with number of elements equal to the number of patients. Each element is a matrix with two rows: ECM p-values for all possible subsets of tumors from this patient, and minP adjusted p-values. The tumors are denoted 1,2,3,... in the order they appear in LOHtable. Any tumor subsets with minP adjusted p-value ≤ 0.05 can be considered clonal.

References

Ostrovnyaya, I. "Testing clonality of three and more tumors using their loss of heterozygosity profiles", Statistical Applications in Genetics and Molecular Biology, 2012

Examples

```
set.seed(25)
LOHtable<-cbind(1:15,matrix(sample(c(0,1,2),15*12,replace=TRUE),ncol=12))
ECMtesting(LOHtable,rep(1:3,each=4),noloh=0,loh1=1,loh2=2,Nsim=100)
```

freqdata

TCGA pancancer mutation frequencies

Description

Mutational data from TCGA

Usage

```
data(freqdata)
```

Details

Object 'freqdata' contains the frequencies of the mutations observed in the exome sequencing data in 3 cancer types, COAD, LUAD, and BRCA, profiled by TCGA.

Full set of 33 cancer types is available by loading a full object in GitHub 'load(url("https://github.com/IOstrovnyaya/MutFreq/b After that proceed with the analysis in the same way.

There are 3 columns for 3 cancer types abbreviated in TCGA as COAD, LUAD, and BRCA (for the full list of abbreviations see <https://gdc.cancer.gov/resources-tcga-users/tcga-code-tables/tcga-study-abbreviations>). The first column is the number of patients profiled in each cancer type. Each subsequent row is a mutation, where number of patients with this particular mutation in each cancer

type is given. The mutation ID, contained in the row names of 'freqdata', is of the following format: "Chromosome Location RefAllele AltAllele", each entry separated by space, where chromosome is a number 1-22 or "X" or "Y"; location is genomic location in GRCh37 build; RefAllele is a reference allele and AltAllele is Alternative allele/Tumor_Seq_Allele2. For example "10 100003849 G A", is the mutation at chromosome 10, genomic location 100003849, where reference allele G is substituted with A, or "10 100011448 - CCGCTGCAAT" is the insertion of "CCGCTGCAAT" at chromosome 10, location 100011448. The ref and alt alleles follow standard TCGA maf file notations.

The code that was used to obtain freqdata is available in the vignette.

References

Ostrovnya, Irina, Venkatraman E. Seshan, and Colin B. Begg. 2015. USING SOMATIC MUTATION DATA TO TEST TUMORS FOR CLONAL RELATEDNESS. *The Annals of Applied Statistics* 9 (3): 1533-48.

Ellrott K, Bailey MH, Saksena G, Covington KR, Kandoth C, Stewart C, Hess J, Ma S, Chiotti KE, McLellan M, Sofia HJ, Hutter C, Getz G, Wheeler D, Ding L; MC3 Working Group; Cancer Genome Atlas Research Network, Scalable Open Science Approach for Mutation Calling of Tumor Exomes Using Multiple Genomic Pipelines. *Cell Syst.* 2018 Mar 28;6(3):271-281.e7. doi: 10.1016/j.cels.2018.03.002. PubMed PMID: 29596782.

genomewidePlots *Plot of the genomewide copy number profiles of a pair of tumors.*

Description

Plot contains genomewide profiles from a pair of tumors. It uses the output from the function `clonality.analysis()`.

Usage

```
genomewidePlots(data.seg1, classall, ptlist, ptpair, ptLR, plot.as.in.analysis = TRUE)
```

Arguments

<code>data.seg1</code>	Output of one-step segmentation - output <code>OneStepSeg</code> of <code>clonality.analysis()</code> . The chromosomes should be in the format "chr01p", "chr01q" etc.
<code>classall</code>	Classifications of the chromosomes - output <code>ChromClass</code> of <code>clonality.analysis()</code>
<code>ptlist</code>	Vector of the patient IDs in the order of the samples appearing in the data.
<code>ptpair</code>	Two sample names for which the plot is desired
<code>ptLR</code>	Matrix with the likelihood ratios - output <code>LR</code> of <code>clonality.analysis()</code>

`plot.as.in.analysis`

If TRUE then the gain/loss patterns will be highlighted in accordance with the chromosome classification. For example, if there are three segments in a chromosome, then the middle one determines the chromosome status. If it is normal, no color will be plotted in the chromosome even if the 1st and 3rd segments are gains or losses. Another example: if there are 2 or 3 different segments of gains, they will be combined and only one segment will be plotted. If `plot.as.in.analysis` is equal to FALSE, the original one-step CBS segmentation will be plotted.

Details

Function produces genomewide plots of a pair of tumors. The log-ratios are plotted in grey in the order of their genomic locations, gains are plotted in blue, and losses are plotted in red.

Examples

```
# See example as in clonality.analysis()
```

`get.mutation.frequencies`

Estimating mutation frequencies based on the TCGA data or a submitted reference mutations file

Description

Function to estimate the mutation frequencies in the specific cancer subtype using TCGA data or using a submitted reference mutations file

Usage

```
get.mutation.frequencies (xmut.ids, tcga.cancer.type=NULL, reference.data=NULL, combine.with.TCGA=FAL
```

Arguments

`xmut.ids` vector of mutation IDs for which frequencies are needed. Usually these are mutation IDs observed in one or two tumors from the same patient that will be tested for clonality. `xmut.ids` should have the following format: Chromosome Location RefAllele AltAllele, each entry separated by space, where chromosome is a number 1-22 or X or Y; location is genomic location in GRCh37 build; RefAllele is a reference allele and AltAllele is Alternative allele/Tumor_Seq_Allele2. For example '10 100003849 G A' is the mutation at chromosome 10, genomic location 100003849, where reference allele G is substituted with A, or '10 100011448 - CCGCTGCAAT' is the insertion of 'CCGCTGCAAT' at chromosome 10, location 100011448. The ref and alt alleles follow standard TCGA maf file notations.

`tcga.cancer.type`

String denoting one of 33 TCGA cancer types: ACC, BLCA, BRCA, CESC, CHOL, COAD, DLBC, ESCA, GBM, HNSC, KICH, KIRC, KIRP, LAML, LGG, LIHC, LUAD, LUSC, MESO, OV, PAAD, PCPG, PRAD, READ, SARC, SKCM, STAD, TGCT, THCA, THYM, UCEC, UCS, or UVM. See <https://gdc.cancer.gov/resources-tcga-users/tcga-code-tables/tcga-study-abbreviations> for details. If `reference.data` is supplied, `cancer.type` can be NULL

`reference.data`

is a matrix that contains the mutations in external dataset from which the frequencies can be estimated. Maf files (Mutation Annotation Format) can be used here. Each line represents a mutation, and the matrix should include the following columns: 'PatientID', 'Tumor_Sample_Barcode' (used as a sample ID), and Chromosome, Start_Position, Reference_Allele, and Tumor_Seq_Allele2 that are used to create mutation IDs.

If the frequencies are estimated for a specific patient with mutations 'xmut.ids', then this matrix should not contain the data from this patient. In this case the denominator for the frequency will be (1+n) where n is number of patients in a reference. Alternatively, if the list of mutations in 'xmut.ids' is the same as list of all mutations in `reference.data` then the denominator will be n. If the reference dataset contains multiple samples from the same patient, union of them will be taken so that clonal mutations are not counted twice - denominator is the number of patients and not the number of samples.

Another word of caution that if you are planning to compare pairs of tumors for clonality profiled by exome sequencing, the reference dataset should also be from exome sequencing (like in TCGA), to avoid counting mutations in the genes that were not sequenced by targeted panel for example as non-observed/never mutated.

The possible reasons to use `reference.data` include: 1) the dataset you analyze is large and combining it with tcga data will improve the frequency; 2) you analyze tumors from cancer type that is not among 33 pancancer TCGA sites; or 3) you want to use specific subtype of cancers, for example Luminal A breast tumors, for which overall TCGA breast cancer frequencies are not applicable.

`combine.with.TCGA`

TRUE if you want to combine mutations in reference dataset with TCGA for frequency calculation. In this case you need to both specify 'cancer.type' and 'reference.data'. This option makes sense when the dataset you assembled for clonality study is large and therefore will improve frequency calculation. Default is FALSE.

Details

For each mutation in 'xmut.ids' we calculate mutation frequency in the following way.

If `tcga.cancer.type` is chosen and `reference.data` is NULL, we calculate X1, number of patients that have that mutation, among n1 TCGA samples of specific subtype. Mutation frequency then is $(1+X1)/(1+n1)$, where 1 is added to denominator and nominator to incorporate data in a current patient.

If `reference.data` is specified and `tcga.cancer.type` is not chosen, we calculate X2 - number of patients with a particular mutation among n2 patients (not samples) in `reference.data`. Then there

are two possibilities. If the list of mutations in 'xmut.ids' is the same as list of all mutations in reference.data, we assume the reference data is the dataset for which clonality is under question, and the frequency is defined as $X2/n2$. Otherwise if 'xmut.ids' is not the same as list of all mutations in reference.data then we assume it's an external data and like in TCGA the frequency then is $(1+X2)/(1+n2)$.

If both tcga.cancer.type and reference.data are specified and combine.with.TCGA is TRUE, then the frequency is defined as either $(1+X1+X2)/(1+n1+n2)$ or $(X1+X2)/(n1+n2)$ depending on whether the list of mutations in 'xmut.ids' is the same as list of all mutations in reference.data.

Value

Vector of frequencies with names same as 'xmut.ids'

Author(s)

Irina Ostrovnaya <ostrovni@mskcc.org>

References

Ostrovnaya, Irina, Venkatraman E. Seshan, and Colin B. Begg. 2015. USING SOMATIC MUTATION DATA TO TEST TUMORS FOR CLONAL RELATEDNESS. The Annals of Applied Statistics 9 (3): 1533-48.

Examples

```
data(lcis)

#Analysis of data from patient 53
mut.matrix<-create.mutation.matrix(lcis )
table(mut.matrix$TK53IDC2,mut.matrix$TK53LCIS2)

freq<-get.mutation.frequencies(rownames(mut.matrix),"BRCA")

SNVtest(mut.matrix$TK53IDC2,mut.matrix$TK53LCIS2,freq)
```

grid.lik

Auxiliary function: Grid of conditional probabilities

Description

This auxiliary function generates the grid of likelihood values for each tumor pair (rows) and each value of xi (columns): $P(\text{observed mutations} \mid x_i)$

Usage

```
grid.lik(xigrid, mutns, probamut)
```

Arguments

xigrid	Grid of the values of xi, corresponding to its domain of definition.
mutns	Matrix of the mutations observed, with all mutations in rows and the cases (tumor pairs) in columns. The data are coded as 0=mutation not observed, 1=shared mutation (observed in both tumors), 2=private mutation (observed in one tumor only).
probamut	Vector of the probabilities of occurrence for each mutation.

Value

Return the matrix of the likelihood values for each tumor pair (rows) and each value of xi (columns). This matrix is called by the auxiliary function grid.lik, returned as a parameter by the function clonal.est, and used as a parameter by the function clonal.proba.

histogramPlot	<i>Histograms of Log-Likelihood Ratios</i>
---------------	--

Description

Function produces the histograms of the within-patient and between-patient log-Likelihood Ratios.

Usage

```
histogramPlot(ptLRvec, refLRvec)
```

Arguments

ptLRvec	Vector with the within-patient likelihood ratios - output LR of clonality.analysis()
refLRvec	Vector with the between-patient likelihood ratios - output refLR of clonality.analysis()

Details

Functions plots two overlapping histograms: within-patient log-likelihood ratios are in red and between-patient log-likelihood ratios (reference distribution under the hypothesis of independence) are in black.

Examples

```
# See example as in clonality.analysis()
```

 lcis

Breast cancer data

Description

Data from the LCIS study, with mutations listed for all pairs of LCIS-invasive tumors in a series of breast cancer cases.

Usage

```
data(lcis)
```

Details

This is exome sequencing data from study of Lobular Carcinoma in Situ (LCIS) and Invasive lobular carcinomas (ILC) or Invasive Ductal Carcinomas (IDC) in the same patients. Each row represents single mutation in a patient. This structure is similar to MAF (mutation annotation format) files.

The columns are "PatientID, Lesion,Hugo_Symbol, Chromosome, Start_Position, Reference_Allele, Tumor_Seq_Allele2 (Alternative allele), Tumor_Sample_Barcode (sample ID). Lesion contains Lesion type, which takes values DCIS (ductal carcinoma in situ), IDC (invasive ductal carcinoma), IDC2 (second profiled IDC), ILC (invasive lobular carcinoma), LCIS and LCIS1,2, etc.

References

Begg CB, Ostrovnaya I, Carniello JV, Sakr RA, Giri D, Towers R, Schizas M, De Brot M, Andrade VP, Mauguen A, Seshan VE, King TA. "Clonal relationships between lobular carcinoma in situ and other breast malignancies.", *Breast Cancer Res.* 2016 Jun 23;18(1):66. doi: 10.1186/s13058-016-0727-z.

 LOHclonality

Clonality testing using LOH (Loss of Heterozygosity) profiles

Description

Function to test clonality of two tumors from the same patient based on their LOH profiles. This function implements Concordant Mutations and Likelihood Ratio tests.

Usage

```
LOHclonality(LOHtable, ptlist, refLOHtable = NULL, pfreq = NULL, noloh, loh1, loh2,method="both")
```

Arguments

LOHtable	Matrix of LOH calls. Each row corresponds to a marker. First column contains the names of the markers. Each other column represents a sample and contains LOH calls.
ptlist	Vector of the patient IDs in the order the samples appear in the data. For example, if the first three tumors (columns 3, 4, 5 of data) belong to patient A, and the following two (columns 6, 7 of data) belong to patient B, then ptlist=c('ptA', 'ptA', 'ptA', 'ptB', 'ptB').
refLOHtable	Matrix of LOH calls that should be used to calculate the LOH frequencies used in Likelihood Ratio calculation. The structure is similar to LOHtable. If refLOHtable is not specified, frequencies are calculated from LOHtable.
pfreq	Vector of LOH frequencies known from the literature. Should be in the same order as the markers in LOHtable. If pfreq is not specified, frequencies are calculated from LOHtable.
no loh	The string or a number that denotes absence of LOH.
loh1	The string or a number that denotes presence of LOH.
loh2	The string or a number that denotes presence of LOH that is discordant from loh1.
method	Takes values "CM", "LR" or "both" if only Concordant Mutations test, or only Likelihood Ratio test, or both should be performed. Default value is "both".

Details

Function tests clonality of LOH profiles of tumors from the same patient using two tests. Concordant Mutations test has number of markers with concordant LOH as its test statistic. Its theoretical reference distribution under independence is calculated assuming that the maternal and paternal alleles are equally likely to be lost and that the frequencies of LOH are about the same across different markers.

Likelihood Ratio test uses pre-specified frequencies of LOH to compute Likelihood Ratio statistic. Its reference distribution is obtained by simulating tumors with the given LOH probabilities, and probability of maternal/paternal mutation estimated from the data. If LOH frequencies are not specified then they are estimated from the data.

Value

The function returns a data frame where each row corresponds to the pair of samples that are compared. Columns are: Sample1 - name of sample 1; Sample2 - name of sample 2; a - number of markers with concordant LOH in both tumors (test statistic for Concordant Mutations test); e - number of markers with LOH in both tumors, concordant or discordant; f - number of markers with LOH in the first tumor and Normal in the 2nd tumor; g - number of markers with LOH in the second tumor and Normal in the first tumor; h - number of markers that are Normal in both tumors; Ntot - total number of informative markers for both tumors; CMpvalue - p-value for Concordant Mutations test; LRpvalue - p-value for Likelihood Ratio test.

References

- Begg CB, Eng KH, Hummer AJ. Statistical tests for clonality. *Biometrics* 2007; 63:522-530
- Ostrovnyaya I, Seshan VE, Begg CB. Comparison of properties of tests for assessing tumor clonality. *Biometrics* 2008; 68:1018-1022.

Examples

```
set.seed(25)
LOHtable<-cbind(1:20,matrix(sample(c(0,1,2),20*20,replace=TRUE),20))
LOHclonality(LOHtable,rep(1:10,each=2),pfreq=NULL,no loh=0,loh1=1,loh2=2)
```

LRtesting3or4tumors *Clonality testing of 3 or 4 tumors using Likelihood model based on LOH (Loss of Heterozygosity) profiles*

Description

Function to test clonality of 3 or 4 tumors from the same patient based on their LOH profiles.

Usage

```
LRtesting3or4tumors(LOHtable,ptlist,refLOHtable=NULL, pfreq=NULL,no loh,loh1,loh2,Nsim=100,m=0.5)
```

Arguments

LOHtable	Matrix of LOH calls. Each row corresponds to a marker. First column contains the names of the markers. Each other column represents a sample and contains LOH calls.
ptlist	Vector of the patient IDs in the order the samples appear in the data. For example, if the first three tumors (columns 2, 3, 4 of data) belong to patient A, and the following two (columns 5, 6 of data) belong to patient B, then ptlist=c('ptA', 'ptA', 'ptA', 'ptB', 'ptB').
refLOHtable	Matrix of LOH calls that should be used to calculate the LOH frequencies used in Likelihood Ratio calculation. The structure is similar to LOHtable. If refLOHtable is not specified, frequencies are calculated from LOHtable.
pfreq	Vector of LOH frequencies known from the literature. Should be in the same order as the markers in LOHtable. If pfreq is not specified, frequencies are calculated from LOHtable.
no loh	The string or a number that denotes absence of LOH.
loh1	The string or a number that denotes presence of LOH.
loh2	The string or a number that denotes presence of LOH that is discordant from loh1.
Nsim	Number of simulations used to calculate minP adjusted p-values
m	Probability that a favored allele is affected given that LOH has occurred. Must be a number above 0.5 (equal probability of maternal and paternal allelic loss)

Details

Likelihood ratio test for 3 and 4 tumors. For 3 tumors there are 3 possible tumor orderings, and for 4 tumors there are 2 topologies with 3 and 12 orderings each. The test calculates maximum likelihood ratio across all possible orderings, and the p-value is calculated using simulated reference distribution.

Value

The function returns a list with number of elements equal to the number of patients. Each element is list with two elements. First contains log maximum likelihood ratio value, p-value, and estimates of parameters c , the topology and tumor ordering that have maximum likelihood ratio. If p-value is significant, then the null hypothesis that all tumors are independent can be rejected. The second element has a matrix with all possible topologies and tumor orderings and their corresponding log likelihood ratios.

References

Ostrovnaya, I. "Testing clonality of three and more tumors using their loss of heterozygosity profiles", Statistical Applications in Genetics and Molecular Biology, 2012

Examples

```
set.seed(25)
LOHtable<-cbind(1:15,matrix(sample(c(0,1,2),15*12,replace=TRUE),ncol=12))
q<-LRtesting3or4tumors(LOHtable,rep(1:4,each=3),refLOHtable=NULL, pfreq=NULL,no loh=0,loh1=1,loh2=2,Nsim=100,m=
```

model.lik

Auxiliary likelihood Function

Description

This function computes the likelihood of the model.

Usage

```
model.lik(para, likmat, out0, xigrid)
```

Arguments

para	Value of the model parameters, in the form $c(\mu, \sigma, \pi)$.
likmat	Grid of conditional probabilities for each tumor pair (rows) and each value of x_i (columns). This matrix is generated by the function <code>grid.lik</code> .
out0	a small value that is used when the likelihood goes to infinite values, posing problem for the maximization. The corresponding combination of the parameters will thus be excluded from the search.
xigrid	Grid of the values of x_i , corresponding to its domain of definition.

Value

Return the likelihood value of the model for the given parameters.

mutation.proba	<i>Probability of being clonal</i>
----------------	------------------------------------

Description

This function uses the results from mutation.rem to estimate the diagnostic probability of clonal relatedness for new cases. It is obtained from Bayes theorem using the prior probability of clonal relatedness (π) and the contributions to the likelihood based on the mutations observed for the case. We recommend to use this function to estimate probabilities of clonality for new subjects, ie who are not used for the model estimation. To obtain estimate for the subjects on which the model estimation is based, the option "proba=TRUE" can be used in the mutation.rem function.

Usage

```
mutation.proba(para, likmat, xigrid = c(0, seq(0.0005, 0.9995, by=0.001)))
```

Arguments

para	Value of the model parameters, in the form $c(\mu, \sigma, \pi)$.
likmat	Grid of conditional probabilities for each tumor pair (rows) and each value of x_i (columns). This matrix is generated by the auxiliary function grid.lik, and returned as a parameter by the principal function mutation.rem.
xigrid	Grid of the values of x_i , corresponding to its domain of definition. The default is $c(0, \text{seq}(0.0005, 0.9995, \text{by}=0.001))$.

Value

Returns the vectors of probability of clonality for each pairs of tumors contained in the matrix likmat (the number of pairs is the number of rows of the matrix).

Author(s)

Audrey Mauguen <mauguena@mskcc.org> and Venkatraman E. Seshan.

References

Mauguen A, Seshan VE, Ostrovnya I, Begg CB. Estimating the Probability of Clonal Relatedness of Pairs of Tumors in Cancer Patients. Submitted.

Examples

```
#__ Analysis of LCIS data
data(lcis)
mut.matrix<-create.mutation.matrix(lcis ,rem=TRUE)
freq<-get.mutation.frequencies(rownames(mut.matrix), "BRCA")

#__ Parameters estimation
mod <- mutation.rem(cbind(freq, mut.matrix))
mod

#__ Probability of being clonal for a new subject
# generate a case with 30 mutations
# probabilities of each observed mutation
pi <- runif(30,0.001,0.13)
# mutation 1=shared or 2=private
newpair <- cbind(pi,rbinom(30,1,1-pi^2)+1)
# generate the matrix of likelihood values
new.likmat <- grid.lik(xigrid=c(0, seq(0.0005, 0.9995, by=0.001)), as.matrix(newpair[,c(-1)]), newpair[,1])
# probability of being clonal using the model previoulsy estimated
proba <- mutation.proba(c(mod$mu, mod$sigma, mod$pi), t(as.matrix(new.likmat)) )
```

mutation.rem

Estimation of the random-effect model for clonality based on mutations.

Description

The model estimates the proportion of clonal cases in a population, and the distribution of the clonality signal.

Usage

```
mutation.rem(mutmat, proba=FALSE, print.proba=FALSE, xigrid = seq(0.0005, 0.9995, by=0.001), init.para
```

Arguments

mutmat	Matrix containing the data, with all mutations in rows and the tumor pairs in columns. The data are coded as 0=mutation not observed, 1=shared mutation (observed in both tumors), 2=private mutation (observed in one tumor only). The first column contains the probabilities of occurrence for each mutation.
proba	Indicates whether to compute the individual probabilities of clonality for each pair. The default is FALSE.
print.proba	Indicates whether the individual probabilities of clonality should be printed in the output. The default is FALSE.
xigrid	Grid of the values of xi used to compute the integration; it corresponds to the domain of definition of xi. The default is seq(0.0005, 0.9995, by=0.001).

init.para	Initial values of the parameters for the optimization. The order of the parameters is $c(\mu, \sigma, \pi)$, where μ and σ are the mean and variance of the lognormal distribution of the random-effect ξ , and π is the proportion of clonal cases. The default is $c(0,1,0.5)$.
conv.crit	Criteria used to defined convergence (on all three parameters). The default is $1e-5$.
niter	Maximum number of iterations used for the EM algorithm. The default is 300.

Details

The function estimates a random effects model in which the random effect (the clonality signal, denoted ξ_i for the i th case) reflects the somatic similarity of the tumors on a scale from 0 to 1, where 0 represents independence and higher values represent clonal tumors that are increasingly similar. The proportion of cases that are clonal is represented by the parameter π . Thus the likelihood is a compound of $(1-\pi)$ cases that have a clonality signal of exactly 0, and π cases that have a clonality signal drawn from a lognormal random effects distribution with mean μ and variance σ^2 . The program estimates the parameters using an EM algorithm to maximize the likelihood. The output provides parameter estimates (μ, σ, π). The example dataset presented contains data from a study in which each patient has both a pre-malignant lobular carcinoma in situ (LCIS) and an invasive breast cancer, and we wish to estimate the proportion of these cases for which the LCIS was a direct precursor to the invasive cancer.

Value

mu	Estimated mean of the random-effect distribution.
sigma	Estimated standard-deviation of the random-effect distribution.
pi	Estimated proportion of clonal pairs in the population.
likmat	Grid of likelihood values for each tumor pair (rows) and each value of ξ (columns) needed for the function <code>clonal.proba</code> that computes the individual probabilities of clonality.
likelihood	Value of the maximized likelihood.
convergence	Convergence status, 0=no convergence, 1=convergence reached.
n.iter	Number of iterations used.
pr.clonal	Individual probabilities of clonality.

Author(s)

Audrey Mauguen <mauguena@mskcc.org> and Venkatraman E. Seshan.

References

Mauguen A, Seshan VE, Ostrovnaya I, Begg CB. Estimating the Probability of Clonal Relatedness of Pairs of Tumors in Cancer Patients. *Biometrics* 2018;74(1):321-330.

Examples

```
#__ Analysis of LCIS data
data(lcis)
mut.matrix<-create.mutation.matrix(lcis ,rem=TRUE)
freq<-get.mutation.frequencies(rownames(mut.matrix),"BRCA")

#__ Parameters estimation
mod <- mutation.rem(cbind(freq, mut.matrix))
mod

#__ Probability of being clonal
mod <- mutation.rem(cbind(freq, mut.matrix), proba=TRUE)
mod
```

print.mutation.proba *Print for the mutation.proba function*

Description

Print a summary of results for the probabilities of clonality estimated by the mutation.proba function.

Usage

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

Arguments

x a mutation.proba object
... Other unused arguments.

Value

Print results for the individual probabilities of clonality.

See Also

mutation.proba

```
print.mutation.rem      Print for the mutation.rem function
```

Description

Print a summary of results for the random-effect model estimation estimated by the `clonal.est` function.

Usage

```
## S3 method for class 'mutation.rem'

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

Arguments

```
x          a mutation.rem object
...        Other unused arguments.
```

Value

Print results for the model estimates.

See Also

`mutation.rem`

```
SNVtest      Testing relatedness (clonality) of two tumors from the same patient
              using profiles of somatic mutations
```

Description

Function to test clonality of two tumors from the same patient based on their mutational profiles. This function calculates conditional likelihood ratio relying only on loci where at least one of the tumors has a mutation, and p-value is calculated under the reference distribution under the hypothesis of independence.

Usage

```
SNVtest(tumor1, tumor2, pfreq, nrep = 1000)
```

Arguments

tumor1	Vector of the binary mutation calls from tumor 1, where 0 denotes no mutation, 1 denotes a mutation. Mutations should be in the same order as frequencies in pfreq.
tumor2	Vector of the binary mutation calls from tumor 2, where 0 denotes no mutation, 1 denotes a mutation. Mutations should be in the same order as frequencies in pfreq.
pfreq	Marginal frequencies of mutations known a priori. These can be obtained from TCGA or similar databases. We recommend setting these frequencies to $(x+y)/(n_x+n_y)$, where x is the number of patients with the mutations in the TCGA (or other databases), and n_x is the total number of the patients in TCGA; y and n_y is number of patients with mutations and total number of patients in this study.
nrep	Number of simulations used for generating the reference distribution under the hypothesis of independence.

Details

Only loci where at least one tumor has a mutation contribute to the model. The null distribution is patient specific since it is generated assuming the same total number of mutations in two tumors.

Value

The output is a vector with 5 values: "n1", "n2", "n_match", "LRstat", "maxKsi", "LRpvalue"

n1	Number of mutations in the first tumor.
n2	Number of mutations in the second tumor
n_match	Number of matches. i.e. loci where both tumors have an identical mutation
LRstat	Likelihood ratio statistic
maxKsi	Maximum likelihood estimate of Ksi, parameter of the likelihood representing clonality strength. Value close to 0 indicates independence, value close to 1 indicates perfect concordance in mutational profiles.
LRpvalue	p-value calculated using the null distribution generated using prespecified mutational frequencies pfreq.

Author(s)

Irina Ostrovnaya <ostrovni@mskcc.org>

References

Ostrovnyaya I, Seshan VE, Begg CB. "USING SOMATIC MUTATION DATA TO TEST TUMORS FOR CLONAL RELATEDNESS.", *Ann Appl Stat.* 2015 Sep;9(3):1533-1548

See Also

clonality.analysis() for test using genomewide copy number profiles; mutation.proba() for bayesian inference of clonality probability.

Examples

```
#___ Analysis of LCIS data from the following paper:
#Begg CB, Ostrovnaya I, Carniello JV, Sakr RA, Giri D, Towers R, Schizas M, De Brot M, Andrade VP, Mauguen A, Seshan V

data(lcis)
data(freqdata)
n<-nrow(lcis)

#Example of artificially generated independent tumor pair with marginal mutation frequencies p
n<-100
p<-runif(n)/10
x1<-as.numeric(runif(n)<=p)
x2<-as.numeric(runif(n)<=p)
SNVtest(x1,x2,p)

#Analysis of data from patient 53
mut.matrix<-create.mutation.matrix(lcis )
table(mut.matrix$TK53IDC2,mut.matrix$TK53LCIS2)

freq<-get.mutation.frequencies(rownames(mut.matrix),"BRCA")

SNVtest(mut.matrix$TK53IDC2,mut.matrix$TK53LCIS2,freq)
```

 SNVtest2

Test for tumors from 2 different sites

Description

This functions performs clonality testing of 2 tumors from the same patient that come from different types or organs. The test uses profiles of somatic mutations. The null hypothesis is that two tumors that come from different sites are independent. There are two alternative hypotheses: that they are clonal and from site 1, and clonal from site 2. This function calculates conditional maximum likelihood ratio relying only on loci where at least one of the tumors has a mutation, and p-value is calculated under the reference distribution under the hypothesis of independence.

Usage

```
SNVtest2(tumor1, tumor2, pfreq, nrep = 999)
```

Arguments

tumor1 Vector of the binary mutation calls from tumor 1, where 0 denotes no mutation, 1 denotes a mutation. Mutations should be in the same order as frequencies in pfreq.

tumor2	Vector of the binary mutation calls from tumor 2, where 0 denotes no mutation, 1 denotes a mutation. Mutations should be in the same order as frequencies in pfreq.
pfreq	Two column matrix of marginal frequencies of mutations in two sites known a priori. These can be obtained from TCGA or similar databases or calculated using function "get.mutation.frequencies" in this package.
nrep	Number of simulations used for generating the reference distribution under the hypothesis of independence.

Details

Test is related to the one described by "SNVtest" but it assumes that under null hypothesis two tumors come from 2 different sites with different marginal probabilities, thus the reference distribution is generated from 2 sets of frequencies. The p-value is significant when null hypothesis is rejected. It might be rejected sometimes in the absence of matches if the observed mutational profiles are unlikely to come from the declared tumor sites.

Value

The output is a vector with 6 values: "n.match", "n.site1only", "n.site2only", "xi.site1", "xi.site2", "p.value"

n.match	Number of matches between two tumors
n.site1only	Number of mutations in the first tumor only
n.site2only	Number of mutations in the second tumor only
xi.site1	Maximum likelihood estimate of Ksi, parameter of the likelihood representing clonality strength, for the alternative hypothesis that two clonal tumors come from site 1. Value close to 0 indicates independence, value close to 1 indicates perfect concordance in mutational profiles.
xi.site2	Maximum likelihood estimate of Ksi, parameter of the likelihood representing clonality strength, for the alternative hypothesis that two clonal tumors come from site 2. Value close to 0 indicates independence, value close to 1 indicates perfect concordance in mutational profiles.
p.value	p-value calculated using the null distribution generated using two independent tumors generated from two different sites

Author(s)

Irina Ostrovnaya <ostrovni@mskcc.org>

References

Ostrovnyaya I, Mauguen A, Seshan V, Begg CB "Testing Tumors from Different Anatomic Sites for Clonal Relatedness Using Somatic Mutation Data", submitted

Examples

```
#Example of artificially generated independent tumors from marginal mutation frequencies p1 and p2
n<-100
p1<-runif(n)/10
p2<-runif(n)/10
x1<-as.numeric(runif(n)<=p1)
x2<-as.numeric(runif(n)<=p2)
SNVtest2(x1,x2,cbind(p1,p2))
```

splitChromosomes	<i>Chromosome splitting</i>
------------------	-----------------------------

Description

Divides the chromosomes into p and q arms.

Usage

```
splitChromosomes(chrom,maploc)
```

Arguments

chrom	Vector of chromosomes. They should be numeric 1 to 22.
maploc	Vector of genomic locations. They should be in Kilobases.

Details

The function returns the vector of chromosome arms labeled "chr01p", "chr01q", etc. The split into arms is accomplished using the following centers (in Kb) for chromosomes 1 through 22: (122356.96, 93189.90, 92037.54 , 50854.87 ,47941.40, 60438.12 , 59558.27, 45458.05 , 48607.50, 40434.94 , 52950.78, 35445.46 , 16934.00, 16570.00, 16760.00 , 36043.30 , 22237.13, 16082.90 , 28423.62 , 27150.40, 11760.00, 12830.00).

Examples

```
#simulated data

set.seed(100)
chrom<-rep(c(1:22),each=100)
maploc<- runif(2200)* 200000
chromarm<-splitChromosomes(chrom,maploc)
```

`xidens`*Auxiliary function computing the density of xi*

Description

Density function for the random variable x_i , using a lognormal density for $\phi_i = -\log(1-x_i)$

Usage

```
xidens(pmu, psig, xigrid)
```

Arguments

<code>pmu</code>	Mean parameter of the distribution.
<code>psig</code>	Variance parameter of the distribution.
<code>xigrid</code>	Grid of the values of x_i , corresponding to its domain of definition.

Value

Returns the density value for the given values of x_i .

Index

ave.adj.probes, 3

chromosomePlots, 4
Clonality (Clonality-package), 2
Clonality-package, 2
clonality.analysis, 5
clonEM, 10
create.mutation.matrix, 11

ECMtesting, 12

freqdata, 13

genomewidePlots, 14
get.mutation.frequencies, 15
grid.lik, 17

histogramPlot, 18

lcis, 19
LOHclonality, 19
LRtesting3or4tumors, 21

model.lik, 22
mutation.proba, 23
mutation.rem, 24

print.mutation.proba, 26
print.mutation.rem, 27

SNVtest, 27
SNVtest2, 29
splitChromosomes, 31

xidens, 32