

Clonality: A Package for Clonality testing

Irina Ostrovnaya

April 24, 2017

Department of Epidemiology and Biostatistics
Memorial Sloan-Kettering Cancer Center
ostrovni@mskcc.org

Contents

1 Overview	1
2 Copy number profiles	1
2.1 Choice of segmentation algorithm	7
3 LOH data	8
4 LOH data for 3 and more tumors	9
5 Inference using profiles of somatic mutations	9
5.1 Likelihood model	9
5.2 Random effects model	11

1 Overview

This document presents an overview of the `Clonality` package. This package can be used to test whether two tumors are clonal (metastases) or independent (double primaries) using their somatic mutations, copy number or loss of heterozygosity (LOH) profiles. For LOH data it implements Concordant Mutations (CM) test (Begg et al., 2007) and Likelihood Ratio (LR) test (Ostrovnaya et al., 2008). For copy number profiles the package implements the methodology based on the likelihood ratio described in (Ostrovnaya et al., 2010). For somatic mutations we included the methods described in (Ostrovnaya et al., 2015) and (Mauguen et al., 2017).

2 Copy number profiles

We will show how to test independence of the copy number profiles from the same patient using simulated data. First we simulate the dataset with 10 pairs of tumors with 22 chromosomes, 100 markers each. Simulated log-ratios are equal to signal + noise. The signal

is defined in the following way: each chromosome has 50% chance to be normal, 30% to be whole-arm loss/gain, and 20% to be partial arm loss/gain, where endpoints are drawn at random, and loss/gain means are drawn from standard normal distribution. There are no chromosomes with recurrent losses/gains. Noise is drawn from normal distribution with mean 0, standard deviation 0.4 and added to the signal. First 9 patients have independent tumors, while last patient has two tumors with identical signal, but independent noise.

```

> library(Clonality)
> 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)
> samnms<-paste("pt",rep(1:10,each=2),rep(1:2,10),sep=".")
>

```

Rows of data correspond to probes (genomic markers). The first column is the chromosome and the second column is probe's genomic position. All subsequent columns correspond to the samples and contain log-ratios. Here the genomic is an index, but normally it would be actual probe's location along the genome, and then 'splitChromosomes' function should be used to divide the chromosome into p and q arms, thus increasing the number of independent units for the analysis.

```

> dim(data)

[1] 2200  20

```

As the next step of data preparation, we have to create a CNA (copy number array) object as described DNACopy.

```

> dataCNA<-CNA(data,chrom=chrom,maploc=maploc,sampleid=samnms)
> as.matrix(dataCNA)[1:5,1:10]

  chrom  maploc pt.1.1      pt.1.2      pt.2.1
1 "chr01p" " 1" " 5.229029e-01" " 0.2959505888" "-3.479070e-01"
2 "chr01p" " 2" " 1.787454e-01" "-0.0747496473" " 3.863461e-01"
3 "chr01p" " 3" "-3.404918e-01" " 0.2797033500" " 1.739630e-01"
4 "chr01p" " 4" "-4.191789e-01" " 0.3877484789" " 2.237324e-01"
5 "chr01p" " 5" " 1.597503e-03" " 0.6996900997" "-1.257982e-01"
  pt.2.2      pt.3.1      pt.3.2      pt.4.1
1 " 3.365784e-01" " 0.5740303360" "-0.138725302" "-5.097874e-01"
2 "-2.887743e-01" " 0.1649959341" " 0.643577307" "-1.060686e-01"
3 " 8.558146e-02" " 0.3676130117" "-0.263964372" "-1.285919e-01"
4 "-3.226487e-01" " 0.3157696773" "-0.936102251" "-6.465105e-01"
5 " 1.276517e-01" " 0.4961418049" "-0.126477964" "-2.908564e-01"
  pt.4.2
1 " 1.196744e-01"
2 "-4.829711e-01"
3 "-1.506162e-01"
4 "-1.710028e-01"
5 "-2.945648e-01"
>

```

Our methodology allows at most one genomic change per chromosome arm, estimated by the one-step Circular Binary Segmentation (CBS) algorithm ((Venkatraman and Olshen, 2007)).

If the data had many more than 15,000 markers, most outstanding, and likely a short change would be picked up, which would not be representative of the chromosome pattern. To avoid this, one can use the following function:

```
> dataAve<- ave.adj.probes(dataCNA,2)
```

Total number of markers after averaging is 1100

Here we have averaged every two consecutive markers. For this dataset, though, averaging is not necessary.

Next we have to create a vector of patient labels that matches the samples.

```
> ptlist<- paste("pt",rep(1:10,each=2),sep=".")
```

Finally, we can run the clonality analysis:

```
> results<-clonality.analysis(dataCNA, ptlist, pfreq = NULL, refdata = NULL, nmad = 1, re
```

Calculating LR.....

Calculating reference LR: %completed 10, 20, 30, 40, 50, 60, 70, 80, 90, 100,

The main information is in the output LR:

```
> results$LR
```

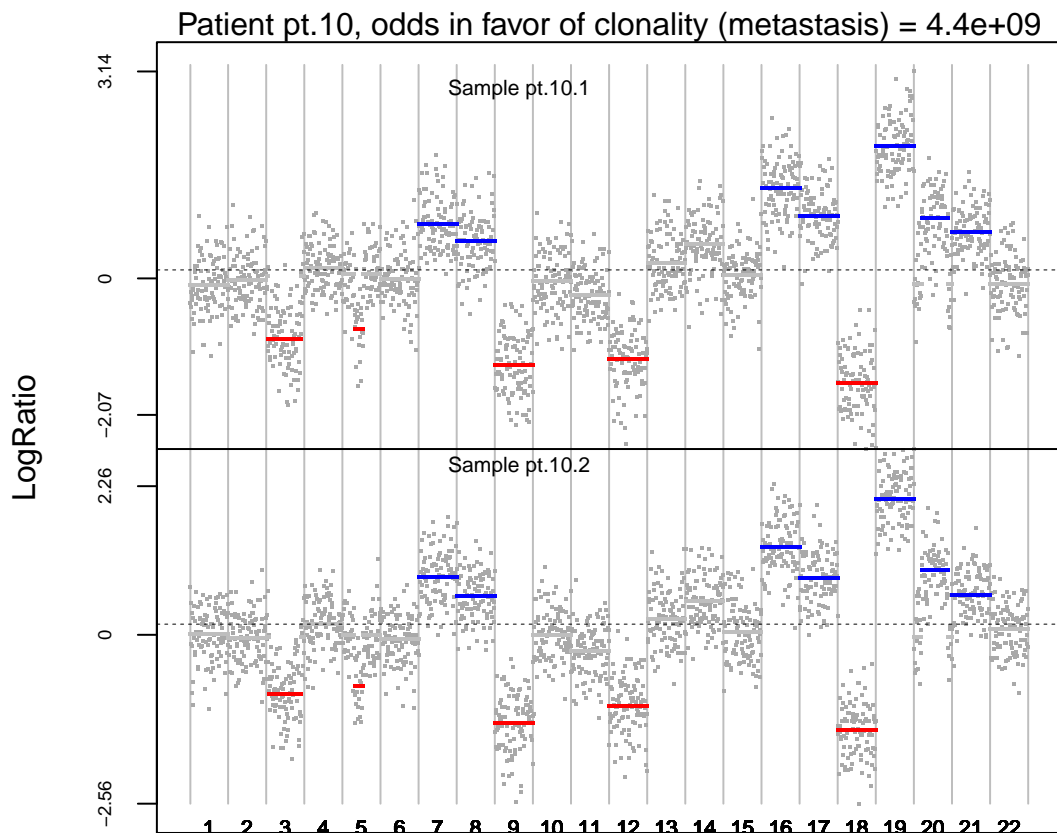
	Sample1	Sample2	LR1	LR2	GGorLL	NN	GL	GNorLN
1	pt.1.1	pt.1.2	2.682477e-02	2.682477e-02	0	12	0	10
2	pt.2.1	pt.2.2	2.830516e-02	4.816474e-03	1	9	0	12
3	pt.3.1	pt.3.2	7.263437e-03	7.263437e-03	0	9	0	13
4	pt.4.1	pt.4.2	1.793088e-01	1.793088e-01	2	10	0	10
5	pt.5.1	pt.5.2	1.897357e-02	2.700118e-03	2	8	3	9
6	pt.6.1	pt.6.2	7.441437e-03	7.441437e-03	0	11	2	9
7	pt.7.1	pt.7.2	1.084280e+00	1.784246e-01	4	8	1	9
8	pt.8.1	pt.8.2	1.350562e-01	1.350562e-01	1	15	3	3
9	pt.9.1	pt.9.2	9.918617e-03	9.918617e-03	1	8	1	12
10	pt.10.1	pt.10.2	5.790525e+04	4.402231e+09	12	10	0	0

	IndividualComparisons	LR2pvalue
1		0.3944444
2	chr20p 0.17	0.7944444
3		0.7388889
4		0.1555556
5	chr15p 0.14	0.8777778
6		0.6833333
7	chr18p 0.16	0.1555556
8		0.2000000
9		0.6388889
10	chr03p 27.5; chr05p 52.81; chr20p 52.34	0.0000000

The likelihood ratios LR2 for patients 1:9 are much smaller than 1, therefore these tumors are independent. Patient 10 has LR2 much higher than one, and we can conclude that this patient's tumors are clonal. The reference distribution for LR2 under the hypothesis of independence is constructed by pairing tumors from different patients that are independent by default. The p-value column reflects the percentiles of a particular patient's LR2 in the reference distribution: clonal tumors would have small p-values.

We can view the genomewide plots of patient 10 using:

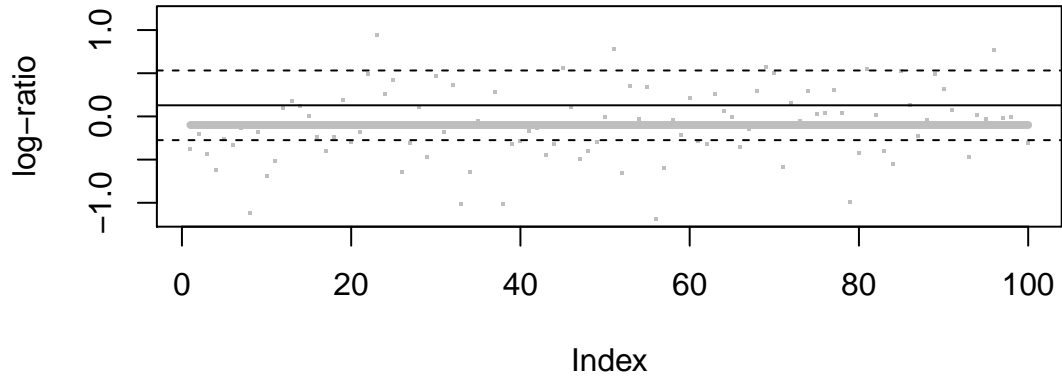
```
> genomewidePlots(results$OneStepSeg, results$ChromClass, ptlist , c("pt.10.1", "pt.10.2"), r
```



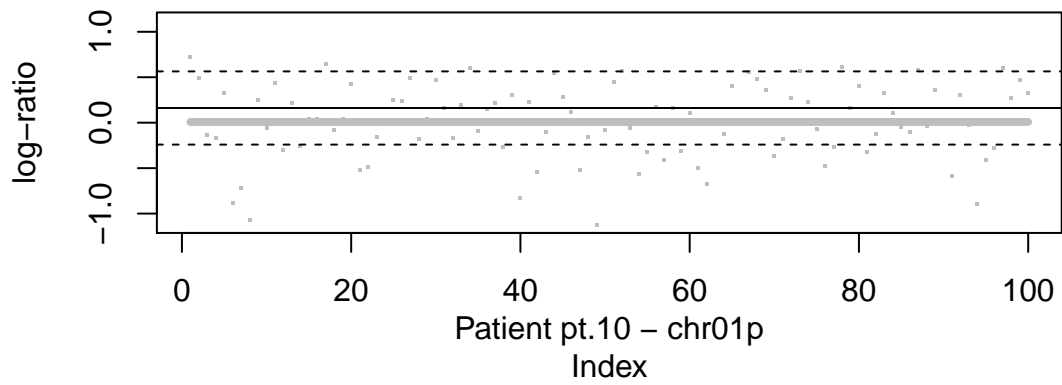
Patterns for each chromosome would be plotted by:

```
> chromosomePlots(results$OneStepSeg, ptlist, ptname="pt.10", nmad=1)
```

pt.10.1



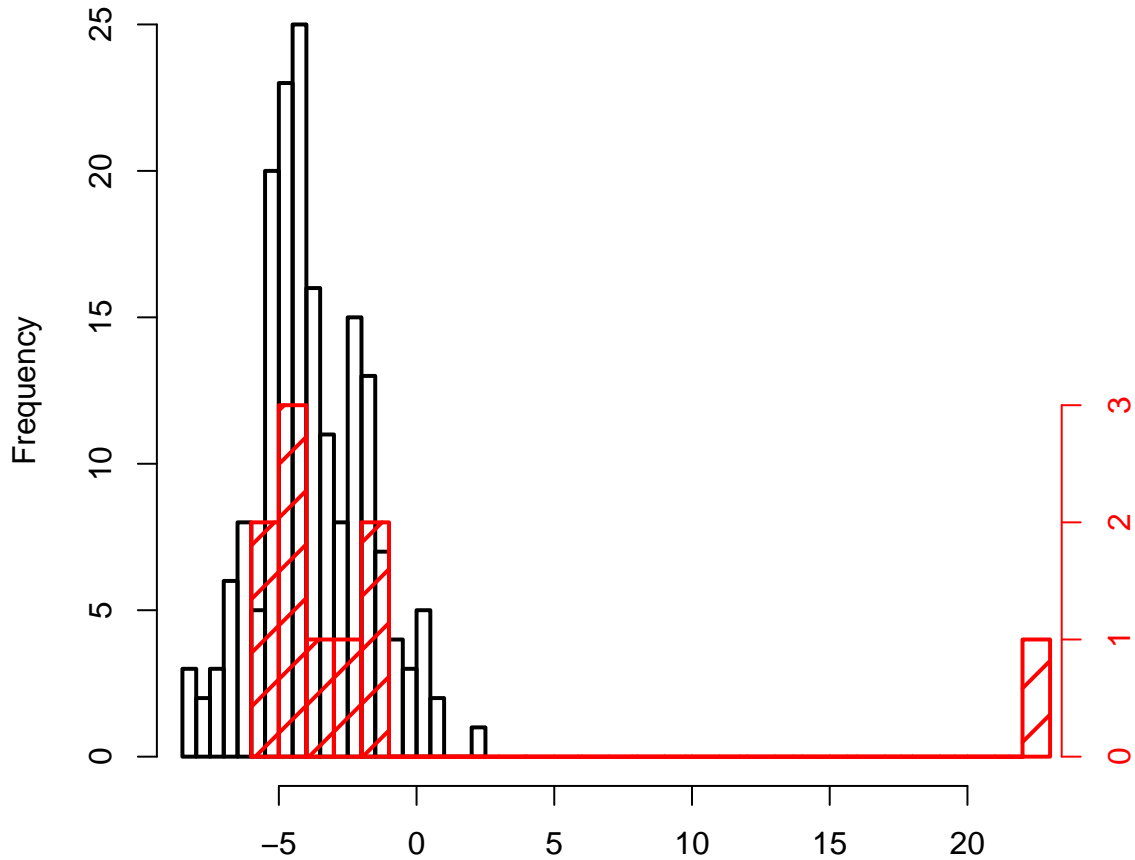
pt.10.2



The overlap between the histograms of LR2 from original pairs of tumors and the reference distribution are produced by:

```
> histogramPlot(results$LR[,4], results$refLR[,4])
```

Reference distribution of logLR (black), tested pairs (red)



2.1 Choice of segmentation algorithm

Note that the user can potentially specify the segmentation method to be used. Currently the default behavior of the `clonality.analysis` function is to use the CBS algorithm to identify the most significant change in each chromosome arm. The internal function for this purpose is "oneseg" called as `oneseg(x, alpha, nperm, sbdry)`

There are 4 arguments to `oneseg`:

- `x`: is the finite logratio data ordered by genomic position.
- `alpha`: the significance level used by CBS.
- `nperm`: the number of permutations for the reference distribution.
- `sbdry`: early stopping boundary for declaring no change (calculated from `alpha` and `nperm`).

The output of this function is a vector of 3 numbers where the first is the number of change-points detected (must be 0, 1 or 2), and the second and the third numbers are the start and end of the left segment if there is only one change-point, and of the middle segment when there are 2 change-points.

The function allows the user to specify alternative alpha and nperm for 'oneseg' as a list using the segpar argument e.g. `segpar=list(alpha=0.05, nperm=1000)`. Since `sbdry` is always calculated in `clonality.analysis` function from alpha and nperm it is not specified.

Alternate segmentation algorithm can be used. It requires the user to create a function that takes the ordered logratio from one chromosome arm as argument "x" as in `oneseg`. The name of this function should not be 'oneseg' and is passed through the 'segmethod' argument and all other necessary arguments that are needed passed as a list through 'segpar' argument.

3 LOH data

The LOH data has to be combined in a matrix where first column has marker names and the following columns have LOH calls for each sample. Here we simulate a dataset with 10 pairs of tumors and 20 markers. First pair of tumor is clonal, and the rest of them are independent. If the marker is heterozygous and there is no LOH, then it is denoted by 0. LOH at maternal or paternal alleles is marked by 1 or 2.

```
> set.seed(25)
> LOHtable<-cbind(1:20,matrix(sample(c(0,1,2),20*20,replace=TRUE),20))
> LOHtable[,3]<-LOHtable[,2]
> LOHtable[1,3]<-0
```

```
> LOHtable[,1:5]
```

	[,1]	[,2]	[,3]	[,4]	[,5]
[1,]	1	1	0	2	0
[2,]	2	2	2	0	0
[3,]	3	0	0	1	0
[4,]	4	2	2	2	2
[5,]	5	0	0	1	1
[6,]	6	2	2	0	2
[7,]	7	1	1	2	1
[8,]	8	1	1	2	2
[9,]	9	0	0	0	1
[10,]	10	0	0	2	0
[11,]	11	0	0	0	2
[12,]	12	1	1	2	0
[13,]	13	2	2	2	0
[14,]	14	1	1	1	2
[15,]	15	2	2	1	0


```
[16,] 16 0 0 0 1
[17,] 17 1 1 0 0
[18,] 18 2 2 1 2
[19,] 19 1 1 0 0
[20,] 20 2 2 0 0
```

```
> LOHclonality(LOHtable,rep(1:10,each=2),pfreq=NULL,noloh=0,loh1=1,loh2=2)
```

```
Testing clonality for patient 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, Done
```

	Sample1	Sample2	a	e	f	g	h	Ntot	CMpvalue	LRpvalue
1	1	1	13	13	0	1	6	20	2.20457220717234e-08	0
2	2	2	3	6	4	6	4	20	0.633257174327221	1
3	3	3	6	9	2	5	4	20	0.13247418005031	0.458
4	4	4	6	9	4	5	2	20	0.271731009940983	0.807
5	5	5	3	6	7	5	2	20	0.768026723950271	1
6	6	6	1	5	8	3	4	20	0.964059678147575	0.442
7	7	7	6	12	3	4	1	20	0.607636663320756	1
8	8	8	5	11	4	2	3	20	0.585520481546597	0.719
9	9	9	4	7	6	5	2	20	0.597049677704141	0.911
10	10	10	6	10	3	6	1	20	0.424944369195046	0.794

First p-value is small, indicating clonality, for both CM and LR tests. The rest of the p-values are not significant.

Markers that are not informative (e.g. homozygous) in a particular tumor should be given NA instead of a call. Such markers will be dropped from the analysis of this specific patient.

4 LOH data for 3 and more tumors

It is possible to test clonality of 3 or more tumors using Extended Concordant Mutations test, implemented in function 'ECMtesting'. The input LOH matrix can be in the same format as for 'LOHclonality' function: first column of a matrix contains marker names, subsequent columns are samples. For each patient all possible subsets of tumors are tested for clonality, with adjustment for multiple comparison performed using permutation MinP procedure.

Likelihood model can be extended for 3 or 4 tumors with function 'LRtesting3or4tumors'. The likelihood function depends on 2 parameters for 3 tumors, and 3 parameters for 4 tumors, allowing for non-symmetric relationship among tumors. Likelihood ratio test is computed and p-value is calculated using permutations.

5 Inference using profiles of somatic mutations

5.1 Likelihood model

In (Ostrovnya et al., 2015) we presented statistical test for evaluating evidence for clonality against null hypothesis that the two tumors are independent using their mutational profiles

obtained by next generation sequencing, such as targeted panel sequencing or whole exome sequencing. It utilizes conditional likelihood model where for each patient only loci where at least one tumor has a mutation are contributing to the test statistic. Marginal frequencies of mutations are assumed to be known and usually can be computed from TCGA data or other similar resources.

Below we download the exome sequencing data from study of Lobular Carcinoma in Situ (LCIS) and Invasive lobular carcinomas (ILC) and Invasive Ductal Carcinomas (IDC) in the same patients ((Begg et al., 2016)). Marginal probabilities in the column *probi* are obtained from breast cancer TCGA data and are not directly applicable to other cancers.

```
> data(lcis)

[1] "lcis"

> n<-nrow(lcis)

[1] 938

> summary(lcis$probi)

      Min.   1st Qu.   Median     Mean   3rd Qu.   Max.
0.0009862 0.0009862 0.0009862 0.0012743 0.0009862 0.1370809

> table(lcis$TK47IDC.TK47LCIS1 )

  0   1   2
880  1  57

> lcis$probi[lcis$TK47IDC.TK47LCIS1==1]

[1] 0.0009861933
```

Here variable TK47IDC.TK47LCIS1 takes values 0 if a mutation is not observed, 1 if shared mutation is observed in both tumors, and 2 if it's a private mutation. We can see that IDC and LCIS tumors in patient 47 have 1 mutation in common, and 57 present in only one of the tumors. The single match is at a locus where the mutations are relatively rare, having probability 0.000986.

Below is the test of clonality of these two tumors. Note that the p-value is calculated using the simulated null distribution, thus setting the random seed is recommended for reproducibility. We will assign private mutations to tumor 1 here since the likelihood doesn't depend on which tumor has the private mutation.

```
> x1<-x2<-rep(0,n)
> x1[lcis$TK47IDC.TK47LCIS1==1]<-x2[lcis$TK47IDC.TK47LCIS1==1]<-1
> x1[lcis$TK47IDC.TK47LCIS1==2]<-1
> set.seed(1)

> SNVtest(x1,x2,lcis$probi)
```

n1	n2	n_match	LRstat	maxKsi	LRpvalue
58.00000000	1.00000000	1.00000000	2.58936185	0.03290253	0.02100000

The p-value of 0.021 confirms that these two tumors are clonal, i.e. originate from the same cell harboring the matching mutation.

5.2 Random effects model

Here we show how to test the independence of the somatic mutation profile, following the random effects model proposed by Mauguen et al (<http://biostats.bepress.com/mskccbiostat/paper33>, (Mauguen et al., 2017)). The example uses the data from 22 cases with both lobular carcinoma in situ (LCIS) and an invasive breast tumor (Begg et al., 2016). Data from whole-exome sequencing were available and used to compare the mutation profile of the two tumors. Those data correspond to the dataset *lcis* included in the package. The random-effect model is estimated on the data using the following code:

```
> data(lcis)
> mod <- mutation.rem(lcis)

> print(mod)

Estimation done on 22 pairs

___ Parameter estimates

Random-effect distribution
mean mu = -2.26
standard-deviation sigma = 1.47

Proportion of clonal pairs
pi = 0.749

___ Model likelihood and convergence

likelihood -282.1317
convergence status 0
convergence message (from optim) CONVERGENCE: REL_REDUCTION_OF_F <= FACTR*EPSMCH
NULL
```

In this example, the estimation converged (convergence status=0). The proportion of clonal cases in the LCIS dataset is estimated to be 75%. The function allows the computation of standard-errors using the option `sd.err=TRUE`. The individual probabilities of clonality for those 22 cases are obtained using:

```
> data(lcis)
> mod <- mutation.rem(lcis, proba=TRUE)
>
```

```
> print(mod)
```

```
Estimation done on 22 pairs
```

```
___ Parameter estimates
```

```
Random-effect distribution
```

```
mean mu = -2.26
```

```
standard-deviation sigma = 1.47
```

```
Proportion of clonal pairs
```

```
pi = 0.749
```

```
___ Model likelihood and convergence
```

```
likelihood -282.1317
```

```
convergence status 0
```

```
convergence message (from optim) CONVERGENCE: REL_REDUCTION_OF_F <= FACTR*EPSMCH
```

```
___ Individual probabilities
```

Tumor pairs	Probability of being clonal
TK24ILC.TK24LCIS1	1.000
TK24ILC.TK24LCIS2	0.997
TK26IDC.TK26LCIS	0.352
TK46ILC.TK46LCIS	0.352
TK46ILC.TK46LCIS3	0.381
TK47IDC.TK47LCIS1	0.943
TK47IDC.TK47LCIS2	0.875
TK47ILC.TK47LCIS1	1.000
TK47ILC.TK47LCIS2	1.000
TK48ILC.TK48LCIS1	1.000
TK48ILC.TK48LCIS2	0.939
TK53IDC2.TK53LCIS1	0.999
TK53IDC2.TK53LCIS2	0.966
TK55ILC.TK55LCIS	1.000
TK68ILC.TK68LCIS1	0.308
TK69ILC.TK69LCIS	1.000
TK73ILC.TK73LCIS1	0.331
TK74IDC.TK74LCIS1	0.331
TK74IDC.TK74LCIS2	0.323
TK74IDC.TK74LCIS3	1.000
TK75IDC.TK75LCIS2	0.384
TK75ILC.TK75LCIS1	1.000

```
Tumor pairs Probability of being clonal
```

1	TK24ILC.TK24LCIS1	1.000
2	TK24ILC.TK24LCIS2	0.997
3	TK26IDC.TK26LCIS	0.352
4	TK46ILC.TK46LCIS	0.352
5	TK46ILC.TK46LCIS3	0.381
6	TK47IDC.TK47LCIS1	0.943
7	TK47IDC.TK47LCIS2	0.875
8	TK47ILC.TK47LCIS1	1.000
9	TK47ILC.TK47LCIS2	1.000
10	TK48ILC.TK48LCIS1	1.000
11	TK48ILC.TK48LCIS2	0.939
12	TK53IDC2.TK53LCIS1	0.999
13	TK53IDC2.TK53LCIS2	0.966
14	TK55ILC.TK55LCIS	1.000
15	TK68ILC.TK68LCIS1	0.308
16	TK69ILC.TK69LCIS	1.000
17	TK73ILC.TK73LCIS1	0.331
18	TK74IDC.TK74LCIS1	0.331
19	TK74IDC.TK74LCIS2	0.323
20	TK74IDC.TK74LCIS3	1.000
21	TK75IDC.TK75LCIS2	0.384
22	TK75ILC.TK75LCIS1	1.000

The individual probability of clonality varies from 31% for case 68 with no shared mutations to >99% for several cases having 2 or more mutations shared between the two tumors.

Finally, once the model is estimated on a given population, it is possible to estimate the probability of clonality of a new case using:

```
> # generate a case with 30 mutations
> # probabilities of each observed mutation
> set.seed(159)
> 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 previously estimated
> proba <- mutation.proba(c(mod$mu, mod$sigma, mod$pi), t(as.matrix(new.likmat)) )

> print(proba)

[1] 0.47
[1] 0.47
```

For this hypothetical case with 30 private mutations, the probability of being clonal is 47%.

Below are the details of the session information:

```
R version 3.4.0 (2017-04-21)
Platform: x86_64-pc-linux-gnu (64-bit)
Running under: Ubuntu 16.04.2 LTS

Matrix products: default
BLAS: /home/biocbuild/bbs-3.5-bioc/R/lib/libRblas.so
LAPACK: /home/biocbuild/bbs-3.5-bioc/R/lib/libRlapack.so

locale:
 [1] LC_CTYPE=en_US.UTF-8      LC_NUMERIC=C
 [3] LC_TIME=en_US.UTF-8      LC_COLLATE=C
 [5] LC_MONETARY=en_US.UTF-8  LC_MESSAGES=en_US.UTF-8
 [7] LC_PAPER=en_US.UTF-8    LC_NAME=C
 [9] LC_ADDRESS=C            LC_TELEPHONE=C
[11] LC_MEASUREMENT=en_US.UTF-8 LC_IDENTIFICATION=C

attached base packages:
[1] stats      graphics  grDevices  utils      datasets  methods   base

other attached packages:
[1] Clonality_1.24.0 DNACopy_1.50.0

loaded via a namespace (and not attached):
[1] compiler_3.4.0 tools_3.4.0
```

References

- Begg, C., Eng, K., and Hummer, A. (2007). Statistical tests for clonality. *Biometrics*, 63:522–530.
- Begg, C., Ostrovnaya, I., Carniello, J., Sakr, R., Giri, D., Towers, R., Schizas, M., DeBrot, M., Andrade, V., Mauguen, A., Seshan, V., and King, T. (2016). Clonal relationships between lobular carcinoma in situ and other breast malignancies. *Breast Cancer Res*, 18(1):66.
- Mauguen, A., Seshan, V., Ostrovnaya, I., and Begg, C. (2017). Estimating the probability of clonal relatedness of pairs of tumors in cancer patients. *Biometrics*, 000:000.
- Ostrovnaya, I., Olshen, A., Seshan, V., Orlow, I., Albertson, D., and Begg, C. (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.

- Ostrovnaya, I., Seshan, V., and Begg, C. (2008). Comparison of properties of tests for assessing tumor clonality. *Biometrics*, 68:1018–1022.
- Ostrovnaya, I., VE, S., and CB, B. (2015). Using somatic mutation data to test tumors for clonal relatedness. *Annals of Applied Statistics*, 9(3):1533–1548.
- Venkatraman, E. and Olshen, A. (2007). A faster circular binary segmentation algorithm for the analysis of array cgh data. *Bioinformatics*, 23:657–663.