

# sequenza usage example

Francesco Favero\*, Aron C. Eklund

October 12, 2013

## Contents

<b>1</b>	<b>Abstract</b>	<b>2</b>
<b>2</b>	<b>Minimum requirements</b>	<b>2</b>
<b>3</b>	<b>Getting ready with Sequenza package/Installing R/Setting up Sequenza</b>	<b>2</b>
<b>4</b>	<b>Preparing inputs for Sequenza</b>	<b>3</b>
<b>5</b>	<b>First the non-R part: pre-processing data</b>	<b>3</b>
<b>6</b>	<b>Read the pre-processed data (<i>abfreq</i> file) into R</b>	<b>4</b>
<b>7</b>	<b>Quality control step? (EXPLAIN)</b>	<b>5</b>
<b>8</b>	<b>GC-normalization</b>	<b>5</b>
<b>9</b>	<b>Create genomic profiles</b>	<b>6</b>
9.1	First, the depth ratio . . . . .	6
9.2	Next, the B-allele frequencies . . . . .	7
<b>10</b>	<b>Allele-specific segmentation</b>	<b>8</b>
10.1	Find genomic breakpoints . . . . .	8
<b>11</b>	<b>Select mutations by mutation frequency</b>	<b>9</b>
<b>12</b>	<b>Plot chromosome view with mutations, BAF, depth ratio and segments</b>	<b>10</b>
<b>13</b>	<b>Inference of cellularity and DNA-index</b>	<b>11</b>

---

\*favero@cbs.dtu.dk

<b>14 Call CNVs and mutations using the estimated parameters</b>	<b>14</b>
14.1 Detect mutated alleles . . . . .	14
14.2 Detect Copy number variation . . . . .	14
<b>15 Visualize detected copy number</b>	<b>16</b>

# 1 Abstract

Deep sequence of tumor DNA along with corresponding normal DNA can provide a rich picture of the mutations and aberrations that characterize the tumor. However, analysis of this data can be impeded by of tumor cellularity and heterogeneity and by unwieldy data. Here we describe the *sequenza* software system, which comprises a fast python-based pre-processor and an R-based analysis package. Sequenza enables the efficient estimation of tumor cellularity and ploidy, and generation of copy number, loss-of-heterozygosity, and mutation frequency profiles.

This document details a typical analysis of matched tumor-normal exome sequence data using *sequenza*.

# 2 Minimum requirements

- Software: R, Python
- Operating system: Linux, OSX
- Memory: Minimum 4GB of RAM. Recommended >8GB.
- Disk space: 1.5 GB for each sample
- R version: 2.15.1
- Python version: 2.7 with itertools, codecs modules; multiprocessing.pool, multiprocessing.queue for parallelization; rpy2 for running sequenza-utils from R.

# 3 Getting ready with Sequenza package/Installing R/Setting up Sequenza

- download from [bitbucket/cbs.dtu.dk](https://bitbucket.org/cbs.dtu.dk)
- how to install it. R CMD INSTALL sequenza\_version.tar.gz

A typical workflow by Sequenza is as follow :

1. Convert pileup to abfreq

2. GC normalization
3. Obtain depth ratio and B allele frequencies
4. Allele-specific segmentation
5. Infer cellularity and DNA-index by model fitting
6. Segmentation and mutations

## 4 Preparing inputs for Sequenza

In order to obtain precise mutational and aberration patterns in a tumor sample, Sequenza requires a matched normal sample from the same patient. In short, the following files are needed to get started with Sequenza.

1. A pileup file from the tumor specimen
2. A pileup file from the normal specimen
3. A FASTA reference genomic sequence file (optional, for GC-content correction)

We recommend using pre-processed and quality filtered BAM files to obtain mpileup calls for both samples.

Pileup files can be generated using `samtools` (ref). The genome sequence file can be obtained from (url).

```
# samtools mpileup -f hg19.fasta -Q 20 normal.bam
```

```
# samtools mpileup -f hg19.fasta -Q 20 tumor.bam
```

## 5 First the non-R part: pre-processing data

For convenience and efficiency we have implemented pre-processing algorithms in an external (not called from R) Python program. The program is provided with the package; it's exact location can be found like this:

```
> system.file("exec", "sequenza-utils.py", package="sequenza")
```

```
[1] ""
```

You may wish to copy this program to a location on your path. NOTE: this script requires several UNIX tools and thus probably not work on Windows (HOW ABOUT CYGWIN ?).

IS THE GC CONTENT FILE INDEPENDENT OF REFERENCE GENOME USED FOR ALIGNMENT AND MPILEUP ? SHOULD NON-CANONICAL CHROMOSOMES BE REMOVED IF ONE WANTS GC CONTENT CALCULATIONS BE DONE ON THEIR REFERENCE GENOME ?

Extract average GC content in 50-base genomic windows:

```
# sequenza-utils.py GC-windows -w 50 hg19.fa | gzip > hg19.gc50Base.txt.gz
```

Process the two pileup files to obtain an "abfreq" file containing alleles and mutation frequency.

```
# sequenza-utils.py pileup2tab -gc hg19.gc50Base.txt.gz -r 0001-normal_blood.pileup.gz  
-s 0001-met2.pileup.gz -q 20 -n 10 -o 0001-met2.abfreq.txt.gz
```

— UPDATE ME UPDATE ME UPDATE ME UPDATE ME UPDATE ME UPDATE ME —

## 6 Read the pre-processed data (*abfreq* file) into R

The remainder of this example takes place in R.

Load the sequenza package:

```
> library("sequenza")
```

Find the example data file:

```
> data.file <- system.file("data", "abf.data.abfreq.txt.gz", package = "sequenza")  
> data.file
```

```
[1] "/usr/local/Cellar/r/3.0.1/R.framework/Versions/3.0/Resources/library/sequenza/data"
```

The abfreq file can be read all at once, but processing one chromosome at a time is less demanding on computational resources and might be preferable. (Note that the demo data included with sequenza is only chromosome 1)

Read only the data corresponding to chromosome 1:

```
> abf.data <- read.abfreq(data.file, chr.name = "1")
```

Alternatively, read all data at once (not run):

```
> abf.data <- read.abfreq(data.file)
```

```
> str(abf.data)
```

```
'data.frame':      45003 obs. of  13 variables:
 $ chromosome      : chr  "1" "1" "1" "1" ...
 $ n.base          : int  132138 331284 881918 884091 897530 902140 909221 9
 $ base.ref        : chr  "C" "C" "G" "C" ...
 $ depth.normal    : int  17 18 55 85 15 30 41 18 45 100 ...
 $ depth.sample    : int  7 5 37 65 7 14 16 6 39 87 ...
 $ depth.ratio     : num  0.412 0.278 0.673 0.765 0.467 0.467 0.39 0.333 0.8
 $ Af              : num  0.571 0.8 0.514 0.597 0.8 0.5 0.6 0.8 0.583 0.506
 $ Bf              : num  0.429 0 0.486 0.387 0 0 0.4 0 0.417 0.481 ...
 $ ref.zygotity    : chr  "het" "hom" "het" "het" ...
 $ GC.percent      : num  62 62 64 70 54 82 64 70 68 64 ...
 $ sample.reads.above.quality: num  1 1 0.95 0.95 0.71 0.71 0.94 0.83 0.92 0.89 ...
 $ AB.germline     : chr  "AT" "C" "AG" "AG" ...
 $ AB.sample       : chr  "." "T0.2" "." "." ...
```

## 7 Quality control step? (EXPLAIN)

Each nucleotide aligned in the sequencing is associated with a quality score. The *sequenza-utils* software is capable of filtering the base with the quality lower than a specified value (default is 20), and returns the rate of reads that have passed the filter in the column *sample.reads.above.quality*, while the *depth.sample* column contains the raw depth calculated in the pileup (from samtools). The product of the rate of bases that have passed the quality check and the total amount of reads aligned at the same nucleotide returns the number of reads that have passed the quality check.

```
> abf.data$good.s.reads <- abf.data$depth.sample *
+                           abf.data$sample.reads.above.quality
```

## 8 GC-normalization

The number of reads at a given genomic position can be affected by the local GC content. We attempt to remove this bias as in (ref).

It is possible to gather gc-content information from the entire file (normally this would be the entire genome, but in our example it contains only chromosome 1):

```
> gc.stats <- gc.sample.stats(data.file)
```

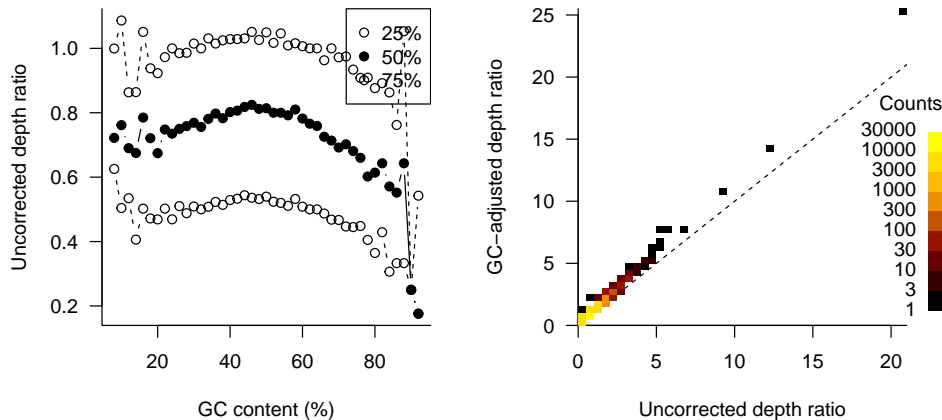
Or alternatively, it is possible to collect the GC-contents information from an object loaded in the environment.

```
> gc.stats <- gc.norm(ratio = abf.data$depth.ratio,
+                     gc = abf.data$GC.percent)
```

Calculate the GC-normalized depth ratio:

```
> gc.vect <- setNames(gc.stats$raw.mean, gc.stats$gc.values)
> abf.data$adjusted.ratio <- abf.data$depth.ratio /
+                             gc.vect[as.character(abf.data$GC.percent)]

> par(mfrow = c(1,2), cex = 1, las = 1, bty = 'l')
> matplot(gc.stats$gc.values, gc.stats$raw,
+         type = 'b', col = 1, pch = c(1, 19, 1), lty = c(2, 1, 2),
+         xlab = 'GC content (%)', ylab = 'Uncorrected depth ratio')
> legend('topright', legend = colnames(gc.stats$raw), pch = c(1, 19, 1))
> hist2(abf.data$depth.ratio, abf.data$adjusted.ratio,
+       breaks = prettyLog, key = vkey, panel.first = abline(0, 1, lty = 2),
+       xlab = 'Uncorrected depth ratio', ylab = 'GC-adjusted depth ratio')
```



## 9 Create genomic profiles

### 9.1 First, the depth ratio

Summarize the depth ratio by binning the data in overlapping genomic windows:

```
> abf.r.win <- windowValues(x = abf.data$adjusted.ratio,
+                             positions = abf.data$n.base,
+                             chromosomes = abf.data$chromosome,
+                             window = 1e6, overlap = 1,
+                             weight = abf.data$depth.normal)
```

```

> plotWindows(abf.r.win[[1]], log2.plot = FALSE,
+           ylab = "Depth ratio", xlab = "Position (bases)",
+           main = names(abf.r.win)[1], las = 1, n.min = 1,
+           ylim = c(0, 2.5))

```

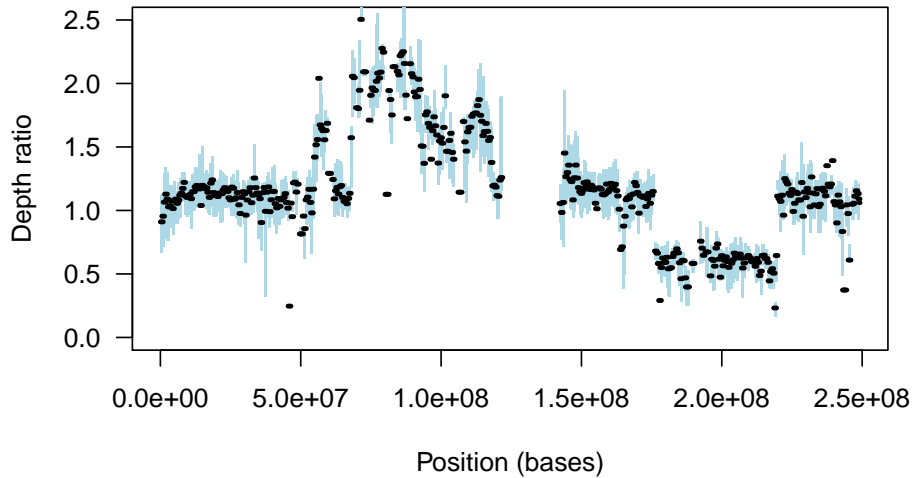


Figure 1: Depth ratio profile visualization over a single chromosome.

## 9.2 Next, the B-allele frequencies

The column *ref.zygosity* contains the zygosity derived from the germline sample. the possible values are *het* for heterozygous positions and *hom* for homozygous positions.

```

> abf.hom <- abf.data$ref.zygosity == 'hom'
> abf.het <- abf.data[!abf.hom, ]

```

Summarize the BAF by binning the data in overlapping genomic windows (including only those positions called heterozygous in the normal sample):

```

> abf.b.win <- windowValues(x = abf.het$Bf,
+ positions = abf.het$n.base,
+ chromosomes = abf.het$chromosome,
+ window = 1e6, overlap = 1,
+ weight = round(x = abf.het$good.s.reads, digits = 0))

```

```
> plotWindows(abf.b.win[[1]], ylim = c(0, 0.5),
+           main = names(abf.r.win)[1], xlab = "Position (bases)",
+           ylab = "B allele frequency", n.min = 10)
```

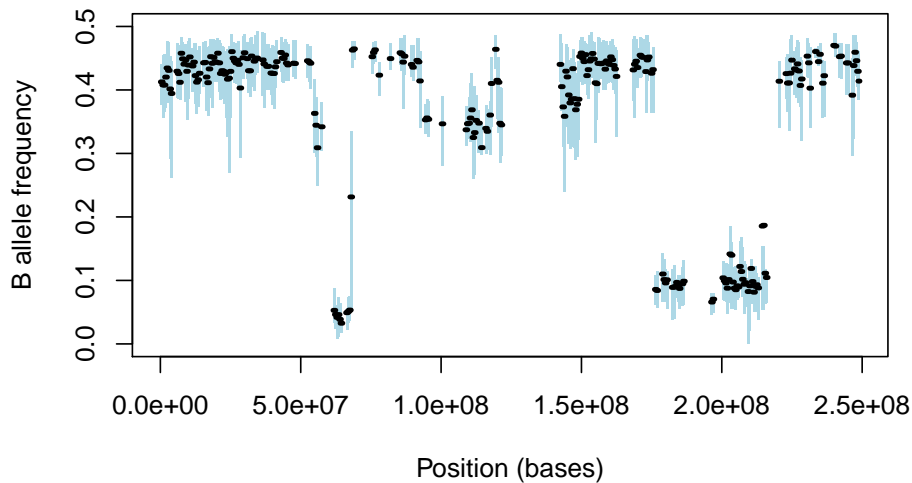


Figure 2: B-allele frequency profile visualization over a single chromosome.

## 10 Allele-specific segmentation

### 10.1 Find genomic breakpoints

To find breakpoints we use the allele-specific segmentation algorithm from the *copynumber* package [1].

```
> breaks <- find.breaks(abf.het, gamma = 80, kmin = 10, baf.thres = c(0, 0.5))
> head(breaks)
```

	chrom	start.pos	end.pos
1	1	132138	6196951
2	1	6202423	54694214
3	1	54700707	59465454
4	1	60381491	67890526
5	1	68151685	92262874
6	1	92445257	118165328

Now obtain the segment values:

```
> seg.s1 <- segment.breaks(abf.data, breaks = breaks)
```



## 11 Select mutations by mutation frequency

In the genotype file (the *abfreq* file) the mutations are detected as homozygous positions with a decreased frequency of the germline nucleotide. A set of nucleotides not present in the germline is present with the relative frequency in the column *AB.sample*. Being a frequency derived by the number of reads covering the position, the accuracy of the measurement is depending on the depth in the considered position. In order to filter the mutations the function *mutation.table* allows to filter the present mutation to a defined level of frequency, a desired number of reads depth, and a desired number of mutated nucleotides per position. Additionally it is possible to swap the *adjusted.ratio* column with the corresponding value after segmentation.

```
> mut.tab <- mutation.table(abf.data, mufreq.threshold = 0.15,  
+                           min.reads = 40, max.mut.types = 1,  
+                           min.type.freq = 0.9, segments = seg.s1)
```

However it is optional, without providing the segmented data the *adjusted.ratio* would remain unchanged.

```
> mut.tab.no.seg <- mutation.table(abf.data, mufreq.threshold = 0.15,  
+                                 min.reads = 40, max.mut.types = 1,  
+                                 min.type.freq = 0.9)
```

```
> dim(mut.tab)
```

```
[1] 185  7
```

```
> head(mut.tab)
```

	chromosome	n.base	GC.percent	good.s.reads	adjusted.ratio	F	mutation
291	1	10436585	50	160.36	1.164266	0.217	C>T
460	1	13112809	48	40.18	1.164266	0.225	C>T
576	1	15821826	54	485.97	1.164266	0.388	G>T
923	1	19983391	72	63.99	1.164266	0.469	G>C
1160	1	26878353	60	50.00	1.164266	0.520	C>A
1280	1	32627966	54	119.56	1.164266	0.403	C>T

```
> head(mut.tab.no.seg)
```

	chromosome	n.base	GC.percent	good.s.reads	adjusted.ratio	F	mutation
291	1	10436585	50	160.36	1.253728	0.217	C>T
460	1	13112809	48	40.18	1.482044	0.225	C>T
576	1	15821826	54	485.97	1.037625	0.388	G>T
923	1	19983391	72	63.99	1.082756	0.469	G>C
1160	1	26878353	60	50.00	1.478172	0.520	C>A
1280	1	32627966	54	119.56	1.257800	0.403	C>T

## 12 Plot chromosome view with mutations, BAF, depth ratio and segments

The visualization can be made by chromosome, using binned data and segmented data. Optionally can be inserted the mutations table as in figure 3 and the estimated parameters to draw the resulting model lines as in figure 6

```
> chromosome.view(mut.tab = mut.tab[mut.tab$chromosome == "1",], baf.windows = abf.b  
+               ratio.windows = abf.r.win[[1]], min.N.ratio = 1,  
+               segments = seg.s1[seg.s1$chromosome == "1",], main = "Chromosome 1
```

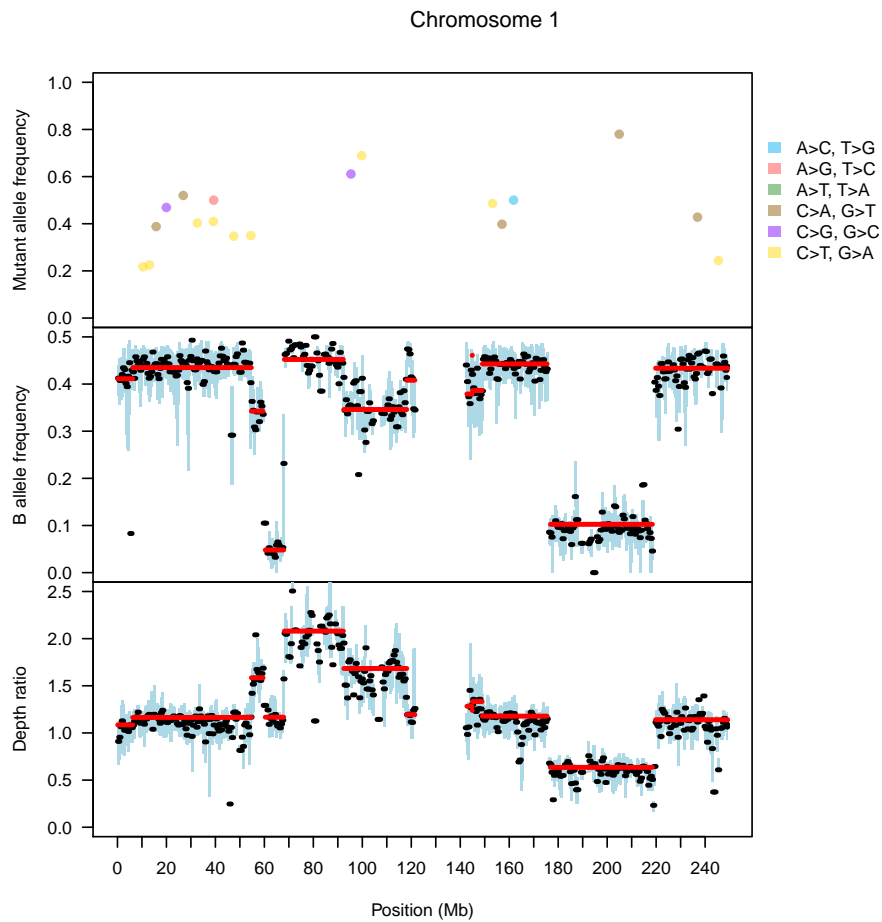


Figure 3: Plots of Mutation (top), B-allele frequencies (middle) and depth ratio (bottom) for chromosome position.

## 13 Inference of cellularity and DNA-index

The parameter estimation is performed on the segmented data, using BAF and depth ratio information. In order to avoid to catch errors due to the heterogeneity of the sample, we can filter the segments smaller then a certain amount. For instance, we can be confident that segments bigger then few megabases have enough data points to usually insure a correct measure. Alternatively it could be possible to exclude chromosomal regions usually know to be problematic, as the region close to the telomeres and also near the centromere

```
> seg.filtered <- seg.s1[(seg.s1$end.pos - seg.s1$start.pos) > 10e6, ]
```

Every segment is evaluated against the model lines using the resulting density of a binomial distribution and the possible values of the model. In order to perform every segment needs to be associate with a sample size, to generate the distribution function. For practical reason the sample size can not be in the order of the size of the segment in nucleotide (millions), So we use the size of the segment in megabases, added to an arbitrary offset (eg 150), to allow each segment to generate a proper distribution.

```
> weights.seg <- 150 + round((seg.filtered$end.pos - seg.filtered$start.pos) / 1e6,
```

The genome wide average depth ratio, should be a value close to 1, after normalization. However, since in this example we only consider chromosome 1, we have a different value (I will add one or more chr to make the avg became 1...).

```
> avg.depth.ratio <- mean(gc.stats$adj[,2])
> avg.depth.ratio
```

```
[1] 1
```

The function *baf.model.fit* evaluate the segmented data to a set of selected value of cellularity and DNA index. Using the implemented model to calculate the theoretic points it returning an x,y,z list containing a matrix z whit the log likelihood for the combinations of the two parameters, a vector x containing all the evaluated DNA index values and a vector y containing all the evaluated cellularity value.

```
> CP <- baf.model.fit(Bf = seg.filtered$Bf, depth.ratio = seg.filtered$depth.ratio,
+                    weight.ratio = weights.seg,
+                    weight.Bf = weights.seg,
+                    avg.depth.ratio = avg.depth.ratio,
+                    cellularity = seq(0.1,1,0.01),
+                    dna.index = seq(0.5,3,0.05), mc.cores = 4)
```

It is possible to calculate the confidence intervals for the two parameters using the function *get.ci*

```
> cint <- get.ci(CP)
```

As well it is possible to plot the likelihood over the combination of the two parameters, highlighting the point estimate and lines delimiting the confidence region.

```
> cp.plot(CP)
> cp.plot.contours(CP, add = TRUE, likThresh = c(0.5, 0.75, 0.95, 0.99))
```

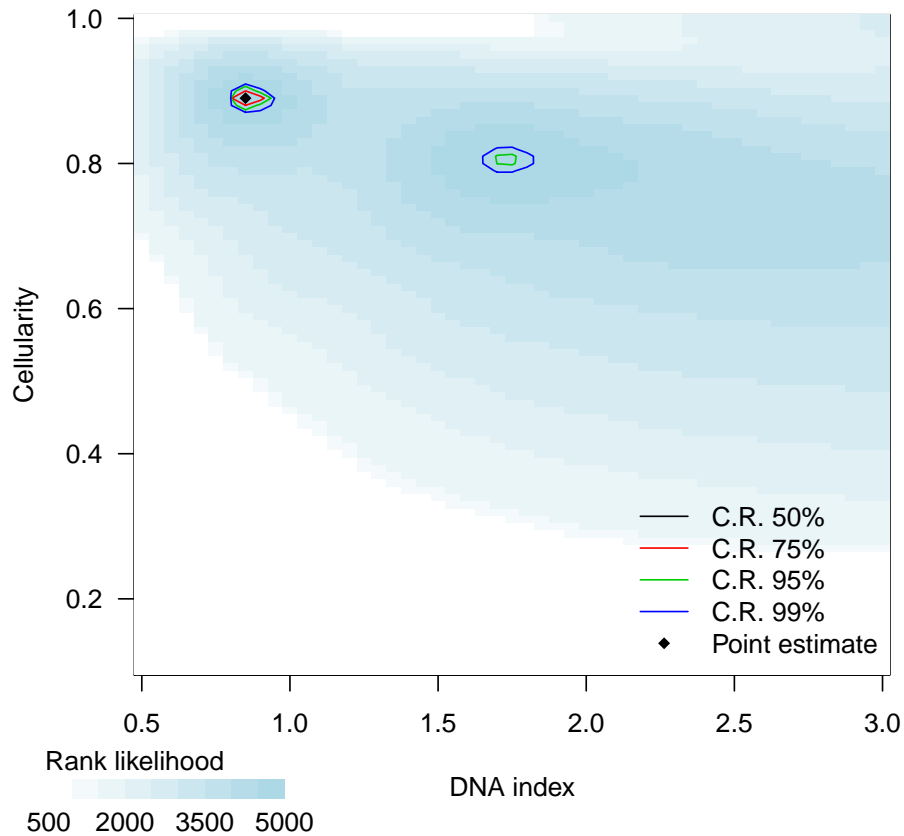


Figure 4: Result from the Bayesian inference over the defined range of cellularity and DNA-index. The color indicates the log-likelihood of the corresponding cellularity/DNA-index values.

Alternatively it is possible to draw the likelihood distribution for each parameter, using the information retrieved by the *get.ci* function.

```

> par(mfrow = c(2,2))
> cp.plot(CP)
> plot(cint$values.y, ylab = "Cellularity",
+       xlab = "likelihood", type = "n")
> select <- cint$confint.y[1] <= cint$values.y[,2] & cint$values.y[,2] <= cint$confi
> polygon(y = c(cint$confint.y[1], cint$values.y[select, 2], cint$confint.y[2]),
+         x = c(0, cint$values.y[select, 1], 0), col='red', border=NA)
> lines(cint$values.y)
> abline(h = cint$max.y, lty = 2, lwd = 0.5)
> plot(cint$values.x, xlab = "DNA index",
+       ylab = "likelihood", type = "n")
> select <- cint$confint.x[1] <= cint$values.x[,1] & cint$values.x[,1] <= cint$confi
> polygon(x = c(cint$confint.x[1], cint$values.x[select, 1], cint$confint.x[2]),
+         y = c(0, cint$values.x[select, 2], 0), col='red', border=NA)
> lines(cint$values.x)
> abline(v = cint$max.x, lty = 2, lwd = 0.5)
>

```

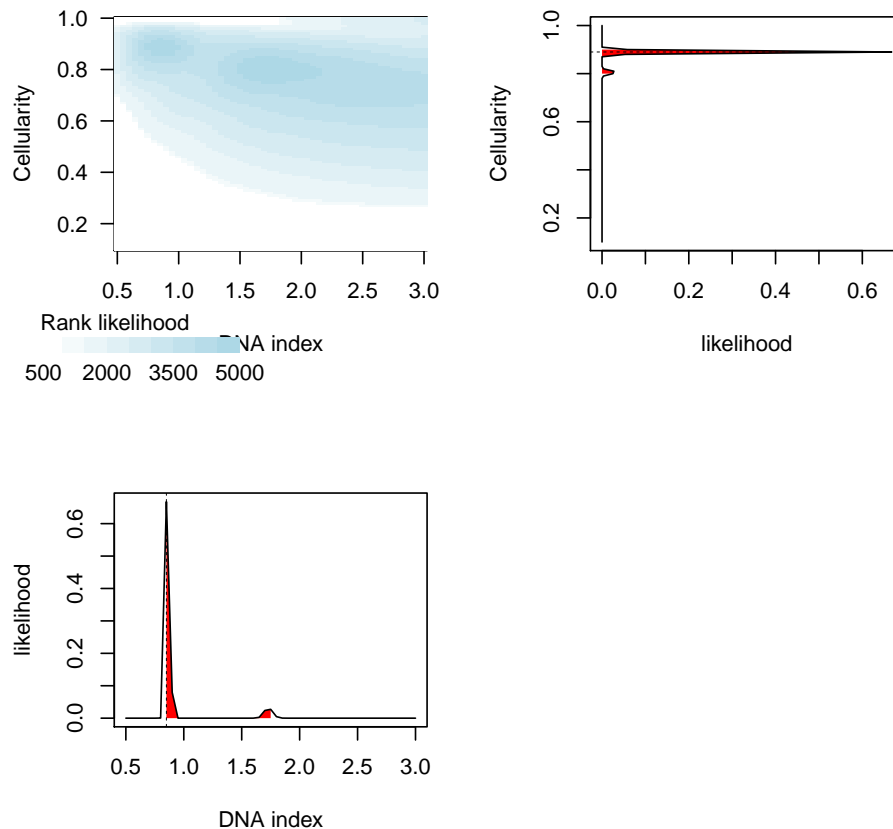


Figure 5: Plot of the log likelihood with respective cellularity and DNA-index probability distribution and confidence intervals.

## 14 Call CNVs and mutations using the estimated parameters

```
> cellularity <- cint$max.y
> cellularity
```

```
[1] 0.89
```

```
> dna.index <- cint$max.x
> dna.index
```

```
[1] 0.85
```

### 14.1 Detect mutated alleles

```
> mut.tab.clean <- na.exclude(mut.tab)
> mut.alleles <- mufreq.bayes(mufreq = mut.tab.clean$F, depth.ratio = mut.tab.clean$
+                               cellularity = cellularity, dna.index = dna.index,
+                               avg.depth.ratio = 1)
> head(mut.alleles)
```

	CNr	CNt	Mt	L
4	2	2	1	-28.10005
41	2	2	1	-28.10005
42	2	2	1	-13.28815
43	2	2	1	-12.61036
44	2	2	1	-14.06081
45	2	2	1	-12.99802

```
> head(cbind(mut.tab.clean[,c("chromosome", "n.base", "F", "adjusted.ratio", "mutation"
```

	chromosome	n.base	F	adjusted.ratio	mutation	CNr	CNt	Mt	L
291	1	10436585	0.217	1.164266	C>T	2	2	1	-28.10005
460	1	13112809	0.225	1.164266	C>T	2	2	1	-28.10005
576	1	15821826	0.388	1.164266	G>T	2	2	1	-13.28815
923	1	19983391	0.469	1.164266	G>C	2	2	1	-12.61036
1160	1	26878353	0.520	1.164266	C>A	2	2	1	-14.06081
1280	1	32627966	0.403	1.164266	C>T	2	2	1	-12.99802

### 14.2 Detect Copy number variation

```
> cn.alleles <- baf.bayes(Bf = seg.s1$Bf, depth.ratio = seg.s1$depth.ratio,
+                           cellularity = cellularity, dna.index = dna.index,
```

```

+                               avg.depth.ratio = 1)
> seg.s1.cn <- cbind(seg.s1, cn.alleles)
> head(seg.s1.cn)

```

	chromosome	start.pos	end.pos	Bf	N.BAF	depth.ratio	N.ratio	CNt	A	B
1	1	132138	6196951	0.4111775	131	1.085438	164	2	1	1
2	1	6202423	54694214	0.4349229	1411	1.164266	1485	2	1	1
3	1	54700707	59465454	0.3424249	60	1.585406	63	3	2	1
4	1	60381491	67890526	0.0480317	101	1.166852	104	2	2	0
5	1	68151685	92262874	0.4521961	216	2.079891	219	4	2	2
6	1	92445257	118165328	0.3459746	261	1.683497	274	3	2	1

L

```

1 -12.05441
2 -11.68261
3 -11.61122
4 -10.48188
5 -11.54389
6 -11.51471

```

## 15 Visualize detected copy number

```
> chromosome.view(mut.tab = mut.tab[mut.tab$chromosome == "1",], baf.windows = abf.b
+           ratio.windows = abf.r.win[[1]], min.N.ratio = 1,
+           segments = seg.s1.cn[seg.s1.cn$chromosome == "1",], main = "Chromo
+           cellularity = cellularity, dna.index = dna.index,
+           avg.depth.ratio = 1)
```

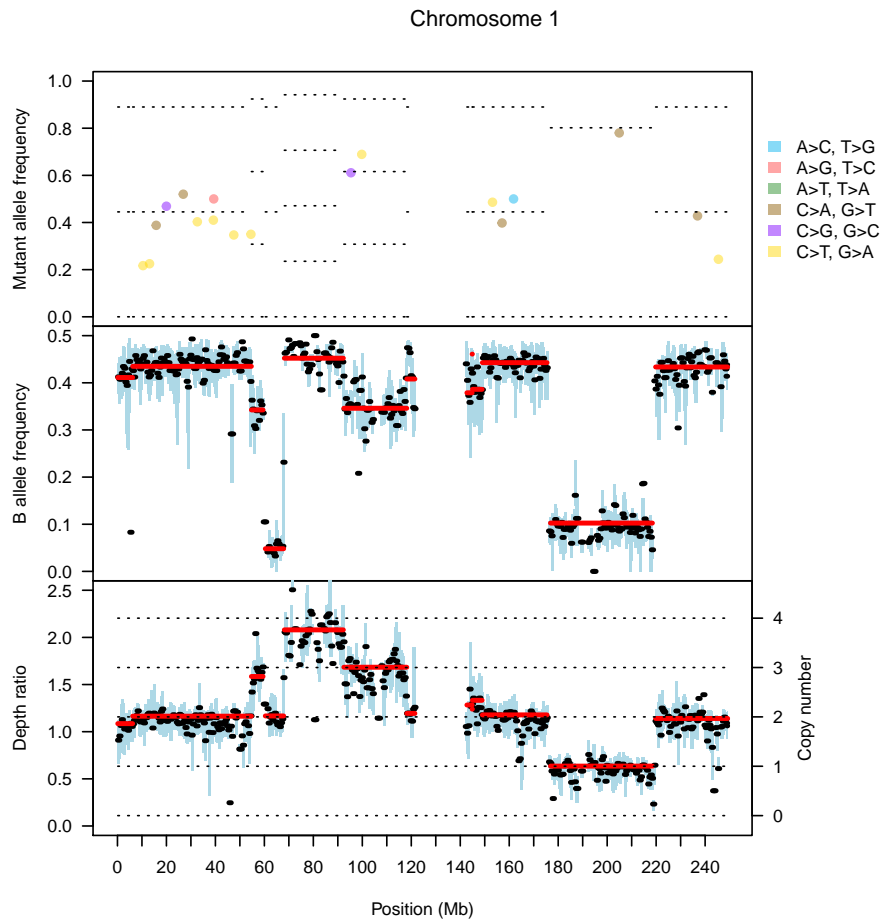


Figure 6: Plots of Mutation (top), B-allele frequencies (middle) and depth ratio (bottom) for chromosome position. Horizontal dotted line indicate different copy number/allelic state.



```

> chromosome.view(mut.tab = mut.tab[mut.tab$chromosome == "17",], baf.windows = abf.
+               ratio.windows = abf.r.win[[17]], min.N.ratio = 1,
+               segments = seg.s1.cn[seg.s1.cn$chromosome == "17",], main = "Chrom
+               cellularity = cellularity, dna.index = dna.index,
+               avg.depth.ratio = 1)

```

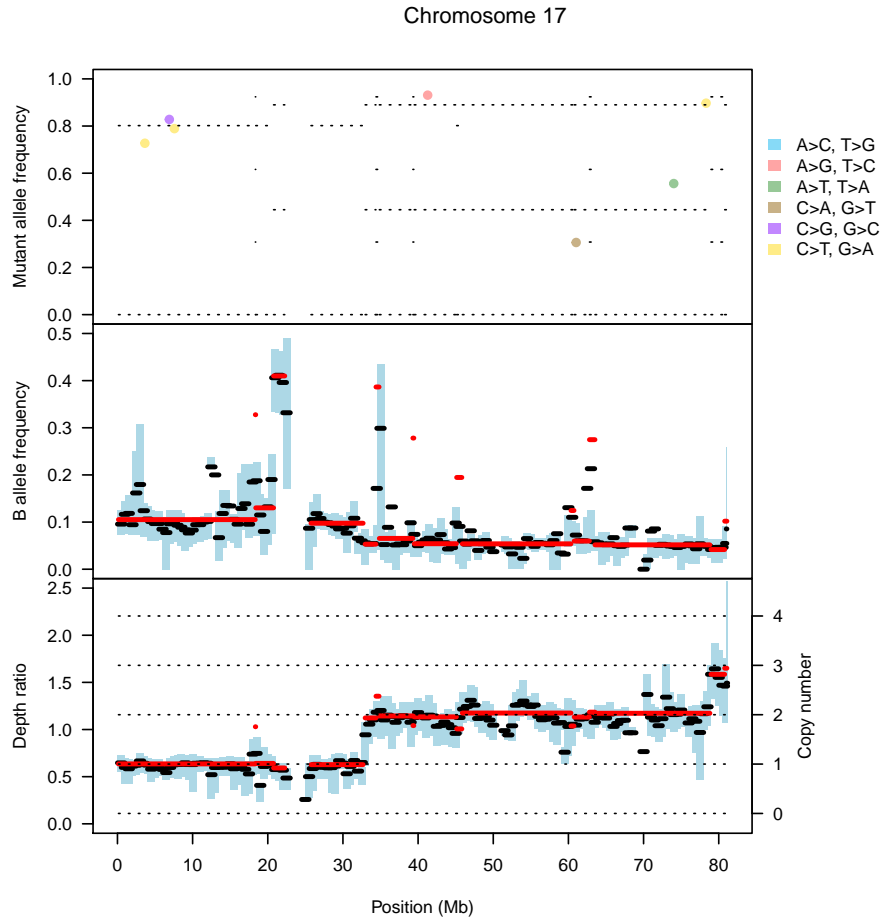


Figure 7: Plots of Mutation (top), B-allele frequencies (middle) and depth ratio (bottom) for chromosome position. Horizontal dotted line indicate different copy number/allelic state.

## References

- [1] Gro Nilsen, Knut Liestøl, Peter Van Loo, Hans Kristian Moen Vollan, Marianne B Eide, Oscar M Rueda, Suet-Feung Chin, Roslin Russell, Lars O Baumbusch, Carlos Caldas, Anne-Lise Børresen Dale, and Ole Christian Lingjaerde. Copynumber:

Efficient algorithms for single- and multi-track copy number segmentation. *BMC genomics*, 13:591, January 2012.