

sequenza usage example

Francesco Favero*, Aron C. Eklund

October 31, 2013

Contents

1	Abstract	1
2	Getting started	2
2.1	Minimum requirements	2
2.2	Installation	2
2.3	Workflow overview	2
2.4	Preparing inputs for Sequenza	2
2.5	Processing the data with <i>sequenza-utils.py</i>	3
3	Exploring the <i>ABfreq</i> file and GC-correction details	3
3.1	Read the file	3
3.2	Quality control and normalization	4
3.3	GC-normalization	4
4	Analyzing sequencing data with <i>sequenza</i>	5
4.1	Extract the information from the <i>ABfreq</i> file.	6
4.1.1	Plot chromosome view with mutations, BAF, depth ratio and segments	6
4.2	Inference of cellularity and ploidy (DNA index)	7
4.2.1	Confidence intervals, confidence region and point estimate	7
4.3	Call CNVs and mutations using the estimated parameters	10
4.3.1	Detect variant alleles (mutations)	10
4.3.2	Detect copy number variations	11
4.4	Visualize detected copy number changes and variant alleles	11
4.4.1	Genome-wide view of the alleles and copy numbers state	13

1 Abstract

Deep sequence of tumor DNA along with corresponding normal DNA can provide a valuable perspective on the mutations and aberrations that characterize the tumor. However, analysis of this data can be impeded by tumor cellularity and heterogeneity and by unwieldy data. Here we describe *Sequenza*, 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*.

*favero@cbs.dtu.dk

2 Getting started

2.1 Minimum requirements

- Software: R, Python
- Operating system: Linux, OSX
- Memory: Minimum 4GB of RAM. Recommended >8GB.
- Disk space: 1.5 GB for sample
- R version: 2.15.1
- Python version: 2.7; rpy2 is required to run sequenza R functions from the python command line programs.

2.2 Installation

In order to install sequenza, you can select download the package from the nearest CRAN mirror doing:

```
> install.packages("sequenza")
```

2.3 Workflow overview

A typical workflow developed with Sequenza on pre-aligned sequencing files(BAM) is structured as follow:

1. Convert pileup to *ABfreq* format
2. GC normalization
3. Allele-specific segmentation using the depth ratio and the B allele frequencies (BAFs)
4. Infer cellularity and DNA-index by model fitting
5. Call CNV and variant alleles

2.4 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 summary, the following files are needed to get started with Sequenza.

1. A bam file or a derived pileup file from the tumor specimen.
2. A bam file or a derived pileup file from the normal specimen.
3. A FASTA reference genomic sequence file (to extract GC-content information, and to transform bam to pileup if needed.)

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

Pileup files can be generated using `samtools`[1]. The genome sequence file can be obtained from (url).

```
1 samtools mpileup -f hg19.fasta -Q 20 normal.bam | gzip > normal.pileup.gz
2 samtools mpileup -f hg19.fasta -Q 20 tumor.bam | gzip > tumor.pileup.gz
```

2.5 Processing the data with *sequenza-utils.py*.

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

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

You may wish to copy this program to a location on your path.

To obtain the GC content information (required to obtain an *ABfreq* file), is possible to use a function from *sequenza-utils*, and extract the average GC content using a fixed genomic windows, or download the gc5Base from golden path (<http://hgdownload-test.cse.ucsc.edu/goldenPath/hg19/gc5Base/>). The following example calculates GC content for 50 nucleotides windows:

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

When the GC content file is available, it is possible to process the two pileup files to obtain an *ABfreq* file containing genotype information, alleles and mutation frequency, and more other features.

```
1 sequenza-utils.py pileup2abfreq -gc hg19.gc50Base.txt.gz \  
2 -r normal.pileup.gz \  
3 -s tumor.pileup.gz | gzip > out.abfreq.gz
```

If you don't have yet the pileup, or you are not interested on storing the pileup for further use, you can use two *FIFO* files, to pipe the samtools output directly to *sequenza-utils*:

```
1 mkfifo normal.fifo tumor.fifo  
2 samtools mpileup -f hg19.fasta -Q 20 normal.bam > normal.fifo &  
3 samtools mpileup -f hg19.fasta -Q 20 tumor.bam > tumor.fifo &  
4 sequenza-utils.py pileup2abfreq -gc hg19.gc50Base.txt.gz \  
5 -r normal.fifo \  
6 -s tumor.fifo | gzip > out.abfreq.gz  
7 rm normal.fifo tumor.fifo
```

3 Exploring the *ABfreq* file and GC-correction details

After the aligned sequence data have been pre-processed, *sequenza* R package would handle all the normalization and analysis steps. So the following part of this vignette will take place in R.

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

3.1 Read the file

In the package we provide an example file, to find the complete path of the example data file:

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

The *ABfreq* file can be read all at once, but processing one chromosome at a time is less demanding on computational resources, especially while processing NGS data and might be preferable in case of limited computational resources.

Read only the data corresponding to chromosome 1:

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

Alternatively, read all data at once:

```

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

> str(abf.data, vec.len = 2)

'data.frame':      45003 obs. of  13 variables:
 $ chromosome   : chr  "1" "1" ...
 $ n.base       : int  133037 330227 883223 884960 896946 ...
 $ base.ref     : chr  "T" "G" ...
 $ depth.normal: int   170  9 126 29 151 ...
 $ depth.sample: int   130 13 96 23 131 ...
 $ depth.ratio  : num   0.765 1.444 ...
 $ Af           : num   0.573 0.615 0.533 0.524 0.557 ...
 $ Bf           : num   0.427 0.385 0.467 0.476 0.443 ...
 $ ref.zygosity: chr   "het" "het" ...
 $ GC.percent   : num   38 56 50 68 66 ...
 $ good.s.reads: num   124 13 92 21 115 ...
 $ AB.germline  : chr   "CT" "CT" ...
 $ AB.sample    : chr   "." "." ...

```

The files can be read even faster, after mapping the chromosomes location in the file, it is possible to select the coordinate (in terms of from line x to line y) of the file to read. see the man page of *read.abfreq* for an example.

3.2 Quality control and normalization

Each aligned base, in the next generation sequencing, is associated with a quality score. The *sequenza-utils* software is capable of filtering the base with a quality score lower than a specified value (default, 20). The number of reads that have passed the filter is returned in the column *good.s.reads*, while the *depth.sample* column contains the raw depth indicated in the pileup (from samtools).

3.3 GC-normalization

The GC content bias is affecting most of the sample, however some samples are more biased than others. We attempt to remove this bias by normalizing with the mean depth ratio value of a corresponding GC content value.

It is possible to gather GC-content information from the entire file and in the meantime map the chromosome position in the file (to fast access chromosome by chromosome later, see *?read.abfreq*):

```

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

> str(gc.stats)

List of 6
 $ raw           : num [1:43, 1:3] 0.625 0.504 0.534 0.406 0.502 ...
 ..- attr(*, "dimnames")=List of 2
 .. ..$ : chr [1:43] "8" "10" "12" "14" ...
 .. ..$ : chr [1:3] "25%" "50%" "75%"
 $ adj           : num [1:43, 1:3] 0.866 0.662 0.775 0.601 0.64 ...
 ..- attr(*, "dimnames")=List of 2
 .. ..$ : chr [1:43] "8" "10" "12" "14" ...
 .. ..$ : chr [1:3] "25%" "50%" "75%"
 $ gc.values     : num [1:43] 8 10 12 14 16 18 20 22 24 26 ...
 $ raw.mean      : Named num [1:43] 0.843 0.829 0.7 0.684 0.821 ...
 ..- attr(*, "names")= chr [1:43] "8" "10" "12" "14" ...

```

```

$ raw.median : Named num [1:43] 0.722 0.762 0.69 0.675 0.785 ...
..- attr(*, "names")= chr [1:43] "8" "10" "12" "14" ...
$ file.metrics:'data.frame':      23 obs. of  4 variables:
..$ chr      : Factor w/ 23 levels "1","10","11",...: 1 12 16 17 18 19 20 21 22 2 ...
..$ n.lines: int [1:23] 4542 2728 2396 1822 1775 2370 2178 1425 2087 2430 ...
..$ start   : num [1:23] 1 4543 7271 9667 11489 ...
..$ end     : num [1:23] 4542 7270 9666 11488 13263 ...

```

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

```

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

```

In either case the the normalization to the depth.ratio is performed in the following way:

```

> 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')

```

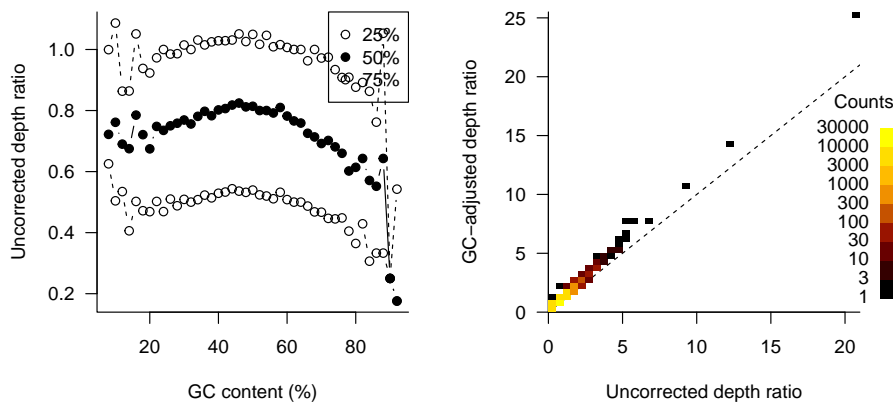


Figure 1: Visualization of depth.ratio bias in relation of GC content (left), and resulting normalization effect (right).

4 Analyzing sequencing data with *sequenza*

The R package *sequenza* offers an ensemble of functions and models that can be used to design customized workflows and analyses. Here we provide generic and most commonly used analyses steps.

- Extract the relevant information from the raw *ABfreq* file.
- Fit the *sequenza* model to infer cellularity and ploidy (*dna.index*).
- Apply the estimated parameter to detect CNV variant alleles

4.1 Extract the information from the *ABfreq* file.

The function *sequenza.extract* is designed to efficiently access the raw ABfreq data and take care of normalization steps. The argument provide the ability of customize a set of action listed below:

- binning depth ratio and B allele frequency in a desired window size (allowing a desired number of overlapping windows);
- performing a fast, allele specific segmentation using the *copynumber* package[2];
- filter mutation by frequency and noise.

```
> test <- sequenza.extract(data.file)
> names(test)
```

After the raw data is processed, the size of the data is considerably reduced. For instance the object resulting from *sequenza.extract* can be stored as a R object within a size contained in few megabytes, even for whole genome sequencing data.

The result of this first step consist in a list of lists. All the sub-lists have a different information subdivided by chromosome. Every list share the same chromosome order.

4.1.1 Plot chromosome view with mutations, BAF, depth ratio and segments

Each chromosome can be visualize using the function *chromosome.view* as in figure 2. The same function can be used to visualize the data after the estimation of *cellularity* and *dna.index* parameters as in Figure 5.

```
> chromosome.view(mut.tab = test$mutations[[1]], baf.windows = test$BAF[[1]],
+                 ratio.windows = test$ratio[[1]], min.N.ratio = 1,
+                 segments = test$segments[[1]], main = test$chromosomes[1])
```

1

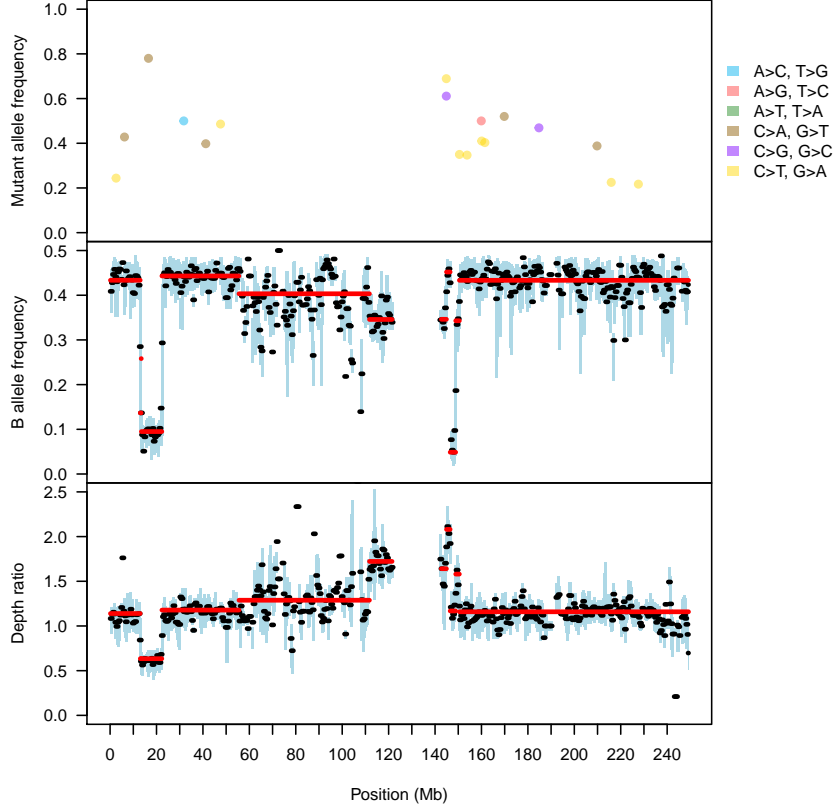


Figure 2: Plots of Mutation (top), B-allele frequencies (middle) and depth ratio (bottom) for chromosome position.

4.2 Inference of cellularity and ploidy (DNA index)

After the raw data is conveniently processed, we can apply the Bayesian inference implemented in the package. The function *sequenza.fit* perform the inference using the calculated B allele frequency and depth ratio of the obtained segments. The method can be explored in more detail by looking the manual pages for the function *baf.model.fit*.

```
> CP.example <- sequenza.fit(test)
```

The result is a list in the format *list(x, y, z)* directly usable by standard graphical function, like *image*. However we provide few customized functions do better display and explore the results, detect the confidence intervals, and finally retrieve the point estimate.

4.2.1 Confidence intervals, confidence region and point estimate

The object resulting from *sequenza.fit* have two vectors, *x* and *y*, indicating respectively the tested values of *dna.index* and *cellularity*, and a matrix *z* with *x* columns and *y* rows, containing the estimated log-likelihood. Confidence intervals for these two parameters can be calculated using the function *get.ci*.

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

It is also possible to plot the likelihood over the combinations of the two parameters, highlighting the point estimate and the confidence region.

```
> cp.plot(CP.example)
```

```
> cp.plot.contours(CP.example, add = TRUE, likThresh = c(0.999))
```

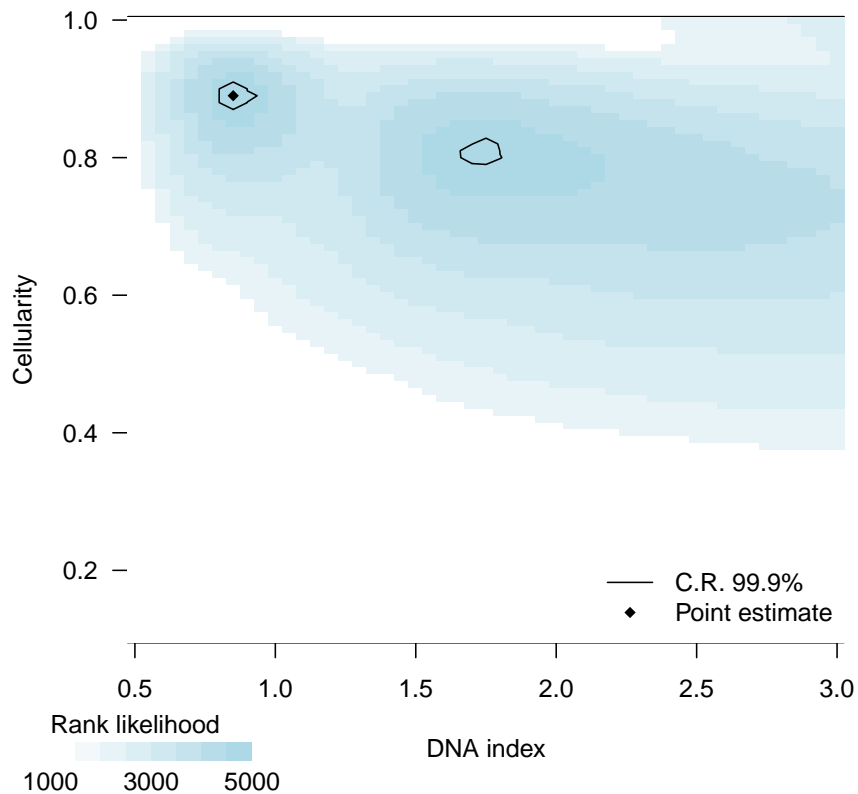


Figure 3: Result from the Bayesian inference over the defined range of cellularity and DNA index. Color intensity indicates the log-likelihood of corresponding cellularity/DNA index values.

By exploring the result for *cellularity* and *dna.index* separately, it is possible to draw the likelihood distribution for each parameter. The information is returned by the *get.ci* function.


```

> par(mfrow = c(2,2))
> cp.plot(CP.example)
> cp.plot.contours(CP.example, add = TRUE)
> plot(cint$values.y, ylab = "Cellularity",
+      xlab = "likelihood", type = "n")
> select <- cint$confint.y[1] <= cint$values.y[,2] &
+      cint$values.y[,2] <= cint$confint.y[2]
> polygon(y = c(cint$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$values.x[,1] <= cint$confint.x[2]
> 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)
>

```

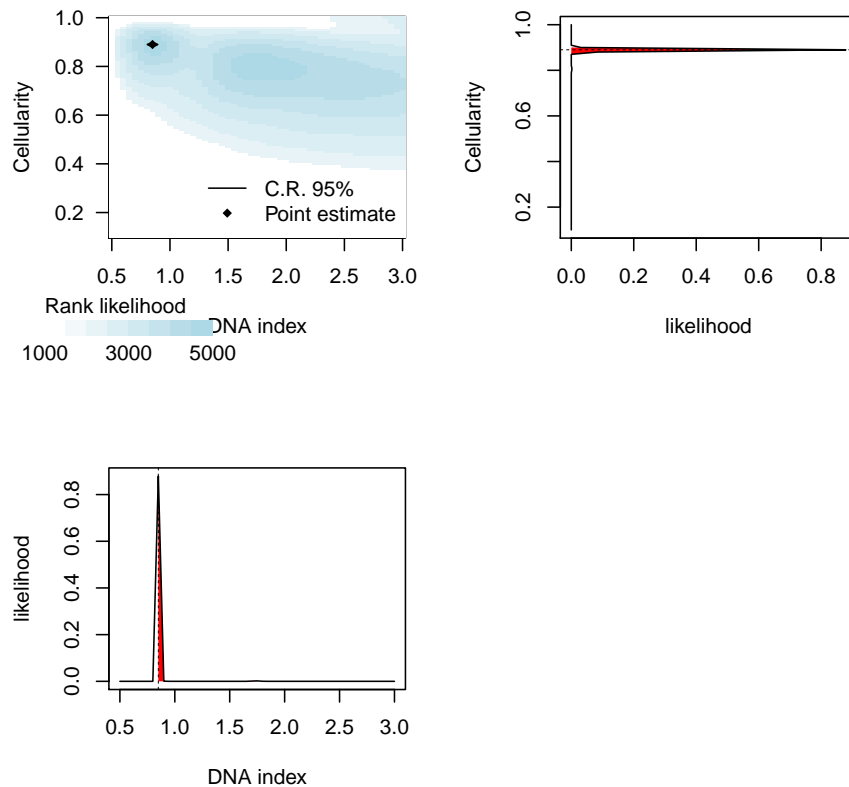


Figure 4: Plot of the log-likelihood with respective cellularity and DNA-index probability distribution and confidence intervals.

4.3 Call CNVs and mutations using the estimated parameters

The point estimate value correspond to the point of maximum likelihood, detected after the confidence interval computation:

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

```
[1] 0.89
```

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

```
[1] 0.85
```

In addition we need to calculate the average normalized depth ratio, used to set a value for the baseline copy number.

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

4.3.1 Detect variant alleles (mutations)

To detect variant alleles, we use a mutation frequency model that is implemented as the *mufreq.bayes* function:

```
> mut.tab      <- na.exclude(do.call(rbind, test$mutations))
> mut.alleles <- mufreq.bayes(mufreq = mut.tab$F,
+                             depth.ratio = mut.tab$adjusted.ratio,
+                             cellularity = cellularity, dna.index = dna.index,
+                             avg.depth.ratio = avg.depth.ratio)
> head(mut.alleles)
```

	CNr	CNt	Mt	L
4	2	2	1	-25.30391
41	2	2	1	-12.49285
2	2	1	1	-12.33863
42	2	2	1	-13.30915
43	2	2	1	-12.99802
44	2	2	1	-13.01904

```
> head(cbind(mut.tab[,c("chromosome", "n.base", "F", "adjusted.ratio", "mutation")],
+               mut.alleles))
```

	chromosome	n.base	F	adjusted.ratio	mutation
1.95	1	2585089	0.244	1.1393516	C>T
1.162	1	6197233	0.428	1.1393516	G>T
1.673	1	16535060	0.780	0.6341823	G>T
1.1253	1	31740870	0.500	1.1779920	T>G
1.1446	1	41234537	0.398	1.1779920	C>A
1.1576	1	47583731	0.486	1.1779920	C>T

	CNr	CNt	Mt	L
1.95	2	2	1	-25.30391
1.162	2	2	1	-12.49285
1.673	2	1	1	-12.33863
1.1253	2	2	1	-13.30915
1.1446	2	2	1	-12.99802
1.1576	2	2	1	-13.01904

The result comprises of four values for every imputed mutation: *CNr* Is the provided copy number of the normal sample at the given position (default = 2); *CNt* is the estimated copy number of the tumor at the given position; *Mt* is the estimated numbers of alleles carrying the mutation; *L* is the log-likelihood of the model fit.

4.3.2 Detect copy number variations

To detect copy number variations we use a B allele frequency model, implemented in the function *baf.bayes*, with the estimated parameters of *cellularity* and *dna.index*:

```
> seg.tab      <- na.exclude(do.call(rbind, test$segments))
> cn.alleles <- baf.bayes(Bf = seg.tab$Bf, depth.ratio = seg.tab$depth.ratio,
+                          cellularity = cellularity, dna.index = dna.index,
+                          avg.depth.ratio = avg.depth.ratio)
> head(cn.alleles)
```

	CNt	A	B	L
[1,]	2	1	1	-11.68157
[2,]	1	1	0	-12.28254
[3,]	1	1	0	-26.57928
[4,]	1	1	0	-10.86406
[5,]	2	1	1	-11.61929
[6,]	2	1	1	-12.33114

```
> seg.tab <- cbind(seg.tab, cn.alleles)
> head(seg.tab)
```

	chromosome	start.pos	end.pos	Bf	N.BAF
1.1	1	133037	12988756	0.43349084	405
1.2	1	13000305	13380725	0.13668009	50
1.3	1	13380510	13448938	0.25821883	9
1.4	1	13450081	22304519	0.09495172	461
1.5	1	22317088	55523989	0.44303403	624
1.6	1	55524484	111741107	0.40325013	471

	depth.ratio	N.ratio	CNt	A	B	L
1.1	1.1393516	430	2	1	1	-11.68157
1.2	0.6345426	57	1	1	0	-12.28254
1.3	0.6341432	10	1	1	0	-26.57928
1.4	0.6341823	541	1	1	0	-10.86406
1.5	1.1779920	646	2	1	1	-11.61929
1.6	1.2879646	538	2	1	1	-12.33114

The result comprises of four values for every imputed segment: *CNt* is the estimated copy number of the tumor of the given segment; *A* is the estimated numbers of A alleles; *B* is the estimated numbers of B alleles; *L* is the log-likelihood of the model fit.

4.4 Visualize detected copy number changes and variant alleles

To visualize the data after detection CNV and variant alleles, it is possible to use the *chromosome.view*. In order to draw the relative model points (and to evaluate how the estimated model fits the real data) more information is needed compared to Figure 2:

- Each segment must have the columns relative to the copy number variation calling.
- *Cellularity* and *dna.index* estimates.
- Average normalized depth ratio.

```

> chromosome.view(mut.tab = test$mutations[[1]], baf.windows = test$BAF[[1]],
+               ratio.windows = test$ratio[[1]], min.N.ratio = 1,
+               segments = seg.tab[seg.tab$chromosome == test$chromosomes[1],],
+               main = test$chromosomes[1],
+               cellularity = cellularity, dna.index = dna.index,
+               avg.depth.ratio = avg.depth.ratio)

```

1

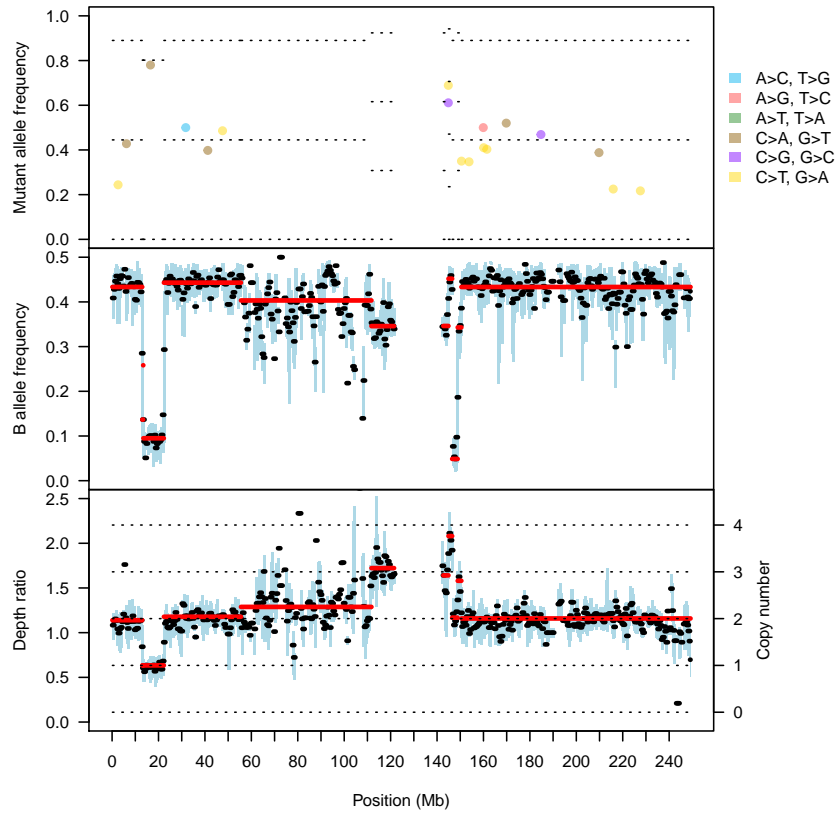


Figure 5: Plots of mutation (top), B-allele frequencies (middle) and depth ratio (bottom) for chromosome position. Horizontal dotted line indicates different copy number/allele state.

4.4.1 Genome-wide view of the alleles and copy numbers state

```
> genome.view(seg.cn = seg.tab, info.type = "CNt")
> legend("bottomright", bty="n", c("Tumor copy number"), col = c("red"),
+       inset = c(0, -0.4), pch=15, xpd = TRUE)
```

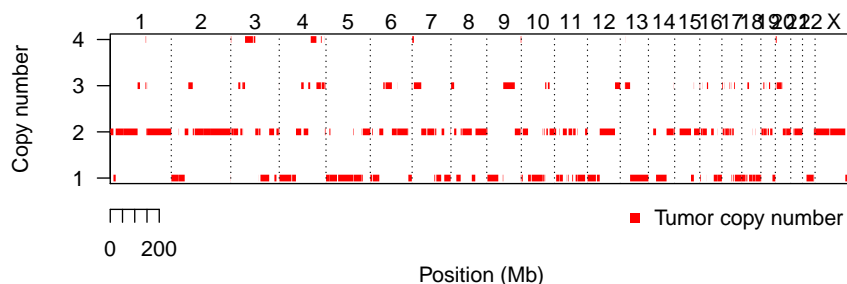


Figure 6: Genome-wide copy number profile obtained from exome sequencing.

```
> genome.view(seg.cn = seg.tab, info.type = "AB")
> legend("bottomright", bty = "n", c("A-allele", "B-allele"), col = c("red", "blue"),
+       inset = c(0, -0.45), pch = 15, xpd = TRUE)
```

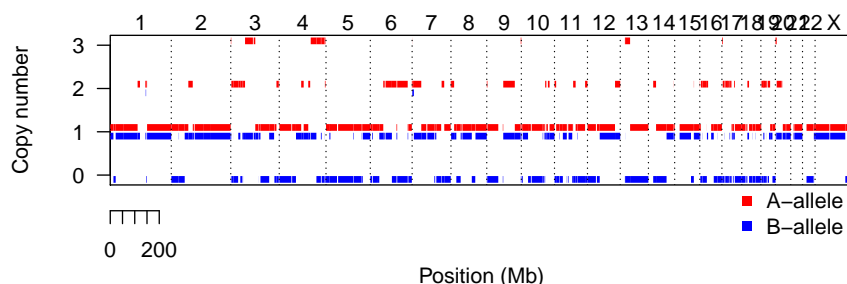


Figure 7: Genome-wide A and B alleles profile obtained from exome sequencing.

References

- [1] Heng Li, Bob Handsaker, Alec Wysoker, Tim Fennell, Jue Ruan, Nils Homer, Gabor Marth, Goncalo Abecasis, and Richard Durbin. The Sequence Alignment/Map format and SAMtools. *Bioinformatics (Oxford, England)*, 25(16):2078–9, August 2009.
- [2] 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.