

# MPAgenomics: An R package for Multi-Patient Analysis of Genomic Markers.

Quentin Grimonprez, Samuel Blanck, Alain Celisse, Guillemette Marot

July, 2014

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Preliminaries</b>	<b>2</b>
2.1	Citing MPAgenomics . . . . .	2
2.2	Installation . . . . .	3
<b>3</b>	<b>Starting example</b>	<b>3</b>
3.1	Download a toy data-set . . . . .	3
3.2	Use of the main functions of the package . . . . .	5
3.3	To go further... . . . .	5
<b>4</b>	<b>Data Normalization</b>	<b>7</b>
4.1	Folder architecture . . . . .	7
4.2	Normalization process . . . . .	9
4.3	Accessing the copy-number, allele B fraction and genotype . . . . .	10
4.4	Usage . . . . .	12
<b>5</b>	<b>Segmentation</b>	<b>13</b>
5.1	From data normalized by aroma . . . . .	13
5.2	For any data provided in matrices . . . . .	15
5.3	To go further . . . . .	16
5.3.1	Choice of the number of segments for cghseg . . . . .	16
5.3.2	Choice of the number of segments for PELT . . . . .	16
<b>6</b>	<b>Calling aberrations in copy-number profiles</b>	<b>18</b>
6.1	From data normalized by aroma . . . . .	18
6.2	For any data provided in matrices . . . . .	19
6.3	filterSeg function . . . . .	21
<b>7</b>	<b>Selection of genomic markers</b>	<b>22</b>
7.1	From data normalized by aroma . . . . .	23
7.2	For any data provided in matrices . . . . .	24
7.3	To go further . . . . .	25

## 1 Introduction

This package provides functions to preprocess and analyze genomic data. Markers refer to SNPs or copy number variations which are designed on the arrays. MPAgenomics is devoted to: *(i)* efficient segmentation and *(ii)* genomic marker selection from multi-patient copy number and SNP data profiles.

For both types of analyses, a pre-processing step (section 4) is proposed as a wrapper of **aroma.\*** packages [Bengtsson, 2004][Bengtsson et al., 2008b][Bengtsson et al., 2010a]. This enables to keep maximum information from the original profiles and improve the multi-patient analysis. In particular, this is useful to keep quantitative data for SNPs rather than usual genotype calls (AA, AB or BB) when these states are not relevant (eg in cancer studies where the number of copies differs from two copies). Note that the use of **aroma.\*** packages [Bengtsson, 2004] implies that the pre-processing step is only available for Affymetrix DNA chips. The other functionalities of the package **MPAgenomics** can however be used independently and such procedures are explained in the sections which detail each step of the proposed analysis.

Efficient segmentation is proposed by the CGHSEG segmentation [Rigaill, 2010][Picard et al., 2005] combined with the slope heuristic [Birgé and Massart, 2007] to choose automatically the number of segments. **MPAgenomics** also suggests an automatic choice of the penalty parameter for the PELT segmentation [Killick et al., 2012] originally implemented in the **changepoint** package [Killick and Eckley, 2013]. **MPAgenomics** offers a pipeline to perform normalization, segmentation and calling at the same time.

Selection of genomic markers relies on the use of penalized regularization techniques implemented in **glmnet** [Friedman et al., 2009] and wrappers are offered to relate this package with the normalisation packages. For the users who would like to go further in the use of penalized regularization techniques, we also provide some guidelines to quickly build inputs for the **HDPenReg** package and interpret some outputs of this latter one.

## 2 Preliminaries

### 2.1 Citing MPAgenomics

Please always cite the following paper when using **MPAgenomics**:

- Grimonprez, Q., Celisse, A., Cheok, M., Figeac, M., and Marot, G. (2014). **MPAgenomics**: An R package for multi-patient analysis of genomic markers. Preprint <http://hal.inria.fr/hal-00933614>

Moreover, **MPAgenomics** wraps and extends functionalities of several packages. Please try to cite the appropriate methodological papers when you use results from the **MPAgenomics** software in a publication, as such citations are the main means by which the authors receive professional credit for their work.

If you use CRMAv2 normalization in **MPAgenomics**, please cite:

- Bengtsson, H., Wirapati, P., and Speed, T. P. (2009). A single-array preprocessing method for estimating full-resolution raw copys from all Affymetrix genotyping arrays including GenomeWideSNP 5 & 6. *Bioinformatics*, 25(17):2149–2156

If you use TumorBoost normalization in **MPAgenomics**, please cite:

- Bengtsson, H., Neuviat, P., and Speed, T. P. (2010b). Tumorboost: Normalization of allele-specific tumor copy numbers from a single pair of tumor-normal genotyping microarrays. *BMC Bioinformatics*, 11

If you use **MPAgenomics** to segment copy number or allele B fraction profiles with the cghseg segmentation method, please cite:

- Picard, F., Hoebeke, M., Lebarbier, E., Miele, V., Rigaill, G., and Robin, S. (2013). cghseg: Segmentation methods for array CGH analysis. R package version 1.0.2

If you use **MPAgenomics** to segment copy number or allele B fraction profiles with the PELT segmentation method, please cite:

- Killick, R., Fearnhead, P., and Eckley, I. A. (2012). Optimal detection of changepoints with a linear computational cost. *Journal of the American Statistical Association*, 107(500):1590–1598

If you use **MPAgenomics** for calling with **CGHcall**, please cite:

- van de Wiel, M. A., Kim, K. I., Vosse, S. J., van Wieringen, W. N., Wilting, S. M., and Ylstra, B. (2007). CGHcall: calling aberrations for array CGH tumor profiles. *Bioinformatics*, 23(7):892–894

## 2.2 Installation

The package can be installed either from a GUI (selecting the R-forge repository) or with the command line (to be typed in the R console)

```
> install.packages("MPAgenomics", repos="http://R-Forge.R-project.org")
```

This vignette can be produced with the command

```
> vignette("MPAgenomics")
```

Since the package provides simple wrappers of many existing packages and all users are not interested by the same wrappers, we did not force **MPAgenomics** to depend on all the packages. We however recommend the user who would like to preprocess data to install **aroma** packages [Bengtsson, 2004] with the following commands:

```
> source("http://www.braju.com/R/hbLite.R")
> hbLite("sfit")
> source("http://bioconductor.org/biocLite.R")
> biocLite("affxparser")
> biocLite("DNAcopy")
> biocLite("aroma.light")
> install.packages("aroma.affymetrix")
> install.packages("aroma.cn")
```

The installation of all these packages at the very beginning of the use of **MPAgenomics** is only for convenience, not mandatory. Indeed, the user who would have forgotten to install some packages would be reminded to install them when applying the corresponding wrapper.

## 3 Starting example

This introductory example aims at helping the user understand the main functions of the package. For more details on each step, we invite the user to read the appropriate section.

### 3.1 Download a toy data-set

The example is based on a free data-set containing 8 .CEL files (251Mo) which can be downloaded from the Broad Institute website, following this link: [http://www.broadinstitute.org/mpg/birdsuite/downloads/birdsuite\\_inputs\\_1.5.3.tgz](http://www.broadinstitute.org/mpg/birdsuite/downloads/birdsuite_inputs_1.5.3.tgz) If you are used to decompress .tgz files, extract this file and copy the following command in the R console, replacing the path to the directory where the raw data (.CEL files) are stored with the appropriate one. Note that we deliberately use “/” and not “\” in `celPATH`. If you copy paste paths from Windows (e.g. looking at properties of the directory), you must change “\” by “/” or “\\”.

```
> celPATH="/home/user/Documents/workdir/CELdata/cel"
```

In this example, the directory pointed by `celPATH` must contain the files shown in Figure 1.

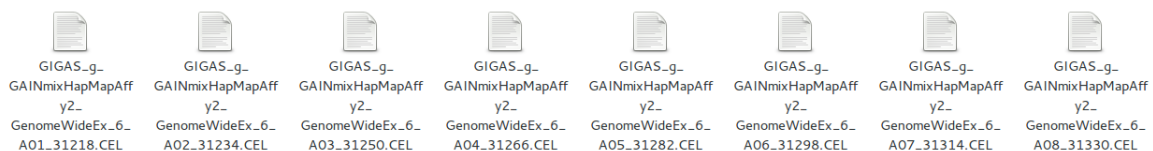


Figure 1: CEL files in `celPATH` (variable containing the path to the downloaded CEL files).

If it is not the case, e.g. if you do not know how to download and extract .tgz files (e.g. via a right click), you can use the following commands in the R console:

```
> #set your working directory (replace with the appropriate path)
> workdir="/home/user/Documents/workdir"
> setwd(workdir)
> #download file
> download.file("http://www.broadinstitute.org/mpg/birdsuite/downloads/birdsuite_inputs_1.5.3.tgz",
+ destfile="./CELdata.tgz")
> #untar the file
> untar("./CELdata.tgz",files="cel",exdir=".")
> #indicate the path containing .cel files
> celPATH="./cel"
```

Once .CEL files are extracted, you must download all the GenomeWideSNP\_6 annotation files which are needed to pre-process the data. These files have extensions .cdf, .ufl, .ugp, .acs. You can download the cdf files (258Mo) from the Affymetrix website (an user account is required), following this link [http://www.affymetrix.com/Auth/support/downloads/library\\_files/genomewidesnp6\\_libraryfile.zip](http://www.affymetrix.com/Auth/support/downloads/library_files/genomewidesnp6_libraryfile.zip)

Extract these files and replace the following path with the appropriate one pointing to the folder containing the extracted .cdf files (only GenomeWideSNP\_6.Full.cdf is needed but you can leave the other files in the same folder):

```
> chipPATH="/home/user/Documents/workdir/CD_GenomeWideSNP_6_rev3/Full/GenomeWideSNP_6/LibFiles/"
```

If you do not know how to decompress .zip files, you can use the following commands in the R console for convenience:

```
> #unzip required files
> unzip("./genomewidesnp6_libraryfile.zip",
+ files=c("CD_GenomeWideSNP_6_rev3/Full/GenomeWideSNP_6/LibFiles/GenomeWideSNP_6.Full.cdf"),exdir=".")
> #indicate the path containing .cdf files
> chipPATH="/home/user/Documents/workdir/CD_GenomeWideSNP_6_rev3/Full/GenomeWideSNP_6/LibFiles/"
```

The other annotation files can be downloaded from [http://www.aroma-project.org/data/annotationData/chipTypes/GenomeWideSNP\\_6/](http://www.aroma-project.org/data/annotationData/chipTypes/GenomeWideSNP_6/). For this example, we download GenomeWideSNP\_6,Full,na31,hg19,HB20110328.ufl.gz (6Mo), GenomeWideSNP\_6,Full,na31,hg19,HB20110328.ugp.gz (7Mo), GenomeWideSNP\_6,HB20080710.acs.gz (37Mo). Extract these files in the folder where your cdf files are in order to have all the annotation files in the same folder (see Figure 2).

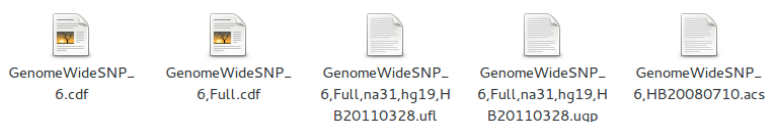


Figure 2: Chip GenomeWideSNP\_6 files in chipPATH (variable containing the path to the downloaded chip files).

If you do not know how to download or decompress .gz files, you can use the following commands in the R console for convenience:

```
> #set the directory where the .cdf files are as your working directory
> setwd(chipPATH)
> ##download the 3 files .ufl, .ugp, .acs
> download.file("http://www.aroma-project.org/data/annotationData/chipTypes/GenomeWideSNP_6/GenomeWideSNP_6,Full,na31,hg19,HB20110328.ufl.gz",
+ destfile="GenomeWideSNP_6,Full,na31,hg19,HB20110328.ufl.gz")
> download.file("http://www.aroma-project.org/data/annotationData/chipTypes/GenomeWideSNP_6/GenomeWideSNP_6,Full,na31,hg19,HB20110328.ugp.gz",
+ destfile="GenomeWideSNP_6,Full,na31,hg19,HB20110328.ugp.gz")
> download.file("http://www.aroma-project.org/data/annotationData/chipTypes/GenomeWideSNP_6/GenomeWideSNP_6,HB20080710.acs.gz",
+ destfile="GenomeWideSNP_6,HB20080710.acs.gz")
> #unzip the gz files
> #install R.utils package containing the gunzip function
> install.packages("R.utils")
> library("R.utils")
> gunzip("GenomeWideSNP_6,Full,na31,hg19,HB20110328.ufl.gz",
+ destname="GenomeWideSNP_6,Full,na31,hg19,HB20110328.ufl")
> gunzip("GenomeWideSNP_6,Full,na31,hg19,HB20110328.ugp.gz",
+ destname="GenomeWideSNP_6,Full,na31,hg19,HB20110328.ugp")
> gunzip("GenomeWideSNP_6,HB20080710.acs.gz",destname="GenomeWideSNP_6,HB20080710.acs")
```

### 3.2 Use of the main functions of the package

In the following example, we suggest to perform segmentation and calling before selection of markers even if only the pre-processing step, which estimates the copy-number and the allele B fraction profile, is required for the last one.

We assume that all the required chip files are in the `chipPATH` folder and that the `.CEL` files are in the `celPATH` folder.

First, we normalize `.CEL` files using the CRMAv2 method [Bengtsson et al., 2009] from `aroma` packages. These packages require a specific architecture containing `.CEL` and annotation files. This normalization step is required before any statistical analysis. The following command performs the creation of architecture and the normalization at the same time:

```
> #list files to check that your path to annotation files (.cdf, .ugp, .ufl, .acs) is correctly set
> dir(chipPATH)
> # list files to check that your path to cel files (.cel) is correctly set
> dir(celPATH)
> #set your working directory (where you have rights to write)
> setwd(workdir)
> #normalize data (might take several hours)
> signalPreProcess(dataSetName="datatest1", chipType="GenomeWideSNP_6",
+ dataSetPath=celPATH, chipFilesPath=chipPATH, path=".",
+ createArchitecture=TRUE, savePlot=TRUE, tags="Full")
```

The first pipeline offered in `MPAgenomics` consists in the segmentation of every copy-number profile of the data. Then, a calling method can be applied on all the segments to label them (e.g. “loss”, “normal”, “gain”). To perform segmentation and calling at the same time, type the following lines in your R console:

```
> segcall=cnSegCallingProcess("datatest1", chromosome=c(1,5), method="PELT")
> #summary of segmentation and calling process
> segcall
```

In this example, chromosomes 1 and 5 for all patients are studied. Figures can be found at `./figures/datatest1/segmentation/`. To better understand the way the results have been produced, we refer to sections 4, 5 and 6.

Results of the calling can be filtered to keep only segments of a minimal size or with a specific label. For example, if you want to only keep segments which represent losses or gains with a minimal size of 10pb and containing at least 2 probes, run the following command:

```
> callfiltered=filterSeg(segcall, minLength=10, minProbes=2, keptLabel=c("gain", "loss"))
> head(callfiltered)
```

To perform selection of markers, a response is associated with each profile and the goal is to find the most relevant markers in the profile according to the response. You can run the following commands:

```
> dataResponse=data.frame(files=getListOfFiles("datatest1"),
+ response=c(2.105092, 1.442868, 1.952103, 1.857819, 2.047897, 1.654766, 2.385327, 2.113406))
> res=markerSelection("datatest1", dataResponse, chromosome=21:22, signal="CN",
+ onlySNP=TRUE, loss="linear")
```

For interpretation of all the results, we refer to section 7.

### 3.3 To go further...

**Use of control samples** In the previous example, there was no control sample. By default, the median of all the samples is used as a reference to estimate the copy-number profile. It is generally better to provide control samples. When this is the case, it is necessary to provide a file or a `data.frame` giving the correspondence between each control and each sample. For this example, we assume that the control is arbitrarily the first file:

```

> #get the file names of our data-set
> files=getListOfFiles("datatest1")
> #create the data.frame linking normal and tumor files
> normalTumorArray=data.frame(normal=rep(files[1],7),tumor=files[2:8])

normal,tumor
GIGAS_g_GAINmixHapMapAffy2_GenomeWideEx_6_A01_31218,GIGAS_g_GAINmixHapMapAffy2_GenomeWideEx_6_A02_31234
GIGAS_g_GAINmixHapMapAffy2_GenomeWideEx_6_A01_31218,GIGAS_g_GAINmixHapMapAffy2_GenomeWideEx_6_A03_31250
GIGAS_g_GAINmixHapMapAffy2_GenomeWideEx_6_A01_31218,GIGAS_g_GAINmixHapMapAffy2_GenomeWideEx_6_A04_31266
GIGAS_g_GAINmixHapMapAffy2_GenomeWideEx_6_A01_31218,GIGAS_g_GAINmixHapMapAffy2_GenomeWideEx_6_A05_31282
GIGAS_g_GAINmixHapMapAffy2_GenomeWideEx_6_A01_31218,GIGAS_g_GAINmixHapMapAffy2_GenomeWideEx_6_A06_31298
GIGAS_g_GAINmixHapMapAffy2_GenomeWideEx_6_A01_31218,GIGAS_g_GAINmixHapMapAffy2_GenomeWideEx_6_A07_31314
GIGAS_g_GAINmixHapMapAffy2_GenomeWideEx_6_A01_31218,GIGAS_g_GAINmixHapMapAffy2_GenomeWideEx_6_A08_31330

```

Figure 3: CSV file linking normal sample with tumoral sample.

Note that the extension .CEL is not provided in the `normalTumorArray` dataframe. In the following of the vignette, a study which provides both control and tumoral samples will be referred as a normal-tumor study.

A new normalization can be run by taking into account information from the control sample. To separate analyses with or without control samples, we will store results from this normal-tumor study under the name “datatest2”. As the architecture has already been created for “datatest1”, it is not necessary to run `signalPreProcess` with `createArchitecture=TRUE`. It is faster to add the CEL files in “datatest2” while leaving the chip definition files as they are by using the function `addData` and then `signalPreProcess` with `createArchitecture=FALSE`.

```

> addData(dataSetName="datatest2",dataPath=celPATH,chipType="GenomeWideSNP_6")
> signalPreProcess(dataSetName="datatest2", chipType="GenomeWideSNP_6",
+ normalTumorArray=normalTumorArray, createArchitecture=FALSE, savePlot=TRUE, tags="Full")

```

**Single-Patient Analysis of allele B fraction** The segmentation example (3.2) was run on the copy-number profile, the allele B fraction can also be segmented after some transformation but only heterozygous SNPs are kept. First, a naive genotype call [Bengtsson et al., 2010b] is performed on each normal sample in order to separate heterozygous SNPs from homozygous SNPs. Naive genotyping method assumes SNPs are bi-allelic and therefore is not recommended for tumor samples. Thus allele B fraction segmentation in `MPAgenomics` requires matched normal-tumor pairs. Then, the resulting signal is symmetrized around 0.5 (a point  $x$  become  $2 * |x - 0.5|$ ) [Staaf et al., 2008](see Figure 4), which makes it similar to the usual copy number.

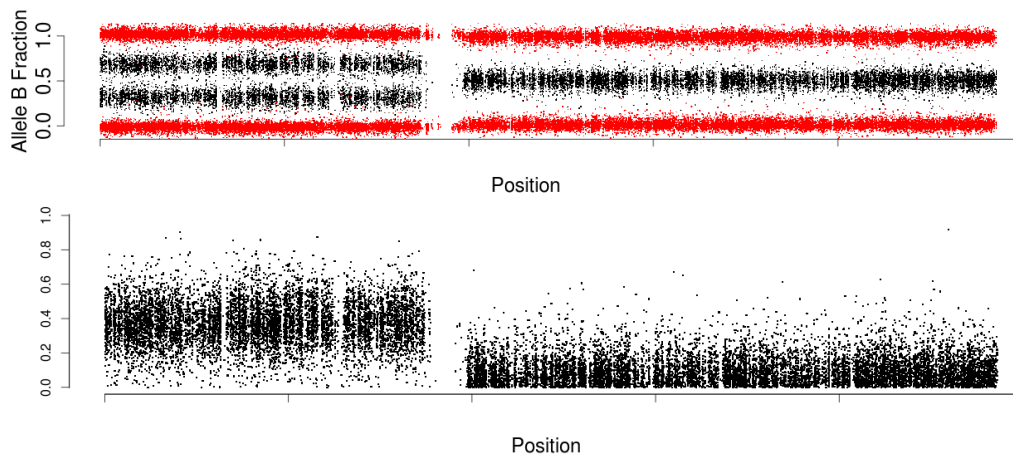


Figure 4: Top: Allele B fraction profile on a chromosome (in red: homozygous SNPs). Bottom: symmetrized profile without homozygous SNPs.

To perform the segmentation on the allele B fraction, run the following command

```

> #run the segmentation
> segfracB=segFracBSignal("datatest1",chromosome=c(1,5))

```

```
> #print summary of segmentation
> segfracB
```

**Selecting steps to be executed to avoid a complete procedure** Each step can be performed with or without the aroma packages if arguments of functions are given in the right format (excepted the normalization procedure described in section 4). Almost each wrapper can be used from matrices independently from the other steps and more details are given in the following of this vignette. The paper is organized as follows: the section 4 presents the preprocessing step and how to access normalized copy-number and allele B fraction. If you already have these data, you can skip this section. Then, an efficient single patient analysis is introduced in the next two sections. The segmentation process is described in section 5. In section 6, we present the calling process of segmented data. The selection of markers is described in section 7.

## 4 Data Normalization

This preprocessing step consists in a correction of biological and technical biases due to the experiment.

MPAgenomics supports Affymetrix<sup>®</sup> arrays (5.0, 6.0 and CytoScanHD). Raw data from Affymetrix DNA chip are provided in different .CEL files. These data need to have some corrections and normalizations before being usable. This normalization process can be performed via functions from **aroma.\*** packages [Bengtsson, 2004] [Bengtsson et al., 2008b]. Here, we provide some user-friendly wrappers to these functions. We assume that all the required chip files are in the **chipPATH** folder and the CEL files are in the **celPATH** folder.

### 4.1 Folder architecture

All the functions from aroma packages have to be used with a particular architecture in the working directory. In MPAgenomics, we use the same architecture and provide some functions to create it.

A function is provided to create the architecture and copy the specified files in the right directory : **createArchitecture**. It also checks that the names of the chip files begin by the specified **chipType** argument.

```
createArchitecture(dataSetName,chipType,dataSetPath,chipFilesPath,path,verbose,tags)
```

Parameters	Description
<b>dataSetName</b>	The name of the data-set folder to create.
<b>chipType</b>	The name of the used chip.
<b>dataSetPath</b>	Path to the folder containing the data CEL files.
<b>chipFilesPath</b>	Path to the folder containing the chip files.
<b>path</b>	Path where the architecture should be created (default=".").
<b>verbose</b>	Print information during the process (default=FALSE).
<b>tags</b>	Common tag which appears in the different file names (cdf, ugp, ufl) of the chip. For no tag, use tags=NULL (default = NULL).

For example, assume that we have the downloaded files of the download section (Fig 5 & 6).

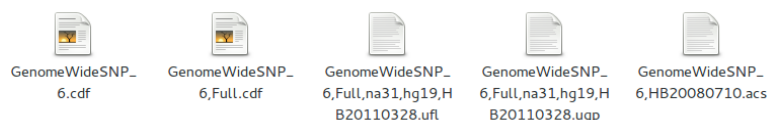


Figure 5: Chip GenomeWideSNP\_6 files in **chipPATH** (variable containing the path to the downloaded chip files).

The **celPATH** and **chipPATH** variables contain the different paths to the downloaded files. The following command enables to both create the architecture and add the previous dataset under the name "datatest1".

```
> createArchitecture("datatest1","GenomeWideSNP_6",celPATH,chipPATH,".",TRUE,"Full")
```

In your working directory, you must find the following architecture :

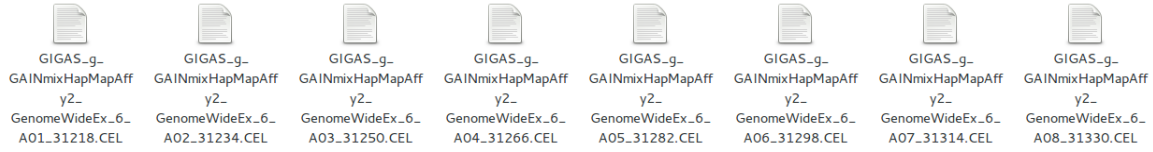


Figure 6: CEL files in `celPATH` (variable containing the path to the downloaded CEL files).

```
<current working directory>
+- annotationData/
| +- chipTypes/
|   +- GenomeWideSNP_6/
|     +- CDF file(s) and other annotation
|
+- rawData/
  +- datatest1/
    +- GenomeWideSNP_6/
      +- CEL files
```

You can see that files are only copied and not moved.

This architecture can also be created with the normalization function (see section 4.2).

The minimal architecture contains 2 folders: `annotationData` and `rawData`. Some supplementary folders will be created during the execution of some functions.

In the `annotationData` folder, each type of chip used for the experiment has to have his own folder containing its definition files (.cdf) and other files (.ugp<sup>1</sup>, .acs<sup>2</sup>, .uff<sup>3</sup>). In this folder, every file name must comply with the following format `<chipType>,<tags>.<extension>` (see Fig. 5).

The `rawData` folder contains the different data-sets with the CEL files. Each data-set folder contains one or more folders (one per different chip type used) containing the CEL files, these files will not be modified during the process. Note that the `chipType` folder in the data-set folder must match exactly a chip type folder under `annotationData`.

In case of new data-sets or new chip types, it is not useful to create a new architecture but just add new folders in the existing architecture as follows:

```
<current working directory>
+- annotationData/
| +- chipTypes/
|   +- <chipTypeA>/ <-- must match exactly the name of the CDF file (fullname minus tags)
|     +- CDF file(s) and other annotation (possibly subdirectories)
|     |
|     +- <chipTypeB>/ <-- must match exactly the name of the CDF file (fullname minus tags)
|       +- CDF file(s) and other annotation (possibly subdirectories)
|       ...
|
+- rawData/
| +- <dataSet1>/
|   +- <chipTypeA>/ <-- must match exactly a chip type folder under annotationData/
|     +- CEL files
|
| +- <dataSet2>/
|   +- <chipTypeB>/ <-- must match exactly a chip type folder under annotationData/
|     +- CEL files
|
| +- <dataSet3>/
```

<sup>1</sup>See <http://www.aroma-project.org/node/43> for more details.

<sup>2</sup>See <http://www.aroma-project.org/node/100> for more details.

<sup>3</sup>See <http://www.aroma-project.org/node/47> for more details.



```
| | +- <chipTypeA>/ <-- must match exactly a chip type folder under annotationData/
| |   +- CEL files
| | +- <chipTypeB>/ <-- must match exactly a chip type folder under annotationData/
| |   +- CEL files
| ...
```

In order to easily update your existing architecture you can use the following `addData` function. It enables to add to your existing architecture a new data-set in the `rawData` folder.

```
addData(dataSetName,dataPath,chipType)
```

Parameters	Description
<code>dataSetName</code>	The name of the data-set folder to create.
<code>dataPath</code>	Path of the folder containing the data CEL files.
<code>chipType</code>	The name of the used chip.

The `addChipType` function enables to add to your existing architecture a new chip type in the `annotationData` folder.

```
addChipType(chipType,chipPath)
```

Parameters	Description
<code>chipType</code>	Name of the new chip type to add.
<code>chipPath</code>	Path to the files to add.

With the function `getListOfFiles`, you can obtain the list of the files contained in the specified data-set.

```
getListOfFiles(dataSetName,chipType)
```

Parameters	Description
<code>dataSetName</code>	The name of a data-set folder.
<code>chipType</code>	The name of the used chip.

If you do not specify a `chipType`, it returns the files for the first chip in the alphabetic order in the `dataSetName` folder.

## 4.2 Normalization process

The main function for the normalization of CEL files is `signalPreProcess`. This function executes 3 different methods on your data:

- **CRMAv2** [Bengtsson et al., 2009]: Normalize profiles to obtain the copy-number and the allele B fraction. This step consists in the correction of biological and technical biases due to the experiment.
- **Genotype calls**: Assign label *AA*, *AB* or *BB* to each SNP.
- **TumorBoost** [Bengtsson et al., 2010b]: Only in a case of normal-tumor study, normalization of the allele B fraction tumor profile using the control profile.

The function `signalPreProcess` executes all these steps and saves the results in different folders in the aroma architecture (`totalAndFracBData,...`), you can also obtain some graphics in the `figures/signal` folder.

The `signalPreProcess` is defined as follows:

```
signalPreProcess(dataSetName, chipType, normalTumorArray, dataSetPath, chipFilesPath, path,
createArchitecture, savePlot)
```

Parameters	Description
<b>dataSetName</b>	The name of the data-set.
<b>chipType</b>	The type of the used chip (e.g. "GenomeWideSNP_6").
<b>normalTumorArray</b>	Only in the case of normal-tumor study. A csv file or a data.frame containing the mapping between normal and tumor files(Fig. 7).
<b>dataSetPath</b>	Only if <b>createArchitecture=TRUE</b> . Path to the folder containing the CEL files of the data-set.
<b>chipFilesPath</b>	Only if <b>createArchitecture=TRUE</b> . Path to the folder containing all the annotations files for the specified chip type.
<b>path</b>	Only if <b>createArchitecture=TRUE</b> . Path where the architecture should be created (default=".").
<b>createArchitecture</b>	if TRUE, the aroma architecture will be automatically created (default=TRUE).
<b>savePlot</b>	if TRUE, graphics of the CN signal and allele B fraction signal will be saved in the <b>figures/signal</b> folder. (default=TRUE).
<b>tags</b>	Common tag which appears in the different filenames (cdf, ugp, ufl) of the chip. For no tag, use tags=NULL (default = NULL).

```

normal,tumor
patient1_normal,patient1_TUMOR
patient2_normal,patient2_tumor

```

Figure 7: Example of the content of a csv file for **normalTumorArray** parameter. The first column contains the name of the different normal files (without the .cel extension) in the data-set folder in rawData. The second column contains the name of the tumor files. The name of the two columns are respectively normal and tumor. The extensions of the files (.CEL for example) should be removed.

If you have specified **createArchitecture=TRUE** and have given the required parameters, the function will create the architecture and copy your files in the right folder. After the creation of the aroma architecture, it runs the normalization process and saves in the aroma architecture all the files necessary to obtain the allele B fraction and the copy-number signal.

In the case of a study without reference, only parameters **dataSetName**, **chipType**, **createArchitecture** and **savePlot** have to be specified. In the case of a normal-tumor study, TumorBoost process is executed. If you want a CEL file to be used as reference for all the tumor data, then fill the normal column of the **normalTumorArray** with the filename of your reference for all files.

### 4.3 Accessing the copy-number, allele B fraction and genotype

After running the **signalPreProcess** function, you can not directly access the copy-number and the allele B fraction signals with the files in the architecture, a treatment is required to read them. Note that the copy-number signal saved is the raw signal, it needs a reference to normalize it. With the **getCopyNumberSignal** and **getFracBSignal** functions, you can access the copy-number and the allele B fraction profiles.

The **getCopyNumberSignal** allows you to extract the copy-number signal a chromosome at a time. In the case of a study without reference, the median of all the signals of the data-set is used to normalize every profile. In the normal-tumor study, the normal signal is used to normalize the tumor signal.

The getter for the copy-number profile is defined as follows:

```
getCopyNumberSignal(dataSetName,chromosome,normalTumorArray,onlySNP,listOfFiles,verbose)
```

Parameters	Description
<b>dataSetName</b>	The name of the data-set folder.
<b>chromosome</b>	A vector containing the chromosomes for which the signal will be extracted.
<b>normalTumorArray</b>	Only in the case of normal-tumor study. A csv file or a data.frame containing the mapping between normal and tumor files(Fig. 7).
<b>onlySNP</b>	If TRUE, only the copy-number for SNPs positions will be returned (default=FALSE).
<b>listOfFiles</b>	A vector containing the names of the files in dataSetName folder for which the copy-number profiles will be extracted (default is all the files).
<b>verbose</b>	If TRUE print some information (default=TRUE).

This function returns a list of size `length(chromosome)` (each element is named `chrX` with `X` the number of the chromosome) containing a data.frame with columns:

- **chromosome** Chromosome of the signal.
- **position** Positions associated with the copy-number.
- **copynumber** Copy number profiles of selected files; the name of each column is the name of the associated data file name.
- **featureNames** Names of the probes.

If you want the allele B fraction, you can use the `getFracBSignal` function. In a normal-tumor study, allele B fraction for both normal and tumoral samples is returned. For the tumoral sample, it is the allele B fraction after the TumorBoost normalization.

```
getFracBSignal(dataSetName,chromosome,normalTumorArray,listOfFiles,verbose)
```

Parameters	Description
<b>dataSetName</b>	The name of the data-set folder.
<b>chromosome</b>	A vector containing the chromosomes for which the allele B fraction signal must be extract.
<b>normalTumorArray</b>	Only in the case of normal-tumor study. A csv file or a data.frame containing the mapping between normal and tumor files(Fig. 7).
<b>listOfFiles</b>	A vector containing the names of the files in dataSetName folder for which the allele B fraction profiles will be extracted (default is all the files).
<b>verbose</b>	If TRUE print some information (default=TRUE).

This function returns a list of size `length(chromosome)` (each element is named `chrX` with `X` the number of the chromosome) containing a data.frame with columns:

- **chromosome** Chromosome of the signal.
- **position** Positions associated with the allele B fraction.
- **fracB** Allele B fraction profiles of selected files; the name of each column is the name of the associated data file name.
- **featureNames** Names of the probes.

If you want the genotype call, you can use the `getGenotypeCalls` function.

```
getGenotypeCalls(dataSetName,chromosome,listOfFiles=NULL,verbose=TRUE)
```

Parameters	Description
<b>dataSetName</b>	The name of the data-set folder.
<b>chromosome</b>	A vector containing the chromosomes for which the genotype call will be extracted.
<b>listOfFiles</b>	A vector containing the names of the files in dataSetName folder for which the genotype signal will be extracted (default is all the files).
<b>verbose</b>	If TRUE print some information (default=TRUE).

This function returns a list of size `length(chromosome)` (each element is named `chrX` with `X` the number of the chromosome) containing a data.frame with columns:

- **chromosome** Chromosome of the signal.
- **position** Positions associated with the genotype.
- **genotype** Genotype calls corresponding to selected files; the name of each column is the name of the associated data file name.
- **featureNames** Names of the probes.

The next getter, `getSymFracBSignal`, returns the symmetrized allele B fraction only for heterozygous positions. To symmetrize the allele B fraction the transformation  $x \mapsto 2 * |x - 0.5|$  is applied. It centers the data in 0.5 corresponding to heterozygous allele B fraction then symmetrize and multiply by 2 to have a signal between 0 and 1.

This method can be run only in a normal-tumor case, and each tumor file must match with a unique normal file.

`getSymFracBSignal(dataSetName,chromosome,normalTumorArray,file,verbose=TRUE)`

Parameters	Description
<b>dataSetName</b>	The name of the data-set folder.
<b>chromosome</b>	A vector with the chromosomes for which the symetrized signal will be extracted .
<b>normalTumorArray</b>	A csv file or a data.frame containing the mapping between normal and tumor files(Fig. 7).
<b>file</b>	The name of the file in dataSetName folder for which the symetrized signal will be extracted .
<b>verbose</b>	If TRUE print some information (default=TRUE).

This function returns a list of size `length(chromosome)` (each element is named `chrX` with `X` the number of the chromosome) containing a data.frame with columns:

- **chromosome** Chromosome of the signal.
- **position** Positions associated with the genotype.
- **fracB** One column named by the data file name. It contains the symmetrized allele B fraction signal for the specified profile.
- **featureNames** Names of the probes.

## 4.4 Usage

We refer to sections 3.2 and 3.3 to normalize input data without or with control samples in the study, respectively.

### Get a signal

The following commands enable to extract the copy-number or the allele B fraction profiles of all the files for the chromosome 5:

```
> #normal-tumor study
> CNdata2=getCopyNumberSignal("datatest2",5,normalTumorArray=normalTumorArray,TRUE)
> fracBdata2=getFracBSignal("datatest2",5,normalTumorArray=normalTumorArray)
> #study without reference
> CNdata1=getCopyNumberSignal("datatest1",5,onlySNP=TRUE)
> fracBdata1=getFracBSignal("datatest1",5)
```

To access the copy-number or allele B fraction profiles, run:

```
> CNdata2$chr5
> fracBdata2$chr5$tumor
> fracBdata2$chr5$normal
> fracBdata1$chr5$tumor
```

## 5 Segmentation

In the package **MPAgenomics**, we provide two segmentation methods. The first one wraps the **cghseg** package which fits a Gaussian maximum-likelihood piecewise constant smoothing model [Picard et al., 2005] using pruned dynamic programming [Rigaill, 2010]. When the user chooses this method, **MPAgenomics** selects the number of segments with the slope heuristic [Birgé and Massart, 2007].

The second segmentation approach uses the PELT method of the **cpt.mean** function from the **changepoint** package [Killick and Eckley, 2013]. In the initial PELT method, the penalty is  $K\rho\log(P)$  with a profile of length  $P$ ,  $K$  the number of segments and a  $\rho > 0$  a constant to optimize. When the user chooses this method, **MPAgenomics** suggests another calibration for the penalty parameter. For more details, see Section 5.3.

### 5.1 From data normalized by aroma

The procedure described above can be run by the **segmentationAroma** function if the working directory respects the aroma architecture.

```
segmentationAroma(dataSetName,normalTumorArray,chromosome,Rho=NULL,Kmax=10,listOfFiles=NULL,onlySNP=TRUE,
savePlot=TRUE,method="cghseg")
```

Parameters	Description
<b>dataSetName</b>	The name of the data-set folder.
<b>normalTumorArray</b>	Only in the case of normal-tumor study. A csv file or a data.frame containing the mapping between normal and tumor files(Fig. 7).
<b>chromosome</b>	A vector with the chromosomes to be segmented.
<b>Rho</b>	Only for method="PELT", a vector containing all the penalization values to test for the segmentation. If no values are provided, default values will be used.
<b>Kmax</b>	Only for method="cghseg", maximal number of segments.
<b>onlySNP</b>	If TRUE, only the copy-number for SNPs positions will be returned (default=TRUE).
<b>listOfFiles</b>	A vector containing the names of the files in dataSetName folder for which the copy number profiles will be segmented (default is all the files).
<b>savePlot</b>	if TRUE, graphics of the segmented CN signal will be saved in the <b>figures/dataSetName/segmentation/CN</b> folder. (default=TRUE).
<b>method</b>	"PELT" or "cghseg". (default="cghseg").

If **savePlot=TRUE**, some graphics will be plotted (see Fig 8).

The graphics of the segmented profiles are saved in the **figures/dataSetName/segmentation** folder. A summary of the segmented profile is saved in a text file (Fig. 9) in the **segmentation/dataSetName/** folder, it contains 5 columns: **chrom**, **chromStart**, **chromEnd**, **probes** and **means** containing the number of the chromosome, the starting position of the segment, the ending position of the segment, the number of probes in the segment and the mean of the segment.

The output of the function is a list where every element of the list contains the segmentation results for a different profile formatted in a list containing:

- **copynumber** A vector containing the copy-number signal.
- **segmented** A vector of the same size as copynumber containing the segmented values.
- **startPos** The position of every probe.
- **chromosome** A vector of the same size as copynumber containing the chromosome number.
- **featureNames** Names of the probes.
- **sampleNames** The name of the signal.
- **segment** A data.frame that sums up the results of the segmentation. Each row is a different segment with the chromosome, start position, end position, number of probes in the segment and the value of the segment.

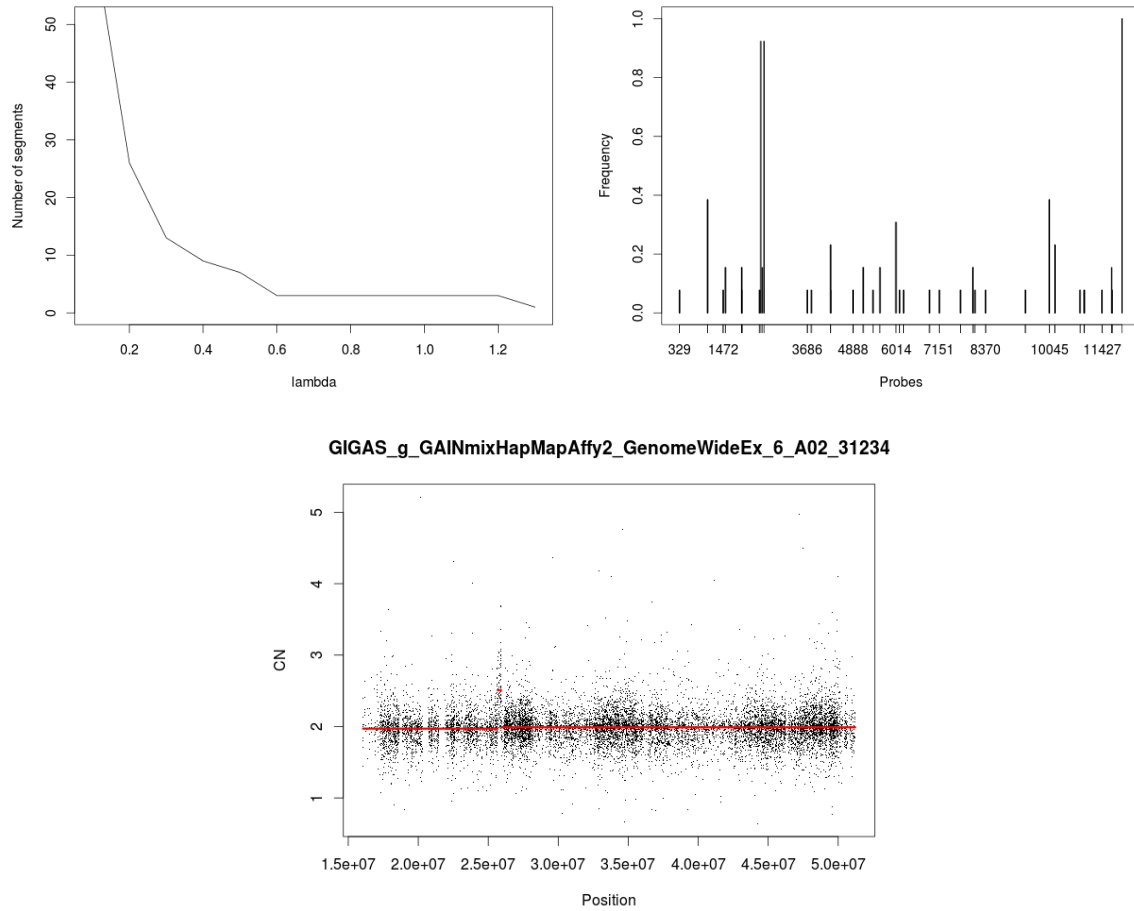


Figure 8: Top left: variation of the number of segments with regards to  $\rho$ . Top right: the frequency of occurrence of breakpoints in all segmentation. Bottom: Segmentation with the optimal  $\rho$  found.

"chrom"	"chromStart"	"chromEnd"	"probes"	"means"
"chr21"	9764384	19525218	1652	1.98922631493452
"chr21"	19529011	19529012	1	8.90689831889179
"chr21"	19529713	48084820	11229	2.00628461788311
"chr22"	16053757	25661698	2465	1.97047216297763
"chr22"	25664248	25918709	85	2.50995556906446
"chr22"	25925077	51219006	9403	1.98876603818565

Figure 9: Summary of a segmentation for a profile.

The output for a profile is formatted in order to be usable in input of the `callingProcess` function (see section 6.2) but if you use the aroma architecture and want to label your segments, you can directly use the `cnSegCallingProcess` function (see section 6.1). This function runs the segmentation and the calling processing.

In the following example, the segmentation processing is run on one file of the downloaded data:

```
> file="GIGAS_g_GAINmixHapMapAffy2_GenomeWideEx_6_A02_31234"
> seg1=segmentationAroma("datatest1", chromosome=21:22, onlySNP=TRUE, plot=TRUE,
+ listOfFiles=file, method="PELT")
```

The results obtained for one profile are shown in figures 8 & 9.

For a normal-tumor study where each tumor file matches with a unique normal file (and only in this case), it is possible to segment the symmetrized allele B fraction with the function `segFracBSignal`:

```
segFracBSignal(dataSetName,normalTumorArray,chromosome=1:22,Rho=NULL,Kmax=10,listOfFiles=NULL,
savePlot=TRUE,verbose=TRUE,method="cghseg")
```

Parameters	Description
<b>dataSetName</b>	The name of the data-set folder.
<b>normalTumorArray</b>	A csv file or a data.frame containing the mapping between normal and tumor files(Fig. 7).
<b>chromosome</b>	A vector with the chromosomes to be segmented.
<b>Rho</b>	Only for method="PELT", a vector containing all the penalization values to test for the segmentation. If no values are provided, default values will be used.
<b>Kmax</b>	Only for method="cghseg", maximal number of segments.
<b>listOfFiles</b>	A vector containing the names of the files in dataSetName folder for which the allele B profile is segmented (default is all the files).
<b>savePlot</b>	if TRUE, graphics of the segmented allele B profile will be saved in the <b>figures/dataSetName/segmentation/fracB</b> folder. (default=TRUE).
<b>verbose</b>	If TRUE, print some information (default=TRUE).

The output of the function is a data.frame containing the segmentation results, each row corresponding to a different segment and with columns:

- **sampleName** Name of the file the segment belongs to.
- **chromosome** Chromosome the segment belongs to.
- **chromStart** Starting position of the segment.
- **chromEnd** Ending position of the segment.
- **probes** Number of probes in the segment.
- **means** Mean of the segment.

## 5.2 For any data provided in matrices

Assuming that normalized data are stored in a matrix, you can use the `segmentation` function to perform segmentations with parameter choices suggested in `MPAgenomics`.

```
segmentation(signal,Rho=NULL,Kmax=10,position=NULL,plot=TRUE,verbose=TRUE,method="cghseg")
```

Parameters	Description
<b>signal</b>	A vector containing the signal.
<b>Rho</b>	Only for method="PELT", a vector containing all the penalization values to test for the segmentation. If no values are provided, default values will be used.
<b>Kmax</b>	Only for method="cghseg", maximal number of segments.
<b>position</b>	A vector containing the position of all elements of the signal (not necessary).
<b>plot</b>	If TRUE, plot some graphics (default=TRUE ).
<b>verbose</b>	If TRUE, print some information (default=TRUE).
<b>method</b>	"PELT" or "cghseg". (default="cghseg").

If `plot=TRUE`, some graphics will be plotted (see Fig 8).

The output of the function is a list containing:

- **signal** A vector containing the signal.
- **segmented** A vector of the same size as the signal containing the segmented values.
- **startPos** The position of each probe.

- **segment** A data.frame that sums up the results of the segmentation. Each row is a different segment with the start position, end position, number of points in the segment and the value of the segment.

For the following example, a matrix is built from the preprocessed files obtained in section 4.2. A copy number profile from one patient chromosome is then isolated:

```
> file="GIGAS_g_GAINmixHapMapAffy2_GenomeWideEx_6_A07_31314"
> CNdata1=getCopyNumberSignal("datatest1",20,onlySNP=TRUE,listOfFiles=file)
> copyNumber=CNdata1$chr20$GIGAS_g_GAINmixHapMapAffy2_GenomeWideEx_6_A07_31314
> position=CNdata1$chr20$position
```

Then, the segmentation process can be run using:

```
> seg=segmentation(copyNumber,position=position,plot=TRUE,verbose=TRUE,method="PELT")
```

This leads to the same plots as in Fig 8 and the results are saved in the right format for the calling function (see section 6.2).

The list of found segments can be accessed with:

```
> seg$segment
```

	start	end	points	means
1	61795	8099663	3573	2.051816
2	8102867	8571315	222	2.634067
3	8571934	62912463	19626	2.047356

## 5.3 To go further

In this part, we describe the different calibration methods associated with cghseg and PELT in **MPAgenomics**. We also illustrate that calibrating the parameter according to a variance criterium can be misleading.

### 5.3.1 Choice of the number of segments for cghseg

cghseg method runs with a predetermined number of segments. In general, users fix a maximal number of segments  $K_{\max}$  and then must choose the best segmentation according to some criteria. **MPAgenomics** implements the slope heuristic [Birgé and Massart, 2007] to choose automatically the number of segments.

Let  $k$  the number of segments, the criterium [Lebarbier, 2005] to minimize is

$$g(k, C) = \frac{1}{P} \left( R(k) + C \times k \times \left( 2.5 + \log\left(\frac{P}{k}\right) \right) \right)$$

where  $R(k)$  is the quadratic loss between the observed and the segmented signal, and  $C$  a constant to optimize.

The first step is to find the optimal constant  $\hat{C}$  to use. Given a list  $\mathcal{C}$  of potential constant, for each element  $C \in \mathcal{C}$ , we find the optimal  $\hat{k}$  minimizing  $g(k, C)$ . Then, the optimal constant  $\hat{C}$  is  $2 \times C_0$ , with  $C_0$  the constant associated with the largest drop of the number of segments in Fig. 10.

Once the optimal constant  $\hat{C}$  found, the optimal number of segments  $\hat{k}$  is the one minimizing  $g(k, \hat{C})$ .

### 5.3.2 Choice of the number of segments for PELT

The penalty of PELT looks like  $\rho \times K \times \log(P)$  with  $K$  the number of segments and  $P$  the length of the profile. Only the parameter  $\rho$  can be chosen by the user. As the default value in PELT ( $\rho = 1$ ) did not provide satisfying segments compared to manually annotated segments, which we observed on 70 profiles, we suggest another method to calibrate the  $\rho$  parameter. To make the choice of the optimal  $\rho$  in **MPAgenomics**, the PELT method is run for a range of  $\rho$  provided by default or given by the user. Then, the evolution of the number of segments is plotted with regards to the  $\rho$  parameter (Fig 11).



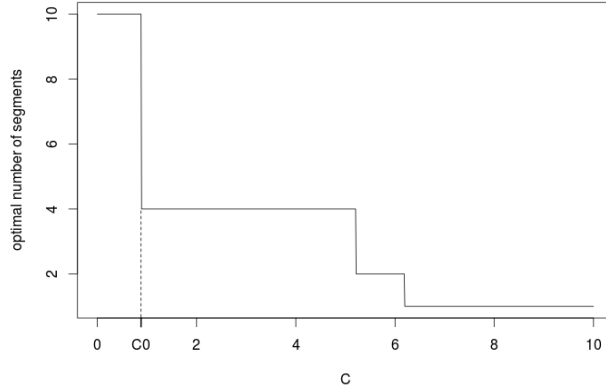


Figure 10: Optimal number of segments minimizing  $g(C, k)$  with regards to  $C$ .  $C_0$  is the constant associated with the largest drop of the number of segments.

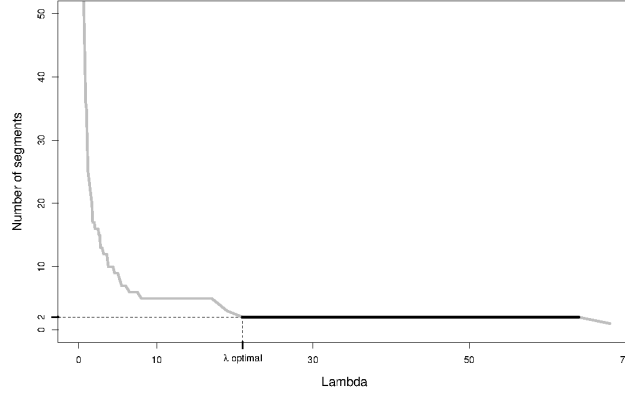


Figure 11: Variation of the number of segment with regards to rho.

We want to adjust the  $\rho$  parameter to have a reasonable number of segments. So, we look for the widest range of  $\rho$  for which the number of segments remains unchanged (and larger than 1). Indeed, our intuition is that we can be confident in the found breakpoints because the penalty has to increase significantly to remove some breakpoints. The optimal data-driven  $\rho$  is the left-most value of the widest range such that the number of segments is larger than 1. Otherwise we consider there is only one segment in the profile.

**Sample-specific versus common parameter** The *sample-specific choice of  $\rho$*  described above was empirically compared with a *common choice of  $\rho$*  depending on the signal-to-noise ratio within each group of given profiles. Seventy profiles from a real dataset [Renneville et al., 2013] were clustered into groups with homogeneous signal-to-noise ratio (SNR) by using a Gaussian mixture model. Three groups were provided by the BIC criterion.

Figure 12 displays results (chromosome 1) for each profile (patient). Following Section 5.3.2 the widest range of  $\rho$  values was plotted for each profile. Colors (light grey, grey, and black) indicate the SNR level in each group (respectively low, middle, and high).

Whereas the lowest SNR group only contains ranges of  $\rho$  with small values, other groups correspond to ranges of both small and large values of  $\rho$ . A common choice of  $\rho$  within each of these two groups would have led to erroneous segmentations. The same conclusion applied to other chromosomes and criteria such as variance.

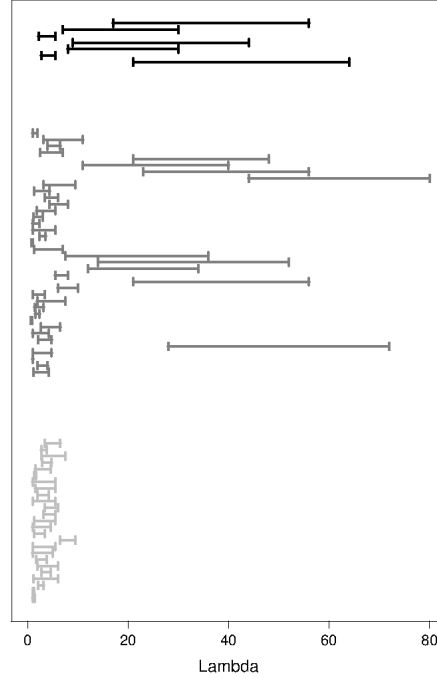


Figure 12: Widest ranges of  $\rho$  ( $x$ -axis) for 70 copy number profiles (chromosome 1) ( $y$ -axis). Colors indicate clusters of signal-to-noise ratios (light grey < grey < black).

## 6 Calling aberrations in copy-number profiles

When performing segmentations with `cghseg` or `PELT`, labels of segments (*loss*, *normal* or *gain*) are not known. It is possible to assign these labels with the `CGHcall` method [van de Wiel et al., 2007]. A wrapper of the `CGHcall` calling process is available in `MPAgenomics`.

Starting with a signal and an associated segmented profile, the method assumes that the segmented values follow a mixture of Gaussians. The parameters of the model are estimated with an EM algorithm with a specific initialization. Then, every segment is called by the most likely label according to the model.

### 6.1 From data normalized by `aroma`

The `cnSegCallingProcess` function executes the segmentation (see section 5) and the calling process. The specified data will be automatically imported by the function.

```
cnSegCallingProcess(dataSetName,normalTumorArray,chromosome,Rho,Kmax,listOfFile,onlySNP,savePlot,
nclass,cellularity,method,...)
```

Parameters	Details
<b>dataSetName</b>	name of the data-set folder in the rawData folder containing the signals to use.
<b>normalTumorArray</b>	Only in the case of normal-tumor study. A csv file or a data.frame containing the mapping between normal and tumor files(Fig. 7).
<b>chromosome</b>	A vector containing the chromosome to segment.
<b>Rho</b>	Only for method="PELT", a vector containing all the penalization values to test for the segmentation. If no values are provided, default values will be used.
<b>Kmax</b>	Only for method="cghseg", maximal number of segments.
<b>listOfFiles</b>	A vector containing the file names from the dataSetName to use.
<b>onlySNP</b>	If TRUE, only the SNP probes will be used.
<b>savePlot</b>	If TRUE, print some graphics (default=TRUE).
<b>nclass</b>	The number of levels to be used for calling. Either 3 (loss, normal, gain), 4 (including amplifications), 5 (including double deletions) (default=3).
<b>cellularity</b>	Percentage of tumor cells in the sample (default=1).
<b>...</b>	Other parameters for the CGHcall function [van de Wiel and Vosse, 2012].
<b>method</b>	"PELT" or "cghseg". (default="cghseg").

This function will save every segment in a text file in the **segmentation** folder in the working directory and return a data.frame with columns:

- **sampleNames** Name of the file.
- **chrom** The chromosome of the segment.
- **chromStart** The starting position (in bp) of the segment. This position is not included in the segment.
- **chromEnd** The ending position (in bp) of the segment. This position is included in the segment.
- **probes** Number of probes in the segment.
- **means** Mean of the segment.
- **calls** The calling of the segment ("double loss", "loss", "normal", "gain" or "amplification").

The graphic of the segmented profile is saved in the **figures** folder of the architecture.

```
> seg2=cnSegCallingProcess("datatest1",chromosome=21:22,method="PELT")
```

The output of the function is shown in the figure 13.

## 6.2 For any data provided in matrices

The function **callingProcess** executes all the calling process and can be used without the aroma architecture. You can use this function with the segmentation method provided in the package or any other methods as long as you give the results of the segmentation in the right format (see **segmentData** argument).

```
callingProcess(segmentData,nclass=5,cellularity=1,...)
```

Parameters	Description
<b>segmentData</b>	A list (see details below).
<b>nclass</b>	Number of labels. 3 corresponds to "loss", "normal" and "gain", 4 adds "amplification" and 5 adds "double loss".
<b>cellularity</b>	Percentage of tumor cells in the sample (default=1).
<b>verbose</b>	If TRUE, print some information.
<b>...</b>	Others parameters for CGHcall function [van de Wiel and Vosse, 2012].

**segmentData** is a list containing:

- **copynumber** A matrix. Every column represents a copy-number signal for a different sample.

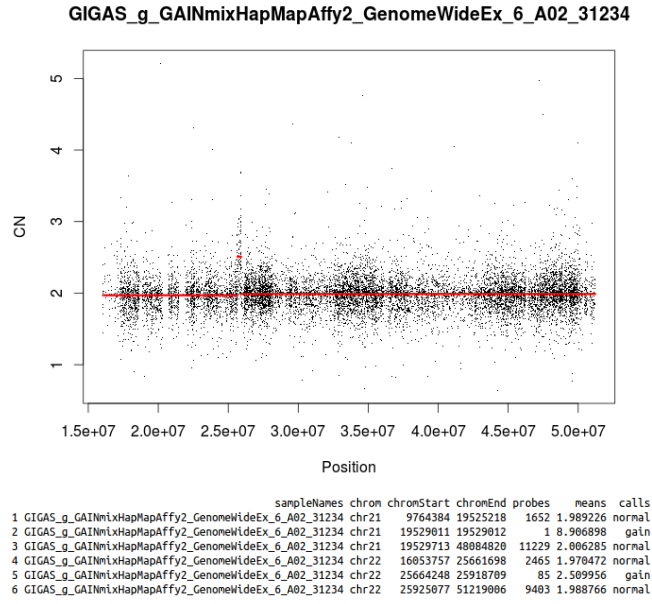


Figure 13: Top: graphics saved in the `figures` folder. Bottom: output of the function (saved in a text file too).

- **segmented** A matrix of the same size as `copynumber` containing the segmented values of the copy-number signal.
- **chromosome** A vector, of length the number of rows of `copynumber` matrix, containing the number of the chromosome for each position.
- **startPos** A vector, of length the number of rows of `copynumber` matrix, containing the genomic position of each probe.
- **featureNames** A vector, of length the number of rows of `copynumber` matrix, containing the name of each probe.
- **sampleNames** A vector, of length the number of columns of `copynumber` matrix, containing the name of each file.

	[,1]	[,2]		[,1]	[,2]
[1,]	2.164182	3.076227	[1,]	2.003522	3.012337
[2,]	1.947533	3.000218	[2,]	2.003522	3.012337
[3,]	2.040600	3.169354	[3,]	2.003522	3.012337
[4,]	2.012960	3.083510	[4,]	2.003522	3.012337
[5,]	1.954961	2.943393	[5,]	2.003522	3.012337
[6,]	2.103451	3.015429	[6,]	2.003522	3.012337
[7,]	1.928264	2.937508	[7,]	2.003522	3.012337
[8,]	2.155039	2.861401	[8,]	2.003522	3.012337
[9,]	1.828346	3.008872	[9,]	2.003522	3.012337
[10,]	1.894207	2.942340	[10,]	2.003522	3.012337
[11,]	2.071411	2.969460	[11,]	2.003522	3.012337
[12,]	2.062989	2.948135	[12,]	2.003522	3.012337
[13,]	2.004539	3.021060	[13,]	2.003522	3.012337
[14,]	1.839818	2.927193	[14,]	2.003522	3.012337
[15,]	2.002250	2.878493	[15,]	2.003522	3.012337
[16,]	2.075132	2.992165	[16,]	2.003522	3.012337
[17,]	2.057482	2.950658	[17,]	2.003522	3.012337
[18,]	1.909619	3.204027	[18,]	2.003522	3.012337
[19,]	2.032840	3.333027	[19,]	2.003522	3.012337
[20,]	1.984808	2.984271	[20,]	2.003522	3.012337

Figure 14: Left: example of `copynumber` matrix for 2 samples. Right: example of the associated `segmented` matrix.

This function returns a list with the same elements as `segmentData` and some supplementary elements:

- **calls** A matrix, of the same size as **copynumber** matrix, containing the label of each point.
- **segment** A data.frame that sums up all the segments found.
- **probdloss** (if you ran **callingProcess** with **nclass=5**) A matrix of the same size as **copynumber** matrix. It contains the probability for each segmented copy-number to be a double loss.
- **probdloss** A matrix of the same size as **copynumber** matrix. It contains the probability for each segmented copy-number to be a loss.
- **probdnorm** A matrix of the same size as **copynumber** matrix. It contains the probability for each segmented copy-number to be normal.
- **probdgain** A matrix of the same size as **copynumber** matrix. It contains the probability for each segmented copy-number to be a gain.
- **probdamp** (if you have run **callingProcess** with **nclass=4** or **5**) A matrix of the same size as **copynumber** matrix. It contains the probability for each segmented copy-number to be an amplification.

A wrapper is provided for creating the **segmentData** argument:

```
callingObject(copynumber, segmented, chromosome, position, featureNames, sampleNames)
```

The parameters **featureNames** and **sampleNames** can be omitted.

The following lines continue the example of the section 5. They create the **segmentData** list and run the calling process:

```
> #create the segmentData object
> callobj= callingObject(copynumber=seg$signal, segmented=seg$segmented,
+ chromosome=rep(20,length(seg$signal)), position=seg$startPos,
+ sampleNames="sample1")
> #run the calling
> call=callingProcess(callobj,nclass=3,cellularity=1,verbose=TRUE)
> call$segment
```

The result of the calling can be seen in figure 15.

	chrom	chromStart	chromEnd	probes	means	calls
1	chr20	61795	8099663	3573	2.051816	normal
2	chr20	8102867	8571315	222	2.634067	gain
3	chr20	8571934	62912463	19626	2.047356	normal

Figure 15: The summation of all the segments after running the function **callingProcess**.

A wrapper **CNAobjectToCGHcallObject** to convert CNA object from **DNAcopy** package [Seshan and Olshen, ] (CBS segmentation) to the desired format (see **segmentData**) is also provided.

```
CNAobjectToCGHcallObject(CNAobject)
```

Parameters	Description
CNAobject	Output of the <b>segment</b> function from <b>DNAcopy</b> package.

### 6.3 filterSeg function

At the end of the calling process, if you have a lot of segments, you might want to filter some uninteresting segments. The **filterSeg** function allows to filter segments from 3 criteria: the length in bp, the size in probes and the calling.

```
filterSeg(segmentList,minLength=1,minProbes=1,normalFilter=TRUE)
```

Arguments	Description
<b>segmentList</b>	A data.frame containing a description of segments, it must have at least columns: "chromStart", "chromEnd", "probes" and "calls" (see the output of <code>cnSegLabelProcess</code> and <code>callingProcess</code> functions).
<b>minLength</b>	The minimum length (in bp) for a segment. All the shorter segments are removed.
<b>minProbes</b>	The minimum number of probes for a segment. All the segments with less probes are removed.
<b>keptLabel</b>	Vector of label to keep. Only segment with one of the specified label will be kept.

It returns a data.frame of the same format as `segmentList` without the filtered segments.

We assume to have a data.frame called `call$segment` containing the columns plot in Fig. 16

```

      chrom chromStart chromEnd probes  means  calls
1 chr20      61795  8099663   3573 2.051816 normal
2 chr20     8102867  8571315    222 2.634067  gain
3 chr20     8571934  62912463   19626 2.047356 normal

```

Figure 16: `call$segment` before filtering.

The following command executes the code to keep only segments presenting a gain.

```

> segmentfilter=filterSeg(call$segment,keptLabel="gain")
> segmentfilter

```

Obtained results are shown in Fig. 17

```

      chrom chromStart chromEnd probes  means  calls
2 chr20     8102867  8571315    222 2.634067  gain

```

Figure 17: Result of the filtering.

## 7 Selection of genomic markers

In this section, the goal is to select some relevant markers according to a response.

It consists in minimizing  $g_\lambda : \beta \in \mathbb{R}^P \mapsto g(\beta)$ , where

$$g_\lambda(\beta) = \sum_{i=1}^I (y_i - (X\beta)_i)^2 + \lambda \sum_{p=1}^P |\beta_p| ,$$

with  $(X\beta)_i = \sum_p x_{i,p}\beta_p$  and  $\lambda > 0$  controlling the number of non-zero coordinates of  $\beta$ . After minimization, non-zero coefficients  $\beta_p$  correspond to influential positions to predict the response.

To solve this problem, we use by default the LARS algorithm [Efron et al., 2004] that solves the lasso problem for all the values of  $\lambda$ . In order to choose the best solution, we use a cross validation to select the best values of  $\lambda$  and return the associated solution.

The spike and slab method [Ishwaran and Rao, 2005][Ishwaran and Rao, 2010] is another alternative in this case of linear regression. It is a three steps procedure including filtering, marker effect estimation and marker selection steps. The filter and marker effect estimation steps use the spike and slab algorithm whereas the marker selection step relies on the elastic net (mixing a  $l_1$  penalty whit a  $l_2$  penalty). The spike and slab algorithm solves a generalized ridge regression problem ( $l_2$  penalty). It is implemented in the R package `spikeslab` [Ishwaran et al., 2013].

## 7.1 From data normalized by aroma

The function `markerSelection` will extract the data from the aroma architecture and run the LARS algorithm and the cross-validation for each chromosome separately.

```
markerSelection(dataSetName,dataResponse,chromosome,signal,normalTumorArray,onlySNP,nbFolds,
loss,plot,pkg,...)
```

Parameters	Description
<b>dataSetName</b>	The name of the data-set folder.
<b>dataResponse</b>	A csv files or a data.frame with 2 columns: "files" and "response". The column "files" contains the filename to extract and the second column the response associated with the file.
<b>chromosome</b>	A vector containing the number of the chromosomes for the SNPs selection.
<b>signal</b>	"CN" or "fracB", corresponding to which signal will be analyzed (default="CN").
<b>normalTumorArray</b>	Only in the case of normal-tumor study. A csv file or a data.frame containing the mapping between normal and tumor files(Fig. 7).
<b>onlySNP</b>	Only if <b>signal</b> ="CN". If TRUE, only the SNPs probes are used (default=FALSE).
<b>nbFolds</b>	Number of folds for the cross-validation (default=10).
<b>loss</b>	either "logistic" (binary response) or "linear" (quantitative response), default is "logistic".
<b>plot</b>	If TRUE, cross-validation mean squared error is plotted (default=TRUE).
<b>pkg</b>	Either "HDPenReg" or "spikeslab". Used package in linear case (default="HDPenReg").
<b>...</b>	supplementary arguments for <code>cv.glmnet</code> function [Friedman et al., 2009] in case of logistic loss or for <code>HDLars</code> function for linear loss.

For **signal**="fracB", the selection of markers is done for the tumor allele B fraction.

The function returns a list containing as many elements as the number of chromosomes specified in parameters. The name **chrX**, where *X* is the number of a chromosome. Each element of this list is a list containing:

- **chr** The chromosome corresponding to the signal.
- **markers.index** A vector containing the index of all selected markers.
- **markers.position** A vector containing the position of all selected markers.
- **markers.names** A vector containing the name of all selected markers.
- **coefficient** A vector containing the coefficients ( $\hat{\beta}$ ) of all selected markers.
- **intercept** Intercept of the model.

For the following example, we assume that the architecture is already created and the data-set is preprocessed as in the previous examples (see section 4.1 and 4.2). The response is stored in a data.frame named **dataResponse** (see Fig. 18).

```

                                files response
GIGAS_g_GAINmixHapMapAffy2_GenomeWideEx_6_A01_31218 2.105092
GIGAS_g_GAINmixHapMapAffy2_GenomeWideEx_6_A02_31234 1.442868
GIGAS_g_GAINmixHapMapAffy2_GenomeWideEx_6_A03_31250 1.952103
GIGAS_g_GAINmixHapMapAffy2_GenomeWideEx_6_A04_31266 1.857819
GIGAS_g_GAINmixHapMapAffy2_GenomeWideEx_6_A05_31282 2.047897
GIGAS_g_GAINmixHapMapAffy2_GenomeWideEx_6_A06_31298 1.654766
GIGAS_g_GAINmixHapMapAffy2_GenomeWideEx_6_A07_31314 2.385327
GIGAS_g_GAINmixHapMapAffy2_GenomeWideEx_6_A08_31330 2.113406
```

Figure 18: Example of **dataResponse** parameters for the **markerSelection** function.

The relevant markers can be selected by executing the following code:

```
> dataResponse=data.frame(files=getListOfFiles("datatest1"),
+ response=c(2.105092,1.442868,1.952103,1.857819,2.047897,1.654766,2.385327,2.113406))
> res=markerSelection(dataSetName="datatest1",dataResponse,chromosome=21:22,signal="CN",
+ onlySNP=TRUE,loss="linear")
```

Figure 19 shows the output of the function.

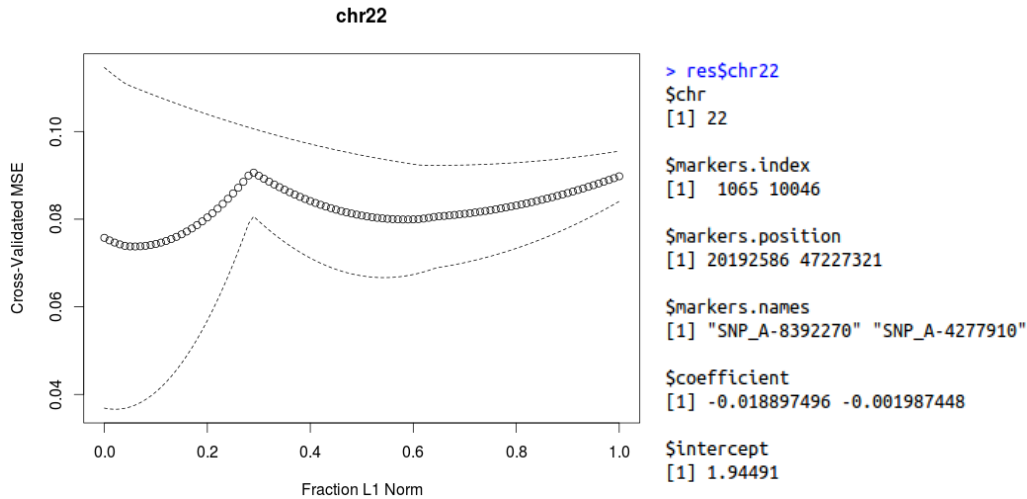


Figure 19: Left: cross-validation mean squared error. Right: results for chromosome 22.

In this example, two markers are selected: the 1065th and 10046th markers of the profile (markers.index) at positions 20192586 pb and 47227321 pb on the chromosome 22. The coefficient associated with each selected marker are given in the coefficient vector.

## 7.2 For any data provided in matrices

A version of the selection of markers without the aroma architecture is available:

```
variableSelection(dataMatrix,dataResponse,nbFolds,loss,plot,pkg)
```

Parameters	Description
<b>dataMatrix</b>	A matrix containing the data, each row is a different signal.
<b>dataResponse</b>	A vector containing the response associated with each signal.
<b>nbFolds</b>	Number of folds for the cross-validation (default=10).
<b>loss</b>	either “logistic” (binary response) or “linear” (quantitative response), default is “logistic”.
<b>plot</b>	If TRUE plot cross-validation mean squared error (default=TRUE).
<b>pkg</b>	Either “HDPenReg” or “spikeslab”. Used package in linear case (default=“HDPenReg”).
<b>...</b>	supplementary arguments for <code>cv.glmnet</code> function [Friedman et al., 2009] in case of logistic loss or for <code>HDlars</code> function for linear loss.

The output of the functions is:

- **variable** A vector containing the index of all selected markers.
- **coefficient** A vector containing the coefficients ( $\hat{\beta}$ ) of all selected variables.
- **intercept** Intercept of the model.

For example, simulate a data-set and the associated response

```
> dataMatrix=matrix(rnorm(5000,0,0.5),nrow=50)
> dataResponse=drop(dataMatrix%%sample(c(rep(0,90),rep(1,10))))
> res=variableSelection(dataMatrix,dataResponse,nbFolds=5,loss="linear",plot=TRUE)
```



### 7.3 To go further

In the current implementation of the **MPAgenomics** R package, two variable selection methods are proposed: (i) the Spike and Slab strategy consisting in first filtering variables to keep most reliable ones, second estimating the effect of remaining variables, and third selecting variables (using AIC), or (ii) Lasso + 10-fold cross validation and elbow search [Salvador and Chan, 2003] to choose the unknown parameter  $\lambda$  (default choice).

In order to compare the two approaches, we have carried out a small simulation study that is shortly described in what follows and mainly summarized by Figure 20. We run the two approaches with 100 samples generated according to the following design:

- $I = 100$  individuals and  $P = 1000, 2000, \dots, 10000$  variables
- $\beta^* \in \mathbb{R}^P$ , the true sparse estimates with only  $K = 10$  non-zero coefficients
- $X \in \mathcal{M}_{I,P}(\mathbb{R})$ ,  $X_1, \dots, X_I \sim \mathcal{N}(0_P, I_P)$
- $\epsilon \in \mathbb{R}^P$  Gaussian noise,  $\epsilon_1, \dots, \epsilon_I \sim \mathcal{N}(0, \sigma^2)$  with  $\sigma^2 = 0.5$
- $y = X\beta^* + \epsilon$

For each approach, the selected variables were compared to the truth (non-zero elements of  $\beta^*$ ), the average number of true and false selected variables were plotted in Figure 20. We can see that the *lars* procedure selects here more true variables – especially when  $P$  grows – and less false variables than *spikeslab*.

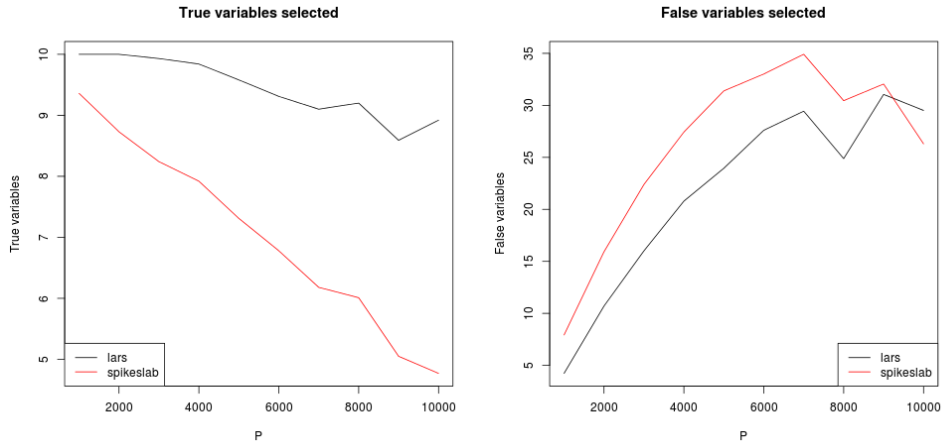


Figure 20: Number of true (Left figure) and false (Right figure) variables selected by two different methods: (i) Lars and 10-fold cross-validation (*lars*), (ii) Spike and Slab (*spikeslab*).

## References

- [Bengtsson, 2004] Bengtsson, H. (2004). aroma - An R Object-oriented Microarray Analysis environment. Preprint in Mathematical Sciences 2004:18, Mathematical Statistics, Centre for Mathematical Sciences, Lund University, Sweden.
- [Bengtsson et al., 2010a] Bengtsson, H., Bullard, J., Hansen, K., Neuviail, P., Purdomand, E., Robinson, M., and Simpson, K. (2010a). <http://www.aroma-project.org/>.
- [Bengtsson et al., 2008a] Bengtsson, H., Irizarry, R., Carvalho, B., and Speed, T. (2008a). Estimation and assessment of raw copy numbers at the single locus level. *Bioinformatics*, 24:759–767.
- [Bengtsson et al., 2010b] Bengtsson, H., Neuviail, P., and Speed, T. P. (2010b). Tumorboost: Normalization of allele-specific tumor copy numbers from a single pair of tumor-normal genotyping microarrays. *BMC Bioinformatics*, 11.
- [Bengtsson et al., 2008b] Bengtsson, H., Simpson, K., Bullard, J., and Hansen, K. (2008b). aroma.affymetrix: A generic framework in R for analyzing small to very large Affymetrix data sets in bounded memory. Technical Report 745, Department of Statistics, University of California, Berkeley.
- [Bengtsson et al., 2009] Bengtsson, H., Wirapati, P., and Speed, T. P. (2009). A single-array preprocessing method for estimating full-resolution raw copys from all Affymetrix genotyping arrays including GenomeWideSNP 5 & 6. *Bioinformatics*, 25(17):2149–2156.
- [Birgé and Massart, 2007] Birgé, L. and Massart, P. (2007). Minimal penalties for gaussian model selection. *Probability Theory and Related Fields*, 138(1-2):33–73.
- [Efron et al., 2004] Efron, B., Hastie, T., Johnstone, I., and Tibshirani, R. (2004). Least angle regression. *Annals of Statistics*, 32:407–499.
- [Friedman et al., 2009] Friedman, J., Hastie, T., and Tibshirani, R. (2009). Regularization paths for generalized linear models via coordinate descent.
- [Grimonprez et al., 2014] Grimonprez, Q., Celisse, A., Cheok, M., Figeac, M., and Marot, G. (2014). MPAge-nomics: An R package for multi-patient analysis of genomic markers. Preprint <http://hal.inria.fr/hal-00933614>.
- [Ishwaran and Rao, 2005] Ishwaran, H. and Rao, J. (2005). Spike and slab variable selection: frequentist and bayesian strategies. *Ann. Statist.*, 33(2):730–773.
- [Ishwaran and Rao, 2010] Ishwaran, H. and Rao, J. (2010). Generalized ridge regression: geometry and computational solutions when p is larger than n. *Manuscript*.
- [Ishwaran et al., 2013] Ishwaran, H., Rao, J., and Kogalur, U. (2013). *spikeslab : Prediction and variable selection using spike and slab regression*. R package version 1.1.5.
- [Killick and Eckley, 2013] Killick, R. and Eckley, I. (2013). *changeoint: An R package for changepoint analysis*. R package version 1.3.
- [Killick et al., 2012] Killick, R., Fearnhead, P., and Eckley, I. A. (2012). Optimal detection of changepoints with a linear computational cost. *Journal of the American Statistical Association*, 107(500):1590–1598.
- [Lebarbier, 2005] Lebarbier, E. (2005). Detecting multiple change-points in the mean of gaussian process by model selection. *Signal Processing*, 85(4):717 – 736.
- [Picard et al., 2013] Picard, F., Hoebeke, M., Lebarbier, E., Miele, V., Rigaiil, G., and Robin, S. (2013). cghseg: Segmentation methods for array CGH analysis. R package version 1.0.2.
- [Picard et al., 2005] Picard, F., Robin, S., Lavielle, M., Vaisse, C., and Daudin, J.-J. (2005). A statistical approach for array CGH data analysis. *BMC Bioinformatics*, 6(1):27.

- [Renneville et al., 2013] Renneville, A., Abdelali, R. B., Chevret, S., Nibourel, O., Cheok, M., Pautas, C., Dulery, R., Boyer, T., Cayuela, J.-M., Hayette, S., Raffoux, E., Farhat, H., Boissel, N., Terre, C., Dombret, H., Castaigne, S., and Preudhomme, C. (2013). Clinical impact of gene mutations and lesions detected by SNP-array karyotyping in acute myeloid leukemia patients in the context of gemtuzumab ozogamicin treatment: Results of the ALFA-0701 trial. *Oncotarget*, 4(9).
- [Rigaill, 2010] Rigaill, G. (2010). Pruned dynamic programming for optimal multiple change-point detection.
- [Salvador and Chan, 2003] Salvador, S. and Chan, P. (2003). Determining the number of clusters/segments in hierarchical clustering/segmentation algorithms. Technical report.
- [Seshan and Olshen, ] Seshan, V. E. and Olshen, A. *DNACopy: DNA copy number data analysis*. R package version 1.38.1.
- [Staaf et al., 2008] Staaf, J., Lindgren, D., Vallon-Christersson, J., Isaksson, A., Göransson, H., Juliusson, G., Rosenquist, R., Höglund, M., Borg, Åke., and Ringnér, M. (2008). Segmentation-based detection of allelic imbalance and loss-of-heterozygosity in cancer cells using whole genome SNP arrays. *Genome Biology*, 9(9).
- [Tibshirani, 1994] Tibshirani, R. (1994). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society, Series B*, 58:267–288.
- [van de Wiel and Vosse, 2012] van de Wiel, M. and Vosse, S. (2012). *CGHcall: Calling aberrations for array CGH tumor profiles*. R package version 2.26.0.
- [van de Wiel et al., 2007] van de Wiel, M. A., Kim, K. I., Vosse, S. J., van Wieringen, W. N., Wilting, S. M., and Ylstra, B. (2007). CGHcall: calling aberrations for array CGH tumor profiles. *Bioinformatics*, 23(7):892–894.