

# A quick introduction to AneuFinder

Aaron Taudt\*

December 8, 2017

## Contents

---

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Quickstart</b>	<b>2</b>
<b>3</b>	<b>A detailed workflow</b>	<b>3</b>
3.1	Mappability correction . . . . .	3
3.2	Blacklisting . . . . .	3
3.3	Running Aneufinder . . . . .	4
3.4	Loading results and plotting single cells . . . . .	5
3.5	Quality control . . . . .	6
3.6	Karyotype measures . . . . .	8
3.7	Principal component analysis . . . . .	9
<b>4</b>	<b>Session Info</b>	<b>10</b>

---

\*aaron.taudt@gmail.com

NOTE: This is the vignette for AneuFinder v1.6.0

## 1 Introduction

*AneuFinder* offers functionality for the study of copy number variations (CNV) in whole-genome single cell sequencing (WGSCS) data. Functionality implemented in this package includes:

- Copy number detection using a Hidden Markov Model or changepoint-algorithm on binned read counts.
- Various plotting capabilities like genomewide heatmaps of copy number state and arrayCGH-like plots.
- Export of copy number calls in BED format for upload to the UCSC genome browser.
- Quality metrics.
- Measures for addressing karyotype heterogeneity.

## 2 Quickstart

The main function of this package is called `Aneufinder()` and performs all the necessary steps to get from aligned reads to interpretable output. *AneuFinder* offers three methods to call copy number variations: (1) A Hidden Markov Model described in [1], (2) an approach based on the *DNACopy* package adopted for single cells from [2], and (3) an approach described in TODO using the edivisive-algorithm from the *ecp* package.<sup>1</sup> We recommend using the "edivisive" method.

```
library(AneuFinder)
Aneufinder(inputfolder='folder-with-BAM-or-BED', outputfolder='output-directory',
           numCPU=2, method=c('edivisive', 'dnacopy', 'HMM'))
```

Although in most cases the above command will produce reasonably good results, it might be worthwhile to adjust the default parameters to improve performance and the quality of the results (see section 3). You can get a description of all available parameters by typing

```
?Aneufinder
```

After the function has finished, you will find the folder **output-directory** containing all produced files and plots. This folder contains the following *files* and **folders**:

- *AneuFinder.config*: This file contains all the parameters that are necessary to reproduce your analysis. You can specify this file as

```
Aneufinder(..., configfile='AneuFinder.config')
```

to run another analysis with the same parameter settings.

- *chrominfo.tsv*: A tab-separated file with chromosome length information.
- **binned**: This folder contains the binned data. If you chose a correction method, you will also see a folder like **binned-GC** in case of GC correction. You can load the data with

```
files <- list.files('output-directory/binned', full.names=TRUE)
binned.data <- loadFromFiles(files)
```

- **BROWSERFILES**: A folder which contains BED files with copy number calls and breakpoint locations that can be uploaded to the UCSC genome browser. If `reads.store=TRUE` it will also contain a subfolder **data** with the mapped reads in BED format.
- **MODELS**: A folder with all produced Hidden Markov Models. You can load the results for further processing, such as quality control and customized plotting. The folder might be named **MODELS\_refined** or **MODELS-StrandSeq**, depending on whether you chose to analyze Strand-seq data or refine breakpoints.

<sup>1</sup>Please cite [1] if you use the *HMM* results and please cite [1] and [2] if you use the *dnacopy* results.

```
files <- list.files('output-directory/MODELS/method-edivisive', full.names=TRUE)
hmms <- loadFromFiles(files)
cl <- clusterByQuality(hmms)
heatmapGenomewide(cl$classification[[1]])
```

- **PLOTS:** All plots that are produced by default will be stored here.
- **data:** Only produced if reads.store=TRUE. This folder stores all the read data as RData objects. This exists mostly for internal usage, namely to estimate confidence intervals and refine breakpoints.

## 3 A detailed workflow

### 3.1 Mappability correction

The first step of your workflow should be the production of a reference file for mappability correction. Mappability correction is done via a variable-width binning approach (as compared to fixed-width bins) and requires a euploid reference. You can either simulate this reference file or take a real euploid reference. For optimal results we suggest to use a real reference, e.g. by merging BAM files of single cells from a euploid reference tissue. This can be achieved with the 'samtools merge' command (not part of R). Be careful: All CNVs that are present in the reference will lead to artifacts in the analysis later. This includes sex-chromosomes that are present in one copy only, so we advice to use a female reference and to exclude the Y-chromosome from the analysis. If you have no reference available, you can simulate one with the simulateReads() command. You can also skip this step and continue without mappability correction. The algorithm will use fixed-width bins in this case.

```
library(AneuFinder)

## Load human genome
library(BSgenome.Hsapiens.UCSC.hg19)

## Get a BAM file for the estimation of quality scores (adjust this to your experiment)
bamfile <- system.file("extdata", "BB150803_IV_074.bam", package="AneuFinderData")

## Simulate reads of length 51bp for human genome
# We simulate reads every 500000bp for demonstration purposes, for a real
# application you should use a much denser spacing (e.g. 500bp or less)
simulatedReads.file <- tempfile() # replace this with your destination file
simulateReads(BSgenome.Hsapiens.UCSC.hg19, readLength=51, bamfile=bamfile,
              file=simulatedReads.file, every.X.bp=500000)
```

This simulated FASTQ file must then be aligned with your aligner of choice (ideally the same that you used for your other samples) and given as reference in Aneufinder(..., variable.width.reference="your-reference-bamfile").

### 3.2 Blacklisting

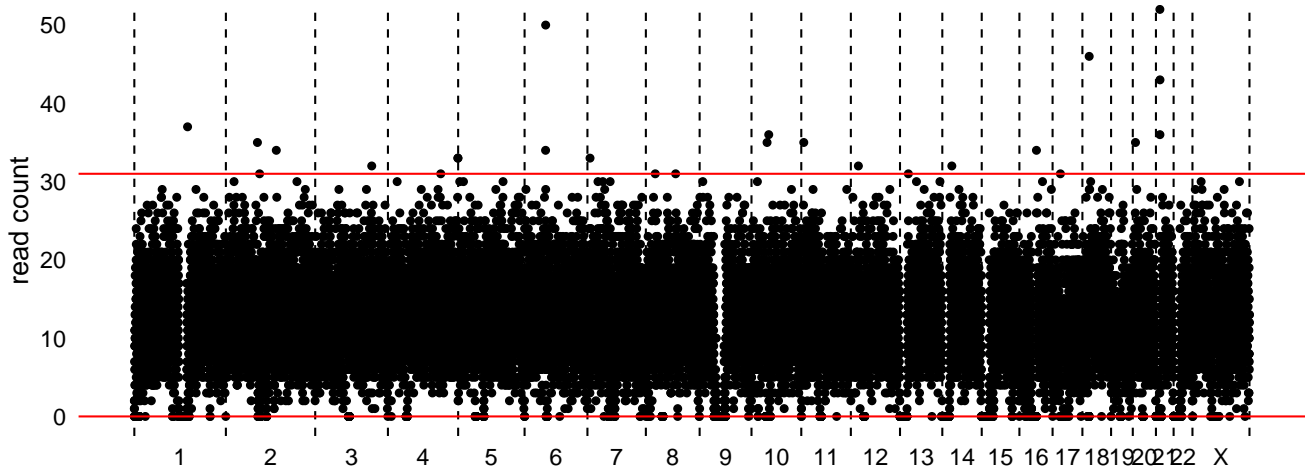
To further improve the quality of the results and remove artifacts caused by high mappability repeat regions, e.g. near centromers, a blacklist can be used with option Aneufinder(..., blacklist). All reads falling into the regions specified by the blacklist will be discarded when importing the read files. You can either download a blacklist from the UCSC genome browser, e.g. the "DAC Blacklisted Regions from ENCODE/DAC(Kundaje)" mappability track, or make your own. For optimal results, we advice to make your own blacklist from a euploid reference. The following code chunk takes a euploid reference and makes fixed-width bins of 100kb. Bins with read count above and below the 0.999 and 0.05 quantile are taken as blacklist:

```
library(AneuFinder)
```

```
## Get a euploid reference (adjust this to your experiment)
bedfile <- system.file("extdata", "hg19_diploid.bam.bed.gz", package="AneuFinderData")

## Make 100kb fixed-width bins
bins <- binReads(bedfile, assembly='hg19', binsizes=100e3,
                 chromosomes=c(1:22,'X'))[[1]]
## Make a plot for visual inspection and get the blacklist
lcutoff <- quantile(bins$counts, 0.05)
ucutoff <- quantile(bins$counts, 0.999)
p <- plot(bins) + coord_cartesian(ylim=c(0,50))
p <- p + geom_hline(aes(yintercept=lcutoff), color='red')
p <- p + geom_hline(aes(yintercept=ucutoff), color='red')
print(p)
```

reads = 0.36M, complexity = NAM, spikiness = 0.4, entropy = 10.17, bhattacharyya = NA, num.segments = 0, loglik = NA, sos = NA



```
## Select regions that are above or below the cutoff as blacklist
blacklist <- bins[bins$counts <= lcutoff | bins$counts >= ucutoff]
blacklist <- reduce(blacklist)
## Write blacklist to file
blacklist.file <- tempfile()
exportGRanges(blacklist, filename=blacklist.file, header=FALSE,
              chromosome.format='NCBI')
```

### 3.3 Running Aneufinder

The function `Aneufinder()` takes an input folder with BAM or BED files and produces an output folder with results, plots and browserfiles. The following code is an example of how to run `Aneufinder()` with variable-width bins (see section 3.1), blacklist (see section 3.2) and GC-correction. Results will be stored in **outputfolder/MODELS** as RData objects for further processing such as quality filtering and customized plotting.

```
library(AneuFinder)

## First, get some data and reference files (adjust this to your experiment)
var.width.ref <- system.file("extdata", "hg19_diploid.bam.bed.gz", package="AneuFinderData")
blacklist <- system.file("extdata", "blacklist-hg19.bed.gz", package="AneuFinderData")
datafolder <- system.file("extdata", "B-ALL-B", package = "AneuFinderData")
```

```
list.files(datafolder) # only 3 cells for demonstration purposes
## [1] "MB140210_I_003.bam.bed.gz" "MB140210_I_004.bam.bed.gz" "MB140210_I_005.bam.bed.gz"
## Library for GC correction
library(BSgenome.Hsapiens.UCSC.hg19)

## Produce output files
outputfolder <- tempdir()
Aneufinder(inputfolder = datafolder, outputfolder = outputfolder, assembly = 'hg19',
            numCPU = 3, binsizes = c(5e5, 1e6), variable.width.reference = var.width.ref,
            chromosomes = c(1:22, 'X', 'Y'), blacklist = blacklist, correction.method = 'GC',
            GC.BSgenome = BSgenome.Hsapiens.UCSC.hg19, refine.breakpoints=FALSE,
            method = 'edivisive')

## Warning in FUN(genome = names(SUPPORTED_UCSC_GENOMES)[idx], circ_seqs = supported_genome$circ_seqs,
: NCBI seqlevel was set to NA for hg19 UCSC seqlevel(s) not in the NCBI assembly:
## chrM
```

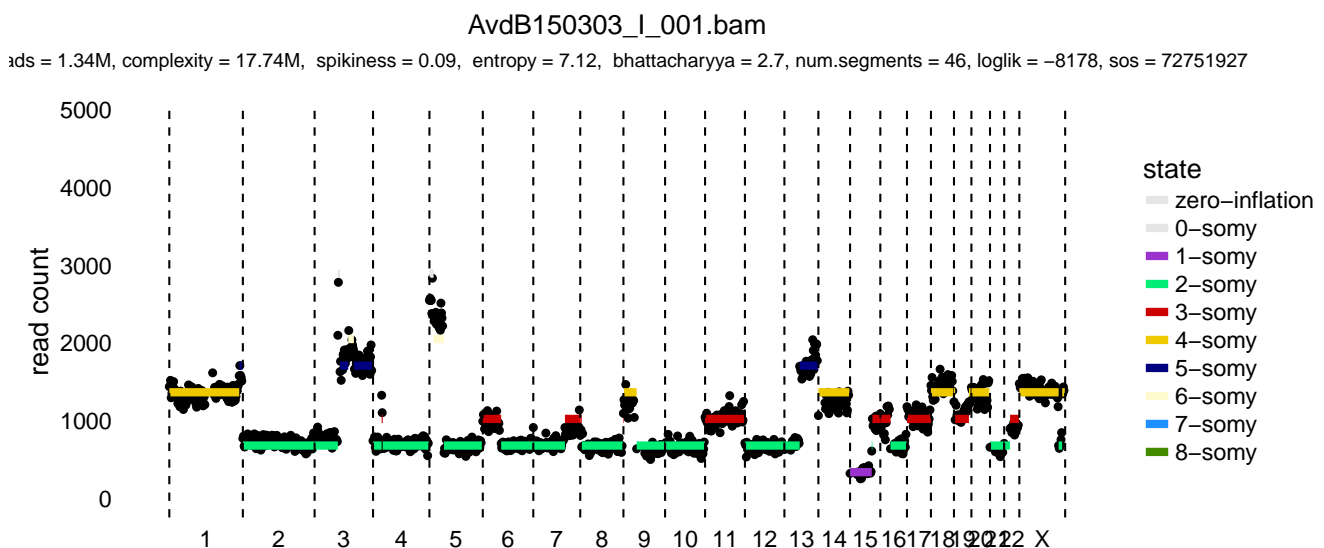
### 3.4 Loading results and plotting single cells

Once the function `Aneufinder()` has completed, results will be accessible as `.RData` files under **outputfolder/MODELS**. You can load the results into R using

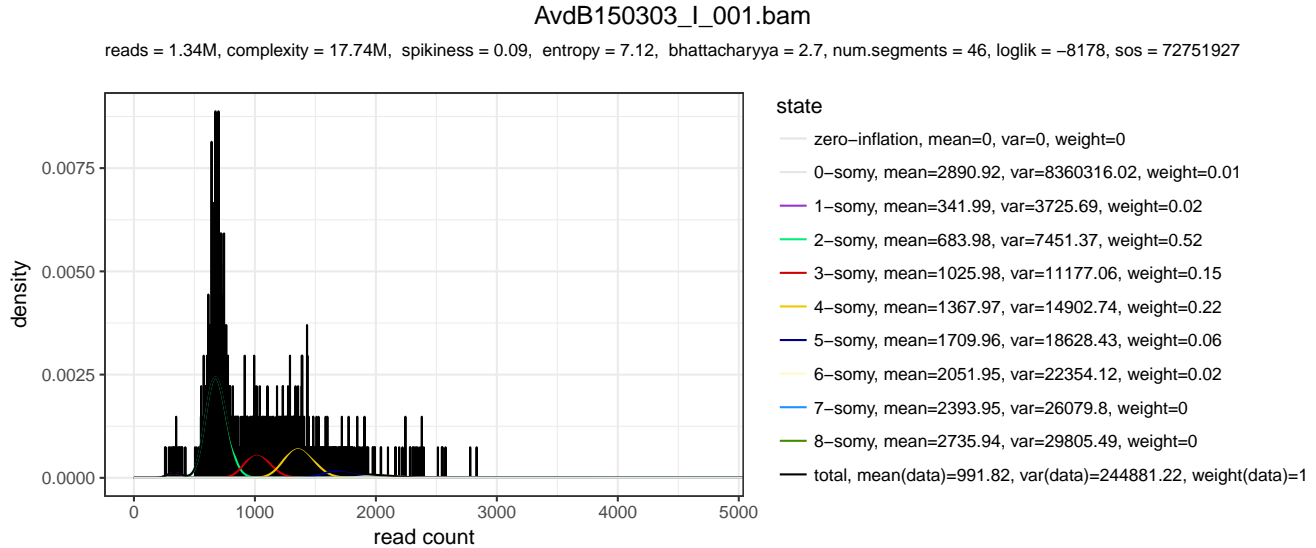
```
library(AneuFinder)
files <- list.files('outputfolder/MODELS/method-edivisive', full.names=TRUE)
models <- loadFromFiles(files)
```

Here are some examples of the plotting functions that are available. Most of these plots are also produced by default by `Aneufinder()` and are available as PDF in **outputfolder/PLOTS**.

```
## Get some pre-produced results (adjust this to your experiment)
results <- system.file("extdata", "primary-lung", "hmms", package="AneuFinderData")
files <- list.files(results, full.names=TRUE)
plot(files[1], type='profile')
```



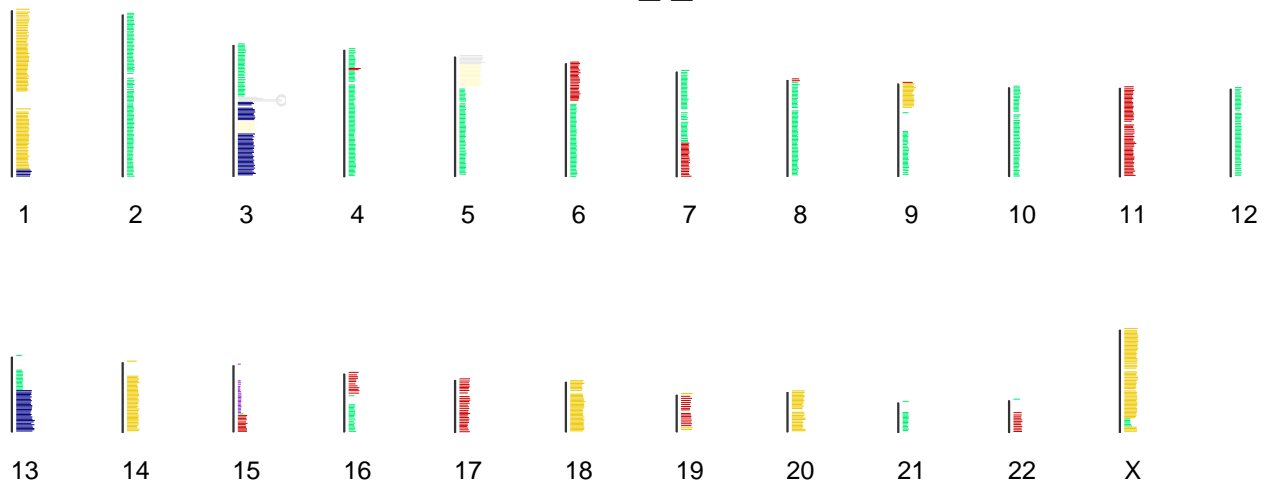
```
plot(files[1], type='histogram')
```



```
plot(files[1], type='karyogram')
```

```
## Warning in plot.karyogram(model, both.strands = both.strands, file = file): Cannot breakpoint coordinates. Please run 'getBreakpoints' first.
```

AvdB150303\_I\_001.bam



1, complexity = 17.74M, spikiness = 0.09, entropy = 7.12, bhattacharyya = 2.7, num.segments = 46, loglik = -8178, so

### 3.5 Quality control

Once the function `Aneufinder()` has completed, results will be accessible as `.RData` files under **outputfolder/MODELS**. Single cell sequencing is prone to noise and therefore it is a good idea to filter the results by quality to retain only high-quality cells. We found that simple filtering procedures such as cutoffs on the total number of reads etc., are insufficient to distinguish good from bad-quality libraries. Therefore, we have implemented a multivariate clustering approach that works on multiple quality metrics (see `?clusterByQuality` for details) to increase robustness of the filtering. Here is an example demonstrating the usage of `clusterByQuality()`.

The figure displays a 5x5 grid of scatter plots, where each row and column corresponds to one of the five variables: spikiness, num.segments, entropy, bhattacharyya, and SOS. The diagonal elements of the grid are empty, showing only the variable name. The off-diagonal plots show the relationship between pairs of variables. Each plot contains data points from various methods, represented by different colors and shapes (e.g., red circles, blue squares, green crosses, orange diamonds, purple pluses, yellow asterisks, cyan crosses). A dashed ellipse in each plot represents the 95% confidence interval for the data distribution. The axes are labeled with the variable names and their corresponding numerical ranges.

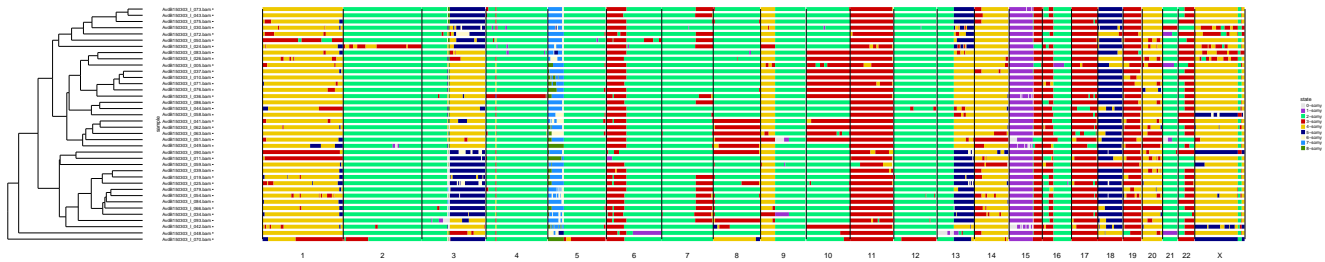
```
print(cl$parameters)
```

##	spikiness	num.segments	entropy	bhattacharyya	sos
## [1,]	0.1060501	64.30637	7.118107	4.52769354	9671255.144
## [2,]	0.1119017	58.86160	7.126918	2.41535861	4722433.201
## [3,]	0.1162552	85.34653	7.124391	3.69585953	3842048.636
## [4,]	0.1182611	72.81237	7.105328	2.64484856	8766630.278
## [5,]	0.1238624	113.06280	7.078982	2.78844568	77991481.726
## [6,]	0.1261168	130.20767	7.095900	3.28808888	14095817.324
## [7,]	0.1611050	140.16982	7.055041	1.63693596	6301025.747

```
## [8,] 1.1625668 458.66667 6.411838 0.03769717 1156.659
```

```
## Apparently, the last cluster corresponds to failed libraries
## while the first cluster contains high-quality libraries
```

```
## Select the two best clusters and plot it
selected.files <- unlist(cl$classification[1:2])
heatmapGenomewide(selected.files)
```



### 3.6 Karyotype measures

This package implements two measures to quantify karyotype heterogeneity, an *aneuploidy* and a *heterogeneity* score. Both measures are independent of the number of cells, the length of the genome, and take into account every position in the genome. The following example compares the heterogeneity and aneuploidy between a primary lung cancer and the corresponding liver metastasis.

```
library(AneuFinder)

## Get some pre-produced results (adjust this to your experiment)
results <- system.file("extdata", "primary-lung", "hmms", package="AneuFinderData")
files.lung <- list.files(results, full.names=TRUE)
results <- system.file("extdata", "metastasis-liver", "hmms", package="AneuFinderData")
files.liver <- list.files(results, full.names=TRUE)

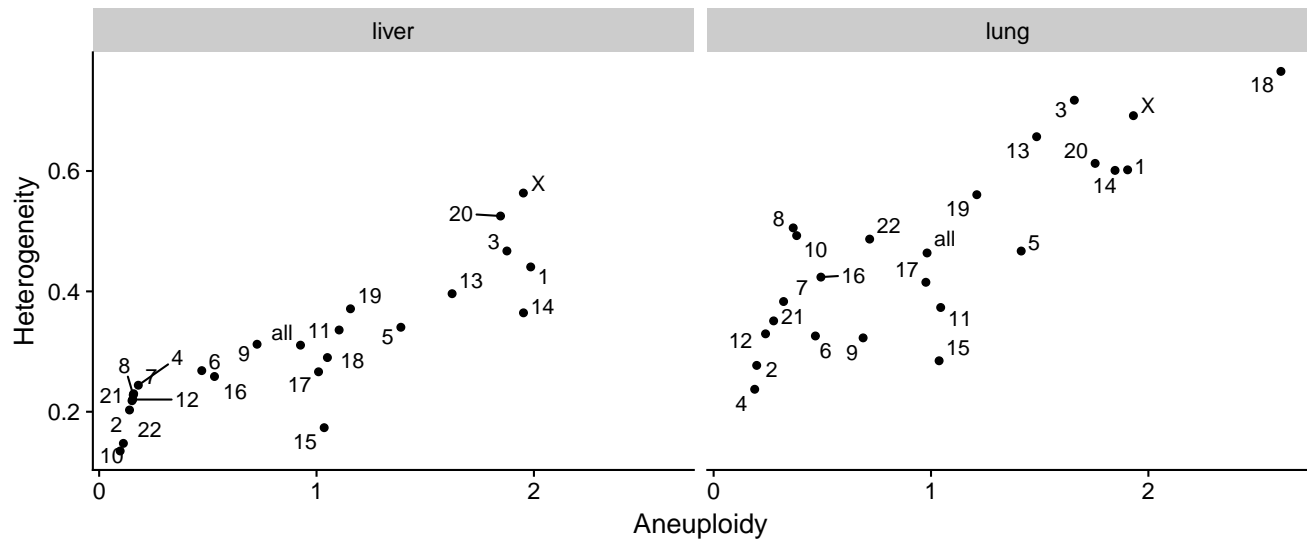
## Get karyotype measures
k.lung <- karyotypeMeasures(files.lung)
k.liver <- karyotypeMeasures(files.liver)

## Print the scores in one data.frame
df <- rbind(lung = k.lung$genomewide, liver = k.liver$genomewide)
print(df)

##      Aneuploidy Heterogeneity
## lung  0.9818370   0.4638431
## liver 0.9262749   0.3106615

## While the aneuploidy is similar between both cancers, the heterogeneity is
## nearly twice as high for the primary lung cancer.
plotHeterogeneity(hmms.list = list(lung=files.lung, liver=files.liver))
```

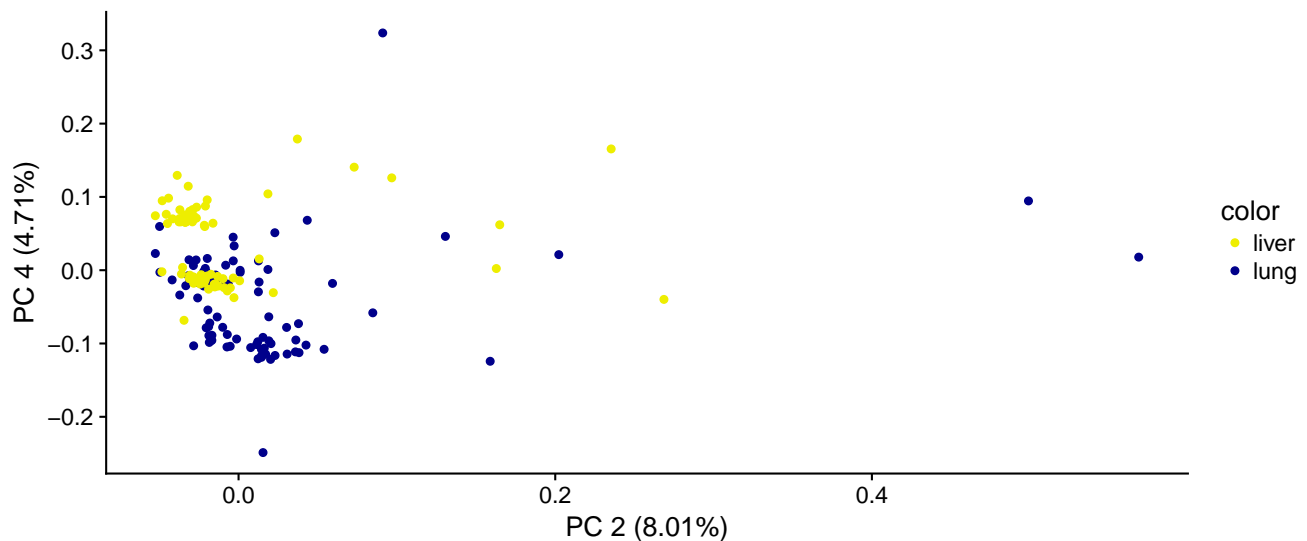




### 3.7 Principal component analysis

The following code let's you plot a PCA for a selection of single cells.

```
## Get results from a small-cell-lung-cancer
lung.folder <- system.file("extdata", "primary-lung", "hmms", package="AneuFinderData")
lung.files <- list.files(lung.folder, full.names=TRUE)
## Get results from the liver metastasis of the same patient
liver.folder <- system.file("extdata", "metastasis-liver", "hmms", package="AneuFinderData")
liver.files <- list.files(liver.folder, full.names=TRUE)
## Plot a clustered heatmap
classes <- c(rep('lung', length(lung.files)), rep('liver', length(liver.files)))
labels <- c(paste('lung', 1:length(lung.files)), paste('liver', 1:length(liver.files)))
plot_pca(c(lung.files, liver.files), colorBy=classes, PC1=2, PC2=4)
```



## 4 Session Info

---

```
toLatex(sessionInfo())
```

- R version 3.4.1 (2017-06-30), x86\_64-pc-linux-gnu
- Locale: LC\_CTYPE=en\_US.UTF-8, LC\_NUMERIC=C, LC\_TIME=de\_DE.UTF-8, LC\_COLLATE=C, LC\_MONETARY=de\_DE.UTF-8, LC\_MESSAGES=en\_US.UTF-8, LC\_PAPER=de\_DE.UTF-8, LC\_NAME=C, LC\_ADDRESS=C, LC\_TELEPHONE=C, LC\_MEASUREMENT=de\_DE.UTF-8, LC\_IDENTIFICATION=C
- Running under: Ubuntu 17.10
- Matrix products: default
- BLAS: /usr/local/lib/R/lib/libRblas.so
- LAPACK: /usr/local/lib/R/lib/libRlapack.so
- Base packages: base, datasets, grDevices, graphics, methods, parallel, stats, stats4, utils
- Other packages: AneuFinder 1.7.2, AneuFinderData 1.4.0, BSgenome 1.44.2, BSgenome.Hsapiens.UCSC.hg19 1.4.0, BiocGenerics 0.22.1, Biostrings 2.44.2, GenomInfoDb 1.12.3, GenomicRanges 1.28.6, IRanges 2.10.5, S4Vectors 0.14.7, XVector 0.16.0, cowplot 0.9.1, ggplot2 2.2.1, rtracklayer 1.36.5
- Loaded via a namespace (and not attached): Biobase 2.36.2, BiocParallel 1.10.1, BiocStyle 2.4.1, DNACopy 1.50.1, DelayedArray 0.2.7, GenomInfoDbData 0.99.0, GenomicAlignments 1.12.2, KernSmooth 2.23-15, MASS 7.3-47, Matrix 1.2-10, RCurl 1.95-4.8, Rcpp 0.12.13, ReorderCluster 1.0, Rsamtools 1.28.0, SummarizedExperiment 1.6.5, XML 3.98-1.9, backports 1.1.1, bamsignals 1.8.0, bitops 1.0-6, caTools 1.17.1, codetools 0.2-15, colorspace 1.3-2, compiler 3.4.1, digest 0.6.12, doParallel 1.0.11, ecp 3.0.0, evaluate 0.10.1, foreach 1.4.3, gdata 2.18.0, gg dendro 0.1-20, ggrepel 0.7.0, gplots 3.0.1, grid 3.4.1, gtable 0.2.0, gtools 3.5.0, highr 0.6, htmltools 0.3.6, iterators 1.0.8, knitr 1.17, labeling 0.3, lattice 0.20-35, lazyeval 0.2.0, magrittr 1.5, matrixStats 0.52.2, mclust 5.4, munsell 0.4.3, plyr 1.8.4, reshape2 1.4.2, rlang 0.1.4, rmarkdown 1.6, rprojroot 1.2, scales 0.5.0, stringi 1.1.5, stringr 1.2.0, tibble 1.3.4, tools 3.4.1, yaml 2.1.14, zlibbioc 1.22.0

## References

---

- [1] Bjorn Bakker, Aaron Taudt, Mirjam E. Belderbos, David Porubsky, Diana C. J. Spierings, Tristan V. de Jong, Nancy Halsema, Hinke G. Kazemier, Karina Hoekstra-Wakker, Allan Bradley, Eveline S. J. M. de Bont, Anke van den Berg, Victor Guryev, Peter M. Lansdorp, Maria Colomé-Tatché, and Floris Foijer. Single-cell sequencing reveals karyotype heterogeneity in murine and human malignancies. *Genome Biology*, 17(1):115, may 2016. [doi:10.1186/s13059-016-0971-7](https://doi.org/10.1186/s13059-016-0971-7).
- [2] Tyler Garvin, Robert Aboukhalil, Jude Kendall, Timour Baslan, Gurinder S Atwal, James Hicks, Michael Wigler, and Michael C Schatz. Interactive analysis and assessment of single-cell copy-number variations. *Nature Methods*, 12(11):1058–1060, sep 2015. [doi:10.1038/nmeth.3578](https://doi.org/10.1038/nmeth.3578).