

# CNVPanelizer: Reliable CNV detection in target sequencing applications

Oliveira, Cristiano  
cristiano.oliveira@med.uni-heidelberg.de

Wolf, Thomas  
thomas\_Wolf71@gmx.de

October 13, 2015

Amplicon based targeted sequencing, over the last few years, has become a mainstay in the clinical use of next generation sequencing technologies. For the detection of somatic and germline SNPs this has been proven to be a highly robust methodology. One area of genomic analysis which is usually not covered by targeted sequencing, is the detection of copy number variations (CNVs). While a large number of available algorithms and software address the problem of CNV detection in whole genome or whole exome sequencing, there are no such established tools for amplicon based targeted sequencing. We introduced a novel algorithm for the reliable detection of CNVs from targeted sequencing.

## 1 Introduction

To assess if a region specific change in read counts correlates with the presence of CNVs, we implemented an algorithm that uses a subsampling strategy similar to Random Forest to reliably predict the presence of CNVs. We also introduce a novel method to correct for the background noise introduced by sequencing genes with a low number of amplicons. To make it available to the community we implemented the algorithm as an R package.

## 2 Using

This section provides an overview of the package functions.

### 2.1 Installing and Loading the package

The package is available through the Bioconductor repository and can be installed and loaded using the following R commands:

```
# To install from Bioconductor  
source("http://bioconductor.org/biocLite.R")  
biocLite("CNVPanelizer")
```

```
# To load the package  
library(CNVPanelizer)
```

## 2.2 Reading data

### 2.2.1 BED file

The BED file is required to obtain amplicon and gene name information associated with the panel.

```
# Bed file defining the amplicons
bedFilepath <- "/somePath/someFile.bed"

# The column number of the gene Names in the BED file.
ampColumnNumber <- 4

# Extract the information from a bed file
genomicRangesFromBed <- BedToGenomicRanges(bedFilepath,
                                             ampliconColumn = ampColumnNumber,
                                             split = "_")

metadataFromGenomicRanges <- elementMetadata(genomicRangesFromBed)
geneNames = metadataFromGenomicRanges["geneNames"][, 1]
ampliconNames = metadataFromGenomicRanges["ampliconNames"][, 1]
```

### 2.2.2 Selecting files

Two sets of data are required. The samples of interest and the set of reference bam files to compare against.

```
# Directory with the test data
sampleDirectory <- "/somePathToTestData"

# Directory with the reference data
referenceDirectory <- "/somePathToReferenceData"

# Vector with test filenames
sampleFilenames <- list.files(path = sampleDirectory,
                              pattern = ".bam$",
                              full.names = TRUE)

# Vector with reference filenames
referenceFilenames <- list.files(path = referenceDirectory,
                                pattern = ".bam$",
                                full.names = TRUE)
```

### 2.2.3 Counting reads

The reads were counted using a wrapper function around the **ExomeCopy** package from the Bioconductor project. All reads overlapping with the region of an amplicon were counted for this amplicon. Only reads with a mapping quality  $\geq 20$  were counted and the function allows to remove PCR Duplicates. For reads with the same start site, end site and chromosomal orientation only one is kept. PCR duplicates might cause a bias for the ratio between reference and the samples of interest. Thus this serves as an additional quality control step in the CNV detection pipeline. It is only recommended for Ion Torrent generated data. For Illumina data this step is not recommended.

```

# Should duplicated reads (same start, end site and strand) be removed
removePcrDuplicates <- FALSE # TRUE is only recommended for Ion Torrent data

# Read the Reference data set
referenceReadCounts <- ReadCountsFromBam(referenceFileNames,
                                         genomicRangesFromBed,
                                         sampleNames = referenceFileNames,
                                         ampliconNames = ampliconNames,
                                         removeDup = removePcrDuplicates)

# Read the sample of interest data set
sampleReadCounts <- ReadCountsFromBam(sampleFileNames,
                                       genomicRangesFromBed,
                                       sampleNames = sampleFileNames,
                                       ampliconNames = ampliconNames,
                                       removeDup = removePcrDuplicates)

```

## 2.2.4 Using Synthetic Data

We also make available synthetic data to test the functions. The following examples make use of two generated data sets, one for the reference and the other as a testing set

```

data(sampleReadCounts)
data(referenceReadCounts)

# Gene names should be same size as row columns
# For real data this is a vector of the genes associated
# with each row/amplicon. For example Gene1, Gene1, Gene2, Gene2, Gene3, ...
geneNames <- row.names(referenceReadCounts)

# Not defined for synthetic data
# For real data this gives a unique name to each amplicon.
# For example Gene1:Amplicon1, Gene1:Amplicon2, Gene2:Amplicon1,
# Gene2:Amplicon2, Gene3:Amplicon1 ...
ampliconNames <- NULL

```

## 2.3 Normalization

To account for sample and sequencing run specific variations the counts obtained for each sample were normalized using a wrapper function around the **tmm** normalization function from the **NOISeq** package.

```

normalizedReadCounts <- CombinedNormalizedCounts(sampleReadCounts,
                                                  referenceReadCounts,
                                                  ampliconNames = ampliconNames)

# After normalization data sets need to be splitted again to perform bootstrap
samplesNormalizedReadCounts = normalizedReadCounts["samples"][[1]]
referenceNormalizedReadCounts = normalizedReadCounts["reference"][[1]]

```

## 2.4 Bootstrap based CNV

This approach is similar to the Random Forest method, which bootstraps the samples and subsamples the features. In our case features would be equivalent to amplicons. The subsampling procedure is repeated  $n$  times to generate a large set of randomized synthetic references  $B = b_1, \dots, b_n$  by selecting

with replacement (bootstrapping) from the set of reference samples. The ratio between the sample of interest and each randomized reference is calculated for each gene, using only a randomly selected subset of amplicons. To get an estimate of significance, the 95 percent confidence interval of the bootstrapping/subsampling ratio distribution was calculated. A significant change is considered as an amplification for a lower bound ratio  $> 1$  and a deletion for an upper bound ratio  $< 1$ . The confidence interval was calculated using the 0.025 and 0.975 percent quantiles of the bootstrapping/subsampling ratio distribution.

```
# Number of bootstrap replicates to be used
replicates <- 10000
```

```
# Perform the bootstrap based analysis
bootList <- BootList(geneNames,
                     samplesNormalizedReadCounts,
                     referenceNormalizedReadCounts,
                     replicates = replicates)
```

## 2.5 Background Estimation

Not all genes have the same number of amplicons  $|A_G|$  and it has been shown that sequencing genes with a higher number of amplicons yields better sensitivity and specificity when detecting putative copy number variations. Still genes sequenced with a low number of amplicons might still show significant changes in observed read counts. While normalization makes the comparison of read counts comparable between samples, genes with a small number of amplicons might still show a bias. To quantify the effect of a low number of amplicons on the calling of CNVs we introduced a background noise estimation procedure. Using the ratio between the mean reference and the sample used for calling we subsample for each unique number of amplicons. In the case of two amplicons we repeatedly sample two random amplicons from the set of all amplicons, and average the ratios. Amplicons that belong to genes showing significant copy number variations  $G_{sig}$  are not included in the subsampling pool. Each amplicon is weighted according to the number of amplicons the respective gene has  $w_A = \frac{1}{|A_g|}$ . Thus the probability of sampling from a gene is the same regardless of the number amplicons. For each number of amplicons a background noise distribution is estimated. The reported background is defined by the lower noise  $meanNoise + qnorm(0.025) * standardDeviationNoise$  and the upper noise  $meanNoise + qnorm(0.975) * standardDeviationNoise$  of the respective distribution. This defines the 95 percent confidence interval. The significance level can also be passed as parameter. Less conservative detection can be achieved by setting the `significanceLevel = 0.1`. To calculate the mean and standard deviation (sd) of ratios the log ratios were used. If setting `robust = TRUE`, median is used instead of mean and mad (median absolute deviation) replaces sd. The upper bound ratio has to be below the lower noise (deletion) or the lower bound ratio above the upper noise (amplification) to be considered reliable.

```
# Estimate the background noise left after normalization
backgroundNoise <- Background(geneNames,
                              samplesNormalizedReadCounts,
                              referenceNormalizedReadCounts,
                              bootList,
                              replicates = replicates,
                              significanceLevel = 0.1,
                              robust = TRUE)
```

## 2.6 Results

To analyse the results we provide two outputs. A plot which shows the detected variations, and a report table with more detailed information about those variations.

### 2.6.1 Report

The report describes the bootstrap distribution for each gene. An example can be found in figure 1. The final report is genewise, and is based on the bootstrapping.

```
# Build report tables
reportTables <- ReportTables(geneNames,
                             samplesNormalizedReadCounts,
                             referenceNormalizedReadCounts,
                             bootList,
                             backgroundNoise)
```

At the figure 1 we can see an example of the report table for a single sample.

##	MeanRatio	LowerBoundBoot	MeanBoot	UpperBoundBoot	LowerNoise	MeanNoise	UpperNoise	Signif.	AboveNoise	Amplicons	Passed
## Gene_01	0.69	0.19	0.69	2.54	0.55	0.91	1.51	FALSE	FALSE	4	0
## Gene_02	0.86	0.48	0.86	1.51	0.59	0.87	1.29	FALSE	FALSE	7	0
## Gene_03	0.98	0.65	0.98	1.49	0.55	0.91	1.51	FALSE	FALSE	4	0
## Gene_04	0.69	0.27	0.70	1.72	0.60	0.86	1.24	FALSE	FALSE	8	0
## Gene_05	0.49	0.23	0.49	0.97	0.60	0.85	1.21	TRUE	FALSE	9	1
## Gene_06	1.31	0.93	1.31	1.90	0.55	1.00	1.82	FALSE	FALSE	2	0
## Gene_07	0.60	0.40	0.60	0.97	0.56	0.89	1.41	TRUE	FALSE	5	1
## Gene_08	0.50	0.16	0.50	1.33	0.60	0.86	1.24	FALSE	FALSE	8	0
## Gene_09	1.02	0.63	1.03	1.58	0.56	0.89	1.41	FALSE	FALSE	5	0
## Gene_10	0.65	0.29	0.66	1.23	0.56	0.89	1.41	FALSE	FALSE	5	0
## Gene_11	0.79	0.43	0.79	1.48	0.60	0.85	1.21	FALSE	FALSE	9	0
## Gene_12	1.40	1.04	1.41	1.88	0.56	0.89	1.41	TRUE	FALSE	5	1
## Gene_13	0.57	0.10	0.57	1.50	0.58	0.88	1.34	FALSE	FALSE	6	0
## Gene_14	0.28	0.07	0.28	0.92	0.56	0.89	1.41	TRUE	FALSE	5	1
## Gene_15	0.47	0.27	0.48	0.81	0.55	1.00	1.82	TRUE	FALSE	2	1
## Gene_16	1.24	0.68	1.23	2.05	0.60	0.86	1.24	FALSE	FALSE	8	0
## Gene_17	2.34	2.10	2.35	2.70	0.55	0.94	1.63	TRUE	TRUE	3	2
## Gene_18	0.53	0.35	0.53	0.82	0.55	1.00	1.82	TRUE	FALSE	2	1
## Gene_19	0.55	0.31	0.55	0.99	0.55	0.91	1.51	TRUE	FALSE	4	1
## Gene_20	1.14	0.74	1.14	1.73	0.60	0.85	1.21	FALSE	FALSE	9	0

Figure 1: Sample report table.

Mean Ratio	The mean amplicon ratio between reference and sample of interest for the gene
LowerBoundBoot   MeanBoot   UpperBoundBoot	The confidence interval of the Bootstrap distribution.
Lower   Mean   Upper Noise	The Lower/Mean/Upper background noise bounds.
Signif.	Boolean value representing if read counts differences are Significant
AboveNoise	Boolean value representing if LowerBoundBoot is above UpperNoise (UpperBoundBoot below LowerNoise)
Amplicons	Number of Amplicons associated with the gene
Passed	if both Signif. and AboveNoise are TRUE then 2, if only one is TRUE then 1 and if both are FALSE then 0

Table 1: Report table Column Description

### 2.6.2 Plots

The generated plots (one per test sample) show the bootstrap distribution for each gene. The function **PlotBootstrapDistributions** generates a list of plots for all test samples (A plot per sample). At the figure 2 we can see an example of the plot for a single sample.

```
PlotBootstrapDistributions(bootList, reportTables)
```

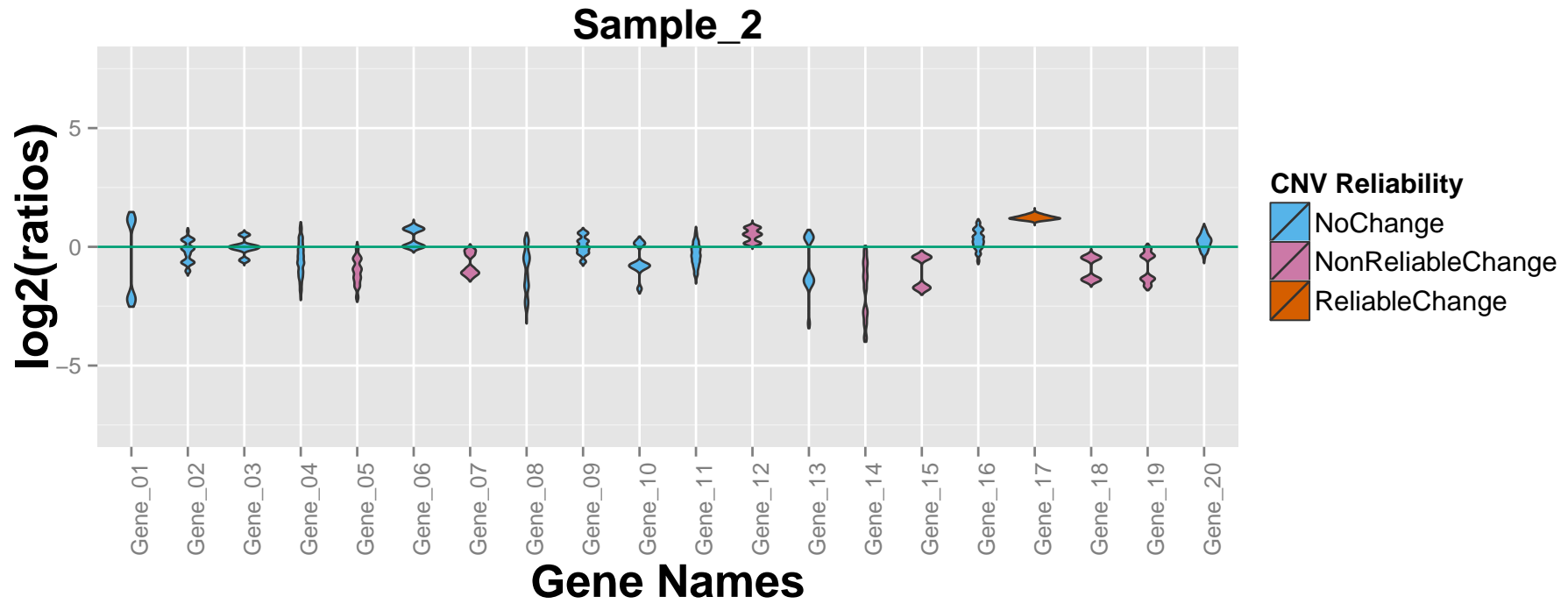


Figure 2: Plot for a single test sample

No Change	There is no change in the read count
Non Reliable Change	The change is significant but below the noise level for the number of amplicons associated with this gene
Reliable Change	The change is above the noise level and significant

Table 2: Description of the CNV detection levels