



# DNA Copy Number and Loss of Heterozygosity Analysis Algorithms

Detection of copy-number variants and chromosomal aberrations in GenomeStudio® software.

## Introduction

Illumina has developed several algorithms for detecting copy number variants (CNVs) and other structural variants from microarray data (Table 1). These algorithms are available as individual software plug-ins for the GenomeStudio Genotyping Module and can be downloaded from the GenomeStudio Portal, from the download section of the Illumina Connect web page<sup>1</sup>, or from MyIllumina (my.illumina.com). The plug-ins are used within the CNV Analysis workbench, and results can be visualized within the GenomeStudio Full Data Table, in the Illumina Genome Viewer (IGV), or in a CNV region display window. This technical note describes the function of these algorithms and how they can be employed to analyze a chromosomal region of interest.

## CNV Region Report

CNV Region Report is a software plug-in for GenomeStudio that generates three separate CNV reports.

- **Standard Report**—Lists each CNV and loss of heterozygosity (LOH) region for each selected sample.
- **Allele-Specific Copy Number Report**—Estimates the allele-specific copy number for each probe entry (e.g., A- or AAB). In the output file, the CN\_GTYPE column is calculated using the CNV Value (as determined by CNVPartition), the B Allele Frequency data, the GTYPE (genotype column), and the theoretical B Allele Frequency normal distributions for each copy number.
- **PLINK CNV Input Report**—Creates input files for some of the CNV features of the PLINK genome-wide association study (GWAS) and CNV analysis application<sup>2</sup>.

## cnvPartition

The goal of the cnvPartition algorithm is to identify regions of the genome that are aberrant in copy number using two Infinium® assay outputs: the log R ratio (LRR) and B allele frequency (BAF). Because LRR is the logged ratio of observed probe intensity to expected intensity, any deviations from zero in this metric are evidence for copy number change. BAF is the proportion of hybridized sample that carries the B allele as designated by the Infinium assay. In a normal sample, discrete BAFs of 0.0, 0.5, and 1.0 are expected for each locus (representing AA, AB, and BB).

Deviations from this expectation are indicative of aberrant copy number. For example, if a locus has a BAF of 0.66, this might indicate that there are two copies of the B allele and one copy of the A allele present in the sample  $\frac{2}{2+1} = 0.66$ . Analyzing both of these metrics provides stronger resolution for detecting true copy number changes.

Table 1: GenomeStudio Copy Number Algorithms

Algorithm	Function
CNV Region Report	Generates three separate CNV reports
cnvPartition	Calculates copy numbers with confidence scores and generates CNV regions
Homozygosity Detector	Autobooks samples with extended tracts of homozygosity (single-sample analysis only)
LOH Score	Estimates the likelihood of a region exhibiting LOH

## Copy Number Estimation

cnvPartition models LRRs and BAFs for each of 14 different copy number scenarios as simple bivariate Gaussian distributions (Table 2).

Modeling copy number in this way allows for computation of a preliminary copy number estimate for each assayed locus by comparing its observed LRR and BAF to values predicted from each of the fourteen models. Specifically, the likelihood of observing a given LRR and BAF under each of the 14 models is calculated. For example, to compute the likelihood of a particular LRR and BAF given a genotype of AAB (L<sub>AAB</sub>), the AAB parameters from Table 2 and the standard normal density are used:

$$L_{AAB} = \frac{1}{0.18\sqrt{2\pi}} \exp -\frac{(LRR - 0.3)^2}{2(0.18^2)} + \frac{1}{0.03\sqrt{2\pi}} \exp -\frac{(BAF - \frac{1}{3})^2}{2(0.03^2)}$$

Likelihoods are also computed for other model genotypes listed in Table 1 with the exception of the homozygous deletion (DD). For homozygous deletions, a very low LRR is expected, but the BAF may be any value between zero and one. Therefore, the likelihood of a double deletion (LDD) is calculated by the equation:

$$L_{DD} = \frac{1}{2 \times \sqrt{2\pi}} \exp -\frac{(LRR - (-5))^2}{2(2^2)}$$

**Table 2: Genotypes Modeled by cnvPartition**

Genotype	CN	LRR Mean	LRR SD	BAF Mean	BAF SD
DD	0	-5	2	NA	NA
A	1	-0.45	0.18	0	0.03
B	1	-0.45	0.18	1	0.03
AA	2	0	0.18	0	0.03
AB	2	0	0.18	0.5	0.03
BB	2	0	0.18	1	0.03
AAA	3	0.3	0.18	0	0.03
AAB	3	0.3	0.18	1/3	0.03
ABB	3	0.3	0.18	2/3	0.03
BBB	3	0.3	0.18	1	0.03
AAAA	4	0.75	0.18	0	0.03
AAAB	4	0.75	0.18	0.25	0.03
ABBB	4	0.75	0.18	0.75	0.03
BBBB	4	0.75	0.18	1	0.03

Parameters for each of the fourteen genotypes considered by cnvPartition are shown. BAFs are modeled as a uniform distribution between zero and one for homozygous deletions (DD). All other distributions are modeled with Gaussian distributions with the given parameters. The genotype AABB is not modeled since this would represent two independent duplication events and rarely occurs in nature. (CN = copy number, DD = double deletion, SD = standard deviation)

These likelihoods are then summarized by four composite copy number likelihoods:

$$\begin{aligned}
 L_0 &= L_{DD} \\
 L_1 &= L_A + L_B \\
 L_2 &= L_{AA} + L_{AB} + L_{BB} \\
 L_3 &= L_{AAA} + L_{AAB} + L_{ABB} + L_{BBB} \\
 L_4 &= L_{AAAA} + L_{AAAB} + L_{ABBB} + L_{BBBB}
 \end{aligned}$$

where  $L_k$  denotes the likelihood of copy number  $k$  for integer values of  $k$  and the likelihood of a genotype for non-numeric values of  $k$ . The preliminary copy number estimate ( $X$ ) is defined as the average of the five modeled copy numbers, weighted by their respective likelihoods:

$$X = \frac{L_1 + 2L_2 + 3L_3 + 4L_4}{L_0 + L_1 + L_2 + L_3 + L_4}$$

### Breakpoint Identification

Preliminary copy number estimates are the inputs to the core partitioning algorithm. The goal of partitioning is to identify regions of the genome where the values of  $X$  are consistently higher or lower than 2, the expected value for a diploid sample. To find an aberrant region, the algorithm orders the  $X$  values by their position along a chromosome and searches for the indexes  $i$  and  $j$  such that the values  $X_i \dots X_j$  are maximally different than those outside this region. Thus, the algorithm seeks to maximize  $|Z_j|$  over all  $i$  and  $j$  with  $i < j$ , defined by the equations:

$$S_i = X_1 + \dots + X_i, 1 \leq i \leq n$$

$$Z_j = \frac{1}{(j-i)} + \frac{1}{(n-j+i)}^{-1/2} \times \frac{(S_j - S_i)}{(j-i)} - \frac{(S_n - S_j + S_i)}{(n-j+i)}$$

where  $n$  is the number of loci assayed on the chromosome.

An exhaustive search through all pairs of  $i$  and  $j$  scales quadratically with  $n$  and is therefore an inefficient process for use with Illumina whole-genome genotyping products<sup>3,4</sup>. To simplify the calculations required, cnvPartition uses a sliding window strategy to maximize  $|Z_j|$ , but where  $j = i + w$ , with  $w$  the defined window size. After the optimal window size value is found, the algorithm attempts to extend the window in both directions to maximize the value of  $|Z_j|$  further. As implemented, the algorithm repeats this procedure for  $w = 4, 8, 16$ , and 32 then reports the  $i$  and  $j$  corresponding to the maximal  $|Z_j|$  found. When a maximally different segment is found,  $|Z_j|$  is compared to a pre-determined threshold (default is 6). If the threshold is exceeded, the boundaries are noted and the algorithm is applied recursively to the regions between 1 and  $i, i+1$  and  $j, j+1$  through  $n$ . The threshold of 6 was chosen as a default because it minimizes false positives, particularly for short aberrations.

### Copy Number Assignment to Partitioned Regions

The partitioning procedure results in a set of putative breakpoints scattered across the genome. The next step is to assign a copy number for each region lying between two consecutive breakpoints. To do this,  $L_0, L_1, L_2, L_3,$  and  $L_4$  for each locus within the region are used. For each putative copy number (0-4), the logarithms of all  $L_k$  for each  $k$  are summed. The  $k$  with the highest sum is the copy number assigned to this region. For regions with copy numbers other than 2, the algorithm also assigns a confidence score for the copy number that is called. The confidence score is defined as the sum of all logged likelihoods in the region for the assigned copy number minus the sum of all  $\log L_2$  values for loci in the region.

### Additional Usage Notes

- Regions with copy number = 1 on the X or Y for males are filtered from the CNV results.
- Probes that are designated as Intensity Only are treated differently than normal probes. The B Allele Frequency is ignored for these probes.
- Y probes are not considered for samples designated as female.

### Homozygous Region Detection

cnvPartition also includes a homozygosity detection algorithm that runs separately from the partitioning algorithm already described. This algorithm only runs on copy number 2 regions by default. Therefore, it is sometimes called a Copy Neutral LOH Detector. The logic is similar to that used in the homozygosity detector autobookmarking plug-in (see next section). However, additional logic has been added to simplify usage for the end user. Instead of adjusting the ChiSquare threshold as in the autobookmarking plug-in, the user can simply adjust the MinHomozygousRegionSize configuration parameter. By default, this is set to 10 Mb based on empirical testing.





## Summary

GenomeStudio provides several methods to analyze SNP and probe intensity data to identify chromosomal regions with LOH and copy number variations. The software plug-ins described in this technical note are freely available to GenomeStudio users to provide extended functionality.

Automated bookmarking algorithms save time by automatically scanning and categorizing samples. Researchers can use *cnvPartition* to find and calculate copy numbers, or *Homozygosity Detector* to identify extended tracts of LOH.

The LOH Score algorithm provides statistical information about chromosomal aberrations of interest. This information includes the probability of LOH existing. This algorithm can be used to identify interesting regions in large sample sets quickly or to analyze a more refined region further.

The open architecture of Illumina GenomeStudio software allows for customized and advanced analysis tools for the downstream analysis of Illumina DNA Analysis BeadChip Genotyping data. The plug-ins described in this document can be downloaded from the GenomeStudio Portal or from the Illumina Connect web page at [www.illumina.com/illuminaconnect](http://www.illumina.com/illuminaconnect). Illumina Connect is a collaborative program for facilitating the development of third-party tools and applications for DNA and RNA analysis of Illumina BeadArray™ products. Users should check the Illumina Connect web page regularly for new or updated plug-ins.

## References

1. [www.illumina.com/Illuminaconnect](http://www.illumina.com/Illuminaconnect)
2. [pngu.mgh.harvard.edu/~purcell/plink/](http://pngu.mgh.harvard.edu/~purcell/plink/)
3. Olshen AB, Venkatraman ES, Lucito R, Wigler M (2004) Circular binary segmentation for the analysis of array-based DNA copy number data. *Biostatistics* 5: 557–572.
4. Venkatraman ES, Olshen AB (2007) A faster circular binary segmentation algorithm for the analysis of array CGH data. *Bioinformatics* 23: 657–663.
5. [my.illumina.com/](http://my.illumina.com/)

AAAGAATGATAACAGTAAACACACTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACGTACCGTAAACGAACGTATCAATTGAGACTAAATATTAACGTACCATTAAAGAGCTACCGTCTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCA  
AATCAACGTACCGTAAACGAACGTATCAATTAAAGATTACTTGATCCACTGATTCAACGTACCGTAAACGAACGTATCAATTGAGACTAAATATTAACGTACCATTAAAGAGCTACCGTCTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCA  
AACGACGAAAGAATGATAACAGTAAACACACTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACGTACCGTAAACGAACGTATCAATTGAGACTAAATATTAACGTACCATTAAAGAGCTACCGTCTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCA  
TTAAAGTACCATTAAAGAGCTACCGTGCAAACAGTAAACACACTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACGTACCGTAAACGAACGTATCAATTGAGACTAAATATTAACGTACCATTAAAGAGCTACCGTGCAAACGAACGAAAGAATGA  
AAAGAATGATAACAGTAAACACACTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCAACGTACCGTAAACGAACGTATCAATTGAGACTAAATATTAACGTACCATTAAAGAGCTACCGTCTTCTGTTAACCTTAAGATTACTTGATCCACTGATTCA  
AAGATTACTTGATCCACTGATTCAACGTAAAGATTACTTGATCCACTGATTCAACGTACCGTAAACGAACGTATCAATTGAGACTACGAAACGATCAATTGAGACTAGCAACGAAACGATCAATTGAGACTAAATATTAACGTACCATTAAAGAGCTACCGTGC  
AACGTATCAATTGAGACTAAATATTAACGTACCATTAAAGATTACTTGATCCACTGATTCAACGTACCGTAAACGAACGTATCAATTGAGACTAAATATTAACGTACCATTAAAGAGCTACCGTGCAAACGAACGAAAGAATGATAAC

