

# Tools for high throughput SNP chip data

Robert Scharpf

April 5, 2010

## Introduction

*SNPchip* defines classes and methods useful for organizing high throughput genomic data. The classes defined here extend the **eSet** class in *Biobase*, utilizing the existing Bioconductor infrastructure for genomic data. This provides a foundation upon which statistical and visualization tools can be further developed. See [2] for additional details.

## 1 Simple Usage

```
> library(SNPchip)

> data(sample.snpset)
> sample.snpset

oligoSnpSet (storageMode: lockedEnvironment)
assayData: 5859 features, 5 samples
  element names: call, callProbability, cnConfidence, copyNumber
experimentData: use 'experimentData(object)'
Annotation: pd.mapping50k.xba240
phenoData
An object of class "AnnotatedDataFrame": none
featureData
An object of class "AnnotatedDataFrame"
  rowNames: SNP_A-1507972, SNP_A-1641761, ..., SNP_A-1759046 (5859 total)
  varLabels and varMetadata description:
    chromosome: chromosome
    position: position
    ...: ...
    enzyme: enzyme
    (9 total)
Annotation [1] "pd.mapping50k.xba240"
```

For Affymetrix platforms, access to SNP-level annotation such as chromosome and physical position is through *RSQLite*. See the *oligo* vignette for additional information regarding available SNP-level Annotation. Accessing RSQLite databases each time chromosome and physical position are needed (such as for making scatterplots of physical position versus copy number) may increase the time of computation, particularly for small objects such as the **sample.snpset**. Therefore, one may wish to add this information to the **featureData** slot. Subsequent calls to chromosome and physical position will look in the **featureData** before searching in the RSQLite database. Because **sample.snpset** is small and adding this annotation to the **featureData** slot will greatly increase the speed for graphing

the data (the graphs make repeated use of chromosome and position calls), we add this annotation to the `featureData` for this vignette.

Note that if one is using a platform for which there is no `pd.mapping` package available, one must provide the chromosome (character string) and physical position (integer) in the `featureData` slot using the column labels “chromosome” and “position”, respectively.

Here we select SNPs on the first three chromosomes. An object of class `ParSnpSet` contains graphical parameters for visualizing the data for this subset:

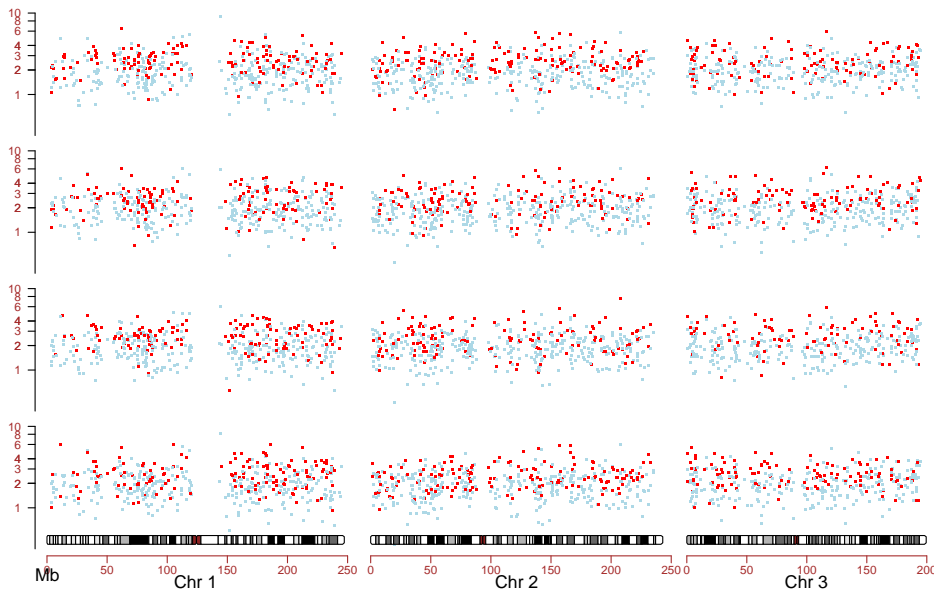
```
> snpset <- sample.snpset[chromosome(sample.snpset) %in%
+   as.character(1:3), 1:4]
> graph.par <- plot(snpset)
> class(graph.par)
```

```
[1] "ParSnpSet"
attr(,"package")
[1] "SNPchip"
```

```
> graph.par$use.chromosome.size <- TRUE
> graph.par$pch <- "."
> graph.par$cex <- 3
> graph.par$oma <- c(3, 3, 3, 0.5)
```

Plot the first few chromosomes for samples 1-4 (note the `print()` command is needed for producing the figures in the vignette, but just typing the object name is sufficient for plotting):

```
> print(graph.par)
```



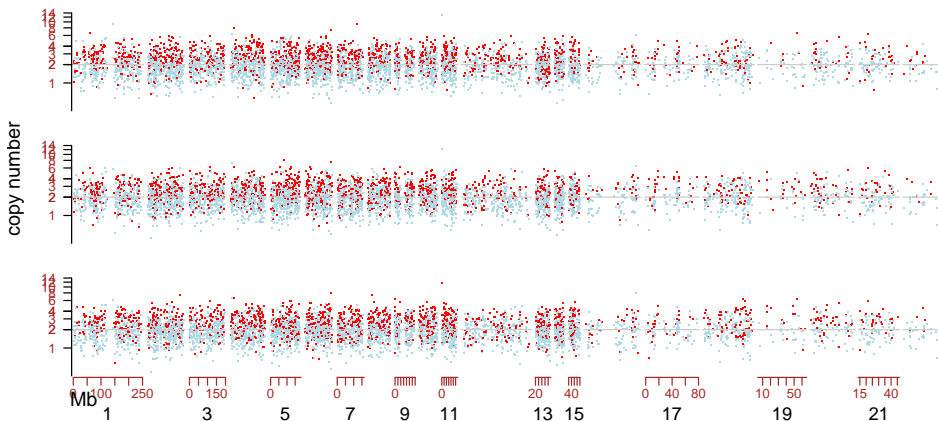
```
> graph.par$cytoband.side <- 3
> graph.par$heights <- rev(graph.par$heights)
```

The samples are plotted by row. For each sample, the copy number (vertical axis) is plotted against the physical position of the SNP in the chromosome. Here, the chromosome labels are plotted beneath the cytobands.

## 2.1 Genome-wide plots for multiple samples

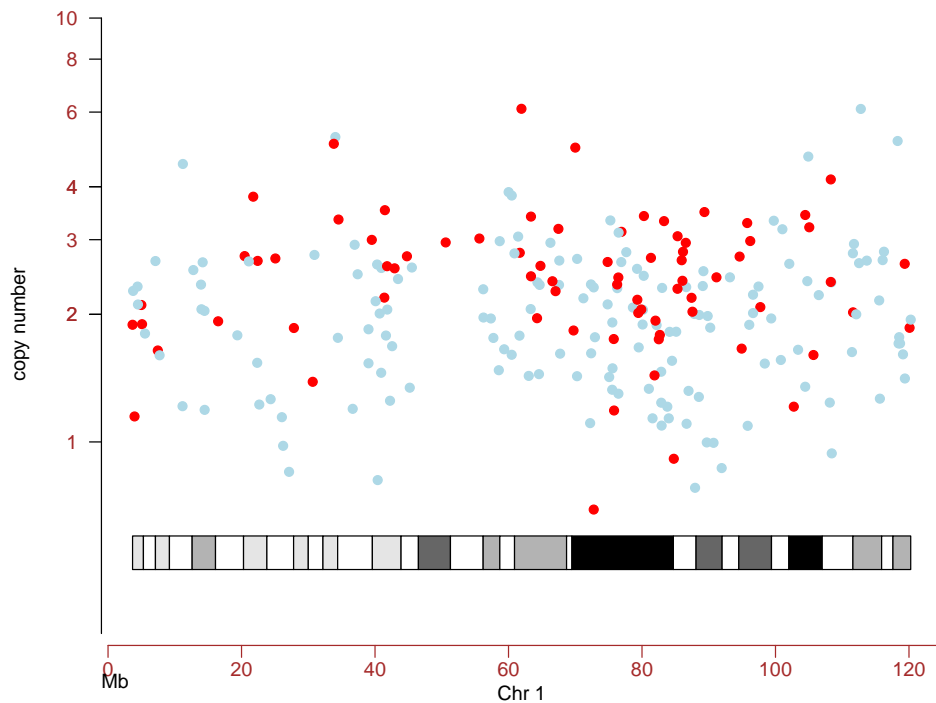
```
> graph.par <- plot(sample.snpset[chromosome(sample.snpset) <
+ 23, 1:3], add.cytoband = FALSE)
> graph.par$one.ylim <- TRUE
> graph.par$mar <- c(0.1, 0.1, 2, 0.1)
> graph.par$oma <- c(3, 4, 2, 1)
> graph.par$cex <- 2
> graph.par$abline <- TRUE
> graph.par$cex.lab <- 0.9
> graph.par$add.cytoband <- FALSE

> print(graph.par)
```



```
> parm <- centromere("1")[1]
> snpset <- sample.snpset[chromosome(sample.snpset) ==
+   "1" & (position(sample.snpset) < parm), 2]
> graph.par <- plot(snpset)
> graph.par$use.chromosome.size <- FALSE
> graph.par$pch <- 21
> graph.par$cex <- 1
> graph.par$ylim <- c(0.4, 9)
> graph.par$cytoband.ycoords <- c(0.5, 0.6)

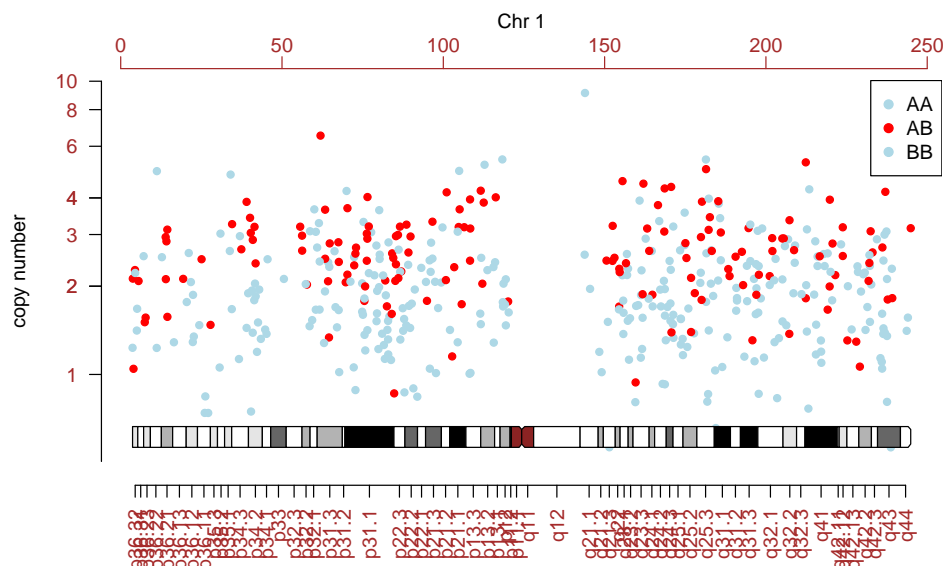
> print(graph.par)
```



Note that the cytoband is automatically subsetting appropriately. Had we instead specified `use.chromosome.size=TRUE`, the x-axis limits would include the entire chromosome (and cytoband) though only the SNPs on the p-arm would be plotted.

Adding a legend for the genotypes

```
> data(sample.snpset)
> x <- sample.snpset[chromosome(sample.snpset) ==
+   "1", 1]
> gp <- plot(x)
> gp$legend <- c("AA", "AB", "BB")
> gp$legend.col <- unique(gp$col)
> gp$legend.bg <- unique(gp$bg)
> gp$pch <- 21
> gp$cex <- 0.8
> gp$label.cytoband <- TRUE
> gp$add.centromere <- FALSE
> gp$xlabel <- ""
> gp$legend.bty = "o"
> gp$ylim[2] <- 9
> print(gp)
```

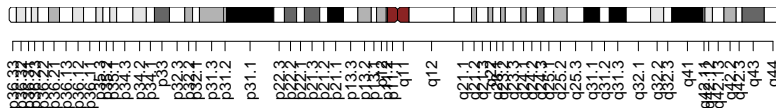


## 2.2 Plotting cytoband

To plot the cytoband of chromosome 1,

```
> plotCytoband("1")
```

NULL



Hidden Markov models for objects of class `SnpSet`, `CopyNumberSet`, and `oligoSnpSet` are available in the package *VanillaICE* available at Bioconductor. A more detailed description of the hidden Markov models fit in the *VanillaICE* package are discussed elsewhere [1]. Here we provide an example of a lowess smoother for copy number estimates. The following code chunk first assigns heterozygous calls to the integer 1 and homozygous calls to the integer zero. It follows that regions of deletions will have homozygous calls of zero. We simulated a deletion of 50 consecutive SNPs and converted the `sample.snpset` to a list where each element in the list is an `oligoSnpSet` object for one chromosome.

## 2.3 Descriptive and statistical summaries

Descriptive statistics for copy number and genotype calls are provided with the `summary` method. For each chromosome in the `oligoSnpSet`, `summary` calculates the average and standard deviation of the copy number estimates, as well as the % homozygous and heterozygous calls. In addition, `summary` calculates the average copy number, standard deviation, % homozygous and heterozygous across all autosomes in the `oligoSnpSet`. The dimensions of the four matrices are  $S \times C + 1$ , where  $S$  is the number of samples and  $C$  is the number of chromosomes in the `oligoSnpSet`.

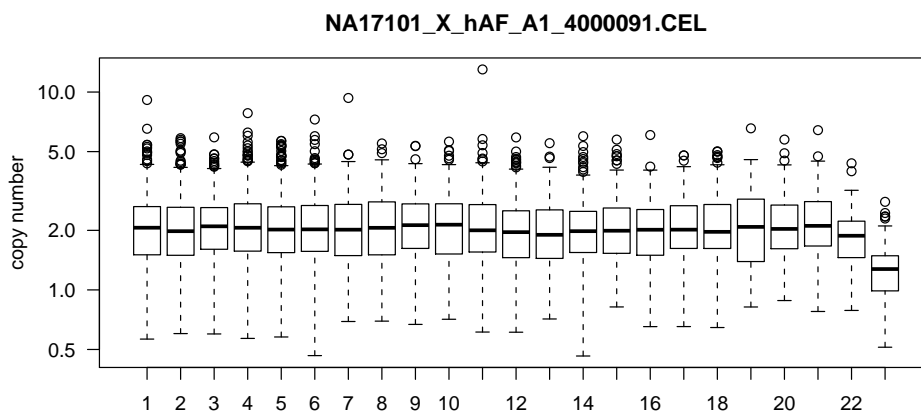
```
> x <- summary(sample.snpset, digits = 1)
> str(x)
```

List of 5

```
$ avg.CN      : num [1:23, 1:5] 2.2 2.1 2.2 2.3 2.2 2.2 2.2 2.2 2.2 2.2 ...
..- attr(*, "dimnames")=List of 2
.. ..$ : chr [1:23] "1" "2" "3" "4" ...
.. ..$ : chr [1:5] "NA17101_X_hAF_A1_4000091.CEL" "NA17102_X_hAF_A2_4000091.CEL" "NA17103_X_hAF_A3_4
$ sd.CN      : num [1:23, 1:5] 1 0.9 0.8 1 1 0.9 0.9 0.9 0.8 0.9 ...
..- attr(*, "dimnames")=List of 2
.. ..$ : chr [1:23] "1" "2" "3" "4" ...
.. ..$ : chr [1:5] "NA17101_X_hAF_A1_4000091.CEL" "NA17102_X_hAF_A2_4000091.CEL" "NA17103_X_hAF_A3_4
$ prop.Hom   : num [1:23, 1:5] 0.7 0.7 0.7 0.7 0.7 0.7 0.7 0.7 0.7 0.7 ...
..- attr(*, "dimnames")=List of 2
.. ..$ : chr [1:23] "1" "2" "3" "4" ...
.. ..$ : chr [1:5] "NA17101_X_hAF_A1_4000091.CEL" "NA17102_X_hAF_A2_4000091.CEL" "NA17103_X_hAF_A3_4
$ prop.Het   : num [1:23, 1:5] 0.3 0.3 0.3 0.3 0.3 0.3 0.3 0.3 0.3 0.3 ...
..- attr(*, "dimnames")=List of 2
.. ..$ : chr [1:23] "1" "2" "3" "4" ...
.. ..$ : chr [1:5] "NA17101_X_hAF_A1_4000091.CEL" "NA17102_X_hAF_A2_4000091.CEL" "NA17103_X_hAF_A3_4
$ prop.NoCall: num [1:23, 1:5] 0 0 0 0 0 0 0 0 0 0 ...
..- attr(*, "dimnames")=List of 2
.. ..$ : chr [1:23] "1" "2" "3" "4" ...
.. ..$ : chr [1:5] "NA17101_X_hAF_A1_4000091.CEL" "NA17102_X_hAF_A2_4000091.CEL" "NA17103_X_hAF_A3_4
```

Boxplot by chromosome:

```
> op <- par(mfrow = c(1, 1), mar = c(4, 4, 3, 1),
+   las = 1)
> boxplot(split(copyNumber(sample.snpset[, 1]),
+   chromosome(sample.snpset)), ylab = "copy number",
+   main = sampleNames(sample.snpset)[1], log = "y")
> par(op)
```



### 3 Annotation

#### 3.1 Chromosome-level

The chromosome-level annotation for build hg18:

```
> list.files(system.file("hg18", package = "SNPchip"))

[1] "centromeres.txt" "chr1_gap.txt"    "chr10_gap.txt"
[4] "chr11_gap.txt"   "chr12_gap.txt"   "chr13_gap.txt"
[7] "chr14_gap.txt"   "chr15_gap.txt"   "chr16_gap.txt"
[10] "chr17_gap.txt"   "chr18_gap.txt"   "chr19_gap.txt"
[13] "chr2_gap.txt"    "chr20_gap.txt"   "chr21_gap.txt"
[16] "chr22_gap.txt"   "chr3_gap.txt"    "chr4_gap.txt"
[19] "chr5_gap.txt"    "chr6_gap.txt"    "chr7_gap.txt"
[22] "chr8_gap.txt"    "chr9_gap.txt"    "chrM_gap.txt"
[25] "chromInfo.txt"   "chrX_gap.txt"    "chrY_gap.txt"
[28] "cytoBand.txt"
```

## 3.2 Feature-level

Feature-level annotation for Affymetrix platforms is available in the `pd.mapping` packages. See the *oligo* vignette for additional information about available feature-level annotation.

# 4 Integration with other Bioconductor packages

## 4.1 *oligo*

For generating `SnpCallSets` from .CEL files, see the R package *oligo*. In particular, the function `crlmm` in *oligo* creates an instance of the class `SnpCallSetPlus`. See the *oligoHowTo* vignette for details on coercing an object of class `SnpCallSetPlus` to `oligoSnpSet`.

## 4.2 *RSNPper*

To retrieve additional annotation on the known SNP's in the region of this simulated deletion, we could use the *RSNPper*.

```
> library(RSNPper)
> (dbId <- dbSnpId(annSnpset)[snps[2] == featureNames(annSnpset)])
> dbId <- strsplit(dbId, "rs")[[1]][2]
> print(SNPinfo(dbId))
```

# 5 Session Information

The version number of R and packages loaded for generating the vignette were:

- R version 2.11.0 Under development (unstable) (2009-11-20 r50517),  
x86\_64-apple-darwin9.8.0
- Locale: en\_US.UTF-8/en\_US.UTF-8/C/C/en\_US.UTF-8/en\_US.UTF-8
- Base packages: base, datasets, graphics, grDevices, methods, stats, tools, utils
- Other packages: Biobase 2.7.4, genefilter 1.29.3, mybase 0.0.1, oligoClasses 1.9.53,  
SNPchip 1.11.1
- Loaded via a namespace (and not attached): affyio 1.15.2, annotate 1.25.0, AnnotationDbi 1.9.2,  
Biostrings 2.15.11, DBI 0.2-4, IRanges 1.5.21, RSQLite 0.7-3, splines 2.11.0, survival 2.35-7,  
xtable 1.5-6

## References

- [1] Robert B Scharpf, Giovanni Parmigiani, Jonathan Pevsner, and Ingo Ruczinski. Hidden Markov models for the assessment of chromosomal alterations using high-throughput SNP arrays. *Annals of Applied Statistics*, 2008. To appear.
- [2] Robert B Scharpf, Jason C Ting, Jonathan Pevsner, and Ingo Ruczinski. SNPchip: R classes and methods for SNP array data. *Bioinformatics*, 23(5):627–628, Mar 2007.