

# Infrastructure for copy number analysis in `crlmm`

Rob Scharpf

March 29, 2011

## 1 Set up

As in the previous vignettes, we load the required libraries and specify a path for storing output files.

```
R> library(ff)
R> library(crlmm)
R> library(cacheSweave)
R> require(DNAcopy)
R> require(VanillaICE)
R> if (getRversion() < "2.13.0") {
  rpath <- getRversion()
} else rpath <- "trunk"
R> outdir <- paste("/thumper/ctsa/snpmicroarray/rs/ProcessedData/crlmm/",
  rpath, "/copynumber_vignette", sep = "")

R> ldPath(outdir)
R> setCacheDir(outdir)
```

We begin by loading the `cnSet` object created by the `AffymetrixPreprocessCN` vignette.

```
R> if (!exists("cnSet")) load(file.path(outdir, "cnSet.rda"))
```

## 2 Supported platforms

The supported Affymetrix and Infinium platforms are those for which a corresponding annotation package is available. The annotation packages contain information on the markers, such as physical position and chromosome, as well as pre-computed parameters estimated from HapMap used during the preprocessing and genotyping steps. For Affymetrix, the 5.0 and 6.0 platforms are supported and the corresponding annotation packages are `genomewidesnp5Crlmm` and `genomewidesnp6Crlmm`. Supported Infinium platforms are listed in the following code chunk.

```
R> pkgs <- annotationPackages()
R> crlmm.pkgs <- pkgs[grepl("Crlmm", pkgs)]
R> crlmm.pkgs[grepl("human", crlmm.pkgs)]

[1] "human370v1cCrlmm"      "human370quadv3cCrlmm"
[3] "human550v3bCrlmm"     "human650v3aCrlmm"
[5] "human610quadv1bCrlmm"  "human660quadv1aCrlmm"
[7] "human1mduov3bCrlmm"   "humanomni1quadv1bCrlmm"
```

**Large data:** In order to reduce `crlmm`'s memory footprint for copy number estimation, we require the `ff`. The `ff` package provides infrastructure for accessing and writing data to disk instead of keeping data in memory. As the functions for preprocessing, genotyping, and copy number estimation do not simultaneously require all samples and all probes in memory, memory-usage by `crlmm` can be fine-tuned by reading in

and processing subsets of the markers and/or samples. The functions `ocSamples` and `ocProbesets` in the `oligoClasses` package can be used to declare how many markers and samples to read at once. In general, specifying smaller values should reduce the RAM required for a particular job. In general, smaller values will increase the run-time. In the following code-chunk, we declare that `crmm` should process 150,000 markers at a time (when possible) and 500 samples at a time. If our dataset contained fewer than 500 samples, the `ocSamples` option would not have any effect. One can view the current settings for these commands, by typing the functions without an argument.

```
R> ocProbesets(50000)
R> ocSamples(200)
```

### 3 The *CNSet* container

The `show` method provides a concise summary of the `cnSet` object.

```
R> cnSet

CNSet (storageMode: lockedEnvironment)
assayData: 1852215 features, 180 samples
  element names: alleleA, alleleB, call, callProbability
protocolData
  rowNames: NA06985_GW6_C.CEL NA06991_GW6_C.CEL ...
            NA19240_GW6_Y.CEL (180 total)
  varLabels: ScanDate
  varMetadata: labelDescription
phenoData
  sampleNames: NA06985_GW6_C.CEL NA06991_GW6_C.CEL ...
              NA19240_GW6_Y.CEL (180 total)
  varLabels: SKW SNR gender
  varMetadata: labelDescription
featureData
  featureNames: SNP_A-2131660 SNP_A-1967418 ... CN_954736
              (1852215 total)
  fvarLabels: chromosome position isSnp
  fvarMetadata: labelDescription
experimentData: use 'experimentData(object)'
Annotation: genomewidesnp6
batch:      C 90, Y 90, grandMean 90
batchStatistics: 29 elements, 1852215 features, 2 batches
```

We briefly outline some of the unique aspects of the *CNSet*-class using in `crmm` that may differ from the more standard extensions of the virtual class *eSet* defined in the `Biobase` package.

**Instantiating a *CNSet*:** An object of class *CNSet* can be instantiated by one of two methods:

1. during the preprocessing of the raw intensities for Illumina and Affymetrix arrays by the the functions `constructInf` and `genotype`, respectively.
2. by subsetting an existing *CNSet* object. As per usual, the `['` method can be used to extract a subset of markers *i* as in `['i, ]'`, a subset of samples *j* as in `['[, j]'`, or a subset of markers *i* and samples *j* as in `['i, j]'`.

There are important differences in the underlying data representation depending on how the object was instantiated. In particular, objects generated by the functions `constructInf` and `genotype` store high-dimensional data on disk rather than in memory through protocols defined in the R package `ff`. For instance,

the normalized intensities and genotype calls in a *CNSet*-instance from approach (1) are ff-derived objects. However, when such an object is subset by the '[' method, the data is pulled from disk to memory and stored as an ordinary matrix in R. To illustrate, the *cnSet* loaded at the start of this session was created by the function `genotype`. As the data is stored on disk rather than in memory, we inspect attributes of the object representing the normalized intensities for the A allele by first opening the file connection and then extracting the ff-derived class of the object as well as where the data is stored on disk.

```
R> invisible(open(cnSet))
R> class(A(cnSet))

[1] "ff_matrix" "ff_array"  "ff"

R> filename(A(cnSet))

[1] "/thumper/ctsa/snpmicroarray/rs/ProcessedData/crlmm/trunk/copynumber_vignette/crlmmA-52fbff7a.ff"
```

Using approach (2), we construct a new *CNSet* object.

```
R> cnset.subset <- cnSet[1:5, 1:10]
R> invisible(close(cnSet))
R> class(A(cnset.subset))

[1] "matrix"
```

The files with .ff extension should not be removed or relocated. The safest way to move these files if necessary is to clone all of the ff objects using the `clone`, followed by the `delete` function to remove the original files on disk. See the documentation in the ff package for additional details.

**Order of operations:** For *CNSet*-instances derived by approach (1), users should be mindful of the substantial I/O when using accessors to extract data from the class. For example, the following 2 methods would extract identical results, with the latter being much more efficient (extra parentheses are added to the second operation to emphasize the order of operations):

```
R> A(cnSet[1:5, 1:5])
```

	NA06985_GW6_C.CEL	NA06991_GW6_C.CEL	NA06993_GW6_C.CEL
SNP_A-2131660	1381	1226	720
SNP_A-1967418	289	223	251
SNP_A-1969580	917	1207	1124
SNP_A-4263484	1502	2424	2319
SNP_A-1978185	1130	1243	1377

  

	NA06994_GW6_C.CEL	NA07000_GW6_C.CEL
SNP_A-2131660	459	582
SNP_A-1967418	193	265
SNP_A-1969580	1084	1420
SNP_A-4263484	448	1653
SNP_A-1978185	1262	1296

  

```
R> (A(cnSet))[1:5, 1:5]
```

	NA06985_GW6_C.CEL	NA06991_GW6_C.CEL	NA06993_GW6_C.CEL
SNP_A-2131660	1381	1226	720
SNP_A-1967418	289	223	251
SNP_A-1969580	917	1207	1124
SNP_A-4263484	1502	2424	2319
SNP_A-1978185	1130	1243	1377

  

	NA06994_GW6_C.CEL	NA07000_GW6_C.CEL
SNP_A-2131660	459	582
SNP_A-1967418	193	265
SNP_A-1969580	1084	1420
SNP_A-4263484	448	1653
SNP_A-1978185	1262	1296

SNP_A-2131660	459	582
SNP_A-1967418	193	265
SNP_A-1969580	1084	1420
SNP_A-4263484	448	1653
SNP_A-1978185	1262	1296

#### featureData

Information on physical position, chromosome, and whether the marker is a SNP can be accessed through accessors defined for the `featureData`.

```
R> fvarLabels(cnSet)
```

```
[1] "chromosome" "position" "isSnp"
```

```
R> position(cnSet)[1:10]
```

```
[1] 1156131 2234251 2329564 2553624 2936870 2951834 3095126 3165267
[9] 3302871 3705226
```

```
R> chromosome(cnSet)[1:10]
```

```
[1] 1 1 1 1 1 1 1 1 1 1
```

```
R> is.snp <- isSnp(cnSet)
```

```
R> table(is.snp)
```

```
is.snp
```

```
FALSE TRUE
945615 906600
```

```
R> snp.index <- which(is.snp)
```

```
R> np.index <- which(!is.snp)
```

```
R> chr1.index <- which(chromosome(cnSet) == 1)
```

#### assayData

The `assayData` elements are of the class `ff_matrix/ffdf` or `matrix`, depending on how the `CNSet` object was instantiated. Elements in the `assayData` environment can be listed using the `ls` function.

```
R> ls(assayData(cnSet))
```

```
[1] "alleleA" "alleleB" "call"
[4] "callProbability"
```

The normalized intensities for the A and B alleles have names `alleleA` and `alleleB` and can be accessed by the methods `A` and `B`, respectively. Genotype calls and confidence scores can be accessed by `snpCall` and `snpCallProbability`, respectively. Note that the confidence scores are represented as integers to reduce the filesize, but can be translated to the probability scale using the R function `i2p`. For example,

```
R> scores <- as.matrix(snpCallProbability(cnSet)[1:5, 1:2])
```

```
R> i2p(scores)
```

	NA06985_GW6_C.CEL	NA06991_GW6_C.CEL
SNP_A-2131660	0.9999949	0.9999995
SNP_A-1967418	0.9999965	0.9999996
SNP_A-1969580	0.9995187	0.9995134
SNP_A-4263484	0.9999999	1.0000000
SNP_A-1978185	1.0000000	1.0000000

Note that for the Affymetrix 6.0 platform the assay data elements each have a row dimension corresponding to the total number of polymorphic and nonpolymorphic markers interrogated by the Affymetrix 6.0 platform. A consequence of keeping the rows of the assay data elements the same for all of the statistical summaries is that the matrix used to store genotype calls is larger than necessary.

Note that NA's are stored in the slot for normalized 'B' allele intensities:

```
R> np.index <- which(!is.snp)
R> stopifnot(all(is.na(B(cnSet)[np.index, ])))
```

#### batch and batchStatistics

As defined in Leek *et al.* 2010, *Batch effects are sub-groups of measurements that have qualitatively different behaviour across conditions and are unrelated to the biological or scientific variables in a study.* The **batchStatistics** slot is an environment used to store SNP- and batch-specific summaries, such as the sufficient statistics for the genotype clusters and the linear model parameters used for copy number estimation. The **batch** slot is used to store the 'batch name' for each array. For small studies in which the samples were processed at similar times (e.g., within a month), all the samples can be considered to be in the same batch. For large studies in which the samples were processed over several months, users should the scan date of the array or the chemistry plate are useful surrogates. The only constraint on the **batch** variable is that it must be a character vector that is the same length as the number of samples to be processed. The **batch** is specified as an argument to the R functions **constructInf** and **genotype** that instantiate *CNSet* objects for the Illumina and Affymetrix platforms, respectively. The **batch** function can be used to access the **batch** information on the samples as in the following example.

```
R> table(batch(cnSet))
```

```
  C  Y
90 90
```

For the **batchStatistics** slot, the elements in the environment have the class *ff.matrix/ffdf* or *matrix*, depending on how the *CNSet* object was instantiated. The dimension of each element is the number of markers (SNPs + nonpolymorphic markers)  $\times$  the number of batches. The names of the elements in the environment can be list using the R function **ls**.

```
R> ls(batchStatistics(cnSet))
```

```
[1] "corrAA"      "corrAB"      "corrBB"      "flags"       "madA.AA"
[6] "madA.AB"     "madA.BB"     "madB.AA"     "madB.AB"     "madB.BB"
[11] "medianA.AA"  "medianA.AB"  "medianA.BB"  "medianB.AA"  "medianB.AB"
[16] "medianB.BB"  "N.AA"        "N.AB"        "N.BB"        "nuA"
[21] "nuB"         "phiA"        "phiB"        "phiPrimeA"   "phiPrimeB"
[26] "tau2A.AA"    "tau2A.BB"    "tau2B.AA"    "tau2B.BB"
```

Currently, the batch-specific summaries are stored to allow some flexibility in the choice of downstream analyses of copy number and visual assessments of model fit. Documentation for such applications will be expanded in future versions of *crIimm*, and are currently not intended to be accessed directly by the user.

#### phenoData

Sample-level summaries obtained during the preprocessing/genotyping steps include skew (SKW), the signal to noise ratio (SNR), and gender (1=male, 2=female) are stored in the **phenoData** slot. As for other *eSet* extensions, the **\$** method can be used to extract these summaries. For *CNSet* objects generated by approach (1), these elements are of the class *ff\_vector*.

```
R> varLabels(cnSet)
```

```
[1] "SKW"      "SNR"      "gender"
```

```
R> class(cnSet$gender)
```

```
[1] "ff_vector" "ff"

R> invisible(open(cnSet$gender))
R> cnSet$gender

ff (open) integer length=180 (180)
  [1]  [2]  [3]  [4]  [5]  [6]  [7]  [8]      [173] [174]
    2    2    1    1    2    2    1    1      :    1    1
[175] [176] [177] [178] [179] [180]
    2    2    1    2    1    2
```

The '[' methods without arguments can be used to coerce to a vector.

```
R> cnSet$gender[]

 [1] 2 2 1 1 2 2 1 1 1 1 2 2 2 2 1 1 2 1 1 2 1 2 1 1 2 1 1 2 1 2 2 1
[33] 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 1 1 2 2 2 2 1 2 1 1
[65] 1 2 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2 1 2 2 1 2 1 1 2 1 1 2
[97] 1 1 2 1 1 2 1 1 2 2 1 1 2 1 1 2 1 1 2 1 1 2 1 1 2 1 1 2 2 1 2
[129] 2 1 2 1 2 1 1 2 1 2 1 2 2 2 1 1 2 1 1 2 1 1 2 1 1 2 1 1 1
[161] 2 1 1 2 2 1 2 1 2 1 1 2 1 1 2 2 1 2 1 2

R> invisible(close(cnSet$gender))

protocolData
The scan date of the arrays are stored in the protocolData.

R> varLabels(protocolData(cnSet))

[1] "ScanDate"

R> protocolData(cnSet)$ScanDate[1:10]

 [1] 2007-03-06 2007-03-06 2007-03-06 2007-03-06 2007-03-06 2007-03-06
 [7] 2007-03-06 2007-03-06 2007-03-06 2007-03-06
Levels: 2007-03-02 2007-03-03 2007-03-06 2007-03-07 2007-03-08
```

## 4 Trouble shooting

### 4.1 Missing values

Most often, missing values occur when the genotype confidence scores for a SNP were below the threshold used by the `crlmmCpynumber` function. For the HapMap analysis, we used a confidence threshold of 0.80 (the default).

```
R> GT.CONF.THR <- 0.8
R> autosome.index <- which(isSnp(cnSet) & chromosome(cnSet) <
  23)
R> sample.index <- which(batch(cnSet) == "C")
R> ca <- CA(cnSet, i = autosome.index, j = sample.index)
R> cb <- CB(cnSet, i = autosome.index, j = sample.index)
R> missing.ca <- is.na(ca)
R> missing.cb <- is.na(cb)
R> (nmissing.ca <- sum(missing.ca))

[1] 0
```

```
R> (nmissing.cb <- sum(missing.cb))
```

```
[1] 0
```

```
R> identical(nmissing.ca, nmissing.cb)
```

```
[1] TRUE
```

If `nmissing.ca` is nonzero, check the genotype confidence scores provided by the `crmm` genotyping algorithm against the threshold specified in `crmmCopynumber`.

```
R> if (nmissing.ca > 0) {
  invisible(open(snpCallProbability(cnSet)))
  gt.conf <- i2p(snpCallProbability(cnSet)[autosome.index,
    sample.index])
  invisible(close(snpCallProbability(cnSet)))
  below.thr <- gt.conf < GT.CONF.THR
  index.allbelow <- as.integer(which(rowSums(below.thr) ==
    length(sample.index)))
  nmissingBecauseOfGtThr <- length(index.allbelow) *
    length(sample.index)
  stopifnot(identical(nmissingBecauseOfGtThr, nmissing.ca))
  length(index.allbelow) * length(sample.index)/nmissing.ca
}
```

One could inspect the cluster plots for the 'low confidence' calls.

We repeat the above check for missing values at polymorphic loci on chromosome X. In this case, we compare the `rowSums` of the missing values to the number of samples to check whether all of the estimates are missing for a given SNP.

```
R> X.index <- which(isSnp(cnSet) & chromosome(cnSet) ==
  23)
```

```
R> ca.X <- CA(cnSet, i = X.index, j = sample.index)
```

```
R> missing.caX <- is.na(ca.X)
```

```
R> (nmissing.caX <- sum(missing.caX))
```

```
[1] 21060
```

```
R> missing.snp.index <- which(rowSums(missing.caX) == length(sample.index))
```

```
R> index <- which(rowSums(missing.caX) == length(sample.index))
```

```
R> length(index) * length(sample.index)/nmissing.caX
```

```
[1] 1
```

From the above codechunk, we see that 234 SNPs have NAs for all the samples. Next, we tally the number of NAs for polymorphic markers on chromosome X that are below the confidence threshold. For the HapMap analysis, all of the missing values arose from SNPs in which either the men or the women had confidence scores that were all below the threshold.

```
R> if (nmissing.caX > 0) {
  invisible(open(snpCallProbability(cnSet)))
  gt.conf <- i2p(snpCallProbability(cnSet)[X.index,
    sample.index])
  invisible(close(snpCallProbability(cnSet)))
  below.thr <- gt.conf < GT.CONF.THR
  F <- which(cnSet$gender[sample.index] == 2)
  M <- which(cnSet$gender[sample.index] == 1)
```

```

    index.allbelowF <- as.integer(which(rowSums(below.thr[,
      F]) == length(F)))
    index.allbelowM <- as.integer(which(rowSums(below.thr[,
      M]) == length(M)))
    index.allbelow <- as.integer(unique(c(index.allbelowF,
      index.allbelowM)))
    sum(index.allbelow %in% index)/length(index)
  }
[1] 0.9957265

```

For nonpolymorphic loci, the genotype confidence scores are irrelevant and estimates are available at most markers.

```

R> np.index <- which(!isSnp(cnSet) & chromosome(cnSet) ==
  23)
R> ca.F <- CA(cnSet, i = np.index, j = F)
R> ca.M <- CA(cnSet, i = np.index, j = M)
R> ca.F <- ca.F[-match("CN_974939", rownames(ca.F)), ]
R> ca.M <- ca.M[-match("CN_974939", rownames(ca.M)), ]
R> sum(is.na(ca.F))
[1] 0
R> sum(is.na(ca.M))
[1] 0

```

In total, there were 234 polymorphic markers on chromosome X for which copy number estimates are not available. Lowering the confidence threshold would permit estimation of copy number at most of these loci. A confidence threshold is included as a parameter for the copy number estimation as an approach to reduce the sensitivity of genotype-specific summary statistics, such as the within-genotype median, to intensities from samples that do not clearly fall into one of the biallelic genotype clusters. There are drawbacks to this approach, including variance estimates that can be a bit optimistic at some loci. More direct approaches for outlier detection and removal may be explored in the future.

Copy number estimates for other chromosomes, such as mitochondrial and chromosome Y, are not currently available in crlmm.

## 5 Session information

```

R> toLatex(sessionInfo())

```

- R version 2.14.0 Under development (unstable) (2011-03-29 r55147), x86\_64-unknown-linux-gnu
- Locale: LC\_CTYPE=en\_US.iso885915, LC\_NUMERIC=C, LC\_TIME=en\_US.iso885915, LC\_COLLATE=en\_US.iso885915, LC\_MONETARY=C, LC\_MESSAGES=en\_US.iso885915, LC\_PAPER=en\_US.iso885915, LC\_NAME=C, LC\_ADDRESS=C, LC\_TELEPHONE=C, LC\_MEASUREMENT=en\_US.iso885915, LC\_IDENTIFICATION=C
- Base packages: base, datasets, graphics, grDevices, methods, stats, tools, utils
- Other packages: Biobase 2.11.9, bit 1.1-6, cacheSweave 0.4-5, crlmm 1.9.22, DNACopy 1.25.0, ff 2.2-1, filehash 2.1-1, IRanges 1.9.27, oligoClasses 1.13.20, stashR 0.3-3, VanillaICE 1.13.4
- Loaded via a namespace (and not attached): affyio 1.19.2, annotate 1.29.3, AnnotationDbi 1.13.18, Biostrings 2.19.12, DBI 0.2-5, digest 0.4.2, ellipse 0.3-5, genefilter 1.33.1, mvtnorm 0.9-96, preprocessCore 1.13.6, RSQLite 0.9-4, splines 2.14.0, survival 2.36-5, xtable 1.5-6