

An Introduction to the Oligo Package

Benilton Carvalho

March, 2007

1 Introduction

The `oligo` package is designed to support all microarray designs provided by Affymetrix and NimbleGen: expression, tiling, SNP and exon arrays. As of now, chip-specific packages are built via `makePlatformDesign` and transitioning to the `pdInfoBuilder` package, which creates the data packages for the Affymetrix SNP arrays.

2 Analyzing Affymetrix SNP Arrays

Genotyping can be performed using `oligo` and you will need:

- `oligo` and its dependencies;
- Chip specific data package, eg. `pd.mapping50k.xba240`: package that contains the array specifications and SNP annotation.
- CEL files.

We will start by loading the `oligo` package and importing the CEL files available on the `hapmap100kxba` package.

```
R> library("oligo")
R> library("hapmap100kxba")
R> pathCelFiles <- system.file("celFiles", package = "hapmap100kxba")
R> fullFilenames <- list.celfiles(path = pathCelFiles,
  full.names = TRUE)
```

The density of the SNP arrays increased in a way that it is not interesting to store the intensity matrix in memory; efficient methods to handle this situations have been developed and batches of SNPs are analyzed one at a time.

The genotyping algorithm implemented in `oligo` is described in [1]. The whole process can be described in three steps:

1. Normalization against a reference distribution;

2. Summarization via SNPRMA;
3. Genotype calling via CRLMM.

The normalization step is done by equalizing the observed intensities according to a reference distribution, based on the 270 Hapmap samples. These samples are available to the public at <http://www.hapmap.org>. The normalization step will remove systematic biases, which are not due to biological factors.

The SNPRMA algorithm is responsible for summarizing the data. The design of the SNP arrays is such that up to the 250K density, there are probes for both alleles on both strands. On the SNP 6.0 platform, given a SNP, there are probes only on one strand.

Therefore, for the designs up to 250K, SNPRMA will create summaries at the SNP-Allele-Strand-level. For each SNP there are four numbers ($\theta_{A-}, \theta_{B-}, \theta_{A+}, \theta_{B+}$), which are proportional to the log-intensities in each of these combinations of allele and strand ($-$: antisense; $+$: sense). They are represented by four matrices: **antisenseThetaA**, **antisenseThetaB**, **senseThetaA** and **senseThetaB**, which are the components of the *SnpQSet* object. One can extract these objects using accessors of the same name.

For the SNP 6.0 array, a similar approach is taken, but the summaries are given at the SNP-Allele-level and there will be only (θ_A, θ_B) estimates for each SNP. An object of class *SnpCnvQSet* is returned and contains two matrices: **thetaA** and **thetaB**. Accessors with the same name are provided.

Average intensities and log-ratios are defined as across allele and within strand, ie:

$$A_s = \frac{\theta_{A,s} + \theta_{B,s}}{2} \quad (1)$$

$$M_s = \theta_{A,s} - \theta_{B,s}, \quad (2)$$

where s defines the strand (antisense or sense). These quantities can be obtained via **getA()** and **getM** methods, which return high-dimensional arrays with dimensions corresponding to SNP's, samples and strands, respectively. These measures are later used for genotyping.

The CRLMM algorithm can be applied on a *SnpQSet* or *SnpCnvQSet* object in order to produce genotype calls. It involves running a mixture of regressions via EM algorithm to adjust for average intensity and fragment length in the log-ratio scale. These adjustments may take long time to run, depending on the combination of number of samples and computer resources available.

A word of warning: the **crlmm()** method searches for a variable **gender** in the *phenoData* slot of the *SnpQSet* object. If it fails to find that variable, it will try estimate the gender from the data. If there is not enough discrimination power to estimate the gender, the following error message will be returned:

empty cluster: try a better set of initial centers

Increasing the sample size is one of the possible solutions, although the preferred one is to have **gender** already defined in the *phenoData* slot.

The *phenoData* slot includes covariates about the samples. Genotyping and copy number analyses often make use of gender information in order to provide more precise inferences. The code below exemplifies the creation of the *phenoData* object.

```
R> aboutSamples <- data.frame(gender = c("female",
    "female", "male"))
R> rownames(aboutSamples) <- basename(fullFileNames)
R> aboutVars <- data.frame(labelDescription = "male/female")
R> rownames(aboutVars) <- "gender"
R> pd <- new("AnnotatedDataFrame", data = aboutSamples,
    varMetadata = aboutVars)

R> crlmmOut <- justCRLMM(fullFileNames, phenoData = pd,
    verbose = FALSE)

    user  system elapsed
0.139   0.244   0.383
```

The *crlmmOut* object above belongs to the *SnpCallSet* class and contains the genotype calls and confidence measures associated to the calls, represented respectively by the *calls* and *callsConfidence* matrices. These matrices can be accessed using the methods of the same name as demonstrated below:

```
R> calls(crlmmOut)[1:5, 1:2]

                NA06985.CEL NA06991.CEL
SNP_A-1507972           3           3
SNP_A-1510136           3           3
SNP_A-1511055           3           3
SNP_A-1518245           2           3
SNP_A-1641749           3           3

R> callsConfidence(crlmmOut)[1:5, 1:2]

                NA06985.CEL NA06991.CEL
SNP_A-1507972  0.9999254  0.9998893
SNP_A-1510136  0.9999254  0.9999254
SNP_A-1511055  0.9999254  0.9999254
SNP_A-1518245  0.9997851  0.9999254
SNP_A-1641749  0.9997838  0.9997551
```

The genotype calls are represented by 1 (AA), 2 (AB) and 3 (BB). The confidence is the log-likelihood ratio of the two most likely calls.

2.1 Exploring the Annotation Package

The user who is willing to make deeper investigation using the annotations provided for each SNP array can use SQL queries to access more other information that might not be directly exposed.

The example below demonstrates how to see the available tables, fields and extract chromosome, physical location and cytoband for the first five SNPs (probes querying specific SNPs have names starting with the string “SNP”).

```
R> conn <- db(crlmmOut)
R> dbListTables(conn)
```

```
[1] "featureSet"    "mmfeature"     "pm_mm"
[4] "pmfeature"     "qcmmfeature"   "qcpm_qcmm"
[7] "qcpmfeature"   "sequence"      "sqlite_stat1"
[10] "table_info"
```

```
R> dbListFields(conn, "featureSet")
```

```
[1] "fsetid"         "man_fsetid"     "affy_snp_id"
[4] "dbsnps_rs_id"   "chrom"          "physical_pos"
[7] "strand"        "cytoband"       "allele_a"
[10] "allele_b"       "gene_assoc"     "fragment_length"
```

```
R> sql <- "SELECT man_fsetid, chrom, physical_pos FROM featureSet WHERE man_fsetid LIKE 'SNP'"
R> dbGetQuery(conn, sql)
```

```
      man_fsetid chrom physical_pos
1 SNP_A-1650338    2    168433267
2 SNP_A-1716667   19    40749462
3 SNP_A-1712945   19    53411226
4 SNP_A-1711654   21    31501701
5 SNP_A-1717655    1    15312743
```

References

- [1] Benilton Carvalho, Henrik Bengtsson, Terence P Speed, and Rafael A Irizarry. Exploration, normalization, and genotype calls of high density oligonucleotide snp array data. *Biostatistics*, Dec 2006.

3 Details

This document was written using:

```
R> sessionInfo()
```

```
R version 2.7.0 Under development (unstable) (2007-10-10 r43139)
x86_64-unknown-linux-gnu
```

```
locale:
```

```
LC_CTYPE=en_US.UTF-8;LC_NUMERIC=C;LC_TIME=en_US.UTF-8;LC_COLLATE=en_US.UTF-8;LC_MONETARY=en_
```

```
attached base packages:
```

```
[1] splines    tools      stats      graphics  grDevices
[6] utils      datasets  methods    base
```

```
other attached packages:
```

```
[1] pd.mapping50k.xba240_0.3.4 hapmap100kxba_1.3
[3] oligo_1.3.10                oligoClasses_1.1.10
[5] affxparser_1.11.3           AnnotationDbi_1.1.8
[7] preprocessCore_1.1.5        RSQLite_0.6-4
[9] DBI_0.2-4                   Biobase_1.17.8
```