

Genotyping with the crlmm Package

Benilton Carvalho

March, 2009

1 Quick intro to crlmm

The `crlmm` package contains a new implementation for the CRLMM algorithm (Carvalho et. al. 2007). Our focus is on efficient genotyping of SNP 5.0 and 6.0 Affymetrix arrays, although extensions of the method are under development for similar platforms.

This implementation, compared to the previous one (in `oligo`), offers improved confidence scores, quality scores for SNP's and batches, higher accuracy on different datasets and better performance.

Additionally, this package does not use the `pd.genomewidesnp` packages created via `pdInfoBuilder` for `oligo`. Instead, it uses different annotation packages (`genomewidesnp.5` and `genomewidesnp.6`), which use simple R objects to store only the information needed for genotyping. This allowed us to improve the speed of the method, as SQL queries are no longer performed here.

It is also our priority to make the package simple to use. Below we demonstrate how to get genotype calls with the 'new' CRLMM. We use 3 samples on SNP 5.0 made available via the `hapmapsnp5` package.

```
R> require(oligoClasses)
R> library(crlmm)
R> library(hapmapsnp6)
R> path <- system.file("celFiles", package="hapmapsnp6")
R> celFiles <- list.celfiles(path, full.names=TRUE)
R> system.time(crlmmResult <- crlmm(celFiles, verbose=FALSE))
```

```
   user  system elapsed
66.64    2.37    78.33
```

The `crlmmResult` is a *SnpSet* (see Biobase) object.

- **calls:** genotype calls (1 - AA; 2 - AB; 3 - BB);
- **confs:** confidence scores, which can be translated to probabilities by using:

$$1 - 2^{-(\text{confs}/1000)},$$

although we prefer this representation as it saves a significant amount of memory;

- SNPQC: SNP quality score;
- SNR: Signal-to-noise ratio.

```
R> calls(crlmmResult)[1:10,]
```

| | NA06985_GW6_C.CEL | NA06991_GW6_C.CEL |
|---------------|-------------------|-------------------|
| SNP_A-2131660 | 2 | 2 |
| SNP_A-1967418 | 3 | 3 |
| SNP_A-1969580 | 3 | 3 |
| SNP_A-4263484 | 2 | 1 |
| SNP_A-1978185 | 1 | 1 |
| SNP_A-4264431 | 1 | 1 |
| SNP_A-1980898 | 3 | 3 |
| SNP_A-1983139 | 1 | 1 |
| SNP_A-4265735 | 2 | 2 |
| SNP_A-1995832 | 2 | 3 |
| | NA06993_GW6_C.CEL | |
| SNP_A-2131660 | 3 | |
| SNP_A-1967418 | 3 | |
| SNP_A-1969580 | 3 | |
| SNP_A-4263484 | 1 | |
| SNP_A-1978185 | 1 | |
| SNP_A-4264431 | 1 | |
| SNP_A-1980898 | 3 | |
| SNP_A-1983139 | 1 | |
| SNP_A-4265735 | 1 | |
| SNP_A-1995832 | 3 | |

```
R> confs(crlmmResult)[1:10,]
```

| | NA06985_GW6_C.CEL | NA06991_GW6_C.CEL |
|---------------|-------------------|-------------------|
| SNP_A-2131660 | 0.9999963 | 0.9999996 |
| SNP_A-1967418 | 0.9999969 | 0.9999997 |
| SNP_A-1969580 | 0.9995187 | 0.9995139 |
| SNP_A-4263484 | 0.9999999 | 1.0000000 |
| SNP_A-1978185 | 1.0000000 | 1.0000000 |
| SNP_A-4264431 | 1.0000000 | 1.0000000 |
| SNP_A-1980898 | 0.9995192 | 0.9995206 |
| SNP_A-1983139 | 1.0000000 | 0.9999878 |
| SNP_A-4265735 | 0.9999827 | 0.9999863 |
| SNP_A-1995832 | 0.9999762 | 1.0000000 |
| | NA06993_GW6_C.CEL | |
| SNP_A-2131660 | 0.9999998 | |
| SNP_A-1967418 | 0.9999969 | |
| SNP_A-1969580 | 0.9995124 | |
| SNP_A-4263484 | 1.0000000 | |

```

SNP_A-1978185      1.0000000
SNP_A-4264431      1.0000000
SNP_A-1980898      0.9995134
SNP_A-1983139      1.0000000
SNP_A-4265735      0.9999998
SNP_A-1995832      0.9999999

```

```
R> crlmmResult[["SNR"]]
```

```
[1] 8.481305 8.446096 7.379559
```

2 Details

This document was written using:

```
R> sessionInfo()
```

```

R version 4.4.1 (2024-06-14 ucrt)
Platform: x86_64-w64-mingw32/x64
Running under: Windows Server 2022 x64 (build 20348)

```

Matrix products: default

locale:

```

[1] LC_COLLATE=C
[2] LC_CTYPE=English_United States.utf8
[3] LC_MONETARY=English_United States.utf8
[4] LC_NUMERIC=C
[5] LC_TIME=English_United States.utf8

```

time zone: America/New_York

tzcode source: internal

attached base packages:

```

[1] stats      graphics  grDevices  utils      datasets
[6] methods    base

```

other attached packages:

```

[1] genomewidesnp6Crlmm_1.0.7 hapmapsnp6_1.47.0
[3] crlmm_1.63.0                preprocessCore_1.67.1
[5] oligoClasses_1.67.0

```

loaded via a namespace (and not attached):

```

[1] Matrix_1.7-1      bit_4.5.0
[3] limma_3.61.12     jsonlite_1.8.9

```

| | | |
|------|-----------------------------|----------------------|
| [5] | compiler_4.4.1 | BiocManager_1.30.25 |
| [7] | crayon_1.5.3 | Rcpp_1.0.13 |
| [9] | SummarizedExperiment_1.35.4 | ellipse_0.5.0 |
| [11] | Biobase_2.65.1 | GenomicRanges_1.57.2 |
| [13] | Biostrings_2.73.2 | parallel_4.4.1 |
| [15] | VGAM_1.1-12 | splines_4.4.1 |
| [17] | IRanges_2.39.2 | statmod_1.5.0 |
| [19] | lattice_0.22-6 | R6_2.5.1 |
| [21] | XVector_0.45.0 | S4Arrays_1.5.11 |
| [23] | RcppEigen_0.3.4.0.2 | GenomeInfoDb_1.41.2 |
| [25] | ff_4.5.0 | BiocGenerics_0.51.3 |
| [27] | iterators_1.0.14 | DelayedArray_0.31.14 |
| [29] | MatrixGenerics_1.17.0 | openssl_2.2.2 |
| [31] | GenomeInfoDbData_1.2.13 | DBI_1.2.3 |
| [33] | illuminaio_0.47.0 | affyio_1.75.1 |
| [35] | base64_2.0.2 | SparseArray_1.5.45 |
| [37] | zlibbioc_1.51.2 | foreach_1.5.2 |
| [39] | grid_4.4.1 | mvtnorm_1.3-1 |
| [41] | askpass_1.2.1 | beanplot_1.3.1 |
| [43] | S4Vectors_0.43.2 | codetools_0.2-20 |
| [45] | abind_1.4-8 | stats4_4.4.1 |
| [47] | httr_1.4.7 | matrixStats_1.4.1 |
| [49] | tools_4.4.1 | UCSC.utils_1.1.0 |