

ceu1kg: resources for exploring the 1000 genomes data on individuals of central European ancestry in Bioconductor

VJ Carey

May 7, 2020

1 Introduction

Using results of next generation sequencing experiments, a consortium of geneticists produced calls for SNP at approximately 8 million loci of the genomes of individuals of central European ancestry.

Full genotype calls are held in a folder of SnpMatrix instances:

```
> library(ceu1kg)
> dir(system.file("parts", package="ceu1kg"))
[1] "chr1.rda"   "chr10.rda"  "chr11.rda"  "chr12.rda"  "chr13.rda"  "chr14.rda"
[7] "chr15.rda"  "chr16.rda"  "chr17.rda"  "chr18.rda"  "chr19.rda"  "chr2.rda"
[13] "chr20.rda" "chr21.rda" "chr22.rda" "chr3.rda"   "chr4.rda"   "chr5.rda"
[19] "chr6.rda"  "chr7.rda"  "chr8.rda"  "chr9.rda"

> lk = load(dir(system.file("parts", package="ceu1kg"), full=TRUE)[1])
> c1gt = get(lk)
> c1gt

A SnpMatrix with 60 rows and 605756 columns
Row names: NA06985 ... NA12874
Col names: chr1:533 ... chr1:247196267
```

Metadata about the loci are provided in GRanges instances available from SNPLocs packages. Here we consider the 2010 November release.

```
> library(SNPLocs.Hsapiens.dbSNP.20101109)
> if (!exists("c1loc")) c1loc = getSNPLocs("ch1", as.GRanges=TRUE)
> c1loc
```

```

GRanges object with 1849438 ranges and 2 metadata columns:
      seqnames      ranges strand |  RefSNP_id alleles_as_ambig
      <Rle>    <IRanges>  <Rle> | <character>    <character>
[1]   ch1     10327    * |  112750067          Y
[2]   ch1     10440    * |  112155239          M
[3]   ch1     10469    * |  117577454          S
[4]   ch1     10492    * |  55998931           Y
[5]   ch1     10519    * |  62636508           S
...
[1849434]   ch1  249232732    * |  80129254          R
[1849435]   ch1  249232742    * |  28850958          S
[1849436]   ch1  249232749    * |  77296965          R
[1849437]   ch1  249232757    * |  28782254          Y
[1849438]   ch1  249232758    * |  28837504          R
-----
seqinfo: 25 sequences from an unspecified genome; no seqlengths

> rsn1 = paste("rs", elementMetadata(c1loc)$RefSNP_id, sep="")
> length(intersect(rsn1, colnames(c1gt)))

[1] 401489

> ext1 = grep("chr", colnames(c1gt))
> ext1 = as.numeric(gsub("chr1:", "", colnames(c1gt)[ext1]))
> length(intersect(ext1, start(c1loc)))

[1] 1608

```

The last computation shows that most of the 1KG locations are not in dbSNP.

The Bioconductor *GGdata* package includes HapMap phase II genotypes on 90 CEU individuals in 30 trios, coupled with expression data as distributed at the Sanger GENEVAR project (<ftp://ftp.sanger.ac.uk/pub/genevar/>). The 1KG genotypes are available for 43 of these 90 and the associated genotype plus expression data for these 43 can be acquired using getSS, for any chromosome or set of chromosomes.

```

> c20 = getSS("ceu1kg", "chr20")
> c20

```

The above code throws warning because the genotype data are present for 60 individuals, but only 43 have expression values. To create the same structure without a warning:

```

> data(eset) # assume ceu1kg is first in line, yields ex in global
> c1m = c1gt[sampleNames(ex),]
> c1ss = make_smlSet( ex, list(chr1=c1m) )
> c1ss

```

```
SnpMatrix-based genotype set:  
number of samples: 43  
number of chromosomes present: 1  
annotation: illuminaHumanv1.db  
Expression data dims: 47293 x 43  
Total number of SNP: 605756  
Phenodata: An object of class 'AnnotatedDataFrame'  
  sampleNames: NA06985 NA06994 ... NA12874 (43 total)  
  varLabels: famid persid ... male (7 total)  
  varMetadata: labelDescription
```

2 Session information

```
> sessionInfo()
```

```
R version 4.0.0 (2020-04-24)  
Platform: x86_64-pc-linux-gnu (64-bit)  
Running under: Ubuntu 18.04.4 LTS
```

```
Matrix products: default  
BLAS: /home/biocbuild/bbs-3.11-bioc/R/lib/libRblas.so  
LAPACK: /home/biocbuild/bbs-3.11-bioc/R/lib/libRlapack.so
```

```
locale:  
[1] LC_CTYPE=en_US.UTF-8      LC_NUMERIC=C  
[3] LC_TIME=en_US.UTF-8       LC_COLLATE=C  
[5] LC_MONETARY=en_US.UTF-8    LC_MESSAGES=en_US.UTF-8  
[7] LC_PAPER=en_US.UTF-8      LC_NAME=C  
[9] LC_ADDRESS=C              LC_TELEPHONE=C  
[11] LC_MEASUREMENT=en_US.UTF-8 LC_IDENTIFICATION=C
```

```
attached base packages:  
[1] stats4     parallel   stats      graphics   grDevices utils      datasets  
[8] methods    base
```

```
other attached packages:  
[1] SNPlocs.Hsapiens.dbSNP.20101109_0.99.7  
[2] ceu1kg_0.26.0  
[3] GGtools_5.24.0  
[4] Homo.sapiens_1.3.1  
[5] TxDb.Hsapiens.UCSC.hg19.knownGene_3.2.2  
[6] org.Hs.eg.db_3.11.1
```

```
[7] GO.db_3.11.1
[8] OrganismDbi_1.30.0
[9] GenomicFeatures_1.40.0
[10] GenomicRanges_1.40.0
[11] GenomeInfoDb_1.24.0
[12] AnnotationDbi_1.50.0
[13] IRanges_2.22.1
[14] S4Vectors_0.26.0
[15] Biobase_2.48.0
[16] BiocGenerics_0.34.0
[17] data.table_1.12.8
[18] GGBase_3.50.0
[19] snpStats_1.38.0
[20] Matrix_1.2-18
[21] survival_3.1-12
```

loaded via a namespace (and not attached):

[1] colorspace_1.4-1	ellipsis_0.3.0
[3] biovizBase_1.36.0	htmlTable_1.13.3
[5] XVector_0.28.0	base64enc_0.1-3
[7] dichromat_2.0-0	rstudioapi_0.11
[9] hexbin_1.28.1	bit64_0.9-7
[11] splines_4.0.0	knitr_1.28
[13] Formula_1.2-3	Rsamtools_2.4.0
[15] annotate_1.66.0	cluster_2.1.0
[17] dplyr_1.4.3	png_0.1-7
[19] graph_1.66.0	BiocManager_1.30.10
[21] compiler_4.0.0	httr_1.4.1
[23] backports_1.1.6	assertthat_0.2.1
[25] lazyeval_0.2.2	acepack_1.4.1
[27] htmltools_0.4.0	prettyunits_1.1.1
[29] tools_4.0.0	gttable_0.3.0
[31] glue_1.4.0	GenomeInfoDbData_1.2.3
[33] reshape2_1.4.4	dplyr_0.8.5
[35] rappdirs_0.3.1	Rcpp_1.0.4.6
[37] biglm_0.9-1	vctrs_0.2.4
[39] Biostrings_2.56.0	rtracklayer_1.48.0
[41] iterators_1.0.12	xfun_0.13
[43] stringr_1.4.0	lifecycle_0.2.0
[45] ensemblldb_2.12.1	XML_3.99-0.3
[47] zlibbioc_1.34.0	scales_1.1.0
[49] BSgenome_1.56.0	VariantAnnotation_1.34.0

```
[51] hms_0.5.3                  ProtGenerics_1.20.0
[53] SummarizedExperiment_1.18.1 RBGL_1.64.0
[55] AnnotationFilter_1.12.0    RColorBrewer_1.1-2
[57] curl_4.3                   memoise_1.1.0
[59] gridExtra_2.3              ggplot2_3.3.0
[61] biomaRt_2.44.0             rpart_4.1-15
[63] latticeExtra_0.6-29       stringi_1.4.6
[65] RSQLite_2.2.0              genefilter_1.70.0
[67] checkmate_2.0.0            BiocParallel_1.22.0
[69] rlang_0.4.6                pkgconfig_2.0.3
[71] matrixStats_0.56.0         bitops_1.0-6
[73] lattice_0.20-41            ROCR_1.0-11
[75] purrr_0.3.4                GenomicAlignments_1.24.0
[77] htmlwidgets_1.5.1          bit_1.1-15.2
[79] tidyselect_1.0.0           plyr_1.8.6
[81] magrittr_1.5                R6_2.4.1
[83] Hmisc_4.4-0                DelayedArray_0.14.0
[85] DBI_1.1.0                  pillar_1.4.4
[87] foreign_0.8-79             RCurl_1.98-1.2
[89] nnet_7.3-14                tibble_3.0.1
[91] crayon_1.3.4               BiocFileCache_1.12.0
[93] jpeg_0.1-8.1               progress_1.2.2
[95] grid_4.0.0                  blob_1.2.1
[97] digest_0.6.25              xtable_1.8-4
[99] ff_2.2-14.2                openssl_1.4.1
[101] munsell_0.5.0              Gviz_1.32.0
[103] askpass_1.1
```