

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

VJ Carey

April 26, 2017

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
	<Rle>	<IRanges>	<Rle>	<character>
[1]	ch1	[10327, 10327]	*	112750067
[2]	ch1	[10440, 10440]	*	112155239
[3]	ch1	[10469, 10469]	*	117577454
[4]	ch1	[10492, 10492]	*	55998931
[5]	ch1	[10519, 10519]	*	62636508
...
[1849434]	ch1	[249232732, 249232732]	*	80129254
[1849435]	ch1	[249232742, 249232742]	*	28850958
[1849436]	ch1	[249232749, 249232749]	*	77296965
[1849437]	ch1	[249232757, 249232757]	*	28782254
[1849438]	ch1	[249232758, 249232758]	*	28837504

alleles_as_ambig
<character>

[1]	Y
[2]	M
[3]	S
[4]	Y
[5]	S
...	...
[1849434]	R
[1849435]	S
[1849436]	R
[1849437]	Y
[1849438]	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("ceukg", "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 ceukg 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 3.4.0 (2017-04-21)
Platform: x86_64-pc-linux-gnu (64-bit)
Running under: Ubuntu 16.04.2 LTS
```

```
Matrix products: default
BLAS: /home/biocbuild/bbs-3.5-bioc/R/lib/libRblas.so
LAPACK: /home/biocbuild/bbs-3.5-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_0.14.0
[3] GGtools_5.12.0
[4] Homo.sapiens_1.3.1
[5] TxDb.Hsapiens.UCSC.hg19.knownGene_3.2.2
[6] org.Hs.eg.db_3.4.1
[7] GO.db_3.4.1
[8] OrganismDbi_1.18.0
[9] GenomicFeatures_1.28.0
[10] GenomicRanges_1.28.0
[11] GenomeInfoDb_1.12.0
[12] AnnotationDbi_1.38.0
[13] IRanges_2.10.0
[14] S4Vectors_0.14.0
[15] Biobase_2.36.0
[16] BiocGenerics_0.22.0
[17] data.table_1.10.4
[18] GGBase_3.38.0
[19] snpStats_1.26.0
[20] Matrix_1.2-9
[21] survival_2.41-3
```

loaded via a namespace (and not attached):

```
[1] ProtGenerics_1.8.0 bitops_1.0-6
[3] matrixStats_0.52.2 RColorBrewer_1.1-2
[5] httr_1.2.1 tools_3.4.0
[7] backports_1.0.5 R6_2.2.0
[9] KernSmooth_2.23-15 rpart_4.1-11
[11] Hmisc_4.0-2 DBI_0.6-1
[13] lazyeval_0.2.0 Gviz_1.20.0
[15] colorspace_1.3-2 nnet_7.3-12
[17] gridExtra_2.2.1 bit_1.1-12
[19] compiler_3.4.0 graph_1.54.0
[21] htmlTable_1.9 biglm_0.9-1
```

[23] DelayedArray_0.2.0	rtracklayer_1.36.0
[25] caTools_1.17.1	scales_0.4.1
[27] checkmate_1.8.2	hexbin_1.27.1
[29] genefilter_1.58.0	RBGL_1.52.0
[31] stringr_1.2.0	digest_0.6.12
[33] Rsamtools_1.28.0	foreign_0.8-68
[35] XVector_0.16.0	base64enc_0.1-3
[37] dichromat_2.0-0	htmltools_0.3.5
[39] ensemblDb_2.0.0	BSgenome_1.44.0
[41] htmlwidgets_0.8	RSQLite_1.1-2
[43] BiocInstaller_1.26.0	shiny_1.0.2
[45] gtools_3.5.0	BiocParallel_1.10.0
[47] acepack_1.4.1	VariantAnnotation_1.22.0
[49] RCurl_1.95-4.8	magrittr_1.5
[51] GenomeInfoDbData_0.99.0	Formula_1.2-1
[53] Rcpp_0.12.10	munsell_0.4.3
[55] stringi_1.1.5	yaml_2.1.14
[57] SummarizedExperiment_1.6.0	zlibbioc_1.22.0
[59] ggplots_3.0.1	plyr_1.8.4
[61] AnnotationHub_2.8.0	grid_3.4.0
[63] gdata_2.17.0	lattice_0.20-35
[65] Biostrings_2.44.0	splines_3.4.0
[67] annotate_1.54.0	knitr_1.15.1
[69] reshape2_1.4.2	biomaRt_2.32.0
[71] XML_3.98-1.6	biovizBase_1.24.0
[73] latticeExtra_0.6-28	httpuv_1.3.3
[75] gtable_0.2.0	ggplot2_2.2.1
[77] mime_0.5	xtable_1.8-2
[79] AnnotationFilter_1.0.0	ff_2.2-13
[81] tibble_1.3.0	iterators_1.0.8
[83] GenomicAlignments_1.12.0	memoise_1.1.0
[85] cluster_2.0.6	interactiveDisplayBase_1.14.0
[87] ROCR_1.0-7	