

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

VJ Carey

May 7, 2016

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 3.3.0 (2016-05-03)

Platform: x86_64-pc-linux-gnu (64-bit)

Running under: Ubuntu 14.04.4 LTS

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_0.10.0
[3] GGtools_5.8.0
[4] Homo.sapiens_1.3.1
[5] TxDb.Hsapiens.UCSC.hg19.knownGene_3.2.2
[6] org.Hs.eg.db_3.3.0
[7] GO.db_3.3.0
[8] OrganismDbi_1.14.0
[9] GenomicFeatures_1.24.0
[10] GenomicRanges_1.24.0
[11] GenomeInfoDb_1.8.0
[12] AnnotationDbi_1.34.0
[13] IRanges_2.6.0
[14] S4Vectors_0.10.0
[15] Biobase_2.32.0
[16] BiocGenerics_0.18.0
[17] data.table_1.9.6
[18] GGBase_3.34.0
[19] snpStats_1.22.0
[20] Matrix_1.2-6
[21] survival_2.39-2
```

loaded via a namespace (and not attached):

```
[1] httr_1.1.0                AnnotationHub_2.4.0
[3] splines_3.3.0             gtools_3.5.0
[5] Formula_1.2-1             shiny_0.13.2
[7] interactiveDisplayBase_1.10.0 latticeExtra_0.6-28
[9] RBGL_1.48.0               BSgenome_1.40.0
[11] Rsamtools_1.24.0         RSQLite_1.0.0
[13] lattice_0.20-33          biovizBase_1.20.0
[15] chron_2.3-47             digest_0.6.9
[17] RColorBrewer_1.1-2       XVector_0.12.0
[19] colorspace_1.2-6         htmltools_0.3.5
[21] httpuv_1.3.3             plyr_1.8.3
[23] XML_3.98-1.4             biglm_0.9-1
[25] biomaRt_2.28.0           genefilter_1.54.0
[27] zlibbioc_1.18.0         xtable_1.8-2
[29] scales_0.4.0             gdata_2.17.0
```

[31] ff_2.2-13	BiocParallel_1.6.0
[33] annotate_1.50.0	ggplot2_2.1.0
[35] SummarizedExperiment_1.2.0	ROCR_1.0-7
[37] hexbin_1.27.1	nnet_7.3-12
[39] Gviz_1.16.0	mime_0.4
[41] magrittr_1.5	gplots_3.0.1
[43] foreign_0.8-66	graph_1.50.0
[45] BiocInstaller_1.22.1	tools_3.3.0
[47] matrixStats_0.50.2	stringr_1.0.0
[49] munsell_0.4.3	cluster_2.0.4
[51] ensemblDb_1.4.0	Biostrings_2.40.0
[53] caTools_1.17.1	grid_3.3.0
[55] RCurl_1.95-4.8	iterators_1.0.8
[57] dichromat_2.0-0	VariantAnnotation_1.18.0
[59] bitops_1.0-6	gtable_0.2.0
[61] DBI_0.4	reshape2_1.4.1
[63] R6_2.1.2	GenomicAlignments_1.8.0
[65] gridExtra_2.2.1	rtracklayer_1.32.0
[67] bit_1.1-12	Hmisc_3.17-4
[69] KernSmooth_2.23-15	stringi_1.0-1
[71] Rcpp_0.12.4.5	rpart_4.1-10
[73] acepack_1.3-3.3	