

Package Ravages (RAre Variant Analysis and GENetic Simulation)

Herve Perdry and Ozvan Bocher

2019-09-17

```
library("knitr")
require("Ravages")
```

Introduction

Ravages was developped to simulate genetic data and to perform rare variant association tests (burden tests and the variance-component test SKAT) on more than two groups of individuals (Bocher et al., 2019, Genetic Epidemiology). Ravages relies on the package Gaston developped by Herve Perdry and Claire Dandine-Roulland. Most functions are written in C++ thanks to the packages Rcpp, RcppParallel and RcppEigen.

Functions of Ravages use bed.matrix to manipulate genetic data as in the package Gaston (see documentation of this package for more details).

In this vignette, we illustrate how to perform rare variant association tests on real data. A second vignette is available showing how to simulate genetic data and how to use it to calculate power of the tests. To learn more about all options of the functions, the reader is advised to look at the manual pages.

Example of analysis using LCT data

Below is an example of an association analysis and previous steps of data filtering using the dataset LCT available with the package Ravages. This dataset contains data from the 1000Genome project in the region containing the Lactase gene. In this example, we look for an association between rare variants and the populations of 1000Genomes EUR. The population of each individual is available in the dataframe LCT.matrix.pop1000G. Details about each function is given right after this example.

```
#Importation of data in a bed matrix
x <- as.bed.matrix(x=LCT.matrix.bed, fam=LCT.matrix.fam,
                  bim=LCT.snps)

#Add population
x@ped[,c("pop", "superpop")] <- LCT.matrix.pop1000G[,c("population", "super.population")]

#Select EUR superpopulation
x <- select.inds(x, superpop=="EUR")
x@ped$pop <- droplevels(x@ped$pop)

# Group variants within know genes by extending their positions
# 500bp upstream and downstream
x <- set.genomic.region(x, flank.width=500)
table(x@snps$genomic.region, useNA = "ifany")
```

```
##
## R3HDM1  UBXN4    LCT    MCM6    DARS    <NA>
##   2047   1207   1454   1149    924   1295
```

```

# Group variants within known genes using their exact positions
x <- set.genomic.region(x, flank.width=0)
table(x@snp$genomic.region, useNA = "ifany")

##
## R3HDM1  UBKN4    LCT    MCM6    DARS    <NA>
##    2038    1175    1429    1123    913    1398

# Filter variants with maf (computed on whole sample) < 0.01
# keeping only genomic region with at least 10 SNPs
x1 <- filter.rare.variants(x, filter = "whole", maf.threshold = 0.01, min.nb.snps = 10)
table(x1@snp$genomic.region, useNA="ifany")

##
## R3HDM1  UBKN4    LCT    MCM6    DARS
##    267    167    202    159    134

# run burden test CAST, using the 1000Genome population as "outcome"
burden.mlogit(x1, group=x1@ped$pop, burden = "CAST", ref.level = "CEU")

##
##          p.value is.err
## R3HDM1 1.300274e-04      0
## UBKN4  4.096613e-05      0
## LCT    3.810119e-09      0
## MCM6   1.202259e-07      0
## DARS   2.036275e-03      0

# run SKAT, using the 1000Genome population as "outcome"
SKAT(x1, group=x1@ped$pop)

##
##          stat          p.perm          p.chi2          p.value
## R3HDM1 6.438772 0.0000399992 5.538767e-06 5.538767e-06
## UBKN4  2.101316 0.0069230242 7.721944e-03 6.923024e-03
## LCT    3.483270 0.0007799844 8.434574e-04 8.434574e-04
## MCM6   2.873867 0.0074259246 6.492139e-03 7.425925e-03
## DARS   2.093018 0.0708274895 1.025923e-01 7.082749e-02

```

Defining genomic regions

For rare variant association tests, the unit of analysis is not a single variant but a genomic region, typically a gene. The first step of the analysis is therefore to group variants into genomic regions. This can be done using the function `set.genomic.region()` and known gene positions. It works on a `bed.matrix` (see Gaston) and simply adds a column “genomic.region” to the slot `x@snp` containing the gene assigned to each variant. By default, any variant being outside the gene positions won’t be annotated. Gene positions can be extended to annotate more variants using the argument `flank.width` corresponding to the number of base pair upstream and downstream the gene. If `flank.width=Inf`, each variant will be assigned to the nearest gene. If two genes overlap, variants in the overlapping zone will be attributed to the second one (in the order given by the position of their starting point on the genome).

The files **genes.b37** and **genes.b38** available in Ravages which contain gene positions from ENSEMBL versions hg19 and hg38 can be used to define gene positions.

Rare variant definition

To perform rare variant analysis, it is also important to define what is a rare variant in order to leave out common ones. The function `filter.rare.variants()` enables to keep only variants with a MAF (Minor Allele Frequency) below a given threshold while leaving out monomorphic variants. This function uses and returns a `bed.matrix` which can be filtered in three different ways:

- If *filter*="whole", all the variants with a MAF lower than the threshold in the entire sample will be kept.
- If *filter*="controls", all the variants with a MAF lower than the threshold in the controls group will be kept. In this situation, the controls group needs to be specified to the argument *ref.level*.
- If *filter*="any", all the variants with a MAF lower than the threshold in any of the groups will be kept.

It is also possible to specify the minimum number of variants needed in a genomic region to keep it using the parameter *min.nb.snps*.

Rare variant association tests

We have implemented two rare variant burden association tests extensions: CAST and WSS. The general idea of burden tests is to compute a genetic score per individual and per genomic region and to test if it differs between the different groups of individuals. To extend these tests to more than two groups of individuals, a non-ordinal multinomial regression is used. The independent variable in this regression is the genetic effect of the gene represented by the genetic score. Potential covariates can be added in the model. In addition to the genetic scores CAST and WSS directly implemented in the package, the user can specify another genetic score for the regression.

Genetic score for burden tests

We have implemented two functions to compute CAST and WSS scores respectively. These functions return a matrix with one row per individual and one column by genomic region. They are directly called in the function `burden.mlogit()` if these scores are used to perform the association tests.

CAST

CAST is based on a binary score which has a value of one if an individual carries at least one variant in the considered genomic region, and 0 otherwise. A MAF threshold for the definition of a rare variant is therefore needed as an argument to *maf.threshold*. This score can be computed using the function `CAST()` as shown here on the LCT data:

```
#Calculation of the genetic score with a maf threshold of 2.5%
CAST.score <- CAST(x = x1, genomic.region = x1@snps$genomic.region, maf.threshold = 0.01)
head(CAST.score)
```

##	R3HDM1	UBXN4	LCT	MCM6	DARS
## HG00096	0	0	1	1	1
## HG00097	1	0	0	0	0
## HG00099	1	0	0	0	0
## HG00100	0	1	0	0	0
## HG00101	0	1	0	0	0
## HG00102	1	0	0	0	1

WSS

WSS (Weighted Sum Statistic) is based on a continuous score giving the highest weights to the rarest variants:

$$WSS_j = \sum_{i=1}^R I_{ij} * w_i$$

with

$$w_i = \frac{1}{\sqrt{(t_i * q_i * 1 - q_i)}}$$

and

$$q_i = \frac{n_i + 1}{2 * t_i + 1}$$

Where n_i is the total number of minor alleles genotyped for variant i , t_i is the total number of alleles genotyped for variant i and I_{ij} is the number of minor alleles of variant i for the individual j . In the original method, each variant is weighted according to its frequency in the controls group. In our version of WSS, the weights depend on allele frequency calculated on the entire sample. The function **WSS()** can be used to compute the WSS score as shown on the LCT data:

```
WSS.score <- WSS(x = x1, genomic.region = x1@snp$genomic.region)
head(WSS.score)
```

```
##          R3HDM1    UBXN4      LCT    MCM6    DARS
## HG00096 0.0000000 0.000000 0.8185268 1.26932 1.418436
## HG00097 0.8185268 0.000000 0.0000000 0.00000 0.000000
## HG00099 1.0019881 0.000000 0.0000000 0.00000 0.000000
## HG00100 0.0000000 1.001988 0.0000000 0.00000 0.000000
## HG00101 0.0000000 1.001988 0.0000000 0.00000 0.000000
## HG00102 1.0019881 0.000000 0.0000000 0.00000 1.001988
```

Regressions

We have extended CAST and WSS using non-ordinal multinomial regression models. Let consider C groups of individuals including a group of controls ($c = 1$) and $C - 1$ groups of cases with different sub-phenotypes of the disease. We can compute $C - 1$ probability ratios, one for each group of cases:

$$\ln \frac{P(Y_j = c)}{P(Y_j = 1)} = \beta_{0,c} + \beta_{G,c} X_G + \beta_{k1,c} K_1 + \dots + \beta_{kl,c} K_l$$

Where Y_j corresponds to the phenotype of the individual j and K_l is a vector for the l th covariate with the corresponding coefficient β_{kl} . The genetic effect is represented by X_G and correspond to the genetic score CAST or WSS with $\beta_{G,c}$ the log-odds ratio associated to this burden score.

The p-value associated to the genetic effect is calculated using a likelihood ratio test comparing this model to the same model without the genetic effect (null hypothesis). If only two groups are compared, a classical logistic regression is performed.

This regression can be performed on a bed.matrix using the function **burden.mlogit()** which relies on the package mlogit. To do so, the user needs to specify a vector with the phenotype of each individual (argument *group*) and the gene associated to each variant (argument *genomic.region*).

The CAST or WSS genetic scores can be directly calculated in the regression (*burden*="CAST" or *burden*="WSS"). The user can also use another genetic score in the regression, which has to be specified as a matrix with one individual per row and one genomic region per column to *burden*. In this situation, no bed matrix is needed. The reference group of individuals should be given to the argument *ref.level*. Potential covariates could also be included in the regression as a matrix with one row per individual and one column per covariate to the argument *data*. If only a subset of covariates from *data* are to be included in the model, a R formula should be given to *formula* with these covariates, otherwise all the covariates will be included. **burden.mlogit()** will return the p-value associated to the regression for each genomic region. If there is a convergence problem with the regression, the function will return 1 in the column *is.err*. The odds ratio

associated to each group of cases compared to the reference group (*ref.level*) with its confidence interval at a given alpha threshold (argument *alpha*) can also be obtained if *get.OR.value=TRUE*.

An example of the p-value and OR calculation with its 95% confidence interval using WSS on the LCT data is shown below with or without the inclusion of covariates. The outcome here corresponds to the population of the individuals from 1000Genome.

```
#WSS
burden.mlogit(x=x1, group=x1@ped$pop, burden="WSS", ref.level="CEU", alpha=0.05,
              get.OR.value=TRUE)

##                p.value is.err  OR.TSI    OR.FIN    OR.GBR    OR.IBS
## R3HDM1 4.876886e-08      0 1.804428 1.0189455 0.8398775 1.313583
## UBXN4 8.341120e-07      0 1.511287 0.4875926 0.5403705 1.033206
## LCT 9.459153e-11      0 2.883109 1.4384475 0.9698497 2.344099
## MCM6 7.339135e-10      0 3.666278 2.1369606 1.1022770 3.006620
## DARS 6.797271e-04      0 1.677731 0.9977584 0.8258296 1.457758
##                1.lower.TSI 1.lower.FIN 1.lower.GBR 1.lower.IBS 1.upper.TSI
## R3HDM1 1.360068 0.7323019 0.5813563 0.9760668 2.393967
## UBXN4 1.042657 0.2882958 0.3212035 0.6919826 2.190546
## LCT 1.881794 0.8967677 0.5715946 1.5278099 4.417231
## MCM6 2.149127 1.2032805 0.5687063 1.7592628 6.254444
## DARS 1.076988 0.5798540 0.4551123 0.9320409 2.613568
##                1.upper.FIN 1.upper.GBR 1.upper.IBS
## R3HDM1 1.4177895 1.2133596 1.767809
## UBXN4 0.8246618 0.9090818 1.542691
## LCT 2.3073213 1.6455868 3.596520
## MCM6 3.7951258 2.1364537 5.138383
## DARS 1.7168491 1.4985194 2.280005

#Simulation of covariates with different probabilities in GBR/CEU/FIN and IBS/TSI
covar <- data.frame( sex = c(sample(0:1, sum(table(LCT.pop)[c("CEU", "GBR", "FIN")])),
                             TRUE, c(0.2,0.8))),
                    sample(0:1, sum(table(LCT.pop)[c("TSI", "IBS")])),
                             TRUE, c(0.8,0.2))),
              u = runif(length(LCT.pop)))

#Regression with the covariate "sex" without OR values
burden.mlogit(x=x1, group=x1@ped$pop, burden="WSS", ref.level="CEU",
              data=covar, formula = ~ sex)

##                p.value is.err
## R3HDM1 1.777007e-08      0
## UBXN4 8.737432e-06      0
## LCT 3.980182e-11      0
## MCM6 1.878952e-10      0
## DARS 1.233003e-04      0

#WSS using directly the score matrix computed previously
burden.mlogit(burden=WSS.score, group=x1@ped$pop, ref.level="CEU")

##                p.value is.err
## R3HDM1 4.876886e-08      0
## UBXN4 8.341120e-07      0
## LCT 9.459153e-11      0
## MCM6 7.339135e-10      0
## DARS 6.797271e-04      0
```

SKAT

We also extended the variance-component test SKAT using a geometric interpretation. Unlike the burden tests, there is no burden calculated in this test: the distribution of the genetic effects in the gene is compared to a null distribution. SKAT is based on a linear mixed model where the random effects correspond to the genetic effects.

Permutations are used to compute the p-values with the arguments *perm.target* and *perm.max*. *perm.target* corresponds to the number of times the observed statistics should be lower than a permuted statistics, while *perm.max* corresponds to the maximum number of permutations to perform. A sequential procedure is used for these permutations: the program stops when any of these two values is reached. Two types of p-values estimations are then performed: if *perm.target* is reached, the p-value is computed as *perm.target* divided by the number of permutations performed to reach this value; if *perm.max* is reached before *perm.target* (that is, for pretty small p-values), the SKAT small sample procedure is used, and p-values are estimated using a chi-square distribution based on statistics moments obtained from the permutations.

In this function, covariates can be included using the argument *Pi*. This argument should be a matrix containing the probabilities that each individual belongs to each group, with one row per individual and one column per group of individuals. It can be directly computed using the function **Pi.matrix()** which has the same arguments *group*, *data*, *formula* and *ref.level* than **burden.mlogit()**. This function uses a regression with the *group* as the dependant variable and the covariates as the independant variables, and returns the adjusted probabilities. *ref.level* won't have any impact on the probability of each individual, but is needed to perform the regression.

An example of this function and how to use it with SKAT is shown below.

```
#Compute the Pi matrix with the covariate sex
Pi.matrix.LCT <- Pi.matrix(group=x1@ped$pop, data=covar, formula= ~ sex, ref.level="CEU")
```

```
#SKAT with the covariates
```

```
SKAT(x1, group=x1@ped$pop, Pi=Pi.matrix.LCT)
```

```
##          stat      p.perm      p.chi2      p.value
## R3HDM1 6.095925 0.0000799984 5.443322e-05 5.443322e-05
## UBXN4  1.892915 0.7214285714 7.177469e-01 7.214286e-01
## LCT    3.427898 0.0308679707 3.458440e-02 3.086797e-02
## MCM6   2.785803 0.0752047655 8.746218e-02 7.520477e-02
## DARS   2.099419 0.2853107345 2.932180e-01 2.853107e-01
```

Data management

Data in the plink format or in the vcf format can be loaded in R using the functions **read.bed.matrix()** and **read.vcf()** respectively from the package *gaston*.

If the data for the controls and the different groups of cases are in different files, they can be loaded separately and then combined using the function **gaston::rbind()** as long as the same variants are present between the different groups of individuals.

An example is given below where the simulated data have been split according to the group of each individual, and then combined in a *bed.matrix*:

```
#Selection of each group of individuals
```

```
CEU <- select.inds(x1, pop=="CEU")
```

```
CEU
```

```
## A bed.matrix with 99 individuals and 929 markers.
```

```
## snps stats are set
```

```
## There are 732 monomorphic SNPs
```

```
## ped stats are set
```

```
FIN <- select.inds(x1, pop=="FIN")
FIN
```

```
## A bed.matrix with 99 individuals and 929 markers.
## snps stats are set
##   There are 755 monomorphic SNPs
## ped stats are set
```

```
GBR <- select.inds(x1, pop=="GBR")
GBR
```

```
## A bed.matrix with 91 individuals and 929 markers.
## snps stats are set
##   There are 778 monomorphic SNPs
## ped stats are set
```

```
#Combine in one file:
```

```
CEU.FIN.GBR <- rbind(CEU, FIN, GBR)
CEU.FIN.GBR
```

```
## A bed.matrix with 289 individuals and 929 markers.
## snps stats are set
##   There are 502 monomorphic SNPs
## ped stats are set
```