

Genetic algorithm using the evola package

Giovanny Covarrubias-Pazaran

2024-08-10

The evola package is nice wrapper of the AlphaSimR package that enables the use of the evolutionary algorithm (EA) to solve complex questions in a simple manner.

The vignettes aim to provide several examples in how to use the evola package under different optimization scenarios. We will spend the rest of the space providing examples for:

- 1) Optimizing the selection of one feature with a constraint in other
- 2) Obtaining a subsample of a population to maximize a feature while constraining the relationship between possibilities 2a) best parents for the next generation 2b) best crosses of the next generation
- 3) Optimizing a subsample of size N to be representative 3a) Of its own 3b) For another sample
- 4) How to specify constraints: 4a) gender 4b) number of times a parent should be used
- 5) How to optimize the number of progeny to produce per cross

Because of CRAN requirements I will only run few generations but please when you run your analysis let it run for many generations.

1) Optimizing the selection of one feature with a constraint in other

The example presented here is a list of gems (Color) that have different weights in Kg (Weight) and a given value (Value).

```
set.seed(1)
# Data
Gems <- data.frame(
  Color = c("Red", "Blue", "Purple", "Orange",
            "Green", "Pink", "White", "Black",
            "Yellow"),
  Weight = round(runif(9,0.5,5),2),
  Value = round(abs(rnorm(9,0,5))+0.5,2)
)
head(Gems)
```

```
##   Color Weight Value
## 1   Red    1.69  8.20
## 2   Blue    2.17  5.14
## 3 Purple    3.08  1.97
## 4 Orange    4.59  0.53
## 5  Green    1.41 12.52
## 6   Pink    4.54  4.32
```

The task to optimize here is to be able to pick in your bag all the possible gems that maximize the value with the constraint that your bag would break after 10Kg. In the evolafit function this would be specified as follows:

```

# Task: Gem selection.
# Aim: Get highest combined value.
# Restriction: Max weight of the gem combined = 10.
res0<-evolafit(cbind(Weight,Value)~Color, dt= Gems,
  # constraints: if greater than this ignore
  constraintsUB = c(10,Inf),
  # constraints: if smaller than this ignore
  constraintsLB= c(-Inf,-Inf),
  # weight the traits for the selection
  traitWeight = c(0,1),
  # population parameters
  nCrosses = 100, nProgeny = 20, recombGens = 1,
  # coancestry parameters
  A=NULL, lambda=c(0,0), nQTLperInd = 1,
  # selection parameters
  propSelBetween = .9, propSelWithin =0.9,
  nGenerations = 30, verbose = FALSE
)

```

```

## Warning in selectFam(pop = pop, nFam = round(nCrosses * propSelBetween), : Suitable
## families smaller than nFam, returning 36 families

```

Noticed that the formula `cbind(Weight,Value)~Color` specified the traits to be considered in the optimization problem are specified in the left side of the formula whereas the right side of the formula specifies the term corresponding to the genes that will form the ‘genome’ of the solution. Each trait in the formula requires a value for the constraints, weights in the fitness function (i.e., selection index) and lambda (weight for group relationship between the genes in the genome). The rest of the parameters are the parameters controlling the evolution of the population of solutions.

When looking at the results of the evolution we can observe that the best solution for the traits under the constraints can be extracted with the `bestSol()` function.

```

best=bestSol(res0)[2]; best # best solution for trait 1

```

```

## Trait2
##      2

```

```

res0$M[best,]

```

```

##      Red   Blue Purple Orange  Green   Pink  White  Black Yellow
##      1     0     0     0     1     0     0     1     1

```

```

xa = res0$M[best,] %*% as.matrix(Gems[,c("Weight", "Value")]); xa

```

```

##      Weight Value
## [1,]    9.9 28.91

```

The best selection of Gems was the one one found in the M element of the resulting object.

2) Obtaining subsample of a population to maximize a feature while constraining the relationship between possibilities

2a) Best parents for the next generation One situation that occurs in plant and animal breeding is the so called ‘optimal contribution’ problem where we want to pick a set of parents that can maximize the gain while managing genetic variance as much as possible. In the following example we take a population of 363 possible parents and pick the best 20 while conserving genetic variance.

```
data(DT_cpdata)
DT <- DT_cpdata
head(DT)
```

```
##      id Row Col Year      color  Yield FruitAver Firmness Rowf Colf occ
## P003 P003  3  1 2014 0.10075269 154.67    41.93  588.917   3   1   0
## P004 P004  4  1 2014 0.13891940 186.77    58.79  640.031   4   1   1
## P005 P005  5  1 2014 0.08681502  80.21    48.16  671.523   5   1   1
## P006 P006  6  1 2014 0.13408561 202.96    48.24  687.172   6   1   1
## P007 P007  7  1 2014 0.13519278 174.74    45.83  601.322   7   1   1
## P008 P008  8  1 2014 0.17406685 194.16    44.63  656.379   8   1   1
```

Our surrogate of fitness will be the Yield trait and we will have a second trait to control the number of individuals we can select. We will set a constraint for the occurrence (occ) trait to 20 but the only trait contributing to fitness will be Yield (traitWeight parameter).

```
# get best 20 individuals weighting variance by 0.5
res<-evolafit(cbind(Yield, occ)~id, dt= DT,
  # constraints: if sum is greater than this ignore
  constraintsUB = c(Inf,20),
  # constraints: if sum is smaller than this ignore
  constraintsLB= c(-Inf,-Inf),
  # weight the traits for the selection
  traitWeight = c(1,0),
  # population parameters
  nCrosses = 100, nProgeny = 10,
  # coancestry parameters
  A=A, lambda=c(0.5,0), nQTLperInd = 2,
  # selection parameters
  propSelBetween = 0.5, propSelWithin =0.5,
  nGenerations = 20, verbose=FALSE)
```

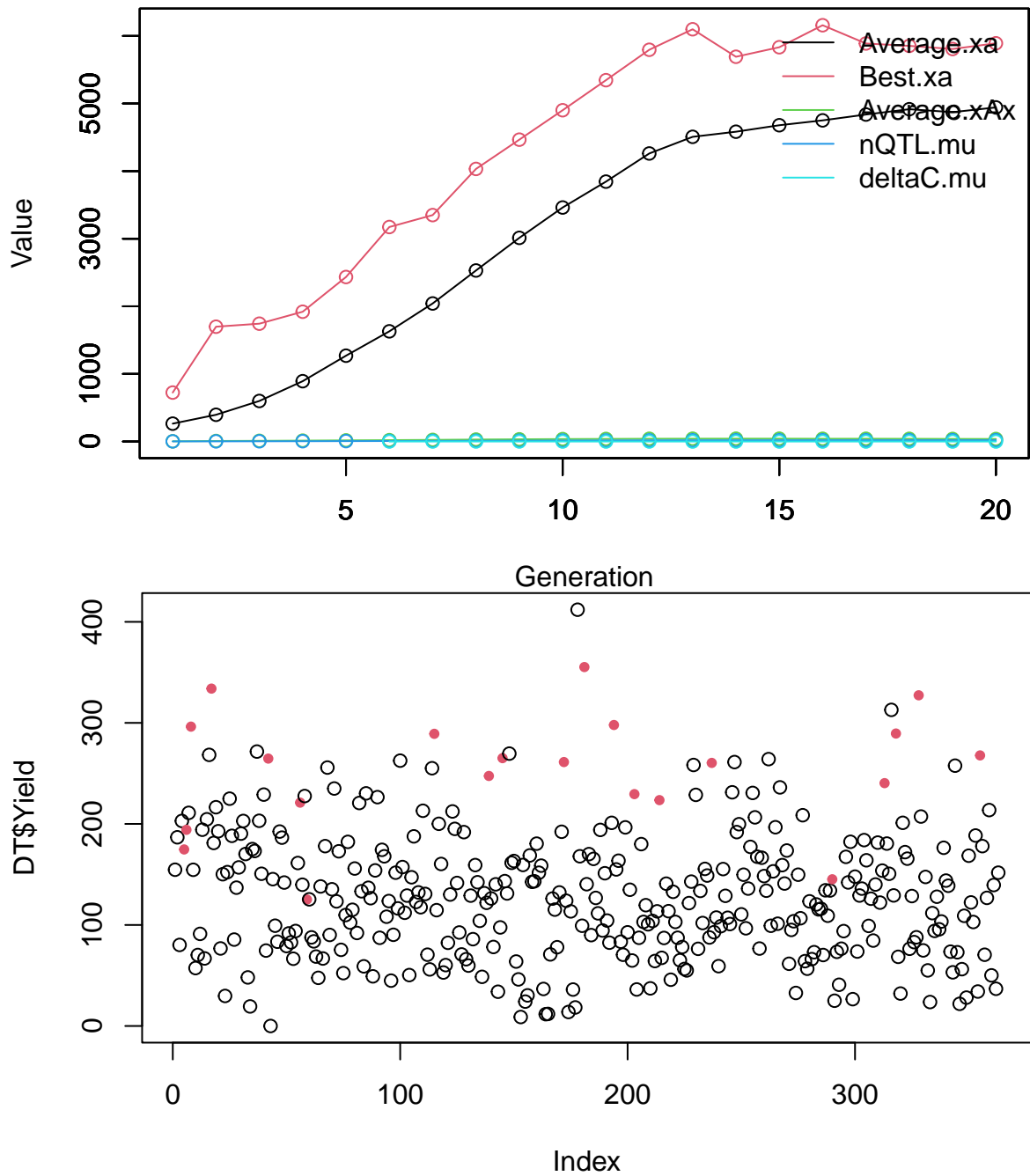
We then use the bestSol() function to extract the solution that maximizes our fitness function and constraints.

```
best = bestSol(res)[1];
sum(res$M[best,]) # total # of inds selected
```

```
## [1] 21
```

We can use the pmonitor() function to see if convergence was achieved between the best and teh average solutions.

```
pmonitor(res)
plot(DT$Yield, col=as.factor(res$M[best,]),
  pch=(res$M[best,]*19)+1)
```



2b) Obtaining optimal N crosses from a population for a given feature A variation of the same problem is when we want to pick the best crosses instead of the best parents to directly find the optimal solution for a crossing block. In the following example we use a dataset of crosses with marker and phenotype information to show how to optimize this problem.

```
data(DT_technow)
DT <- DT_technow
DT$occ <- 1; DT$occ[1]=0
M <- M_technow
A <- A.mat(M)
head(DT)
```

```
##   hybrid dent flint    GY    GM    hy occ
## 1 518.298 518   298  -8.04 -0.85 518:298  0
## 2 518.305 518   305 -11.10  1.70 518:305  1
## 3 518.306 518   306 -16.85  2.24 518:306  1
## 4 518.316 518   316  2.08 -1.33 518:316  1
## 5 518.323 518   323  5.65 -2.71 518:323  1
## 6 518.327 518   327 -16.95 -0.52 518:327  1
```

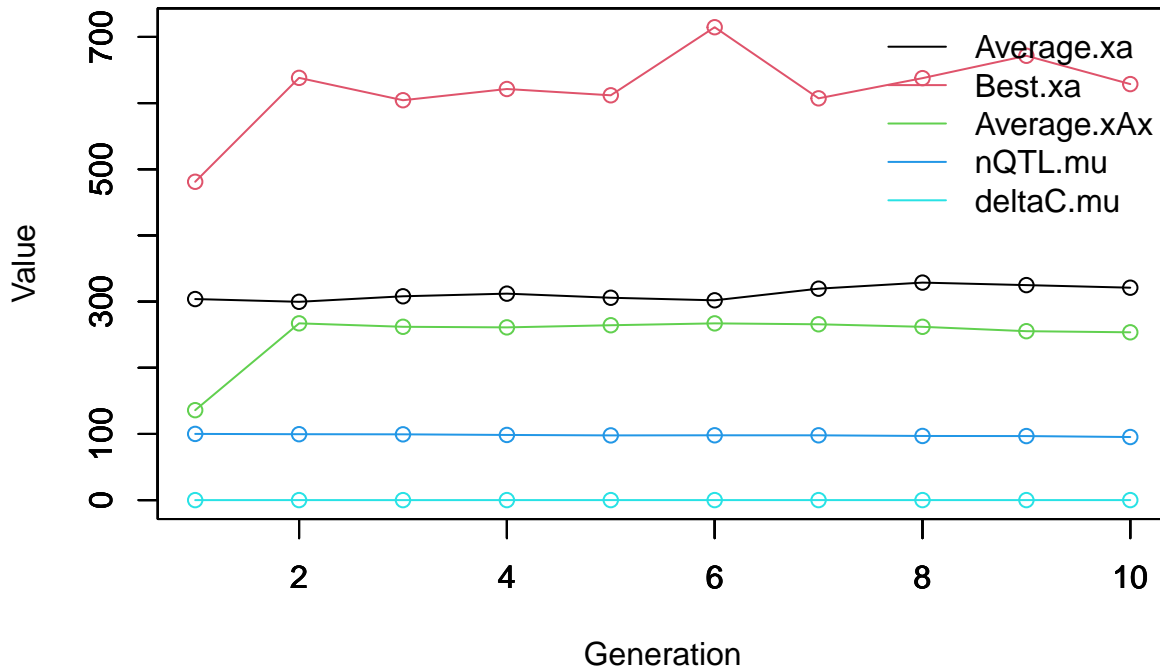
The way to specify this problem is exactly the same than with the optimization of parents but the input information is at the level of predicted crosses instead of individuals.

```
# run the genetic algorithm
res<-evolafit(formula = c(GY, occ)~hy, dt= DT,
  # constraints: if sum is greater than this ignore
  constraintsUB = c(Inf,100),
  # constraints: if sum is smaller than this ignore
  constraintsLB= c(-Inf,-Inf),
  # weight the traits for the selection
  traitWeight = c(1,0),
  # population parameters
  nCrosses = 100, nProgeny = 10,
  # coancestry parameters
  A=A, lambda=c(0.3,0), nQTLperInd = 100,
  # selection parameters
  propSelBetween = 0.9, propSelWithin =0.9,
  nGenerations = 10, verbose=FALSE)
best = bestSol(res)[1]
sum(res$M[best,]) # total # of inds selected
```

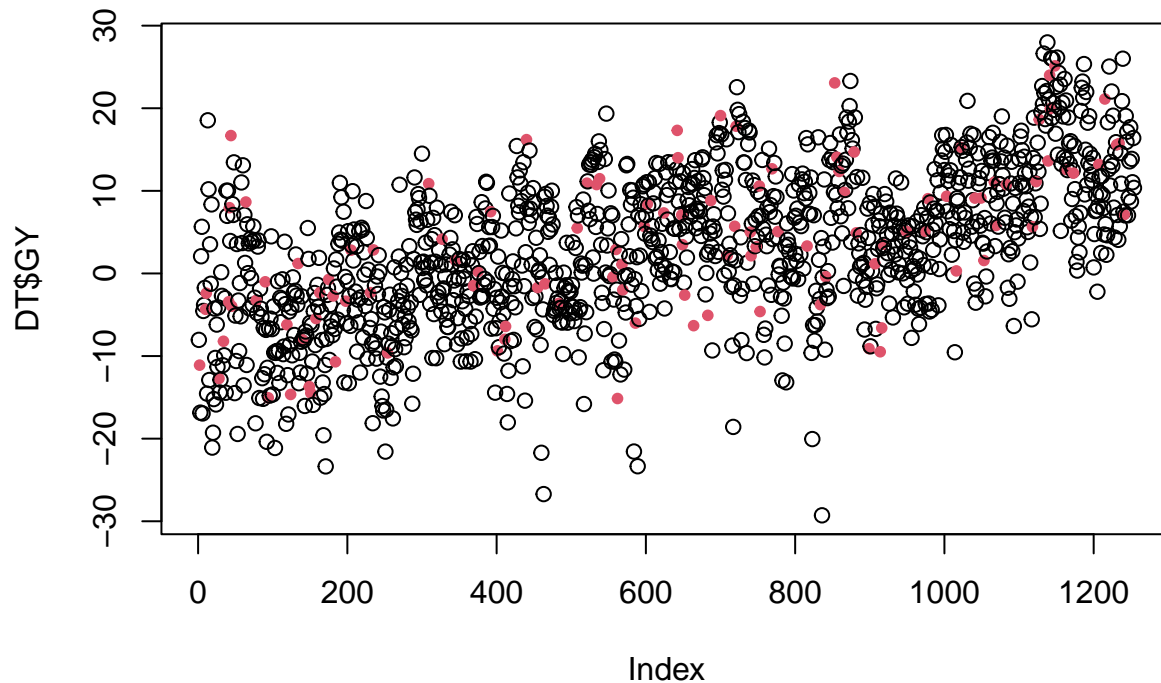
```
## [1] 115
```

You can use the pmonitor() or pareto() functions to see the evolution of the solution and see the performance of the solution selected.

```
pmonitor(res)
```



```
plot(DT$GY, col=as.factor(res$M[best,]),
     pch=(res$M[best,]*19)+1)
```



3) Optimizing a subsample of size N to be representative

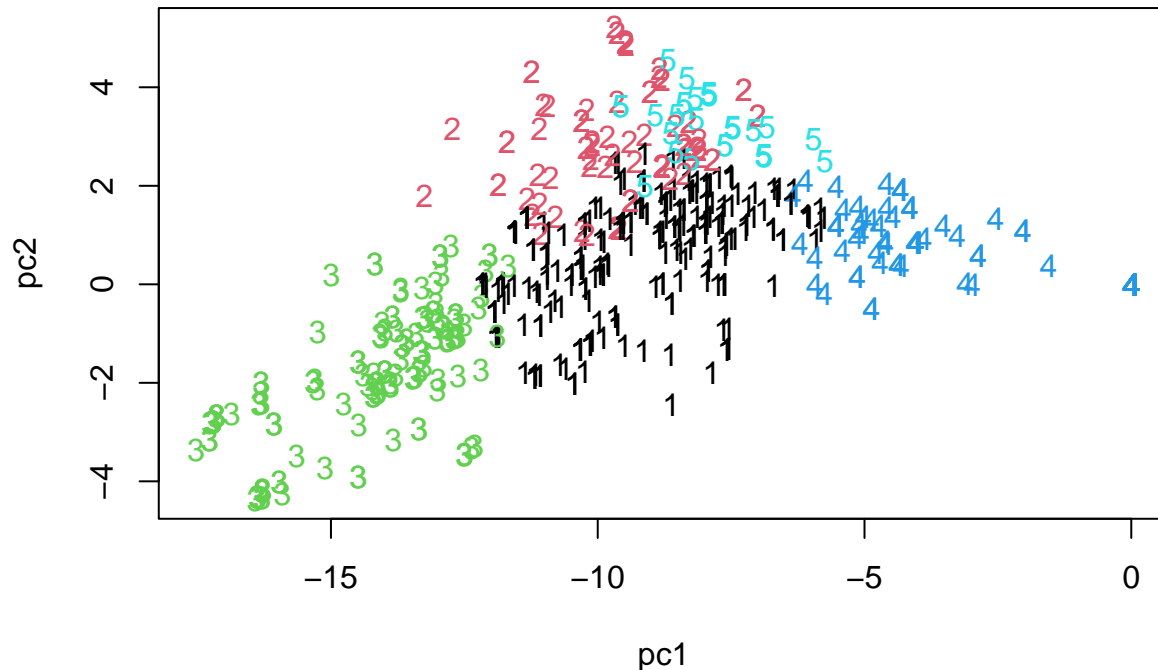
One particular case when we want to pick a representative subsample is when we don't have the resources to test everything (e.g., in the field/farm). In this example we use the information from 599 wheat lines to pick a subsample that maximizes the prediction accuracy for the entire sample. We start loading the data, in particular the phenotypes (DT) and the pedigree relationship matrix (A).

```
data(DT_wheat)
DT <- as.data.frame(DT_wheat)
DT$id <- rownames(DT) # IDs
DT$occ <- 1; DT$occ[1]=0 # to track occurrences
DT$dummy <- 1; DT$dummy[1]=0 # dummy trait
# if genomic
# GT <- GT_wheat + 1; rownames(GT) <- rownames(DT)
# A <- GT%*%t(GT)
# A <- A/mean(diag(A))
# if pedigree
A <- A_wheat
```

Now in order to pick a structured sample we will do a PCA and pick the cluster number 3 to be a subset to predict later (vp), while we will focus in rest of the population as candidates for the training set (tp).

```
##Perform eigenvalue decomposition for clustering
##And select cluster 5 as target set to predict
pcNum=25
svdWheat <- svd(A, nu = pcNum, nv = pcNum)
PCWheat <- A %*% svdWheat$v
rownames(PCWheat) <- rownames(A)
DistWheat <- dist(PCWheat)
```

```
TreeWheat <- cutree(hclust(DistWheat), k = 5 )
plot(PCWheat[,1], PCWheat[,2], col = TreeWheat,
     pch = as.character(TreeWheat), xlab = "pc1", ylab = "pc2")
```



```
vp <- rownames(PCWheat)[TreeWheat == 3]; length(vp)
```

```
## [1] 159
```

```
tp <- setdiff(rownames(PCWheat), vp)
```

3a) Optimizing a subsample of size N to be representative of its own Since the objective is to select a set of 100 lines that represent best the training set (tp) of ~400 lines we will subset a relationship matrix for that training set (As).

```
As <- A[tp, tp]
DT2 <- DT[rownames(As), ]
```

For this particular case there is no trait to optimize ($x'a$) but we just want to make sure that we maintain as much variation in the sample as possible ($x'Ax$). We then just create a dummy trait in the dataset (dummy) to put all the weight into the group relationship ($x'Ax$) using the lambda argument. The trait for occurrence we will use it as before to control the number of individuals to be in the sample.

```
res<-evolafit(cbind(dummy, occ)~id, dt= DT2,
              # constraints: if sum is greater than this ignore
              constraintsUB = c(Inf, 100),
              # constraints: if sum is smaller than this ignore
              constraintsLB= c(-Inf, -Inf),
              # weight the traits for the selection
              traitWeight = c(1,0),
              # population parameters
              nCrosses = 100, nProgeny = 10,
              # coancestry parameters
              A=As,
```

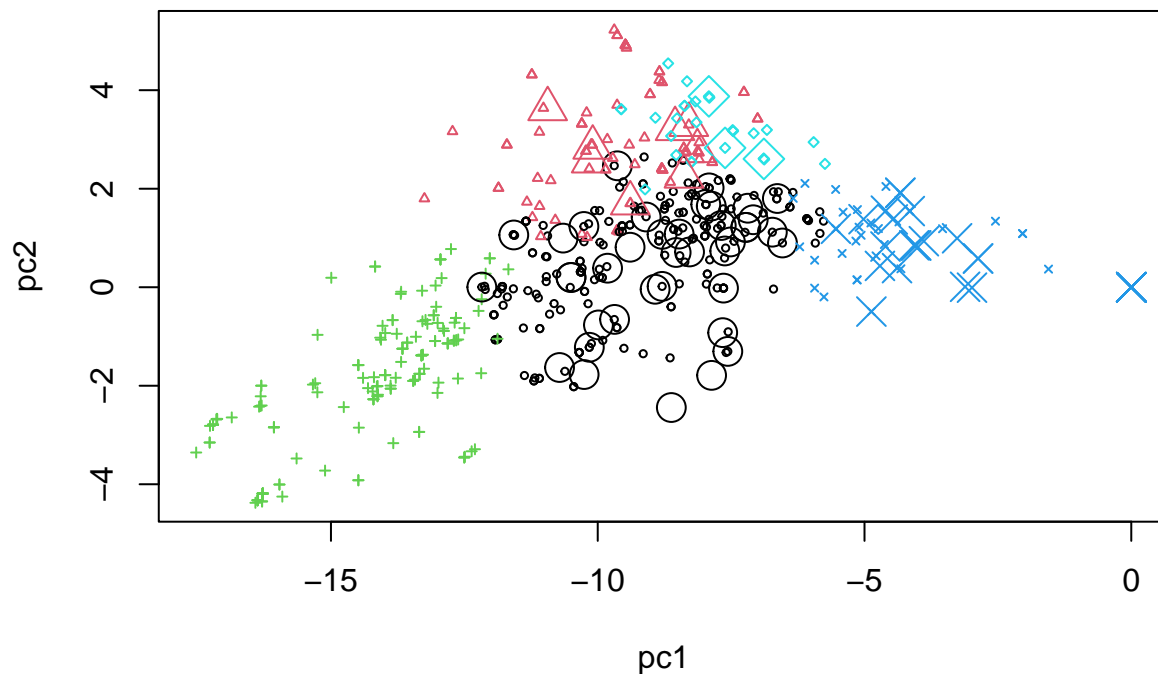
```
lambda=c(1,0), nQTLperInd = 80,
# selection parameters
propSelBetween = 0.5, propSelWithin =0.5,
nGenerations = 15, verbose = FALSE)
```

```
best = bestSol(res)[1]
sum(res$M[best,]) # total # of inds selected
```

```
## [1] 89
```

You can see which individuals were selected.

```
cex <- rep(0.5,nrow(PCWheat))
names(cex) <- rownames(PCWheat)
cex[names(which(res$M[best,]==1))]=2
plot(PCWheat[,1], PCWheat[,2], col = TreeWheat, cex=cex,
     pch = TreeWheat, xlab = "pc1", ylab = "pc2")
```



3b) Optimizing a subsample of size N to be representative of another population

we can use the covariance between the training population and the validation population to create a new trait ($x'a$) that can be used in addition to the group relationship ($x'Ax$).

```
DT2$cov <- apply(A[tp,vp],1,mean)
```

The model can be specified as before with the subtle difference that the covariance between the training and validation population contributes to the fitness function.

```
res<-evolafit(cbind(cov, occ)~id, dt= DT2,
# constraints: if sum is greater than this ignore
constraintsUB = c(Inf, 100),
# constraints: if sum is smaller than this ignore
constraintsLB= c(-Inf, -Inf),
```



```

# weight the traits for the selection
traitWeight = c(1,0),
# population parameters
nCrosses = 100, nProgeny = 10,
# coancestry parameters
A=As,
lambda=c(1,0), nQTLperInd = 80,
# selection parameters
propSelBetween = 0.5, propSelWithin =0.5,
nGenerations = 15, verbose = FALSE)
best = bestSol(res)[1]
sum(res$M[best,]) # total # of inds selected

```

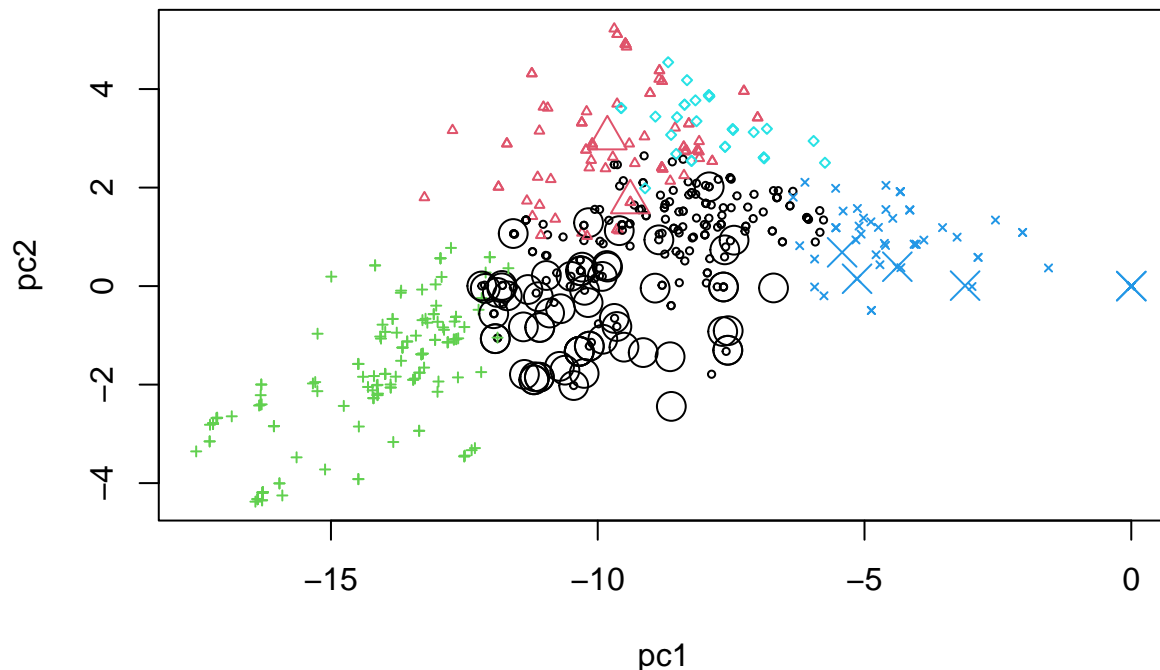
```
## [1] 72
```

You can plot the final results and see which individuals were picked.

```

cex <- rep(0.5,nrow(PCWheat))
names(cex) <- rownames(PCWheat)
cex[names(which(res$M[best,]==1))]=2
plot(PCWheat[,1], PCWheat[,2], col = TreeWheat, cex=cex,
     pch = TreeWheat, xlab = "pc1", ylab = "pc2")

```



6) How to specify constraints

Gender In this case is better if you only create the cross combinations that are possible (e.g., where male and female can couple) and you handed them to the evolutionary algorithm. That means, the rows of the crosses to be in the searching space only include the realistic ones.

Number of times a parent should be used In this case you would need to do 2 things:

- A) Create two new column with the mother and father of the cross

B) Create a new trait that is computed on the fly while running the analysis where the trait is a representation of how many times the parent is present in the current contribution. For example the trait function can be defined for mother as:

```
# data$motherN <- as.numeric(as.factor(data$mother))
# data$fatherN <- as.numeric(as.factor(data$father))
# fitnessf <- list(motherN= function(motherN, fatherN){
# v = table(c(motherN, fatherN))
# res <- ifelse(any(v > 4), 0, 1)
# return(res)
# })
```

Where motherN is a new trait that is just a numerical recoding of the mothers (e.g., `as.numeric(as.factor(dt$mother))`). This new fitness function for the trait motherN can be added to the ‘fitnessf’ argument together with a constraint (e.g., `constraintsLB=1`) to ensure that only crosses that meet the number of times showing in a cross are among the selected candidates.

7) How to optimize the number of progeny to produce per cross

The advice here is to upload directly the phased genotypes (haplotypes) to the AlphaSimR machinery and simulate the possible crosses to explore how many individuals are required to sample a given trait (oligogenic or polygenic) with a given probability. No need to use the `evola` package.

Literature

Giovanny Covarrubias-Pazaran (2024). `evola`: a simple evolutionary algorithm for complex problems. To be submitted to *Bioinformatics*.

Gaynor, R. Chris, Gregor Gorjanc, and John M. Hickey. 2021. AlphaSimR: an R package for breeding program simulations. *G3 Gene|Genomes|Genetics* 11(2):jkaa017. <https://doi.org/10.1093/g3journal/jkaa017>.

Chen GK, Marjoram P, Wall JD (2009). Fast and Flexible Simulation of DNA Sequence Data. *Genome Research*, 19, 136-142. <http://genome.cshlp.org/content/19/1/136>.