

boom-bust-meta-sim

Peta Hill & Bernd Gruber

2023-03-08

Forward

The code below allows you to rerun the simulations as described in the Manuscript of Hill et al. (2023): Fragmentation and gene flow in a boom-bust metapopulation reveal a trade-off between heterozygosity and allelic richness.

As the actual simulation and analysis takes some time to run (about 2 days using 40 cores) we also provide means to run the simulation with fewer parameter combinations.

Please note the simulation is only tested for the stated versions of the R packages and programs and if the reader tries to run it, it may break due to incompatibilities of packages

Required packages and programs

To be able to run the simulation the following programs are needed (links are given which details how to install on your system):

- slim (V 3.7.1) (<https://github.com/MesserLab/SLiM>)
- plink (V 1.9) (<https://www.cog-genomics.org/plink/>)

For the R installation you need to have a recent version of R (>4.0) and the following packages. Please note the version is important as some of the packages are known not to work if versions are different.

- seqinr version ‘4.2.8’
- slimr version ‘0.1.0’
- dartR version ‘2.7.2’
- ggplot2 version ‘3.3.6’
- vcfR version ‘1.12.0’
- purrr version ‘0.3.4’
- furrr version ‘0.2.3’
- parallel version ‘4.1.3’
- doParallel version 1.0.17

Setting the environment

Loading the packages

```
library(seqinr)
library(slimr)
library(dartR)
library(ggplot2)
library(vcfR)
library(purrr)
library(furrr)
library(parallel)
library(doParallel)
```

Once loaded package versions can be checked via: `packageVersion("packagename")`

Setting paths

Next it is important to define the paths where you have installed the external programs (plink and slim) and the folder where all output files are stored. Please be aware in case you run the full simulation output, files need a considerable amount of disc space (~2TB)

```
plinkdir <- "C:/slimtest"
slimdir <- "C:/msys64/mingw64/bin"
Sys.setenv(SLIM_HOME=slimdir)

#folder where all simulated files are stored
outdir <- "c:/TOY/test"
```

Setting numbers of cores

You can check the number of processors (cores) your system has to define how many cores you want to use for the simulation. As a default we use the number of cores on your system -1, to allow for smooth input/outputs, but obviously this could be changed to a lower number. It does not make sense to increase the number to more than the number of cores that are available.

```
cat(paste0("You have ", detectCores()), " cores available on your
  system.\nThe maximum recommended number of cores you may use
  is ", detectCores()-1, "."))
```

```
## You have 8 cores available on your
##   system.
## The maximum recommended number of cores you may use
##   is 7.
```

```
#number of cores to be used
n.cores <- detectCores()-1
```

Setting the parameters for the simulation

The parameter space that is simulated is configured by a data.frame. This allows to scale the simulation to a short (toy) run and also for a run of all parameter combinations as in the manuscript.

The parameter are:

- npops: number of subpopulations in the metapopulation. Specified as a vector, e.g. npops <- c(2,4,8) runs simulation with 2,4 and 8 subpopulation. Numbers have to be powers of 2.
- n.loc: number of loci to be simulated.
- mig: migration rates in percentage, e.g. c(25,50,100) runs simulation with 25, 50 and 100% migration rates
- mev: [migration every x generation], specifies the interval of migration events, e.g. c(2,5,10) creates migration events every second, fifth and 10th generation.
- rep: the number of repeats for a parameter combination

Please note the simulation runs all possible parameter combinations above, so the number of simulations increases quite fast if the number of levels of each parameter is high. The default setting below is to run a toy example of only 72 combinations, we also provide the settings for the simulations as run in the manuscript, but be aware that will take quite some time to run. For a more detailed explanation please refer to the manuscript.

```
#number of subpopulations
pops <- c(2,4,8)

# number of loci (SNPs)
n.loc <- 50
# toy example
df <- expand.grid(rep=1:10, pNum = pops, mig = c(25,100), mev = c(2,10))

#full run as in the manuscript
#n.loc <- 5000
#df <- expand.grid(rep=1:50, pNum = c(1,2,4,8,16,32,64),
#mig = c(1,5,10,25,50,100), mev = c(2,5,10,20))

#number of paramter combinations
nrow(df)
```

```
## [1] 120
```

Creating initial genotype frequencies for each subpopulation

We create the metapopulations as specified with initial genotype frequencies of 0.5 for the number of loci specified and using the script below. The initial populations will be saved as vcf files in the folder as specified in outdir above. The naming of the files is:

- gl_n.loc_subpop_i.vcf

```
ref<-sample(c("A","G","C","T"),n.loc,replace=T, prob=c(0.25,0.25,0.25,0.25))
R<-paste(ref, collapse = "")
#saveRDS(R,file.path(outdir,"TOYreference_allele.rds"))
alt<-comp(ref, forceToLower = FALSE)
locAl<-paste0(ref,"/",alt)
p<-list()
```

```
for (ii in 1:length(pops))
{
  n.pop <- pops[ii]
```



```
## Starting gl2plink
## Processing genlight object with SNP data
## Completed: gl2plink
## Starting gl2plink
## Processing genlight object with SNP data
## Completed: gl2plink
```

Slim templates

Each parameter file will be used for a simulation executed within slim using the file for the corresponding parameter combination. To prepare all the simulations in slim we use `slimr` and therein the `template` function. This function allows to setup templates for the slim runs leaving place holders for the parameter values as set in the `data.frame` (`df`).

```
slim_script(
  slim_block(initialize(),
    { initializeSLiMModelType("nonWF");
      defineConstant("alpha1", 1e-9); #mutation rate
      defineConstant("metasz",256);# total number of individuals in the metapopulation
      defineConstant("subp_num",slimr_template("pNum"));#number of subpopulations
      defineConstant("rep",slimr_template("rep")); #number of replicate simulations
      defineConstant("mig", slimr_template("mig")); #migration rate as percent
      defineConstant("migRate",asInteger(mig/100*metasz)); #migration rate as number of indivi
      defineConstant("subp_sz", asInteger(metasz/subp_num)); #number of individuals in each su
      defineConstant("outpath", slimr_inline(outdir)); #output directory
      defineConstant("mev",slimr_template("mev"));# How often to mix
      initializeSLiMOptions(nucleotideBased=T);
      defineConstant("L",slimr_inline(n.loc))
      initializeAncestralNucleotides(slimr_inline(R)); #reference sequence
      initializeMutationTypeNuc("m1", 0.5, "f", 0.0);
      m1.convertToSubstitution = T;
      m1.mutationStackPolicy = "l";
      mutmat = matrix(c(0.0, alpha1, 0.0,0.0,alpha1,0.0,0.0,0.0,0.0,0.0,0.0,
                        alpha1,0.0,0.0,alpha1,0.0), nrow=4, ncol=4)
      initializeGenomicElementType("g1", m1, 1.0,mutmat);
      initializeGenomicElement(g1, 0, L-1);
      initializeRecombinationRate(0.5);
    },
    slim_block(reproduction(), {# everybody mates, selfing allowed
      for (j in seqLen(subp_num+1)){
        inds = sim.subpopulations[j];
        if (inds.individualCount > 0)
        {
          parents = inds.sampleIndividuals(inds.individualCount, replace = T);
          for (i in seqLen(inds.individualCount))
          {
            mum = parents[i];
            dad = inds.sampleIndividuals(1);
            child = inds.addCrossed(mum, dad);
          }
        }
      }
    })
  )
)
```

```

    else {j=j+1;
    }
    self.active = 0;}
}),
slim_block(1, {
  for (i in seqLen(subp_num)){#set up subpopulations, each of size 'subpop_sz', reading in genotypes
    sim.addSubpop(i,subp_sz);#from .vcf files we generated
    subpops = sim.subpopulations;
    f = paste0(outpath+"/gl_50_" + subp_num + "_" + (i+1) + ".vcf");

    subpops[i]%.%genomes.readFromVCF(f, m1);
    subpops[i]%.%individuals.tag = ((subp_sz*i)+1):((i + 1) * subp_sz);
  }
  sim.addSubpop(100,0);# this is for the migrant pool
}),
slim_block(2,501, early(), {
  # parents die in line with WF assumptions
  inds = sim.subpopulations.individuals;
  inds[inds.age > 0]%.%fitnessScaling = 0.0;
}),
slim_block(2,499, first(), {
  if(sim.generation %% mev == 0){#if the generation is one in which gene flow
    migrants = sample(sim.subpopulations.individuals,migRate);#is to occur, take migrants into
    p100.takeMigrants(migrants);#the migrant pool
  }
}),
slim_block(2,500, first(), {
  if(sim.generation %% mev == 1){#if the generation is the one after migration
    for (id in seqLen(subp_num)){# randomly return migrants to the subpopulations
      subpops = sim.subpopulations[sim.subpopulations.id == id];
      migrants = sample(p100.individuals, asInteger(subp_sz-subpops.individualCount));
      dest = subpops;
      dest.takeMigrants(migrants);
    }
  }
}),
slim_block(470,500, late(), {
  subpops = sim.subpopulations;
  if(500-sim.generation <= (mev+1)){
    subpops.individuals.genomes.outputVCF(paste0(outpath,"/FNAME_",
slimr_template("pNum"),"_mig_",slimr_template("mig"),"_rep_",
slimr_template("rep"),"_mev",slimr_template("mev"),"_gen",
sim.generation,".vcf"),simplifyNucleotides=T, outputNonnucleotides=F);}
  },
  slim_block(500,late(), {
    sim.simulationFinished();
  })
})
)->BBS3

#generate a simulation for all combinations of parameters
script_temp_df <- slim_script_render(BBS3, template =df)

```

You can check the content of the first template via `script_template_df[[1]]`

Running simulations in slim

After this setup, the runs in slim are then very simply executed via two commands. The `plan()` function prepares the parallel execution using the number of cores requested (`n.cores`), and the `slim_run` function executes the simulations in parallel with the provided number of cores.

```
plan(future::multisession(workers = n.cores))

resultsNonWF <- slim_run(script_temp_df, parallel = T, throw_error = TRUE)
```

Load and analysis of the simulations

The next step is to load and summarise all simulation output files. Depending on ‘mev’ the simulations output the mev last generations of each run and then calculate heterozygosity and number of alleles for those last generations and records those in a table (`toy.csv`). To speed this step we also use a parallel implementation with the same number of cores as specified above (`n.cores`).

```
#load all files starting with FNAME
filesM <- list.files(path=outdir,pattern="FNAME" )
gls <- list()

ext <- function(i) {
  dumvcf<- read.vcfR(file.path(outdir,filesM[i]),verbose=FALSE)
  gls<- vcfR2genlight(dumvcf, n.cores = 1)
  gls<-gl.compliance.check(gls, verbose=0)
  namesgls <- filesM[i]
  pNum<-as.numeric(gsub("FNAME_([0-9]+).*$", "\\1",namesgls))
  rep<-as.numeric(gsub(".*?_rep_([0-9]+).*$", "\\1", namesgls))
  mig<-as.numeric(gsub(".*?_mig_([0-9]+).*$", "\\1", namesgls))
  mev<-as.numeric(gsub(".*?_mev([0-9]+).*$", "\\1", namesgls))
  gen<-as.numeric(gsub(".*?_gen([0-9]+).*$", "\\1", namesgls))
  het<-sum(as.matrix(gls)==1)/(nInd(gls)*1*n.loc)
  glout<-gl.filter.monomorphs(gls, verbose = 0)
  Lnum<-(2*nLoc(glout)) + (n.loc-nLoc(glout))
  locd<-rbind(c(pNum, rep, mig, mev,gen, Lnum, het))
  return(locd)
}

#analyse in parallel using the ext(raction) function above.
registerDoParallel(n.cores)
system.time(locd <- foreach(i=1:length(filesM),.combine = rbind,
  .packages = c("vcfR", "dartR")) %dopar% ext(i))

##      user  system elapsed
##      0.25    0.06   61.39

colnames(locd)<- c("pNum", "rep","mig", "mev","gen", "Lnum", "het")

#create the final table (toy)
toy <- data.frame(locd)
write.csv(toy,file.path(outdir,"toy.csv"))
```

Plotting results

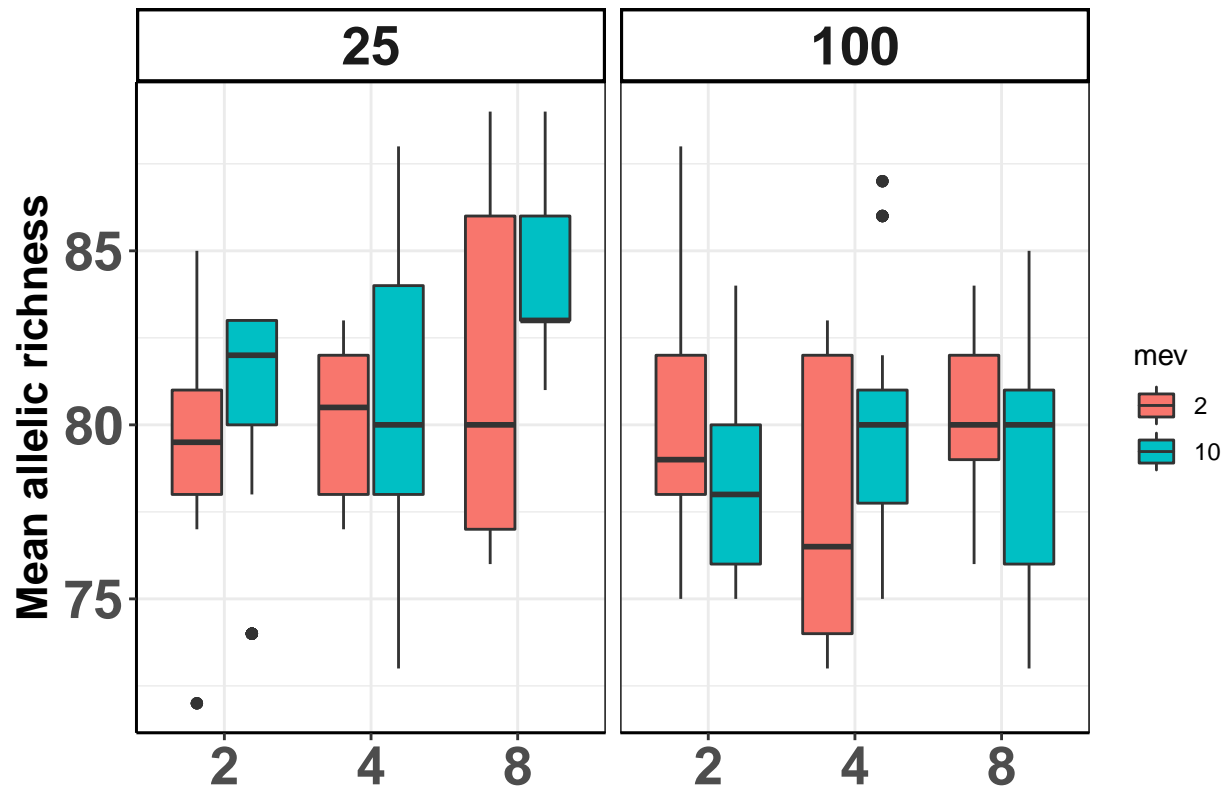
The results of the simulations are recorded in the toy object. We use ggplot to produce the results and concentrate here on the main to plots, namely the allelic richness and heterozygosity for each of the parameter combinations (e.g., Figure 2 A and B of the manuscript)

```
#prepare for plotting
toy$pNum <- factor(toy$pNum ,pops)
toy$mig <- factor(toy$mig)
toy$mev<-factor(toy$mev)
```

Allelic richness

```
nma<-ggplot(aes(pNum,Lnum), data=toy)+geom_boxplot(aes(fill=mev))+theme_bw()+
  theme(strip.background = element_rect(color="black",
    fill="white",size=1, linetype="solid"),strip.text.x=element_text
    (size=20,face="bold"), strip.text.y=element_blank())+
  theme(axis.line = element_line(colour ="black"))+xlab("")+
  ylab("Mean allelic richness")+
  theme(axis.title.x=element_text(size=16,
    face="bold"),axis.title.y=element_text(size=16,face="bold"),
    axis.text.x=element_text(size=20,face="bold"),
    axis.text.y=element_text(size=20,face="bold"))+facet_wrap(~mig)

nma
```

Heterozygosity

```
nmh<-ggplot(aes(pNum,het), data=toy)+geom_boxplot(aes(fill=mev))+
  theme_bw()+theme(strip.background =
    element_rect(color="black",
      fill="white", size=1, linetype="solid"),
    strip.text.x =element_text(size=20, face="bold"),
    strip.text.y = element_blank()+theme(axis.line =
    element_line(colour ="black"))+xlab("")+
    ylab("Mean heterozygosity")+ theme(axis.title.x
    =element_text(size=16,face="bold"),axis.title.y
    =element_text(size=16,face="bold"),axis.text.x=
    element_text(size=20,face="bold"),axis.text.y=
    element_text(size=20,face="bold"))+facet_wrap(~mig)

nmh
```

