

# NB Population Genomic Analysis on Non-Haplotig Masked Neutral SNPs

Amy Zyck

5/27/2020

## Population analyses on Neutral SNP dataset

I first created a new directory NB\_PopGen\_Neutral within PATH: /home/azyck/NB\_capture/NB\_ddocent. I linked neutralloci.recode.vcf to this working directory. The RMarkdown file should be saved in this same working directory.

```
#Loading all the necessary packages
library(adegenet)

## Loading required package: ade4

## Registered S3 method overwritten by 'spdep':
##   method    from
##   plot.mst ape

##
##   /// adegenet 2.1.3 is loaded ///////////
##
##   > overview: '?adegenet'
##   > tutorials/doc/questions: 'adegenetWeb()'
##   > bug reports/feature requests: adegenetIssues()

library(vcfR)

##
##      ****   ***  vcfR   ***      ****
## This is vcfR 1.12.0
##   browseVignettes('vcfR') # Documentation
##   citation('vcfR') # Citation
##      ****   ***      ****      ****

library("radiator") # Conversion from vcf to a lot of other formats

## ***** IMPORTANT NOTICE *****
## radiator v.1.0.0 was modified heavily.
## Read functions documentation and available vignettes.
##
## For reproducibility:
##   radiator version: 1.0.0
##   radiator build date: R 3.5.1; ; 2019-03-20 13:31:04 UTC; unix
##   Keep zenodo DOI.
## *****
```

```

library("dplyr")

##
## Attaching package: 'dplyr'
## The following objects are masked from 'package:stats':
##
##     filter, lag
## The following objects are masked from 'package:base':
##
##     intersect, setdiff, setequal, union
library("hierfstat")

##
## Attaching package: 'hierfstat'
## The following object is masked from 'package:adegenet':
##
##     read.fstat
library("ggplot2") #For plotting
library("reshape2") #For plotting
library("plyr")

## -----
## You have loaded plyr after dplyr - this is likely to cause problems.
## If you need functions from both plyr and dplyr, please load plyr first, then dplyr:
## library(plyr); library(dplyr)

## -----
## 
## Attaching package: 'plyr'
## The following objects are masked from 'package:dplyr':
##
##     arrange, count, desc, failwith, id, mutate, rename, summarise,
##     summarise
library("cowplot") #For plotting manuscript figs
library(PCAviz) #Visualizing output of PCA
library("stringr")
library("bigsnpr") # package for Linkage Disequilibrium pruning

## Loading required package: bigstatsr

```

## All Neutral SNPs

### Making files

#### Make genind object

`neutralloci.recode.vcf` contains neutral SNPs for populations PVD, GB, BIS, and NAR. Steps for generating this VCF file are located in `EecSeq_Cvirginica_dDocent.md` and `EecSeq_Cvirginica_Filtering.md`. Population NIN was removed from the VCF file following steps in `EecSeq_Cvirginica_OutlierDetection.md`.

`strata` contains population, environmental, and library information for each sample - can be accessed here.

```
my_vcf <- read.vcfR("neutralloci.recode.vcf")
```

```
## Scanning file to determine attributes.  
## File attributes:  
##   meta lines: 63  
##   header_line: 64  
##   variant count: 74520  
##   column count: 49  
##  
Meta line 63 read in.  
## All meta lines processed.  
## gt matrix initialized.  
## Character matrix gt created.  
##   Character matrix gt rows: 74520  
##   Character matrix gt cols: 49  
##   skip: 0  
##   nrows: 74520  
##   row_num: 0  
##  
Processed variant 1000  
Processed variant 2000  
Processed variant 3000  
Processed variant 4000  
Processed variant 5000  
Processed variant 6000  
Processed variant 7000  
Processed variant 8000  
Processed variant 9000  
Processed variant 10000  
Processed variant 11000  
Processed variant 12000  
Processed variant 13000  
Processed variant 14000  
Processed variant 15000  
Processed variant 16000  
Processed variant 17000  
Processed variant 18000  
Processed variant 19000  
Processed variant 20000  
Processed variant 21000  
Processed variant 22000  
Processed variant 23000  
Processed variant 24000  
Processed variant 25000  
Processed variant 26000  
Processed variant 27000  
Processed variant 28000  
Processed variant 29000  
Processed variant 30000  
Processed variant 31000  
Processed variant 32000  
Processed variant 33000  
Processed variant 34000
```

```

Processed variant 35000
Processed variant 36000
Processed variant 37000
Processed variant 38000
Processed variant 39000
Processed variant 40000
Processed variant 41000
Processed variant 42000
Processed variant 43000
Processed variant 44000
Processed variant 45000
Processed variant 46000
Processed variant 47000
Processed variant 48000
Processed variant 49000
Processed variant 50000
Processed variant 51000
Processed variant 52000
Processed variant 53000
Processed variant 54000
Processed variant 55000
Processed variant 56000
Processed variant 57000
Processed variant 58000
Processed variant 59000
Processed variant 60000
Processed variant 61000
Processed variant 62000
Processed variant 63000
Processed variant 64000
Processed variant 65000
Processed variant 66000
Processed variant 67000
Processed variant 68000
Processed variant 69000
Processed variant 70000
Processed variant 71000
Processed variant 72000
Processed variant 73000
Processed variant 74000
Processed variant: 74520
## All variants processed

strata <- read.table("strata", header=TRUE)

radfilt <- vcfR2genind(my_vcf, strata = strata, pop = c(rep("BIS", 10),rep("GB", 10),rep("NAR", 10), r

radfilt

## /// GENIND OBJECT ///////////
##
## // 40 individuals; 74,520 loci; 148,794 alleles; size: 67.9 Mb
##
## // Basic content
## @tab: 40 x 148794 matrix of allele counts

```

```

##  @loc.n.all: number of alleles per locus (range: 1-2)
##  @loc.fac: locus factor for the 148794 columns of @tab
##  @call.names: list of allele names for each locus
##  @ploidy: ploidy of each individual (range: 2-2)
##  @type: codom
##  @call: adegenet::df2genind(X = t(x), sep = sep, pop = ..2, strata = ..1)
##
## // Optional content
##  @pop: population of each individual (group size range: 10-10)
##  @strata: a data frame with 12 columns ( Individual, Population, Latitude, Longitude, Distance, SE
#Providing population names for plotting
pop_order <- c("BIS", "GB", "NAR", "PVD")

```

Read in the other info from .strata file and extract information such as locality, latitude, and longitude.

```

info <- as.data.frame(read.table("strata", header = T, sep = "\t", stringsAsFactors = F))
mystrats <- as.data.frame(matrix(nrow = length(indNames(rad.filt)), ncol=10))
colnames(mystrats) <- c("Population", "Latitude", "Longitude", "Distance", "SE", "Temperature", "Salinity", "pH", "Chlorophyll", "Other")
just.strats <- select(info, c("Population"))
stratted.filt <- strata(rad.filt, formula= Population, combine = TRUE, just.strats)
stratted.filt@other <- select(info, Latitude, Longitude, Distance, SE, Temperature, Salinity, pH, Chlorophyll, Other)
stratted.filt

## /// GENIND OBJECT ///////////
##
## // 40 individuals; 74,520 loci; 148,794 alleles; size: 67.9 Mb
##
## // Basic content
##  @tab: 40 x 148794 matrix of allele counts
##  @loc.n.all: number of alleles per locus (range: 1-2)
##  @loc.fac: locus factor for the 148794 columns of @tab
##  @call.names: list of allele names for each locus
##  @ploidy: ploidy of each individual (range: 2-2)
##  @type: codom
##  @call: adegenet::df2genind(X = t(x), sep = sep, pop = ..2, strata = ..1)
##
## // Optional content
##  @pop: population of each individual (group size range: 10-10)
##  @strata: a data frame with 1 columns ( Population )
##  @other: a list containing: Latitude Longitude Distance SE Temperature Salinity pH Chlorophyll

```

### Repeat for quasi-independent set of SNPs

Steps for filtering are documented for OutFLANK in EecSeq\_Cvirginica\_OutlierDetection.md. However, I am repeating the steps here.

```
my_vcf <- read.vcfR("neutralloci.recode.vcf")
```

```

## Scanning file to determine attributes.
## File attributes:
##  meta lines: 63
##  header_line: 64
##  variant count: 74520
##  column count: 49
##
```

```
Meta line 63 read in.  
## All meta lines processed.  
## gt matrix initialized.  
## Character matrix gt created.  
##   Character matrix gt rows: 74520  
##   Character matrix gt cols: 49  
##   skip: 0  
##   nrows: 74520  
##   row_num: 0  
##  
Processed variant 1000  
Processed variant 2000  
Processed variant 3000  
Processed variant 4000  
Processed variant 5000  
Processed variant 6000  
Processed variant 7000  
Processed variant 8000  
Processed variant 9000  
Processed variant 10000  
Processed variant 11000  
Processed variant 12000  
Processed variant 13000  
Processed variant 14000  
Processed variant 15000  
Processed variant 16000  
Processed variant 17000  
Processed variant 18000  
Processed variant 19000  
Processed variant 20000  
Processed variant 21000  
Processed variant 22000  
Processed variant 23000  
Processed variant 24000  
Processed variant 25000  
Processed variant 26000  
Processed variant 27000  
Processed variant 28000  
Processed variant 29000  
Processed variant 30000  
Processed variant 31000  
Processed variant 32000  
Processed variant 33000  
Processed variant 34000  
Processed variant 35000  
Processed variant 36000  
Processed variant 37000  
Processed variant 38000  
Processed variant 39000  
Processed variant 40000  
Processed variant 41000  
Processed variant 42000  
Processed variant 43000  
Processed variant 44000
```

```

Processed variant 45000
Processed variant 46000
Processed variant 47000
Processed variant 48000
Processed variant 49000
Processed variant 50000
Processed variant 51000
Processed variant 52000
Processed variant 53000
Processed variant 54000
Processed variant 55000
Processed variant 56000
Processed variant 57000
Processed variant 58000
Processed variant 59000
Processed variant 60000
Processed variant 61000
Processed variant 62000
Processed variant 63000
Processed variant 64000
Processed variant 65000
Processed variant 66000
Processed variant 67000
Processed variant 68000
Processed variant 69000
Processed variant 70000
Processed variant 71000
Processed variant 72000
Processed variant 73000
Processed variant 74000
Processed variant: 74520
## All variants processed

geno <- extract.gt(my_vcf) # Character matrix containing the genotypes
position <- getPOS(my_vcf) # Positions in bp
chromosome <- getCHROM(my_vcf) # Chromosome information

G <- matrix(NA, nrow = nrow(geno), ncol = ncol(geno))

G[geno %in% c("0/0", "0|0")] <- 0
G[geno %in% c("0/1", "1/0", "1|0", "0|1")] <- 1
G[geno %in% c("1/1", "1|1")] <- 2

# NA should be replaced with "9" to work with the functions in the OutFLANK package
G[is.na(G)] <- 9

```

Chromosomes need to be of class integer for this to work.

```

# Visualizing chromosomes
chrom_unique <- unique(chromosome)
print(chrom_unique)

## [1] "NC_035780.1" "NC_035781.1" "NC_035782.1" "NC_035783.1" "NC_035784.1"
## [6] "NC_035785.1" "NC_035786.1" "NC_035787.1" "NC_035788.1" "NC_035789.1"

```

```
# Removing "NC_" from the chromosome name so it can be converted to an integer
chrom_new <- chromosome%>%str_replace("NC_","")

# Converting the character string into an integer
chrom1 <- as.integer(chrom_new)
```

chromosome and position need to be sorted for imputation to work.

```
chrom_sort <- sort(chrom1)
pos_sort <- sort(position)
```

*Note: This filtering program does not allow for missing genotype values.*

I can either remove the missing genotype values or impute missing genotypes. I'm going to try both and see how the two methods differ during the trimming step.

## Remove missing genotypes

```
# removing missing data
G_miss <- matrix(NA, nrow = nrow(geno), ncol = ncol(geno))

G_miss[geno %in% c("0/0", "0|0")] <- 0
G_miss[geno %in% c("0/1", "1/0", "1|0", "0|1")] <- 1
G_miss[geno %in% c("1/1", "1|1")] <- 2

# removing missing data
G_miss[na.omit(G_miss)]
```

```
##      [1] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
##     [37] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
##     [73] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
##    [109] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
##    [145] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
##    [181] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
##    [217] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
##    [253] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
##    [289] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
##    [325] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
##    [361] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
##    [397] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
##    [433] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
##    [469] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
##    [505] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
##    [541] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
##    [577] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
##    [613] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
##    [649] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
##    [685] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
##    [721] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
##    [757] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
##    [793] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
##    [829] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
##    [865] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
##    [901] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
```

```

# Converting the genotype matrix to class FBM.code256 using the matrix where missing data was removed
G1_miss <- add_code256(big_copy(t(G_miss), type="raw"), code=bigsnpr:::CODE_012)

## Warning in replaceMat(x$address_rw, i, j, value): At least one value changed (nan -> 0)
##   while converting from R type 'double' to C type 'unsigned char (raw)'.

newpc_miss <- snp_autoSVD(G1_miss, infos.chr = chrom_sort, infos.pos = pos_sort)

##
## Phase of clumping (on MAF) at r^2 > 0.2.. keep 18462 SNPs.
## Discarding 5790 variants with MAC < 10.
##
## Iteration 1:
## Computing SVD..
## 0 outlier variant detected..
##
## Converged!

which_pruned_miss <- attr(newpc_miss, which="subset") # Indexes of remaining SNPs after pruning
length(which_pruned_miss)

## [1] 12672

```

### Using snp\_fastImputeSimple to impute missing genotypes

This requires the genotype matrix to be converted to class FBM.code256.

A reference class for storing and accessing up to 256 arbitrary different values using a Filebacked Big Matrix of type unsigned char.

```

# Converting the genotype matrix to class FBM.code256
G1 <- add_code256(big_copy(t(G), type="raw"), code=bigsnpr:::CODE_012)

```

The chunk below imputes missing genotypes in order for the following code to run. Fast imputation via mode, mean, sampling according to allele frequencies, or 0.

Method argument: Either “random” (sampling according to allele frequencies), “mean0” (rounded mean), “mean2” (rounded mean to 2 decimal places), “mode” (most frequent call).

Depending on the method I choose, the # of SNPs kept after running the following chunk changes, so I’m not sure which method to use.

```

# Imputing missing genotypes
G_missSimple <- snp_fastImputeSimple(G1, method = "mode") # Using mode or no method call results in the same number of SNPs

```

Truncated SVD while limiting LD.

```

newpc<-snp_autoSVD(G_missSimple, infos.chr = chrom_sort, infos.pos = pos_sort)

```

```

##
## Phase of clumping (on MAF) at r^2 > 0.2.. keep 18348 SNPs.
## Discarding 5815 variants with MAC < 10.
##
## Iteration 1:
## Computing SVD..
## 0 outlier variant detected..
##
## Converged!

```

```
which_pruned <- attr(newpc, which="subset") # Indexes of remaining SNPs after pruning
length(which_pruned)
```

```
## [1] 12533
```

Either removing missing data or imputation results in a similar number of trimmed SNPs, so I will continue forward with missing data removed.

```
invisible(lapply(which_pruned_miss, write, "pruned_data.txt", append=TRUE))
```

### In terminal

```
$ mawk '!/#/' neutralloci.recode.vcf | cut -f1,2 > totalloci
```

```
$ NUM=(`cat totalloci | wc -l`)
```

```
$ paste <(seq 1 $NUM) totalloci > loci.plus.index
```

```
$ cat pruned_data.txt | parallel "grep -w ^{} loci.plus.index" | cut -f2,3> pruned_data_neutral.loci.txt
```

```
$ head pruned_data_neutral.loci.txt
```

```
output:
```

```
NC_035780.1 573785
```

```
NC_035780.1 573788
```

```
NC_035780.1 574229
```

```
NC_035780.1 574307
```

```
NC_035780.1 574313
```

```
NC_035780.1 574321
```

```
NC_035780.1 574444
```

```
NC_035780.1 574580
```

```
NC_035780.1 574616
```

```
NC_035780.1 574657
```

Create VCF file with just the pruned\_data loci

```
$ vcftools --vcf neutralloci.recode.vcf --recode --recode-INFO-all --positions pruned_data_neutral.loci
```

```
output:
```

```
fter filtering, kept 40 out of 40 Individuals
```

```
Outputting VCF file...
```

```
After filtering, kept 14880 out of a possible 74520 Sites
```

```
Run Time = 3.00 seconds
```

```
my_vcf_u <- read.vcfR("pruned_data_neutral.recode.vcf")
```

```
## Scanning file to determine attributes.
```

```
## File attributes:
```

```
##   meta lines: 63
```

```
##   header_line: 64
```

```
##   variant count: 14880
```

```
##   column count: 49
```

```
##
```

```
Meta line 63 read in.
```

```
## All meta lines processed.
```

```
## gt matrix initialized.
```

```
## Character matrix gt created.
```

```
##   Character matrix gt rows: 14880
```

```
##   Character matrix gt cols: 49
```

```
##   skip: 0
```

```
##    nrows: 14880
##    row_num: 0
##
Processed variant 1000
Processed variant 2000
Processed variant 3000
Processed variant 4000
Processed variant 5000
Processed variant 6000
Processed variant 7000
Processed variant 8000
Processed variant 9000
Processed variant 10000
Processed variant 11000
Processed variant 12000
Processed variant 13000
Processed variant 14000
Processed variant: 14880
## All variants processed
strata<- read.table("strata", header=TRUE)

rad.u <- vcfR2genind(my_vcf_u, strata = strata, pop = c(rep("BIS", 10),rep("GB", 10),rep("NAR", 10), rep("AUS", 10),rep("CHN", 10),rep("COL", 10),rep("ESP", 10),rep("FIN", 10),rep("GBR", 10),rep("LIO", 10)), sex=sex, na.rm=TRUE)

## // GENIND OBJECT ///////////
##
## // 40 individuals; 14,880 loci; 29,760 alleles; size: 13.6 Mb
##
## // Basic content
## @tab: 40 x 29760 matrix of allele counts
## @loc.n.all: number of alleles per locus (range: 2-2)
## @loc.fac: locus factor for the 29760 columns of @tab
## @all.names: list of allele names for each locus
## @ploidy: ploidy of each individual (range: 2-2)
## @type: codom
## @call: adegenet::df2genind(X = t(x), sep = sep, pop = ..2, strata = ..1)
##
## // Optional content
## @pop: population of each individual (group size range: 10-10)
## @strata: a data frame with 12 columns ( Individual, Population, Latitude, Longitude, Distance, SE
stratated.u <- strata(rad.u, formula= Population, combine = TRUE, just.strats)
stratated.u@other <- select(info, Latitude, Longitude, Distance, SE, Temperature, Salinity, pH, Chlorophylla, Chl.a, Turbidity)

stratated.u

## // GENIND OBJECT ///////////
##
## // 40 individuals; 14,880 loci; 29,760 alleles; size: 13.6 Mb
##
## // Basic content
## @tab: 40 x 29760 matrix of allele counts
## @loc.n.all: number of alleles per locus (range: 2-2)
## @loc.fac: locus factor for the 29760 columns of @tab
## @all.names: list of allele names for each locus
```

```

##  @ploidy: ploidy of each individual (range: 2-2)
##  @type: codom
##  @call: adegenet::df2genind(X = t(x), sep = sep, pop = ..2, strata = ..1)
##
## // Optional content
##  @pop: population of each individual (group size range: 10-10)
##  @strata: a data frame with 1 columns ( Population )
##  @other: a list containing: Latitude Longitude Distance SE Temperature Salinity pH Chlorophy
Make hierfstat object
hf.filt <- genind2hierfstat(rad.filt, pop = c(rep("BIS", 10),rep("GB", 10),rep("NAR", 10), rep("PVD",10))
hf.u <- genind2hierfstat(rad.u, pop = c(rep("BIS", 10),rep("GB", 10),rep("NAR", 10), rep("PVD",10)))
hf.u <- hf.u$hierfstat.no.imputation

```

## Estimating effective migration surfaces (EEMS)

The program `runeems_snps` implements the EEMS method for analyzing spatial population structure. This version uses the pairwise genetic dissimilarity matrix computed from SNP data.

Here is the code used to run `runeems_snps` in both RStudio and command-line. Detailed steps for running this program can be accessed here. Input files are created in RStudio or manually (see below) and saved to the directory where `runeems_snps` will be run. The program is then run in the command-line. The outputs are then plotted in RStudio.

`runeems_snps` requires three data input files that have the same file name but different extension. The description below assumes that `datopath` is the full path + the file name (but without the extension).

### 1. `datopath.diffs`

`datopath.diffs` is the matrix of average pairwise genetic dissimilarities. This can be computed with `bed2diffs` from genetic data in plink binary format.

The dissimilarity matrix is nonnegative, symmetric, with 0s on the main diagonal. These conditions are necessary but not sufficient for `diffs` to be a valid dissimilarity matrix. Mathematically, `diffs` should be conditionally negative definite.

Steps for generating `datopath.diffs` are completed in RStudio.

```

# V1 methotd to get diffs matrix, preferred
bed2diffs_v1 <- function(Geno) {
  nIndiv <- nrow(Geno)
  nSites <- ncol(Geno)
  Diffs <- matrix(0, nIndiv, nIndiv)

  for (i in seq(nIndiv - 1)) {
    for (j in seq(i + 1, nIndiv)) {
      x <- Geno[i, ]
      y <- Geno[j, ]
      Diffs[i, j] <- mean((x - y)^2, na.rm = TRUE)
      Diffs[j, i] <- Diffs[i, j]
    }
  }
  Diffs
}

```

```

# V2 method to get .diffs matrix, only if V1 doesn't work
bed2diffs_v2 <- function(Geno) {
  nIndiv <- nrow(Geno)
  nSites <- ncol(Geno)
  Miss <- is.na(Geno)
  ## Impute NAs with the column means (= twice the allele frequencies)
  Mean <- matrix(colMeans(Geno, na.rm = TRUE), ## a row of means
                  nrow = nIndiv, ncol = nSites, byrow = TRUE) ## a matrix with nIndiv identical rows of means
  Mean[Miss == 0] <- 0 ## Set the means that correspond to observed genotypes to 0
  Geno[Miss == 1] <- 0 ## Set the missing genotypes to 0 (used to be NA)
  Geno <- Geno + Mean
  ## Compute similarities
  Sim <- Geno %*% t(Geno) / nSites
  SelfSim <- diag(Sim) ## self-similarities
  vector1s <- rep(1, nIndiv) ## vector of 1s
  ## This chunk generates a `diffs` matrix
  Diffs <- SelfSim %*% t(vector1s) + vector1s %*% t(SelfSim) - 2 * Sim
  Diffs
}

geno <- stratted.filt@tab

# Get rid of non-biallelic loci
multi.loci <- names(which(stratted.filt@loc.n.all != 2))
multi.cols <- which(grepl(paste0("^", multi.loci, "\\\\.\\\\d+$", collapse = "|"), colnames(geno)))
if (length(multi.cols)) geno <- geno[, -multi.cols]
nloci <- dim(geno)[2] / 2
dim(geno)

## [1]      40 148548
stopifnot(identical(stratted.filt@type, 'codom'))

# bed2diffs functions
diffs.v1 <- bed2diffs_v1(geno)
diffs.v2 <- bed2diffs_v2(geno)
# Round to 6 digits
diffs.v1 <- round(diffs.v1, digits = 6)
diffs.v2 <- round(diffs.v2, digits = 6)

```

Check that the dissimilarity matrix has one positive eigenvalue and nIndiv-1 negative eigenvalues, as required by a full-rank Euclidean distance matrix. If the V1 method does not make a Euclidean matrix, you must use V2.

```

tail(sort(round(eigen(diffs.v1)$values, digits = 2)))

## [1] -0.49 -0.48 -0.48 -0.48 -0.47 21.06
tail(sort(round(eigen(diffs.v2)$values, digits = 2)))

## [1] -0.48 -0.48 -0.47 -0.47 -0.46 20.73
# Set suffix for EEMS input files
suf <- "neutraldata-filt"

# This saves the file to directory
write.table(diffs.v1, paste(suf, ".v1.diffs", sep = ""),
            col.names = FALSE, row.names = FALSE, quote = FALSE)

```

## 2. datapath.coord

datapath.coord are the sample coordinates, two coordinates per sample, one sample per line. The sampling locations should be given in the same order as the rows and columns of the dissimilarity matrix.

Steps for generating datapath.coord are completed in RStudio.

```
## Get gps coordinates from previously created info matrix
xOR.info <- dplyr::filter(info)
gps_matrix <- select(xOR.info,c("Longitude","Latitude"))

#write .coord file
write.table(gps_matrix, paste(suf,".v1.coord",sep=""),col.names = FALSE, row.names = FALSE,quote = FALSE)
```

## 3. datapath.outer

datapath.outer are the habitat coordinates, as a sequence of vertices that form a closed polygon. The habitat vertices should be listed counterclockwise and the first vertex should also be the last vertex, so that the outline is a closed ring. Otherwise, EEMS attempts to “correct” the polygon and prints a warning message.

datapath.outer is created manually in Excel,based on site coordinates gathered from Google Maps, copied into terminal using nano, and saved as the file with the appropriate extension.

\*\*runeems\_snps is then run in command-line following the steps documented in NB\_EEMS\_OutlierHap.md

Back in RStudio to plot runeems\_snps outputs

```
# Install rEEMSpots
library(rEEMSpots)

# Plotting EEMS after running runeems_snps
path = "./NB_EEMS_Neutral/"
dirs = c(paste0(path,"neutraldata-D200-chain1"), paste0(path,"neutraldata-D300-chain1"), paste0(path,"neutraldata-D400-chain1"))

eems.plots(mcmcpath = c(paste0(path,"./neutraldata-D200-chain1"), paste0(path,"neutraldata-D300-chain1"),
                        longlat = T,add.grid=F,add.outline = T,add.demes = T,
                        projection.in = "+proj=longlat +datum=WGS84",projection.out = "+proj=merc +datum=WGS84",
                        add.map = T,add.abline = T, add.r.squared = T))

## Input projection: +proj=longlat +datum=WGS84
## Output projection: +proj=merc +datum=WGS84
## Loading rgdal (required by projection.in)
## Loading rworldmap (required by add.map)
## Loading rworldxtra (required by add.map)

## Using the default DarkOrange to Blue color scheme, with 'white' as the midpoint color.
## It combines two color schemes from the 'dichromat' package, which itself is based on
## a collection of color schemes for scientific data graphics:
## Light A and Bartlein PJ (2004). The End of the Rainbow? Color Schemes for Improved Data
## Graphics. EOS Transactions of the American Geophysical Union, 85(40), 385.
## See also http://geog.uoregon.edu/datagraphics/color_scales.htm

## Using 'euclidean' distance to assign interpolation points to Voronoi tiles.

## Processing the following EEMS output directories :
## ./NB_EEMS_Neutral./neutraldata-D200-chain1./NB_EEMS_Neutral/neutraldata-D300-chain1./NB_EEMS_Neutral/neutraldata-D400-chain1

## Plotting effective migration surface (posterior mean of m rates)
```

```

## ./NB_EEMS_Neutral./neutraldata-D200-chain1
## ./NB_EEMS_Neutral/neutraldata-D300-chain1
## ./NB_EEMS_Neutral/neutraldata-D600-chain1

## Using the default DarkOrange to Blue color scheme, with 'white' as the midpoint color.
## It combines two color schemes from the 'dichromat' package, which itself is based on
## a collection of color schemes for scientific data graphics:
## Light A and Bartlein PJ (2004). The End of the Rainbow? Color Schemes for Improved Data
## Graphics. EOS Transactions of the American Geophysical Union, 85(40), 385.
## See also http://geog.uoregon.edu/datagraphics/color_scales.htm

## Plotting effective diversity surface (posterior mean of q rates)

## ./NB_EEMS_Neutral./neutraldata-D200-chain1
## ./NB_EEMS_Neutral/neutraldata-D300-chain1
## ./NB_EEMS_Neutral/neutraldata-D600-chain1

## Using the default DarkOrange to Blue color scheme, with 'white' as the midpoint color.
## It combines two color schemes from the 'dichromat' package, which itself is based on
## a collection of color schemes for scientific data graphics:
## Light A and Bartlein PJ (2004). The End of the Rainbow? Color Schemes for Improved Data
## Graphics. EOS Transactions of the American Geophysical Union, 85(40), 385.
## See also http://geog.uoregon.edu/datagraphics/color_scales.htm

## Plotting posterior probability trace

## ./NB_EEMS_Neutral./neutraldata-D200-chain1
## ./NB_EEMS_Neutral/neutraldata-D300-chain1
## ./NB_EEMS_Neutral/neutraldata-D600-chain1

## Plotting average dissimilarities within and between demes

## ./NB_EEMS_Neutral./neutraldata-D200-chain1
## ./NB_EEMS_Neutral/neutraldata-D300-chain1

## Using 'euclidean' distance to assign interpolation points to Voronoi tiles.
##
##
##
## Using 'euclidean' distance to assign interpolation points to Voronoi tiles.

## EEMS results for at least two different population grids
```

The plots are saved to the same directory where `runeems_snps` was run.

## Pairwise Fst

```

fst.mat <- pairwise.WCfst(hf.filt)

gindF.fst.mat.triN <- as.matrix(fst.mat)
colnames(gindF.fst.mat.triN) <- pop_order
rownames(gindF.fst.mat.triN) <- pop_order

meltedN <- melt(gindF.fst.mat.triN, na.rm =TRUE)
round(gindF.fst.mat.triN,4)
```

```

##          BIS      GB      NAR      PVD
## BIS      NA  0.0005  0.0088 -0.0004
## GB       0.0005      NA  0.0105 -0.0001
## NAR     0.0088  0.0105      NA  0.0102
## PVD    -0.0004 -0.0001  0.0102      NA
summary(meltedN$value)

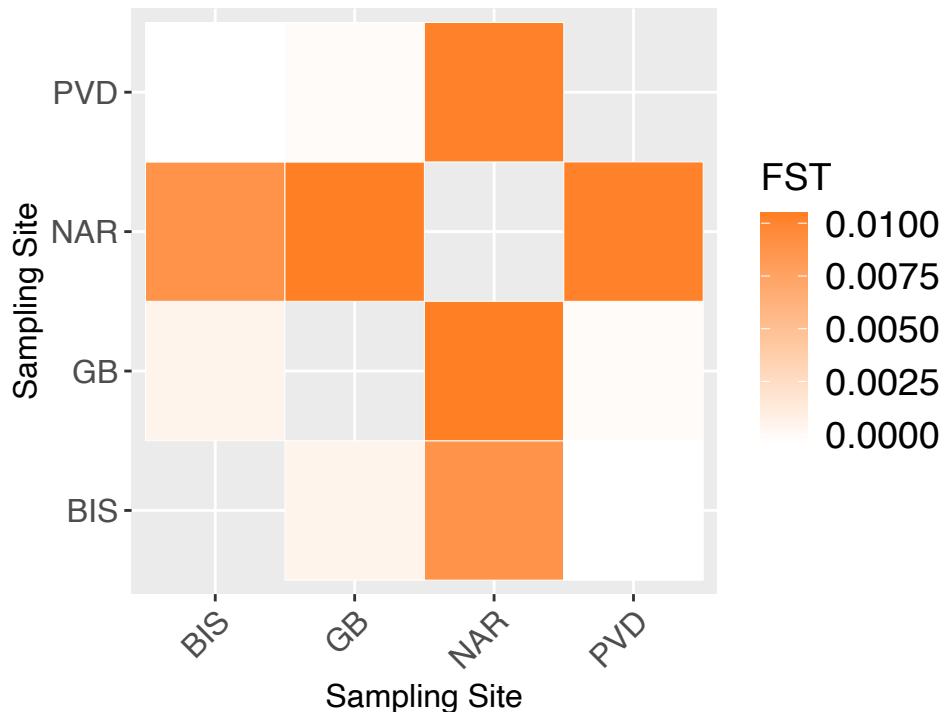
##           Min.    1st Qu.   Median    Mean   3rd Qu.    Max.
## -4.100e-04 -7.949e-05  4.658e-03  4.919e-03  1.022e-02  1.047e-02

#Plotting Pairwise fst
neutral <- ggplot(data = meltedN, aes(Var2, Var1, fill = value)) + geom_tile(color = "white") +
  scale_fill_gradient(low = "white", high = "chocolate1", name="FST") +
  ggtitle(expression(atop("Pairwise FST, WC (1984) Neutral SNPs", atop(italic("N = 40, L = 74,520"), ""))) +
  labs(x = "Sampling Site", y = "Sampling Site") +
  theme(axis.text.x = element_text(angle = 45, vjust = 1, size = 12, hjust = 1), axis.text.y = element_text(size = 12), legend.text = element_text(size = 14), legend.title = element_text(size = 14)) +
  theme(plot.title = element_text(size = 14)) +
  coord_fixed()
neutral

```

Pairwise FST, WC (1984) Neutral SNPs

$N = 40, L = 74,520$



### Genetic diversity (observed and expected heterozygosity)

```

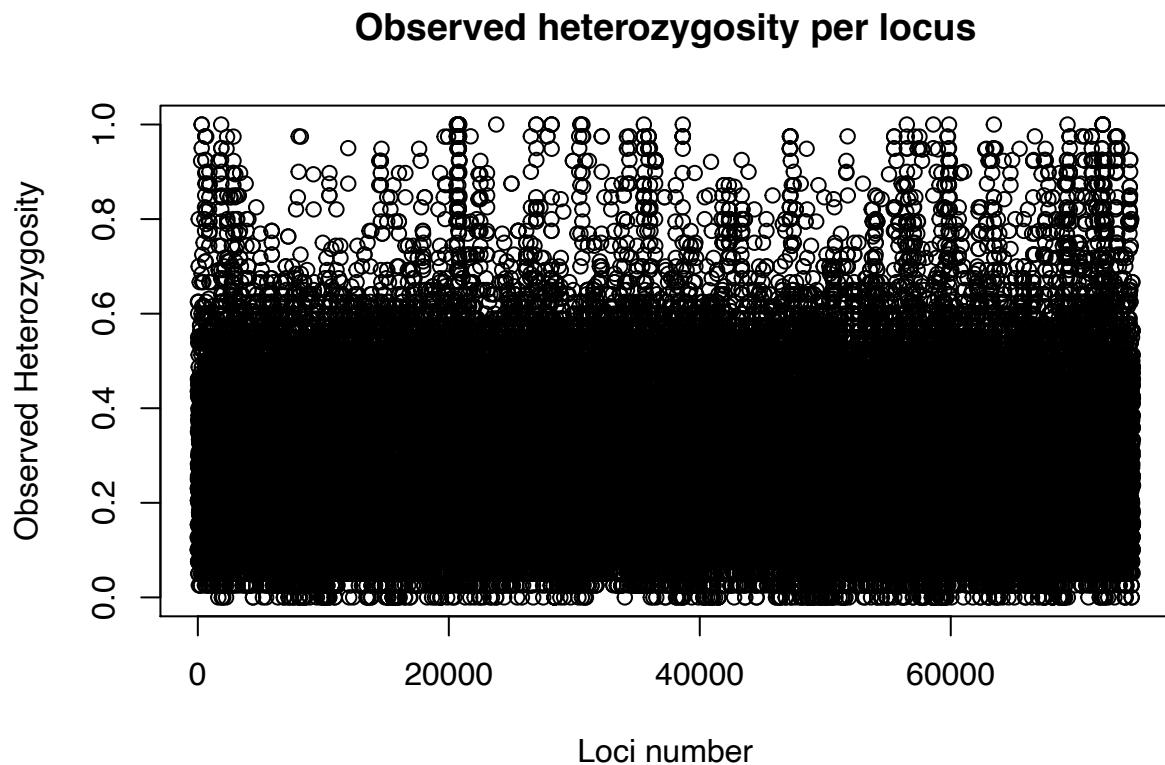
comb <- summary(stratated.filt)
names(comb)

```

```

## [1] "n"          "n.by.pop"   "loc.n.all"  "pop.n.all"  "NA.perc"    "Hobs"
## [7] "Hexp"
plot(comb$Hobs, xlab="Loci number", ylab="Observed Heterozygosity",
     main="Observed heterozygosity per locus")

```

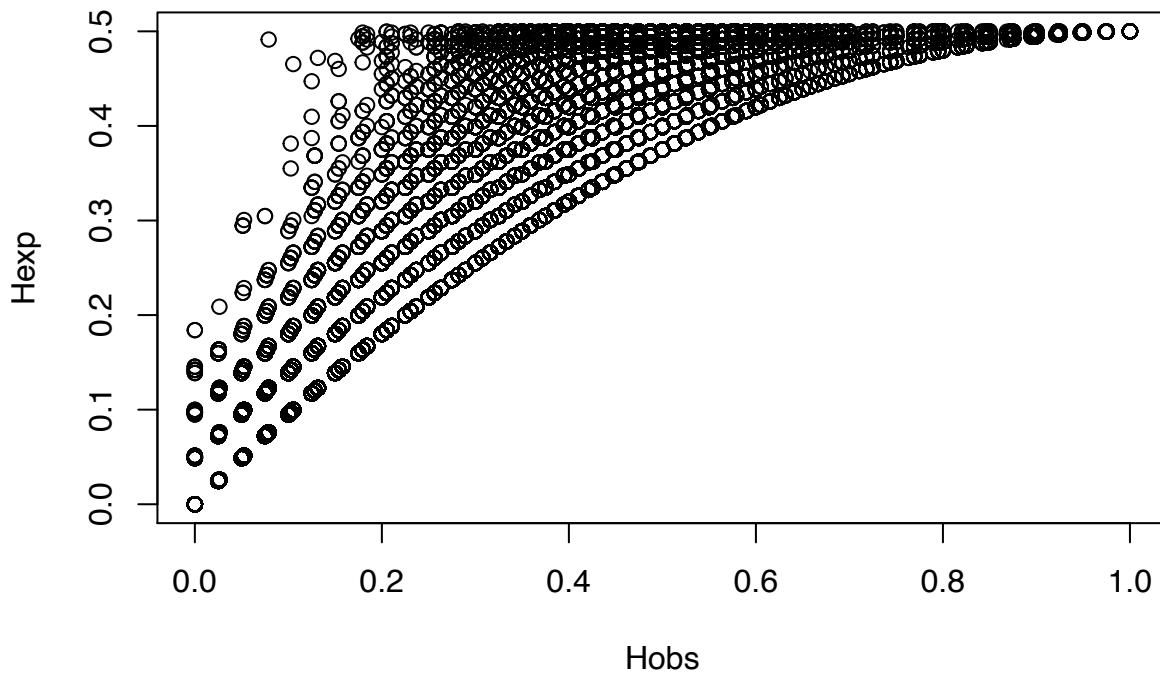


```

plot(comb$Hobs, comb$Hexp, xlab="Hobs", ylab="Hexp",
     main="Expected heterozygosity as a function of observed heterozygosity per locus")

```

## Expected heterozygosity as a function of observed heterozygosity per I



```
bartlett.test(list(comb$Hexp, comb$Hobs)) # a test : H0: Hexp = Hobs
```

```
##  
##  Bartlett test of homogeneity of variances  
##  
## data: list(comb$Hexp, comb$Hobs)  
## Bartlett's K-squared = 2099.3, df = 1, p-value < 2.2e-16
```

*Significant difference between Observed and expected heterozygosity.*

```
basicstat <- basic.stats(hf.filt, diploid = TRUE, digits = 3)
```

```
as.data.frame(basicstat$overall)
```

```
##      basicstat$overall  
##  Ho          0.279  
##  Hs          0.273  
##  Ht          0.274  
##  Dst         0.001  
##  Htp         0.274  
##  Dstp        0.001  
##  Fst         0.004  
##  Fstp        0.005  
##  Fis         -0.021  
##  Dest        0.002
```

```
# get bootstrap confidence values for Fis  
boot <- boot.ppfis(hf.filt, nboot = 1000)
```

```

boot5 <- boot.ppfis(hf.filt,nboot = 1000,quant = 0.5)

# add latitude for each population
latitude = c(41.545, 41.654, 41.505, 41.816)

# add longitude for each population
longitude = c(-71.431, -71.445, -71.453, -71.391)

# add distance for each population
distance = c(4.76, 0.47, 15.41, 1.49)

# add sewage effluent for each population
sewage = c(8.82, 14.60, 2.03, 59.86)

# add temperature for each population
temperature = c(23, 24, 25, 23)

# add salinity for each population
salinity = c(30, 28, 18, 25)

# add pH for each population
pH = c(7.9, 7.4, 7.6, 7.4)

# add Chlorophylla for each population
Chlor_a = c(4.9, 18.8, 4.6, 8.1)

# add DO for each population
DO = c(8.2, 5.7, 7, 4.9)

# combine all pop statistics
colnames(basicstat$Ho) <- pop_order
Ho <- colMeans(basicstat$Ho,na.rm = T)
He <- colMeans(basicstat$He,na.rm = T)
Fis<- boot5$fis.ci$ll
y <- cbind(pop_order, Ho, He, Fis, boot$fis.ci, latitude, longitude, distance, sewage, temperature, salinity, pH, Chlor_a, DO)

##      pop_order      Ho      He      Fis      ll      hl latitude longitude
## BIS        BIS 0.2792641 0.2730862 -0.0226 -0.0257 -0.0196 41.545 -71.431
## GB         GB 0.2831958 0.2745212 -0.0316 -0.0344 -0.0287 41.654 -71.445
## NAR        NAR 0.2779964 0.2703649 -0.0283 -0.0313 -0.0250 41.505 -71.453
## PVD        PVD 0.2737460 0.2729991 -0.0029 -0.0057 0.0000 41.816 -71.391
##      distance sewage temperature salinity     pH Chlor_a     DO
## BIS       4.76   8.82        23      30 7.9    4.9 8.2
## GB        0.47  14.60        24      28 7.4   18.8 5.7
## NAR      15.41   2.03        25      18 7.6    4.6 7.0
## PVD      1.49  59.86        23      25 7.4    8.1 4.9

summary(He)

##      Min. 1st Qu. Median Mean 3rd Qu. Max.
##  0.2704  0.2723  0.2730  0.2727  0.2734  0.2745

summary(Fis)

##      Min. 1st Qu. Median Mean 3rd Qu. Max.
## -0.03160 -0.02913 -0.02545 -0.02135 -0.01767 -0.00290

```

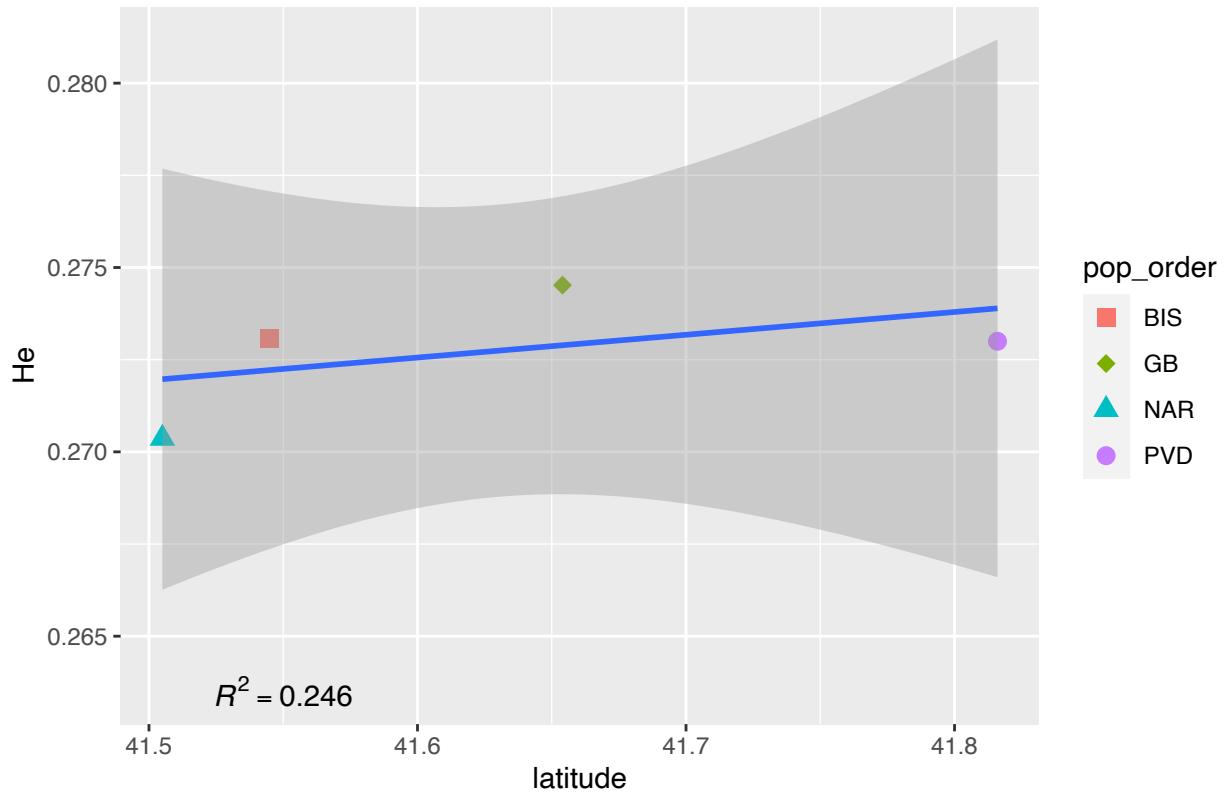
```

# Plot He vs Latitude
R2 = round(summary(lm(y$He ~ y$latitude))$r.squared, 4)
ggplot(y, aes(x = latitude, y = He)) + geom_point(aes(shape=pop_order, color=pop_order), size = 3) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggtitle("Expected heterozygosity vs Latitude, Neutral SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R2), x=41.55, y=0.2635, parse=T) +
  scale_x_continuous()

## `geom_smooth()` using formula 'y ~ x'

```

Expected heterozygosity vs Latitude, Neutral SNPs



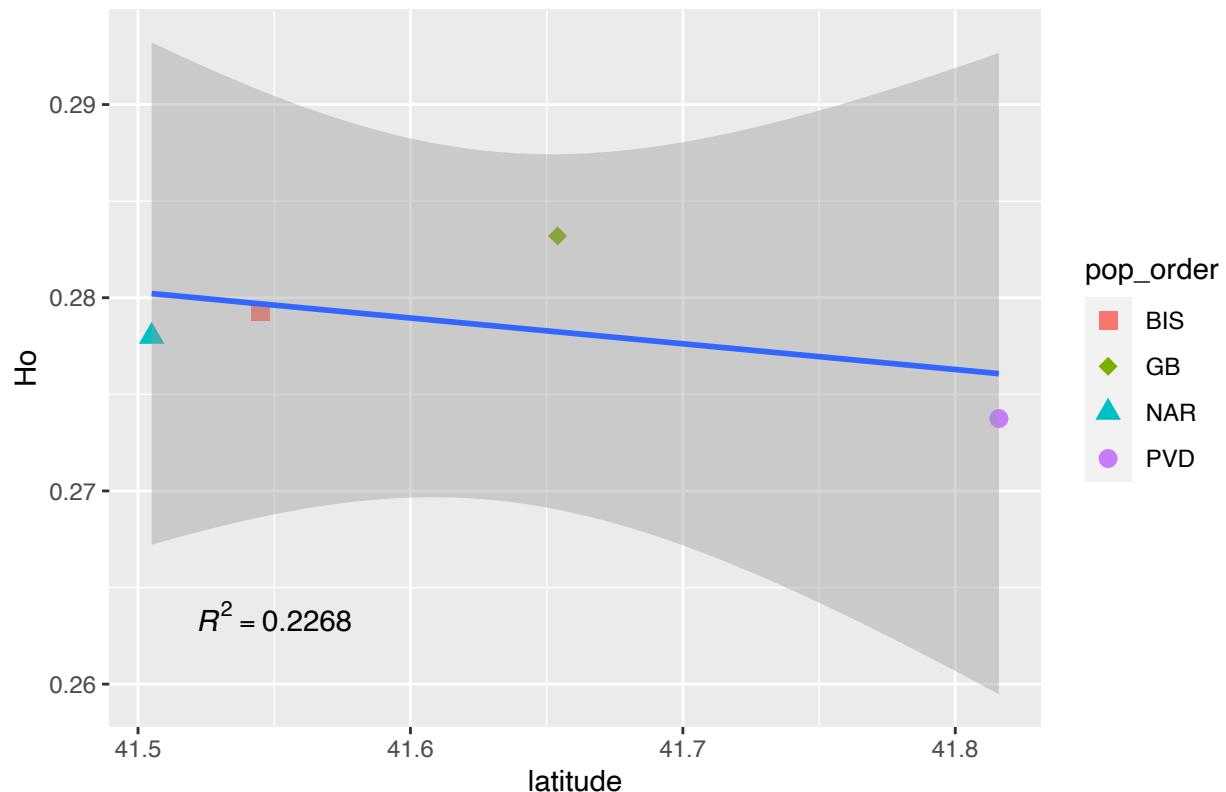
```

# Plot Ho vs Latitude
R2 = round(summary(lm(y$Ho ~ y$latitude))$r.squared, 4)
ggplot(y, aes(x = latitude, y = Ho)) + geom_point(aes(shape=pop_order, color=pop_order), size = 3) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggtitle("Observed heterozygosity vs Latitude, Neutral SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R2), x=41.55, y=0.2635, parse=T) +
  scale_x_continuous()

## `geom_smooth()` using formula 'y ~ x'

```

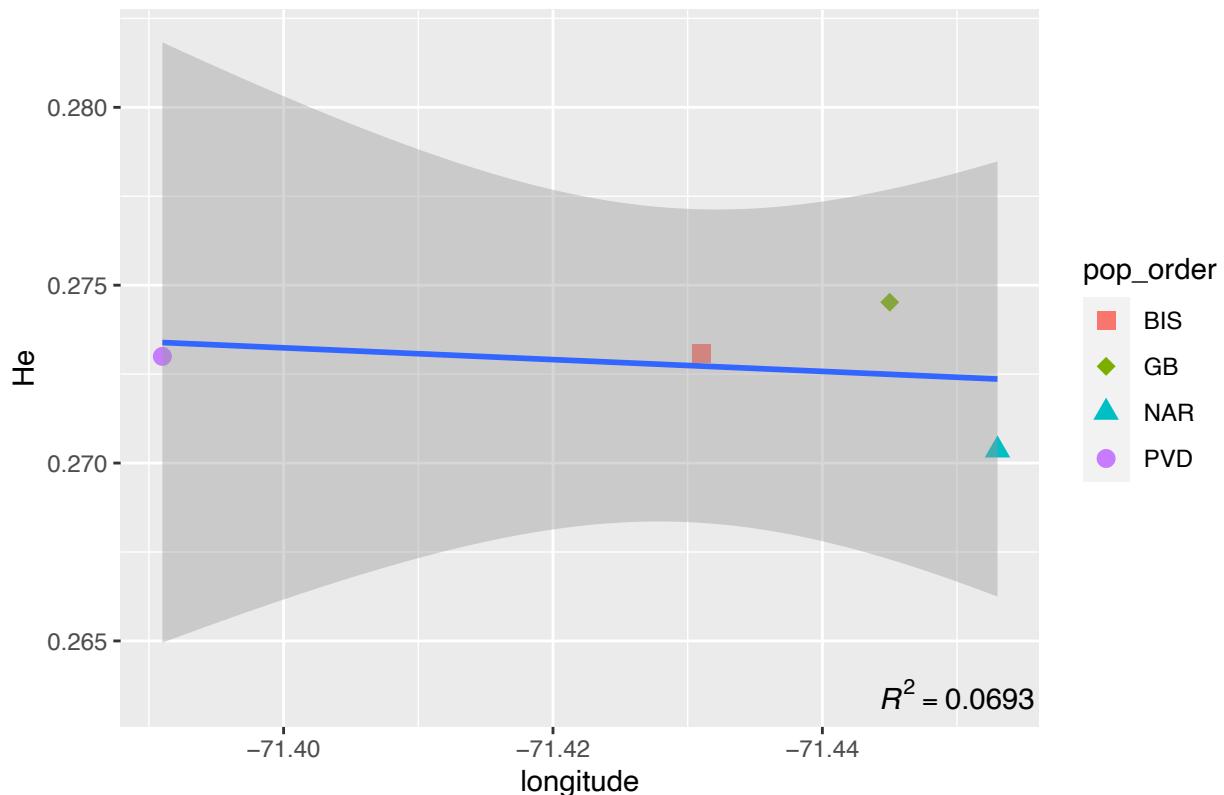
### Observed heterozygosity vs Latitude, Neutral SNPs



```
#Plot He vs Longitude
R3 = round(summary(lm(y$He ~ y$longitude))$r.squared, 4)
ggplot(y, aes(x = longitude, y = He)) + geom_point(aes(shape=pop_order, color=pop_order), size = 3) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggttitle("Expected heterozygosity vs Longitude, Neutral SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R3), x=-71.45, y=0.2635, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

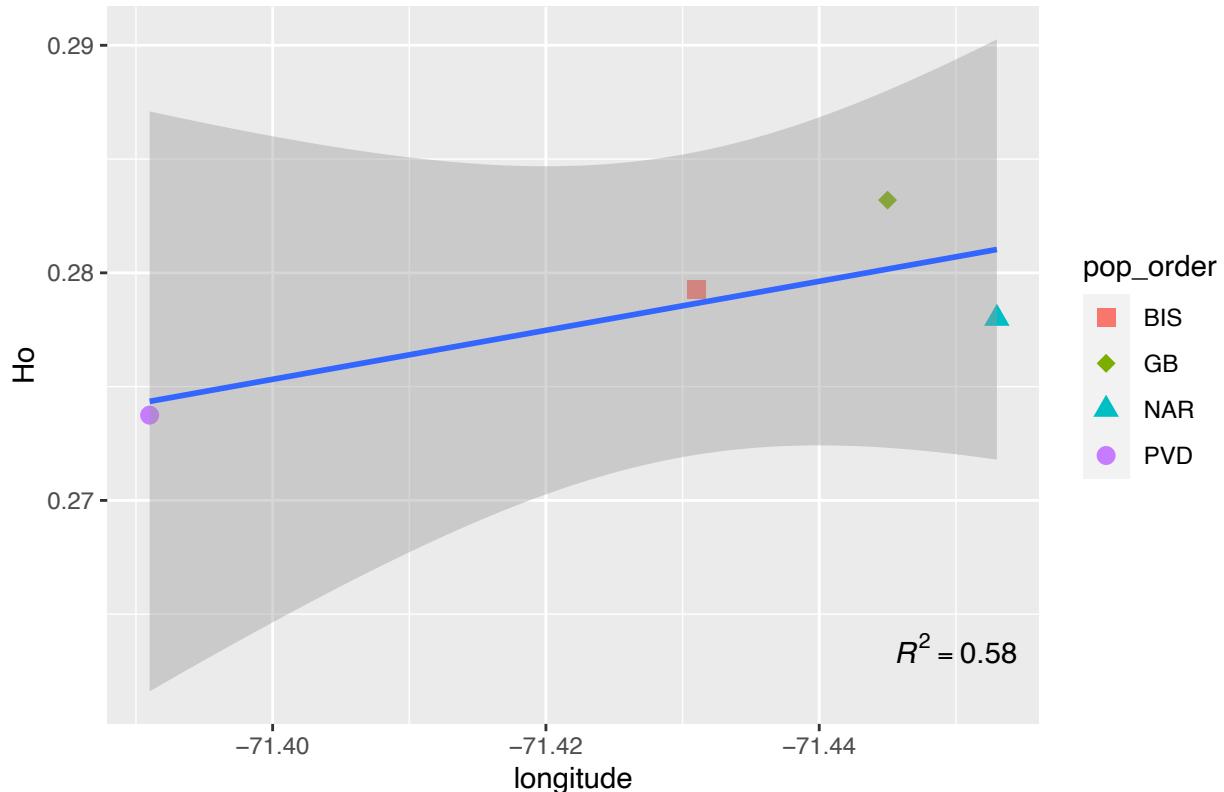
## Expected heterozygosity vs Longitude, Neutral SNPs



```
#Plot Ho vs Longitude
R3 = round(summary(lm(y$Ho ~ y$longitude))$r.squared, 4)
ggplot(y, aes(x = longitude, y = Ho)) + geom_point(aes(shape=pop_order, color=pop_order), size = 3) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  gtitle("Observed heterozygosity vs Longitude, Neutral SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R3), x=-71.45, y=0.2635, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

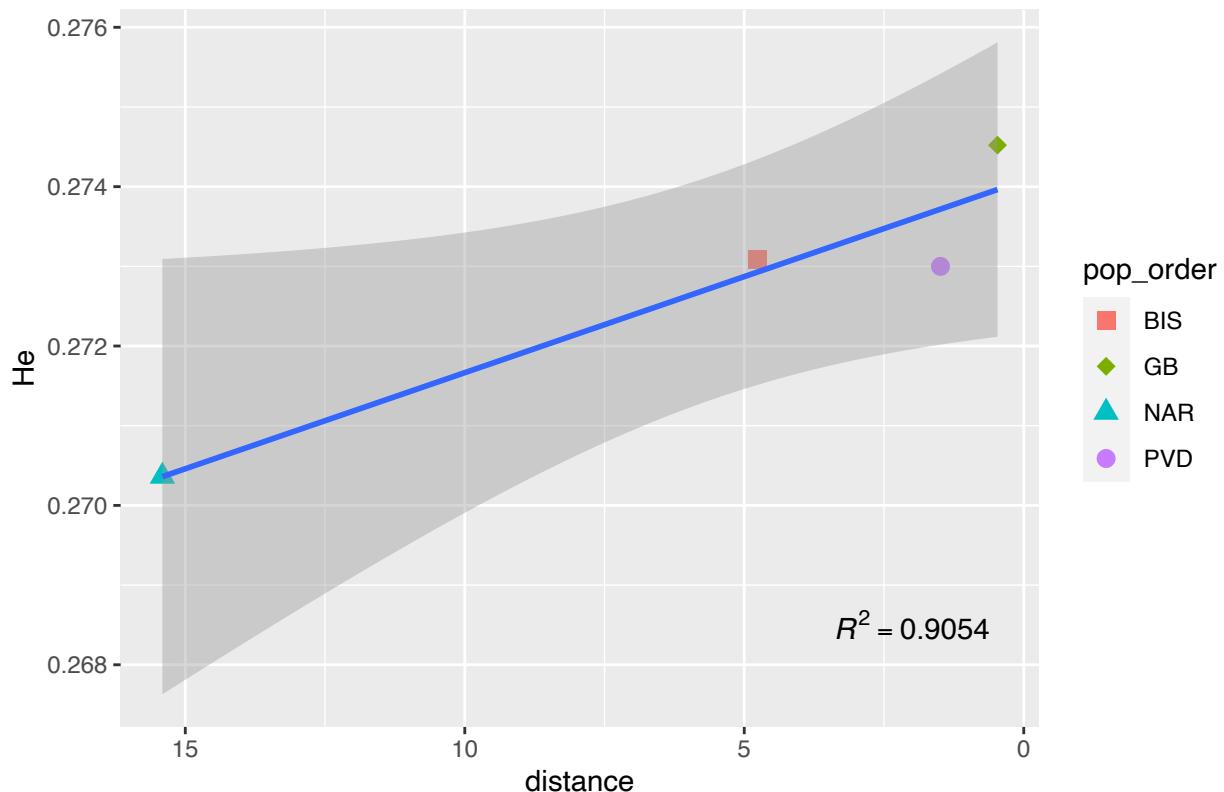
### Observed heterozygosity vs Longitude, Neutral SNPs



```
#Plot He vs Distance from sewage outflow
R4 = round(summary(lm(y$He ~ y$distance))$r.squared, 4)
ggplot(y, aes(x = distance, y = He)) + geom_point(aes(shape=pop_order, color=pop_order), size = 3) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggtitle("Expected heterozygosity vs Distance, Neutral SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R4), x=2, y=0.2685, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

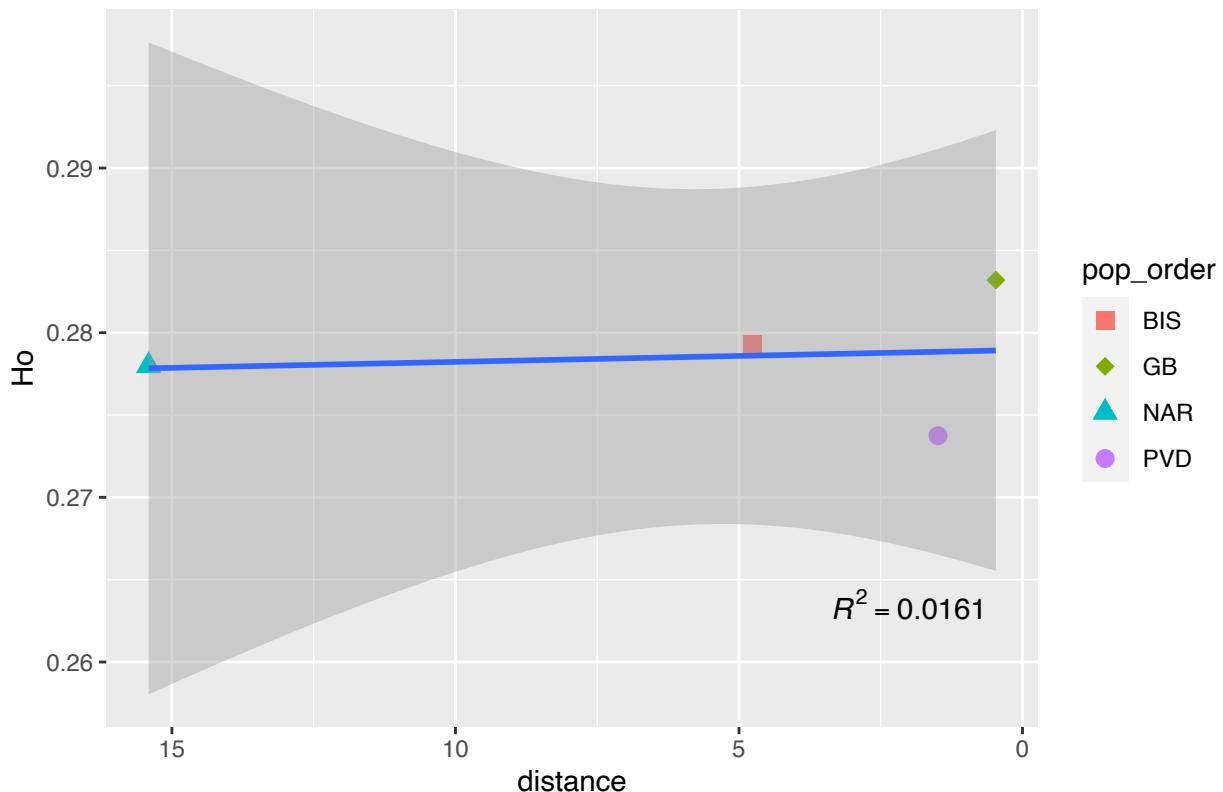
### Expected heterozygosity vs Distance, Neutral SNPs



```
#Plot Ho vs Distance from sewage outflow
R4 = round(summary(lm(y$Ho ~ y$distance))$r.squared, 4)
ggplot(y, aes(x = distance, y = Ho)) + geom_point(aes(shape=pop_order, color=pop_order), size = 3) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggtitle("Observed heterozygosity vs Distance, Neutral SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R4), x=2, y=0.2635, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

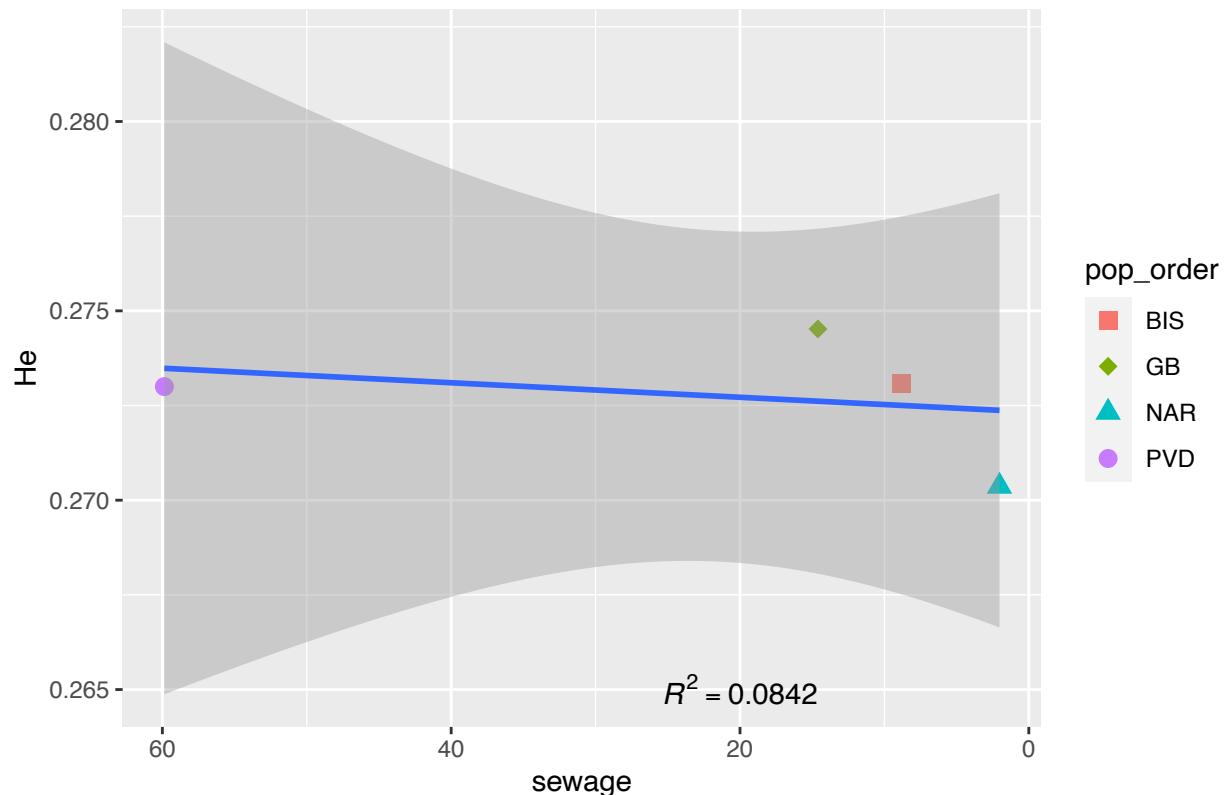
## Observed heterozygosity vs Distance, Neutral SNPs



```
#Plot He vs Sewage Effluent
R4 = round(summary(lm(y$He ~ y$sewage))$r.squared, 4)
ggplot(y, aes(x = sewage, y = He)) + geom_point(aes(shape=pop_order, color=pop_order), size = 3) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggttitle("Expected heterozygosity vs Sewage Effluent, Neutral SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R4), x=20, y=0.265, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

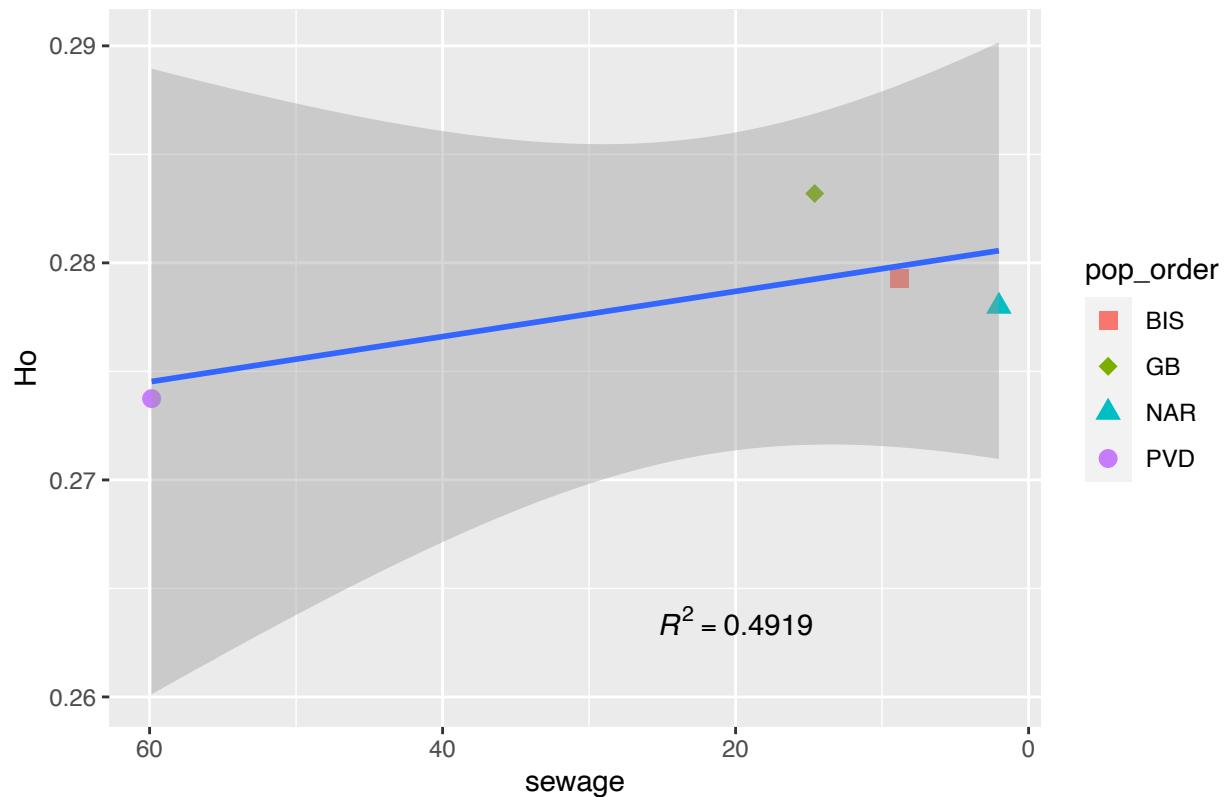
## Expected heterozygosity vs Sewage Effluent, Neutral SNPs



```
#Plot Ho vs Sewage Effluent
R4 = round(summary(lm(y$Ho ~ y$sewage))$r.squared, 4)
ggplot(y, aes(x = sewage, y = Ho)) + geom_point(aes(shape=pop_order, color=pop_order), size = 3) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggttitle("Observed heterozygosity vs Sewage Effluent, Neutral SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R4), x=20, y=0.2635, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

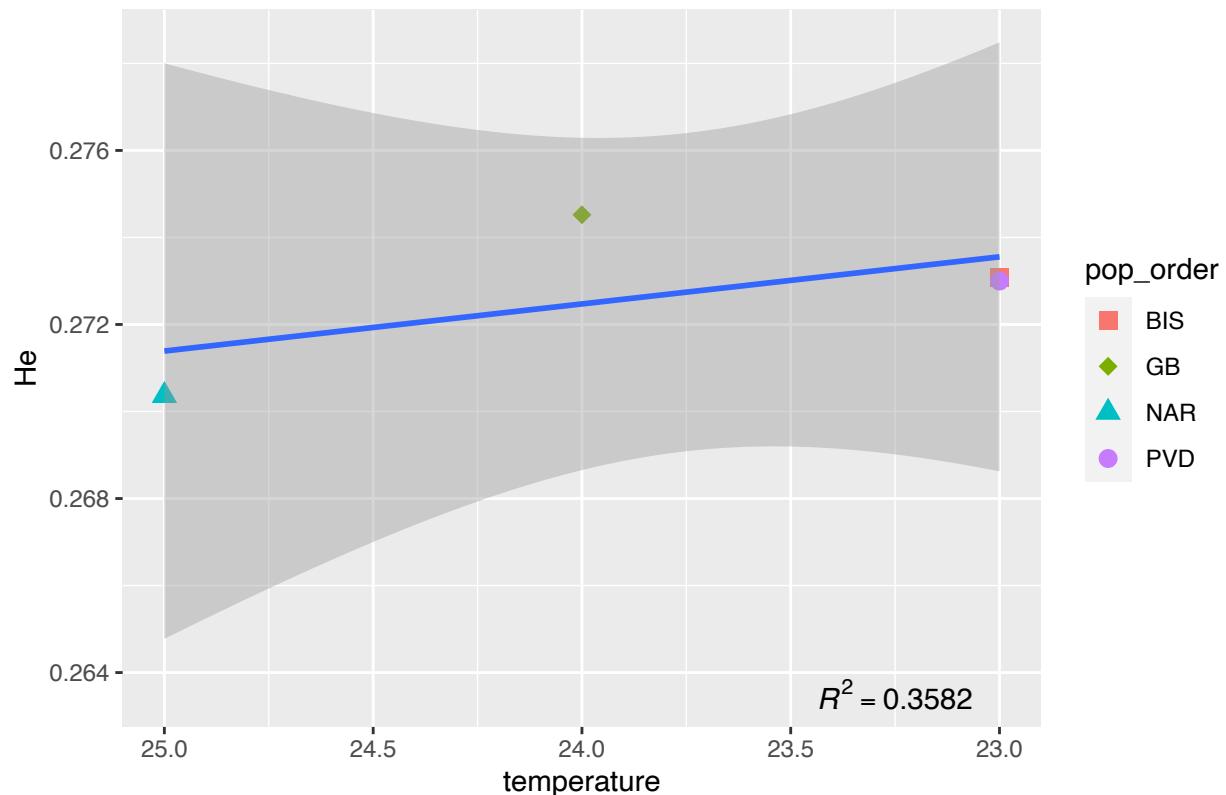
### Observed heterozygosity vs Sewage Effluent, Neutral SNPs



```
#Plot He vs Temperature
R4 = round(summary(lm(y$He ~ y$temperature))$r.squared, 4)
ggplot(y, aes(x = temperature, y = He)) + geom_point(aes(shape=pop_order, color=pop_order), size = 3) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggtitle("Expected heterozygosity vs Temperature, Neutral SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R4), x=23.25, y=0.2635, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

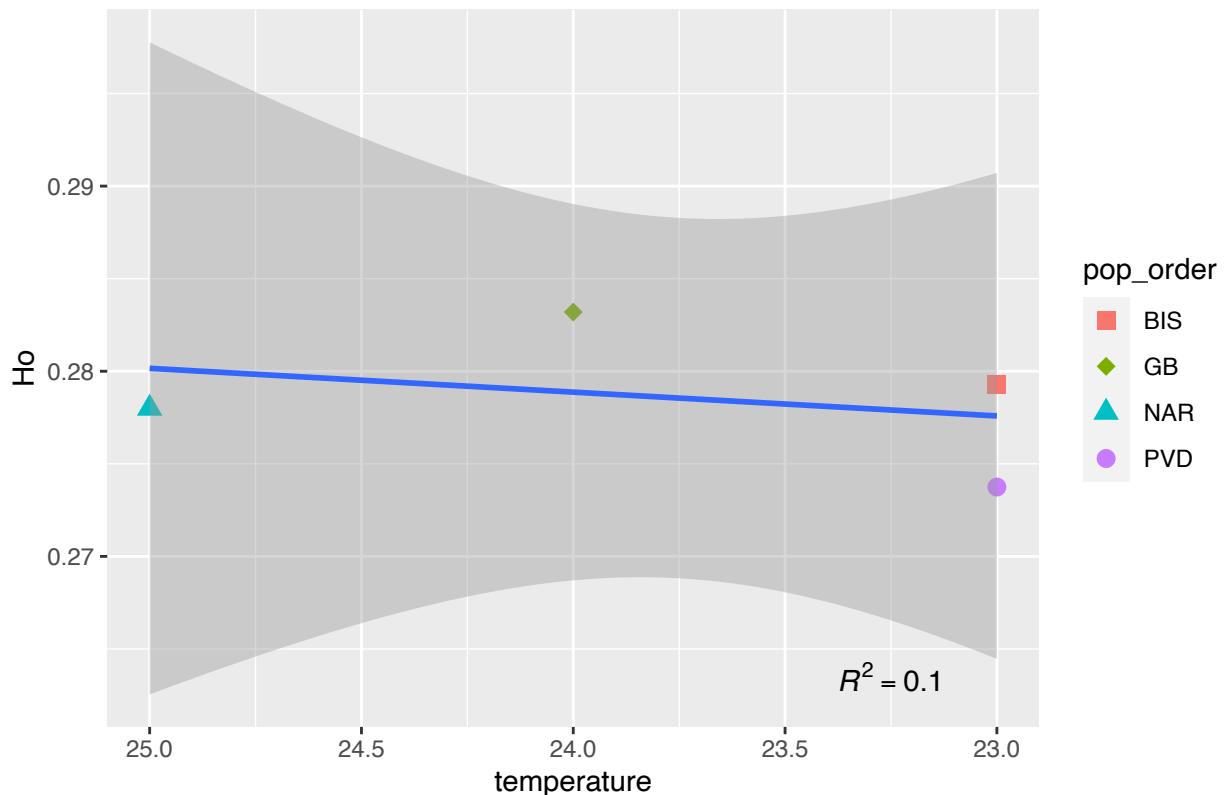
## Expected heterozygosity vs Temperature, Neutral SNPs



```
#Plot Ho vs Temperature
R4 = round(summary(lm(y$Ho ~ y$temperature))$r.squared, 4)
ggplot(y, aes(x = temperature, y = Ho)) + geom_point(aes(shape=pop_order, color=pop_order), size = 3) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggtitle("Observed heterozygosity vs Temperature, Neutral SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R4), x=23.25, y=0.2635, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

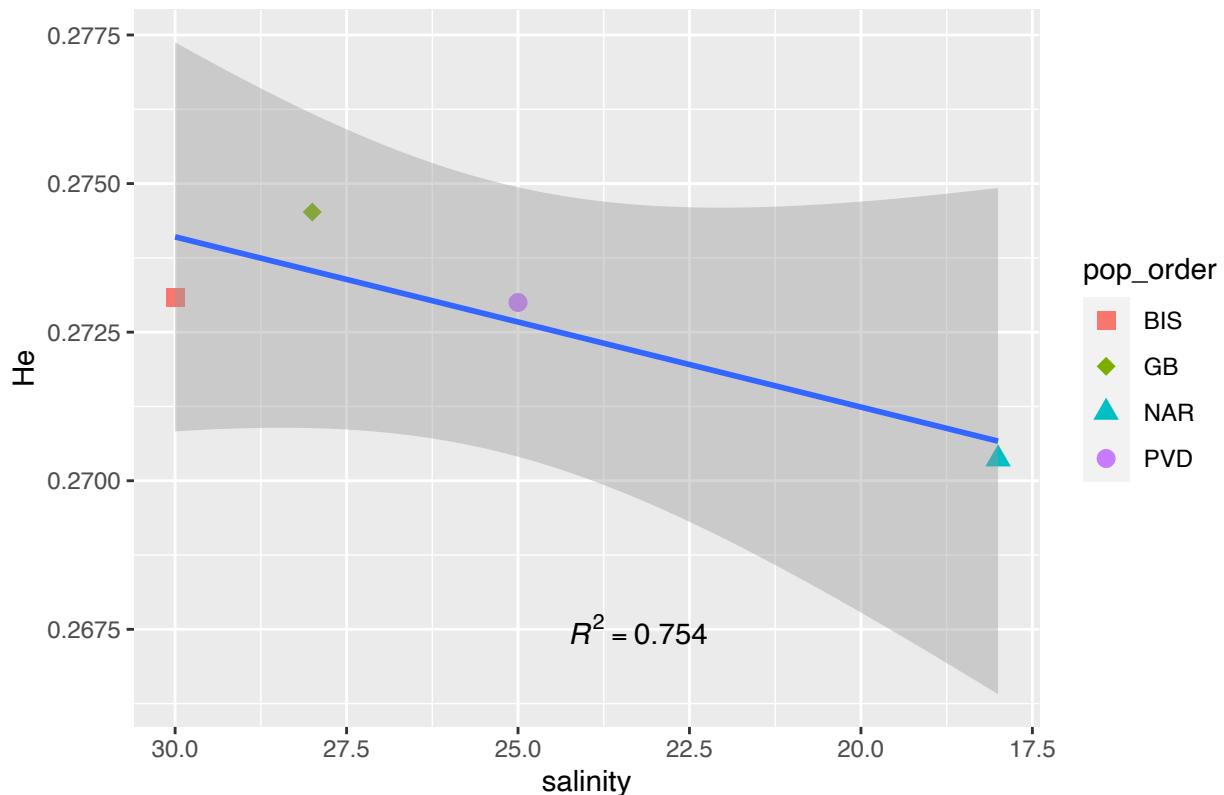
## Observed heterozygosity vs Temperature, Neutral SNPs



```
#Plot He vs Salinity
R4 = round(summary(lm(y$He ~ y$salinity))$r.squared, 4)
ggplot(y, aes(x = salinity, y = He)) + geom_point(aes(shape=pop_order, color=pop_order), size = 3) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggtitle("Expected heterozygosity vs Salinity, Neutral SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R4), x=23.25, y=0.2675, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

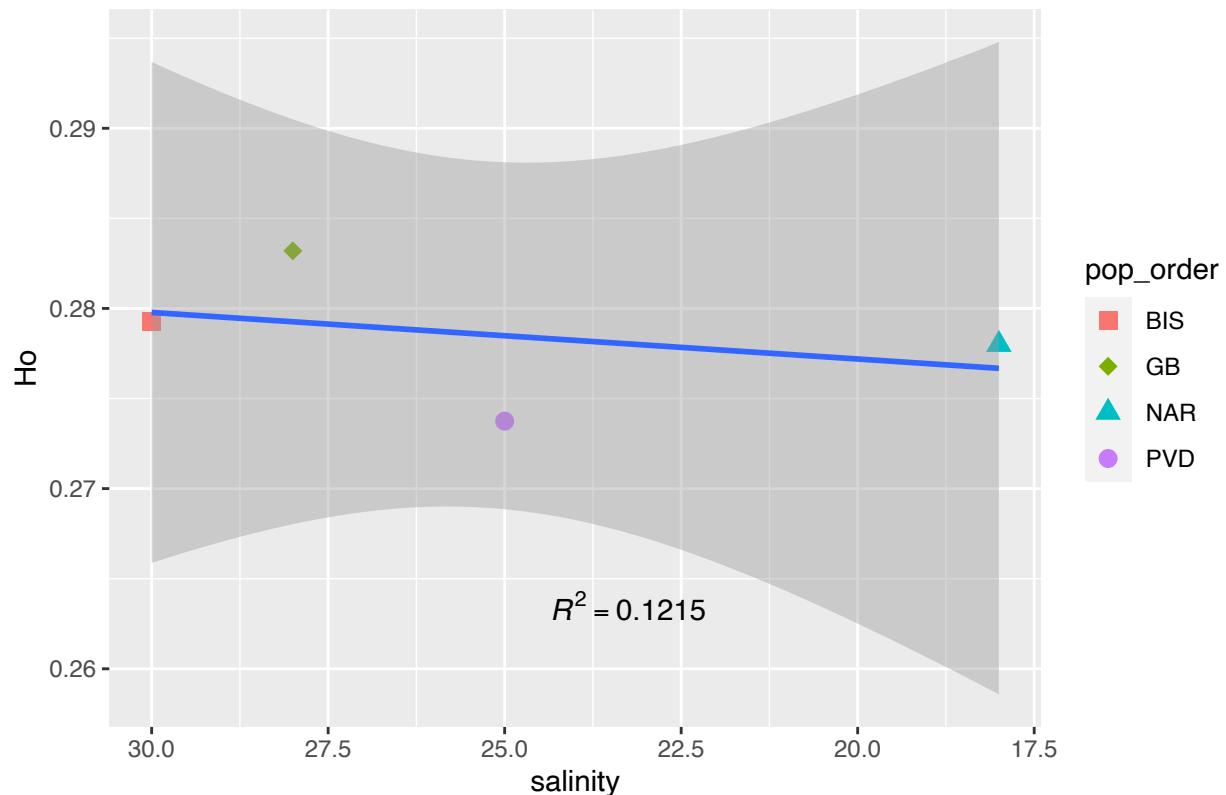
### Expected heterozygosity vs Salinity, Neutral SNPs



```
#Plot Ho vs Salinity
R4 = round(summary(lm(y$Ho ~ y$salinity))$r.squared, 4)
ggplot(y, aes(x = salinity, y = Ho)) + geom_point(aes(shape=pop_order, color=pop_order), size = 3) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggtitle("Observed heterozygosity vs Salinity, Neutral SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R4), x=23.25, y=0.2635, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

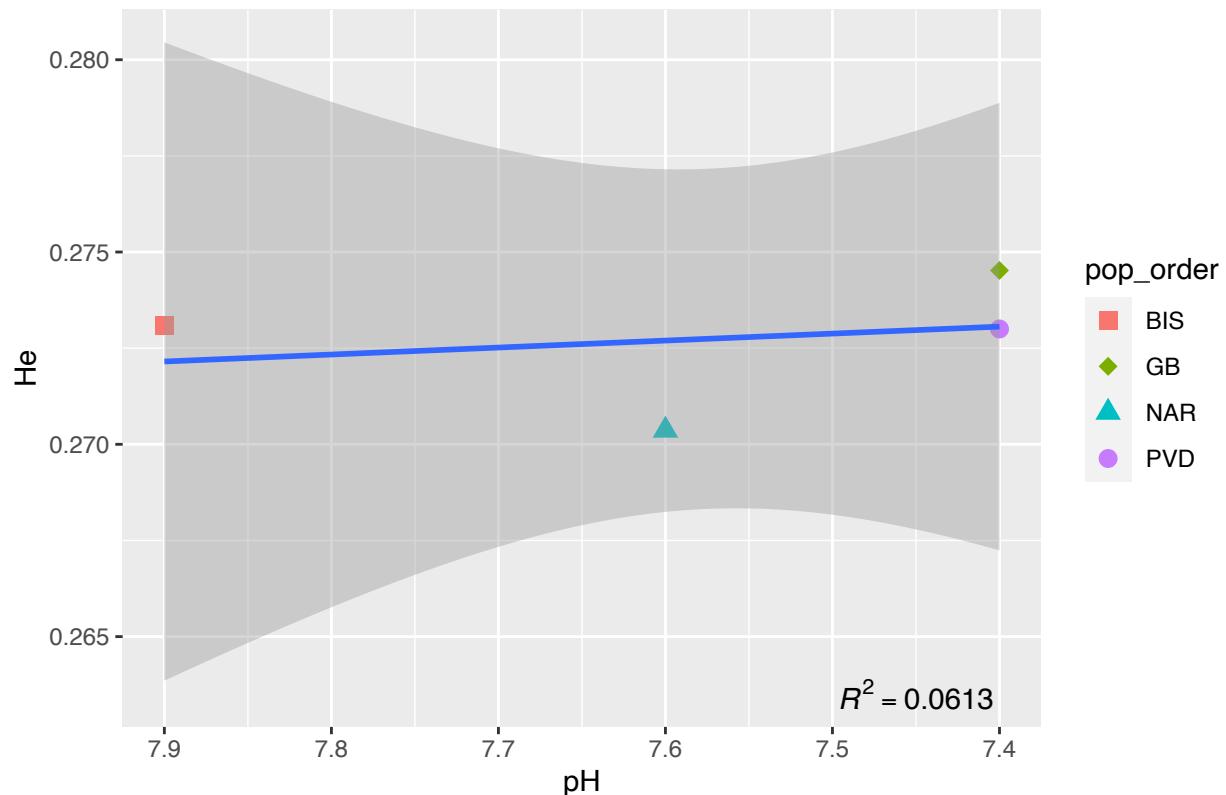
## Observed heterozygosity vs Salinity, Neutral SNPs



```
#Plot He vs pH
R4 = round(summary(lm(y$He ~ y$pH))$r.squared, 4)
ggplot(y, aes(x = pH, y = He)) + geom_point(aes(shape=pop_order, color=pop_order), size = 3) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggttitle("Expected heterozygosity vs pH, Neutral SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R4), x=7.45, y=0.2635, parse=T) +
  scale_x_reverse()

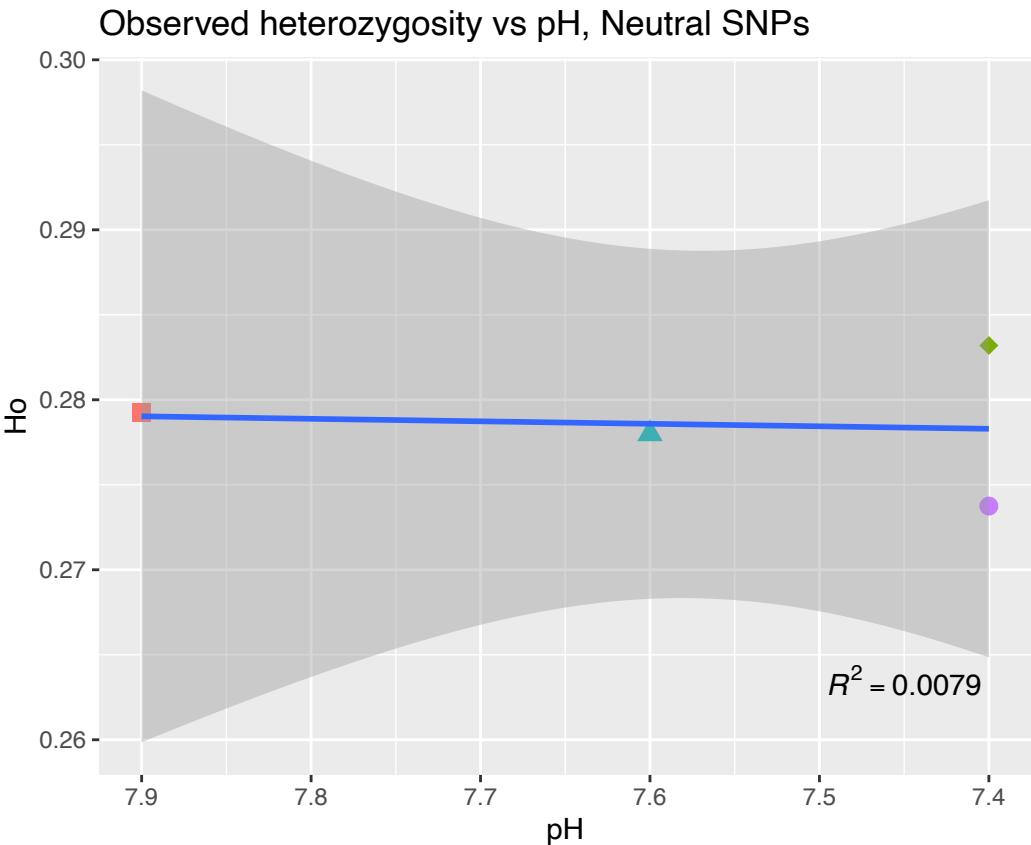
## `geom_smooth()` using formula 'y ~ x'
```

### Expected heterozygosity vs pH, Neutral SNPs



```
#Plot Ho vs pH
R4 = round(summary(lm(y$Ho ~ y$pH))$r.squared, 4)
ggplot(y, aes(x = pH, y = Ho)) + geom_point(aes(shape=pop_order, color=pop_order), size = 3) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggtitle("Observed heterozygosity vs pH, Neutral SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R4), x=7.45, y=0.2635, parse=T) +
  scale_x_reverse()

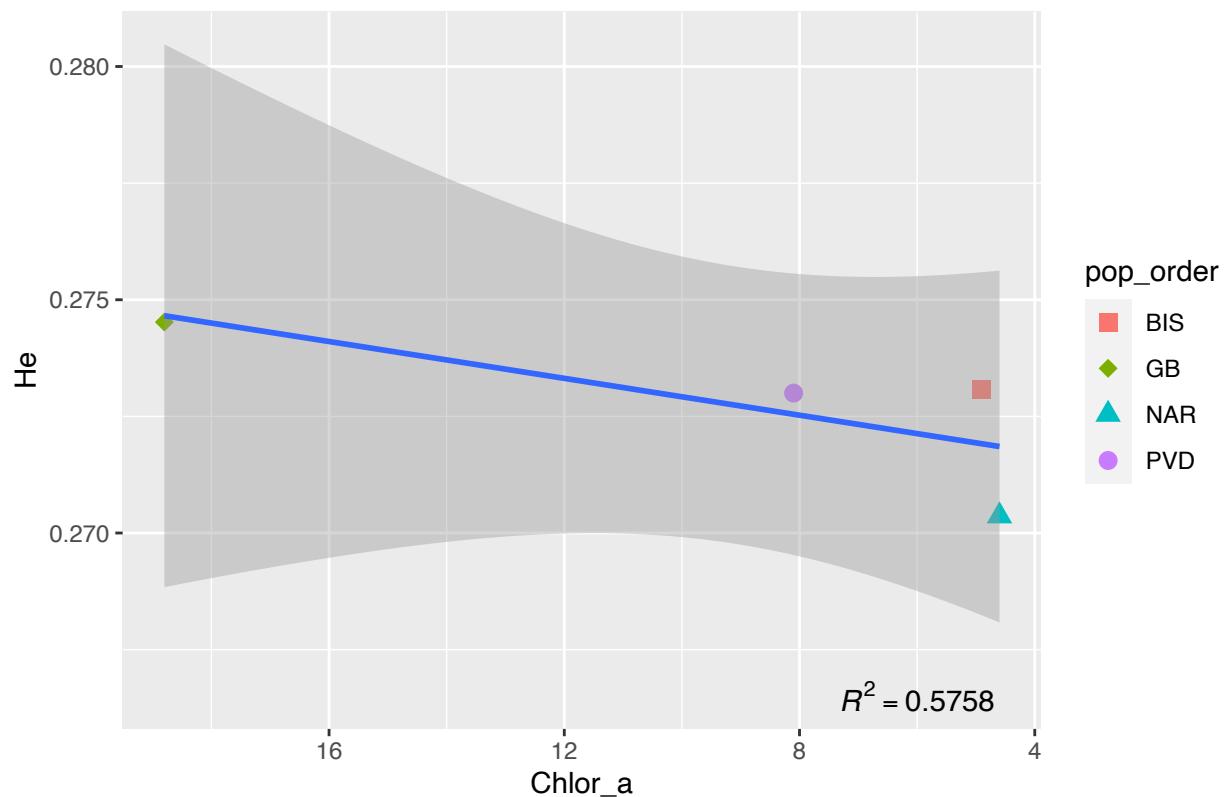
## `geom_smooth()` using formula 'y ~ x'
```



```
#Plot He vs Chlorophyll a
R4 = round(summary(lm(y$He ~ y$Chlor_a))$r.squared, 4)
ggplot(y, aes(x = Chlor_a, y = He)) + geom_point(aes(shape=pop_order, color=pop_order), size = 3) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggtitle("Expected heterozygosity vs Chlorophyll a, Neutral SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R4), x=6, y=0.2665, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

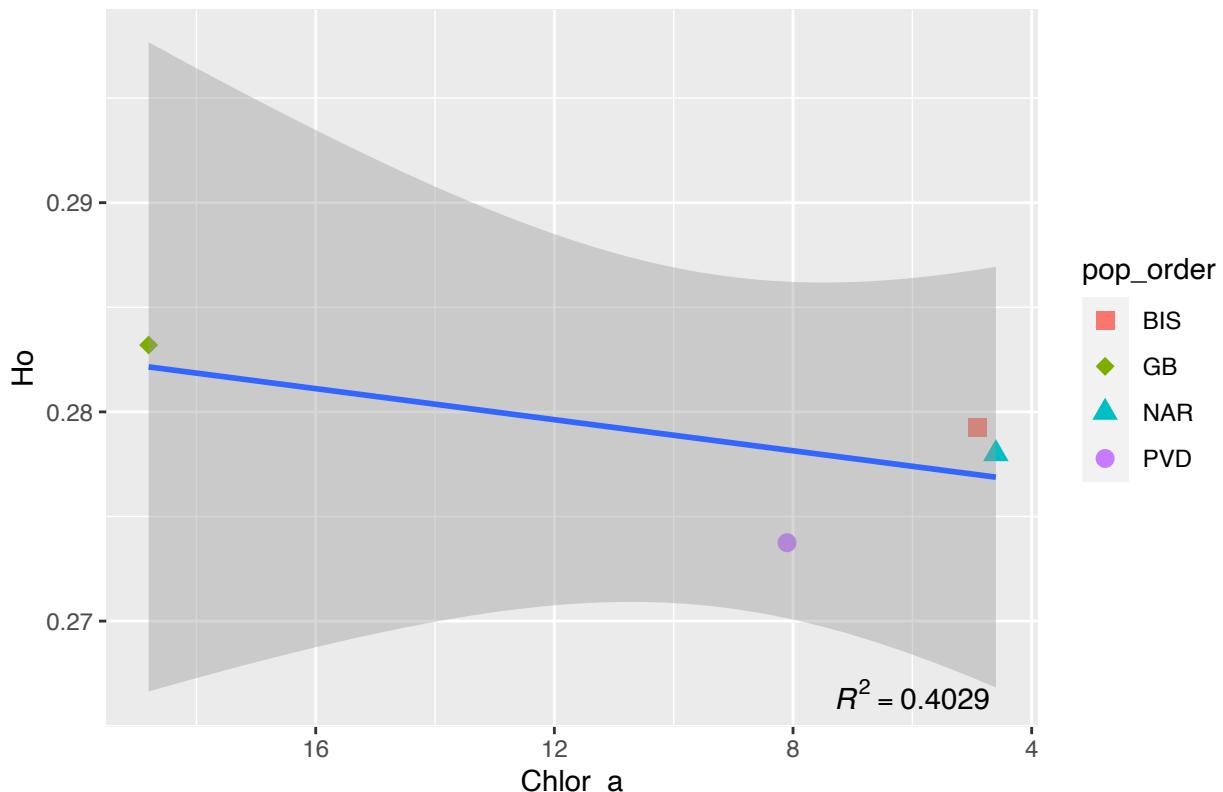
### Expected heterozygosity vs Chlorophyll a, Neutral SNPs



```
#Plot Ho vs Chlorophyll a
R4 = round(summary(lm(y$Ho ~ y$Chlor_a))$r.squared, 4)
ggplot(y, aes(x = Chlor_a, y = Ho)) + geom_point(aes(shape=pop_order, color=pop_order), size = 3) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggttitle("Observed heterozygosity vs Chlorophyll a, Neutral SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R4), x=6, y=0.2665, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

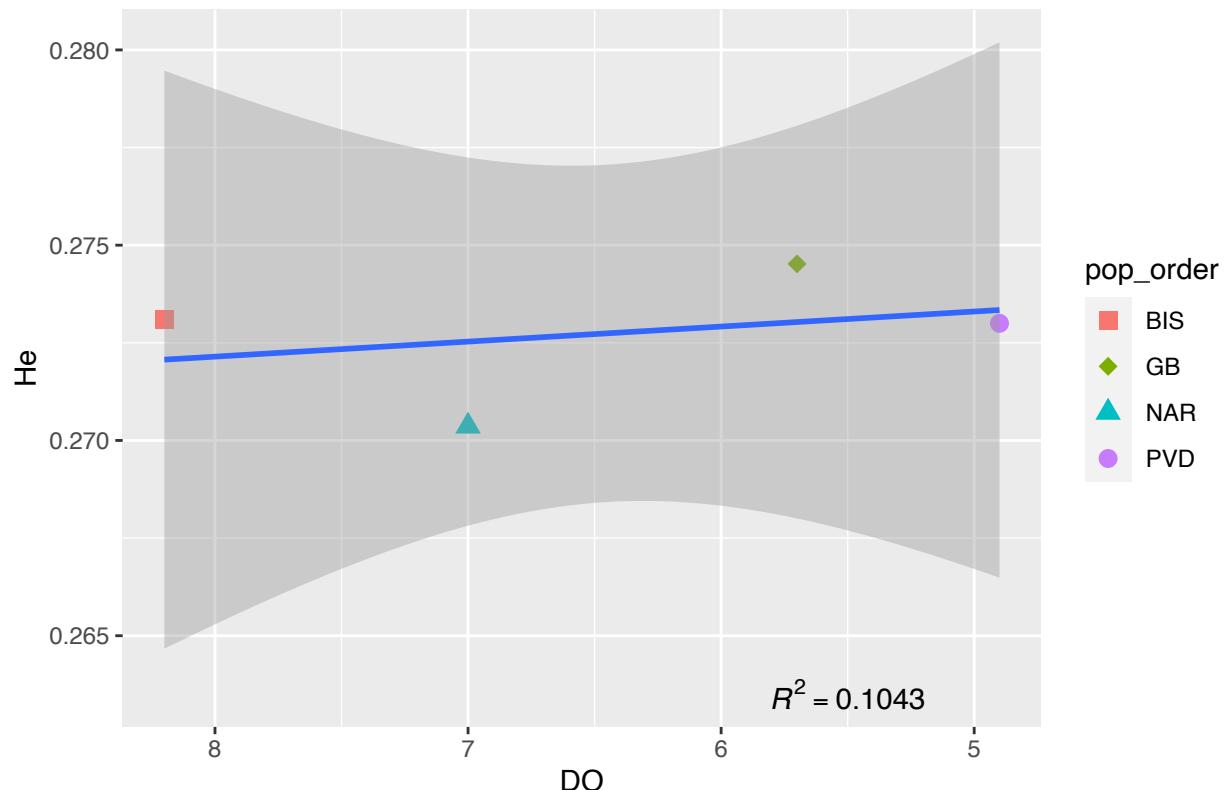
## Observed heterozygosity vs Chlorophyll a, Neutral SNPs



```
#Plot He vs Dissolved Oxygen
R4 = round(summary(lm(y$He ~ y$D0))$r.squared, 4)
ggplot(y, aes(x = D0, y = He)) + geom_point(aes(shape=pop_order, color=pop_order), size = 3) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggttitle("Expected heterozygosity vs Dissolved Oxygen, Neutral SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R4), x=5.5, y=0.2635, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

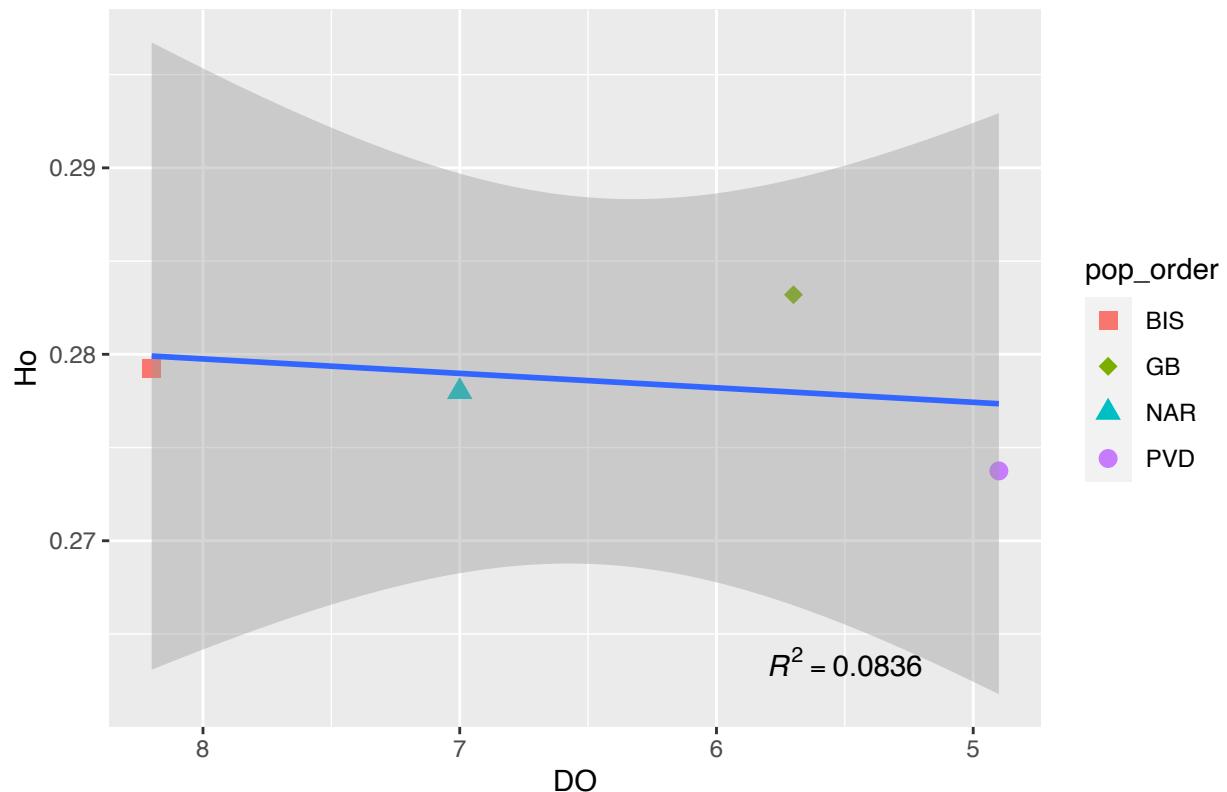
### Expected heterozygosity vs Dissolved Oxygen, Neutral SNPs



```
#Plot Ho vs Dissolved Oxygen
R4 = round(summary(lm(y$Ho ~ y$DO))$r.squared, 4)
ggplot(y, aes(x = DO, y = Ho)) + geom_point(aes(shape=pop_order, color=pop_order), size = 3) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggtitle("Observed heterozygosity vs Dissolved Oxygen, Neutral SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R4), x=5.5, y=0.2635, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

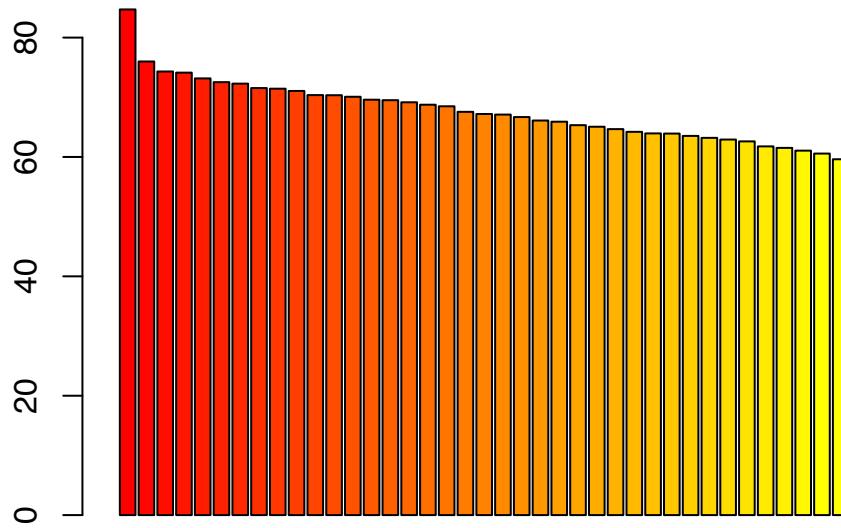
### Observed heterozygosity vs Dissolved Oxygen, Neutral SNPs



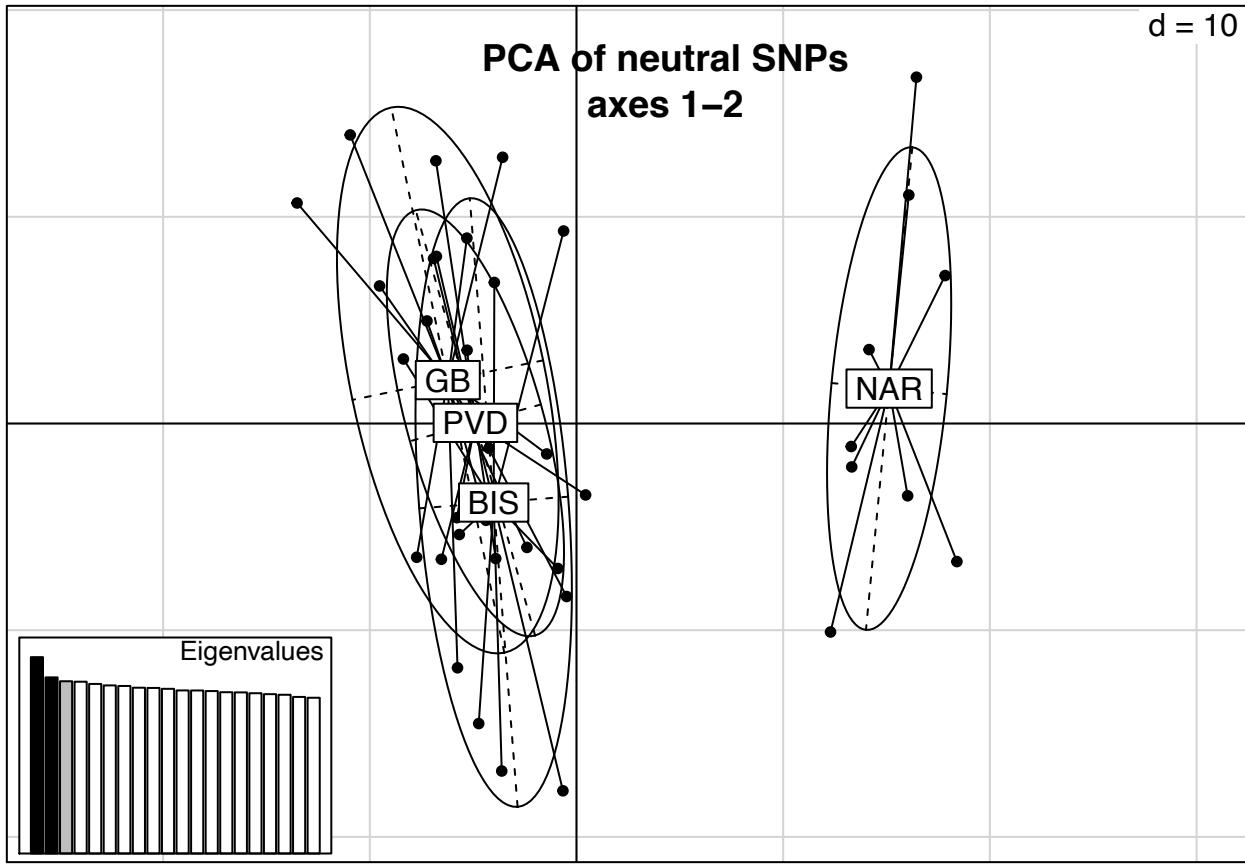
### PCA

```
X <- tab(stratated.u, freq = TRUE, NA.method = "mean")
pca1 <- dudi.pca(X, scale = FALSE, scannf = FALSE, nf = 3)
barplot(pca1$eig[1:50], main = "PCA eigenvalues", col = heat.colors(50))
```

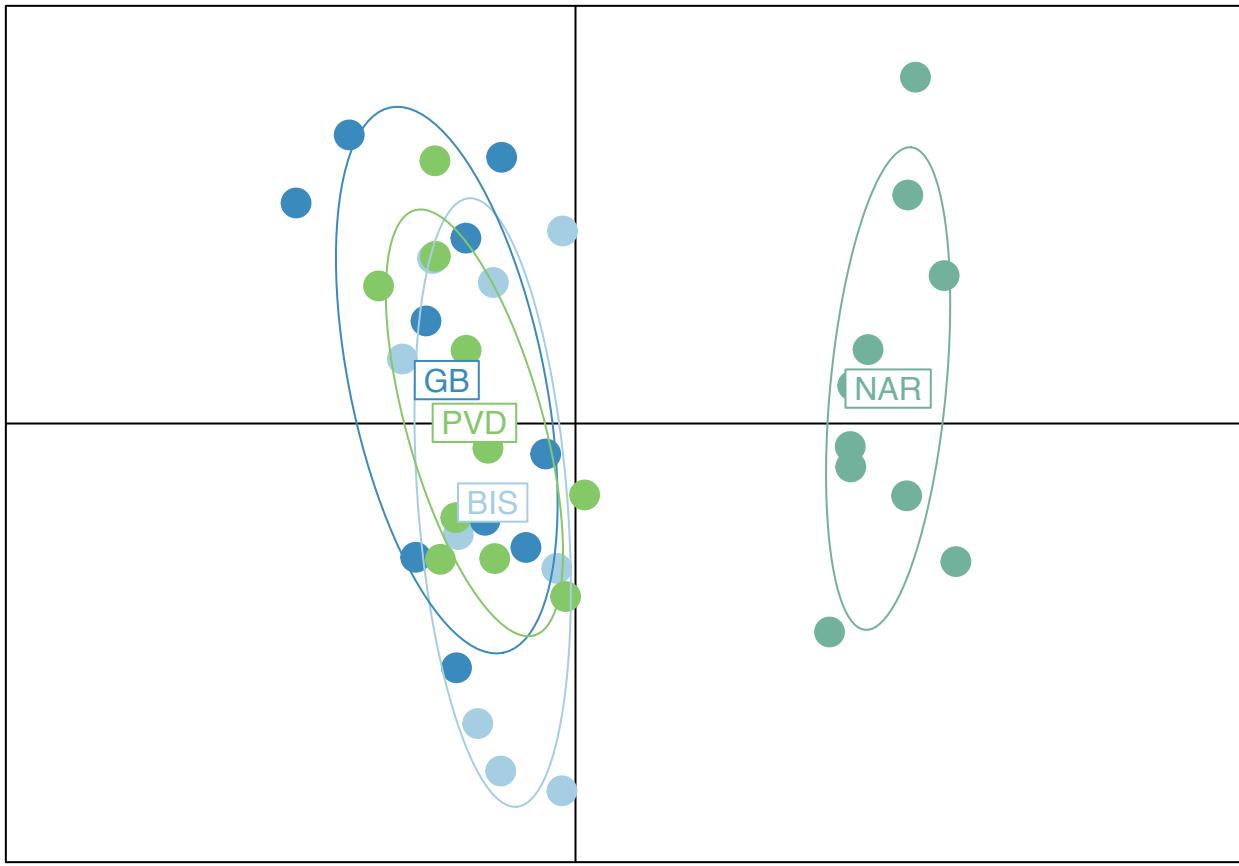
## PCA eigenvalues



```
s.class(pca1$li, pop(stratated.filt))
title("PCA of neutral SNPs\naxes 1-2")
add.scatter.eig(pca1$eig[1:20], 3,1,2)
```



```
col <- funky(15)
s.class(pca1$li, pop(stratified.filt), xax=1, yax=2, col=col, axesell=FALSE, cstard=0, cpoint=3, grid=FALSE)
```



## Discriminant Analysis of Principal Components (DAPC)

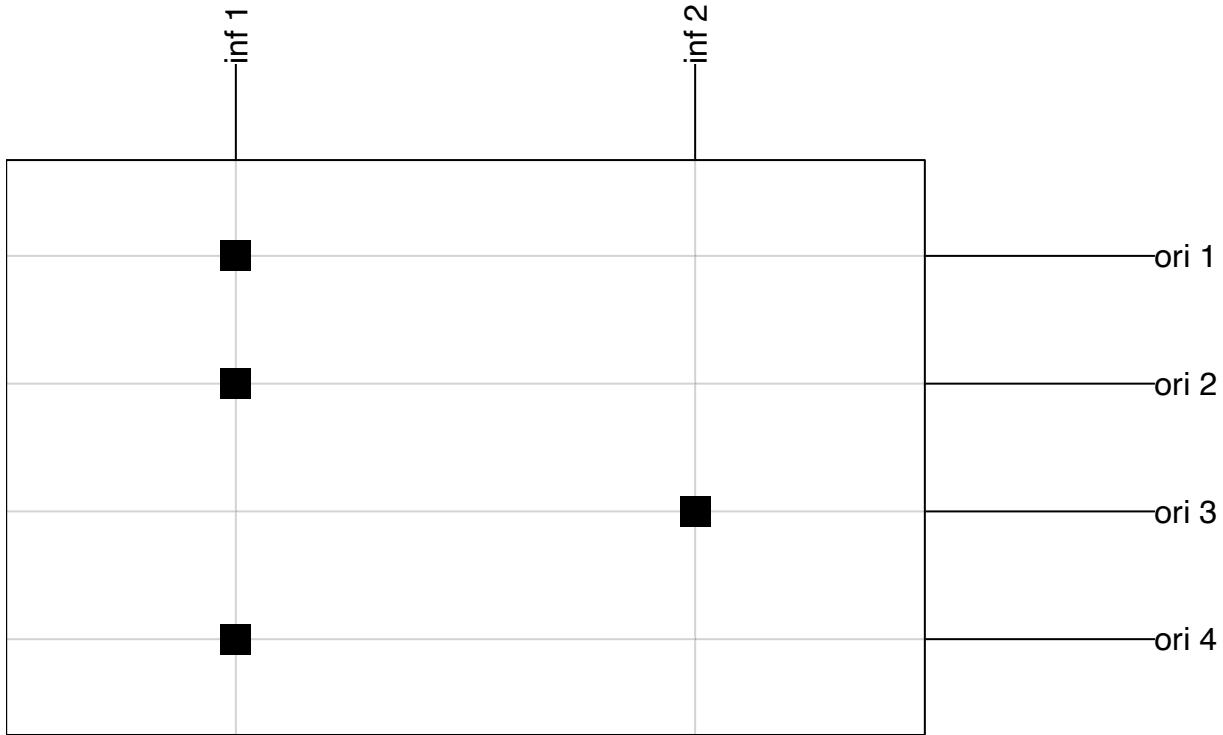
The first line of the code block below is interactive. In the console, it will tell you:

1. Choose the number PCs to retain ( $\geq 1$ ): I picked 1 2. Choose the number discriminant functions to retain ( $\geq 1$ ): I picked 2

```
grp <- find.clusters(stratated.u, n.pca = 1, choose.n.clust = FALSE)
table(pop(stratated.u), grp$grp)

##
##      1 2
##  BIS 10 0
##  GB  10 0
##  NAR  0 10
##  PVD 10 0

table.value(table(pop(stratated.u), grp$grp), col.lab=paste("inf", 1:2), row.lab=paste("ori", 1:4))
```

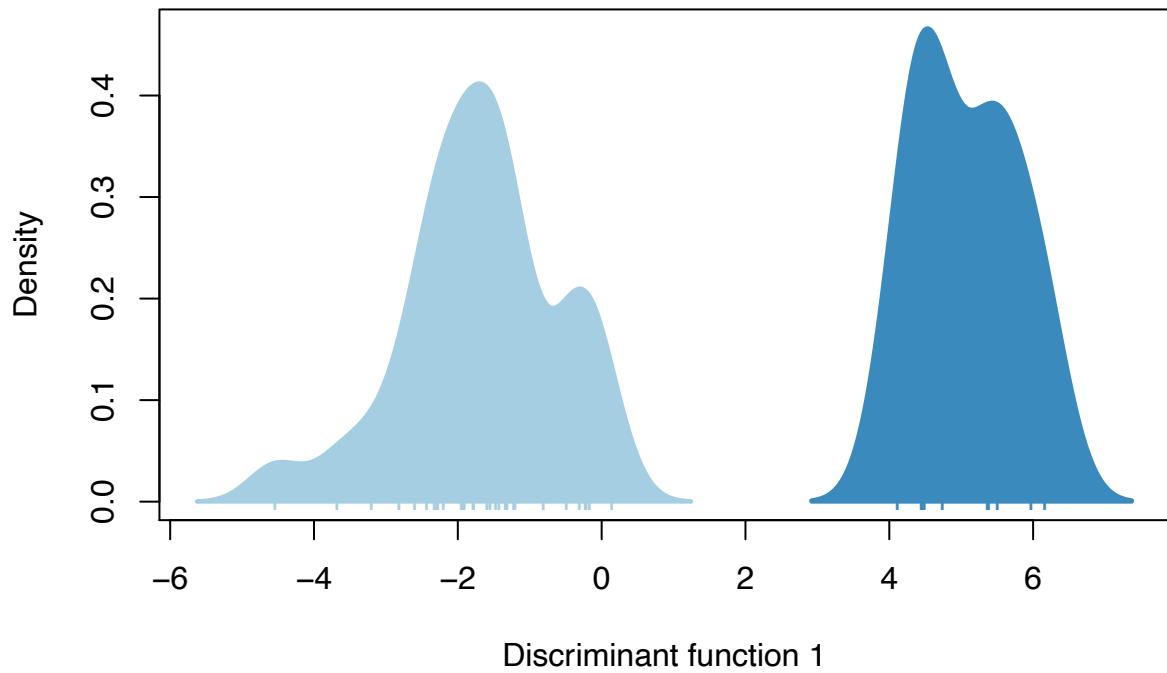


**1 • 3 ■ 5 ■ 7 ■ 9 ■**

Again, the first line of the code block below is interactive. In the console, it will tell you:

1. Choose the number PCs to retain ( $\geq 1$ ): I picked 1 2. Choose the number discriminant functions to retain ( $\geq 1$ ): I picked 1

```
dapc1 <- dapc(stratated.u, grp$grp, n.pca = 1, n.da= 1)
scatter(dapc1,col=col,bg="white", solid=1)
```

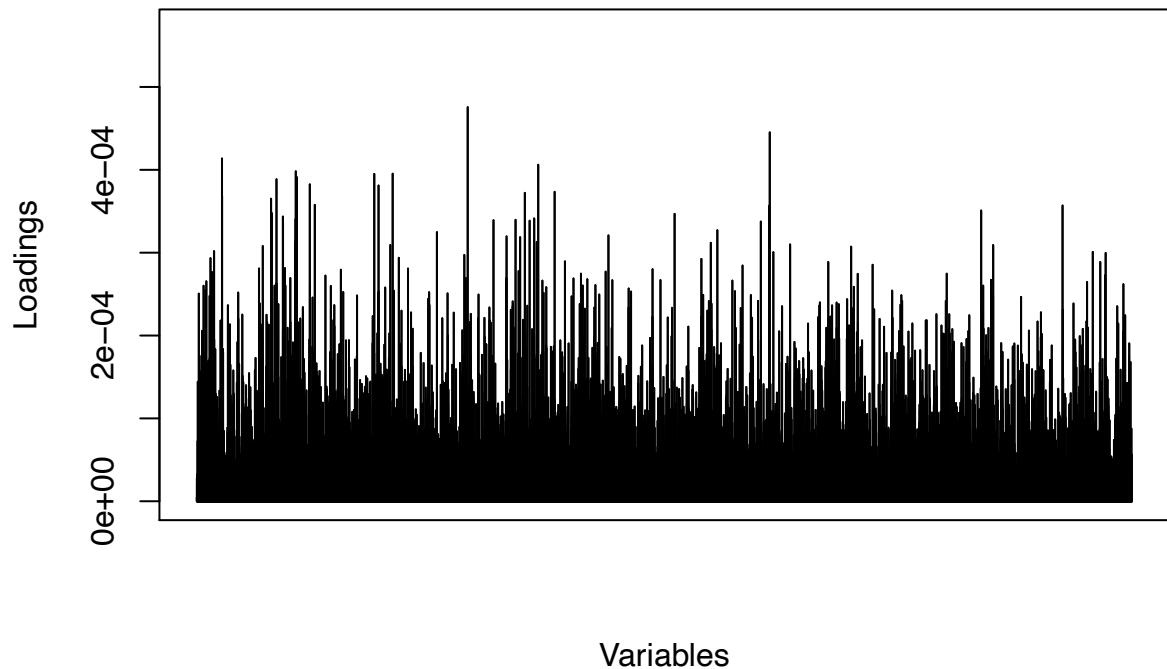


Again, the first line of the code block below is interactive. In the console, it will tell you:

1. Choose the number PCs to retain ( $\geq 1$ ): I picked 1 2. Choose the number discriminant functions to retain ( $\geq 1$ ): I picked 1

```
contrib <- loadingplot(dapc1$var.contr, axis=1, thres=.01, lab.jitter=1)
```

## Loading plot

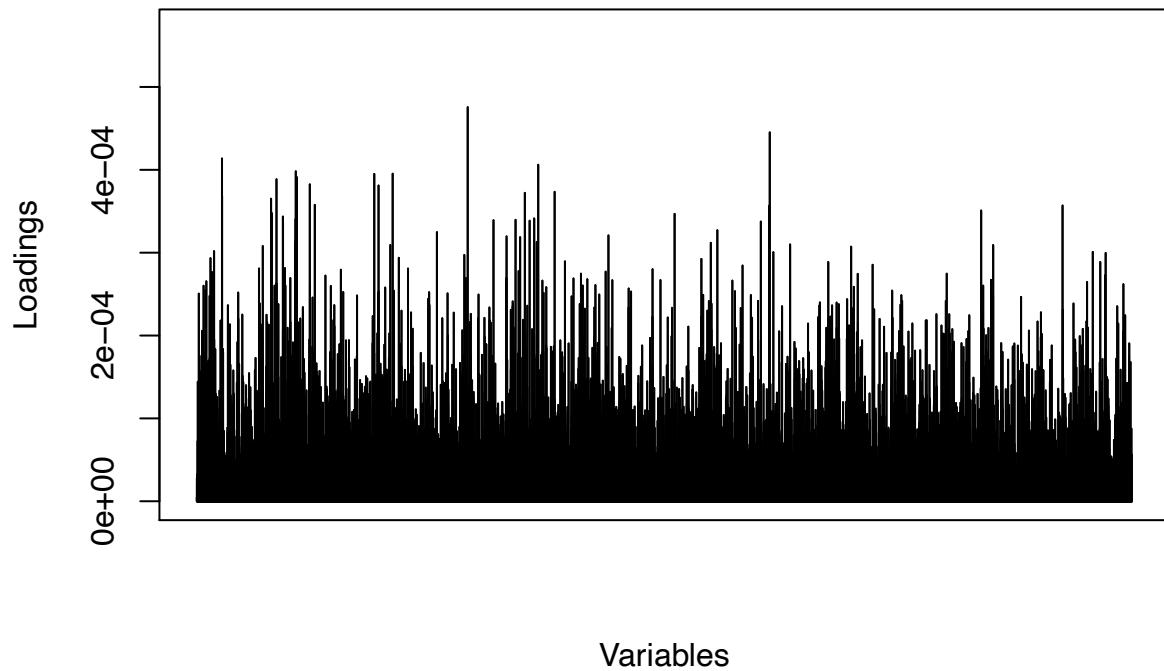


```
contrib

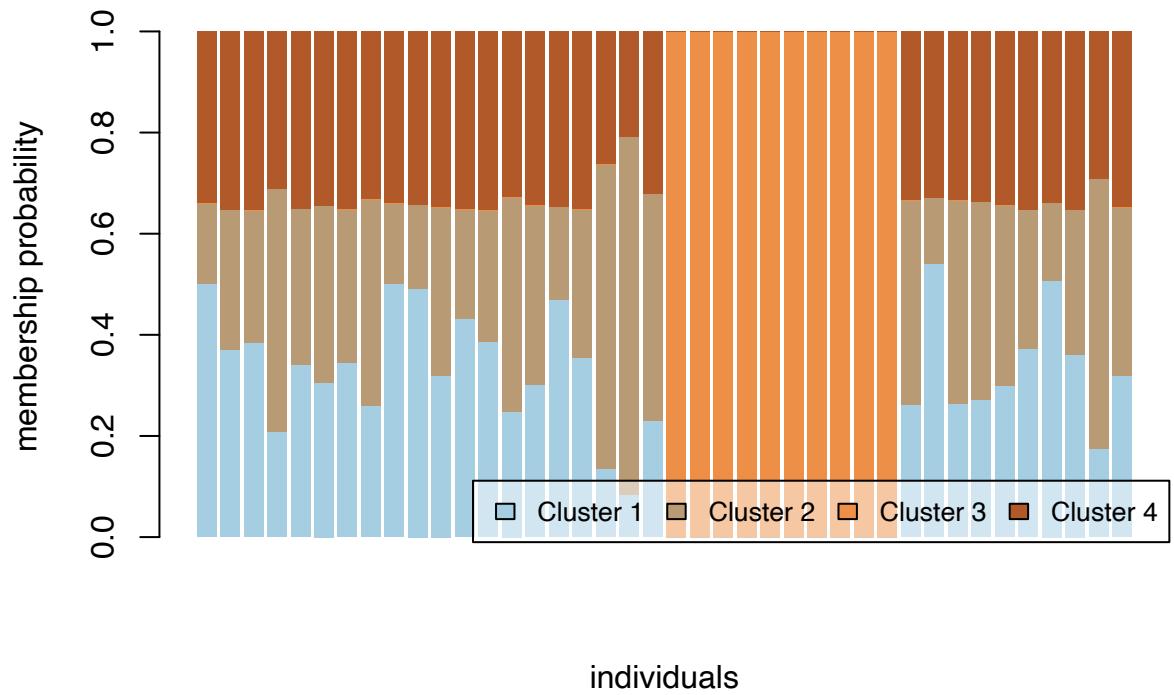
## NULL

setPop(rad.filt) <- ~Library
dapc1 <- dapc(stratified.u, pop(stratified.u), n.pca = 1, n.da = 1)
contrib <- loadingplot(dapc1$var.contr, axis=1, thres=.05, lab.jitter=1)
```

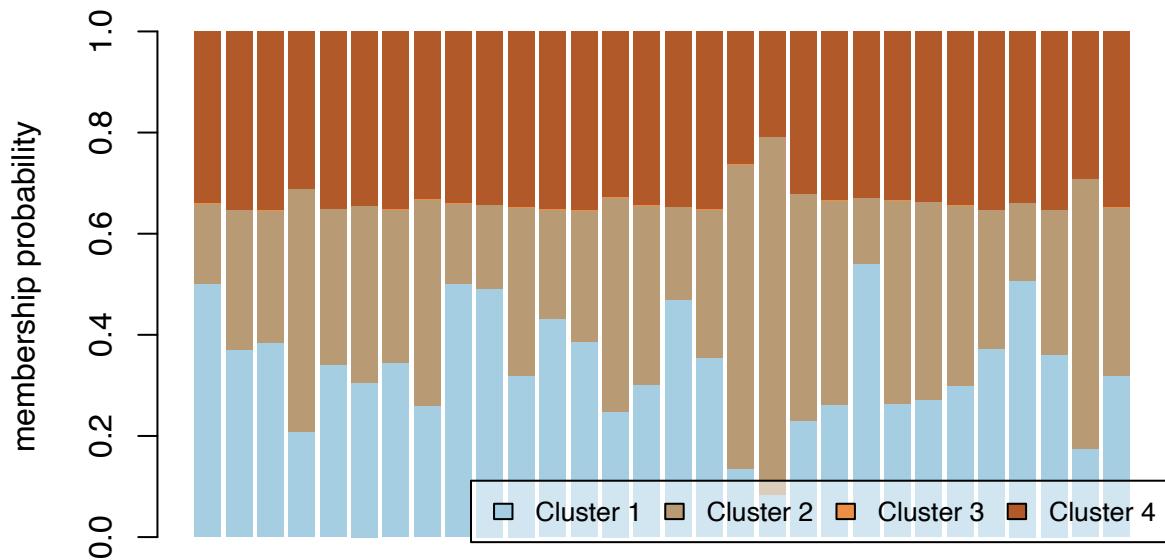
## Loading plot



```
#Structure Like
compoplot(dapc1, posi="bottomright", txt.leg=paste("Cluster", 1:4), lab="", ncol=1, xlab="individuals")
```



```
temp <- which(apply(dapc1$posterior, 1, function(e) all(e<0.9)))
compoplot(dapc1, subset=temp, posi="bottomright", txt.leg=paste("Cluster", 1:4), ncol=2)
```



## PCAviz

```

NA.afDraw<- function(ind){
  ind.mat <- ind@tab
  new.mat <- ind.mat
  af = colSums(ind.mat[,seq(1,ncol(ind.mat)-1,2)],na.rm = TRUE) /
    (2*apply(ind.mat[,seq(1,ncol(ind.mat)-1,2)],2,function(x) sum(!is.na(x))))
  af.Draw <- function(geno, af){
    new <- function(geno,af){
      if(is.na(geno)){
        newA = rbinom(1,2,af)
      } else {newA <- geno}
      return(newA)
    }
    new.row <- mapply(geno,af,FUN = new)
    return(new.row)
  }

  new.mat[,seq(1,ncol(ind.mat)-1,2)] <- t(apply(ind.mat[,seq(1,ncol(ind.mat)-1,2)],1,af.Draw,af))
  new.mat[,seq(2,ncol(ind.mat),2)] <- 2-new.mat[,seq(1,ncol(ind.mat)-1,2)]
  new.ind <- ind
  new.ind@tab <- new.mat
  return(new.ind)
}

```

```

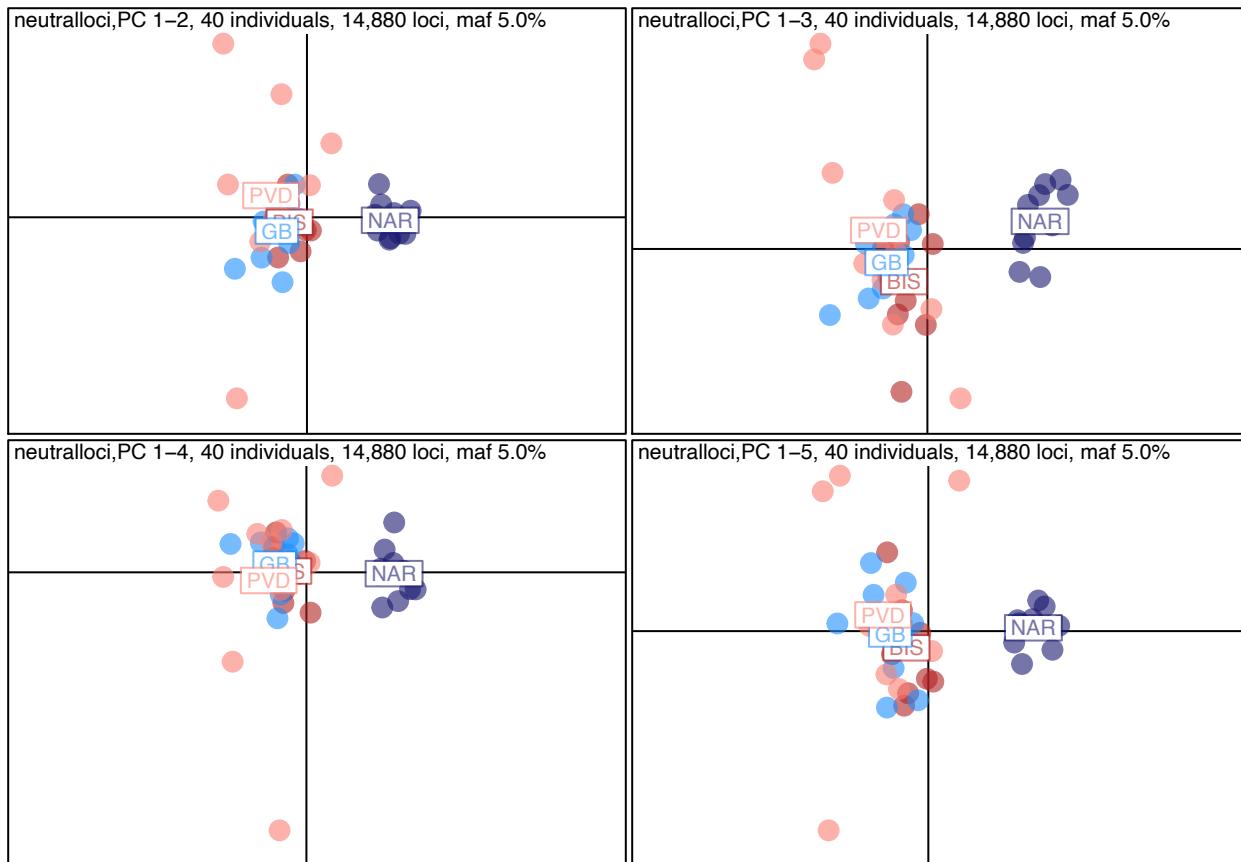
u.na <- NA.afDraw(stratified.u)

pca <- dudi.pca(u.na,center=TRUE,scale=TRUE,scannf = F, nf = 30)

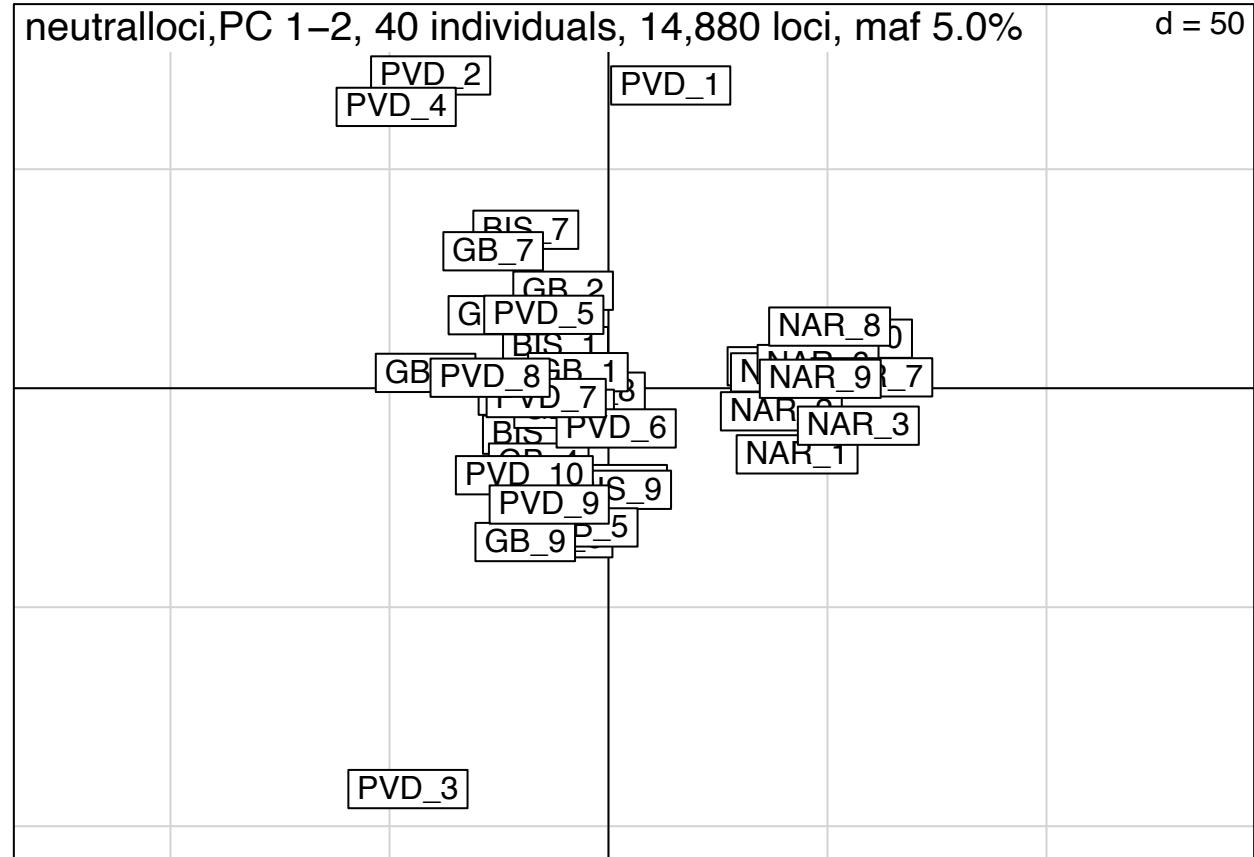
col18 <- funky(length(unique(u.na@strata$Population)))
#Colors that match the neutral Structure results
col6 <- c("firebrick","dodgerblue","midnightblue","salmon")
par(mfrow=c(2,2))

s.class(pca$li, strata(u.na)$Population,xax=1,yax=2,
        sub = "neutralloci,PC 1-2, 40 individuals, 14,880 loci, maf 5.0%",
        possub = "topleft",col=transp(col6,.6),axesell=FALSE,
        cstardot=0, cpoint=3, grid=FALSE, cellipse = 0)
s.class(pca$li, strata(u.na)$Population,xax=1,yax=3,
        sub = "neutralloci,PC 1-3, 40 individuals, 14,880 loci, maf 5.0%",
        possub = "topleft",col=transp(col6,.6),axesell=FALSE,
        cstardot=0, cpoint=3, grid=FALSE, cellipse = 0)
s.class(pca$li, strata(u.na)$Population,xax=1,yax=4,
        sub = "neutralloci,PC 1-4, 40 individuals, 14,880 loci, maf 5.0%",
        possub = "topleft",col=transp(col6,.6),axesell=FALSE,
        cstardot=0, cpoint=3, grid=FALSE, cellipse = 0)
s.class(pca$li, strata(u.na)$Population,xax=1,yax=5,
        sub = "neutralloci,PC 1-5, 40 individuals, 14,880 loci, maf 5.0%",
        possub = "topleft",col=transp(col6,.6),axesell=FALSE,
        cstardot=0, cpoint=3, grid=FALSE, cellipse = 0)

```



```
s.label(pca$li, xax=1,yax=5,
       sub = "neutralloci,PC 1-2, 40 individuals, 14,880 loci, maf 5.0%",
       possub = "topleft")
```



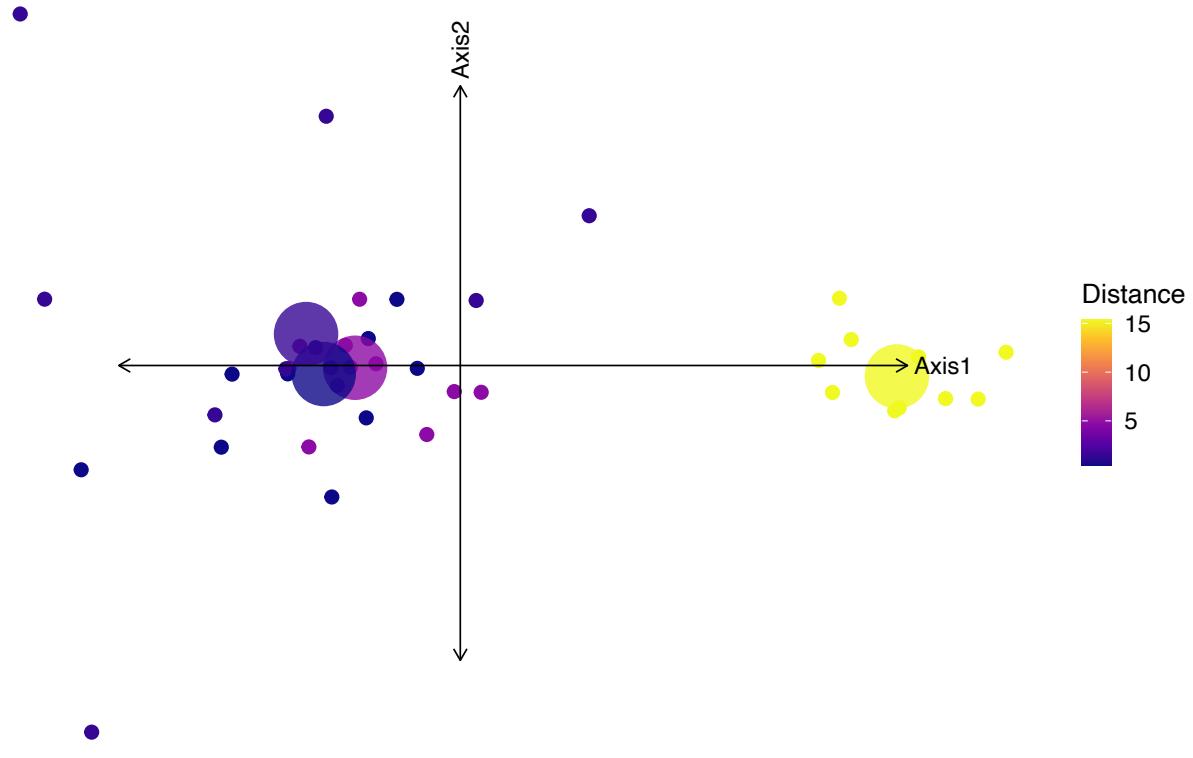
```
eig.perc <- 100*pca$eig/sum(pca$eig)
head(eig.perc)
```

```
## [1] 3.146839 2.947990 2.863383 2.851769 2.815016 2.793945
```

```
li <- pca$li
c1 <- pca$c1
#Create dataframe of info like latitude and population for each individual
info_mat <- as.data.frame(cbind(u.na@strata, u.na@other$Latitude, u.na@other$Longitude, u.na@other$Distance))
colnames(info_mat) <- c("Population", "Latitude", "Longitude", "Distance", "Temperature", "Salinity", "pH",
colnames(c1) <- colnames(li)
#create pcaviz object
pviz <- pcaviz(x=li, rotation=c1, dat=info_mat)

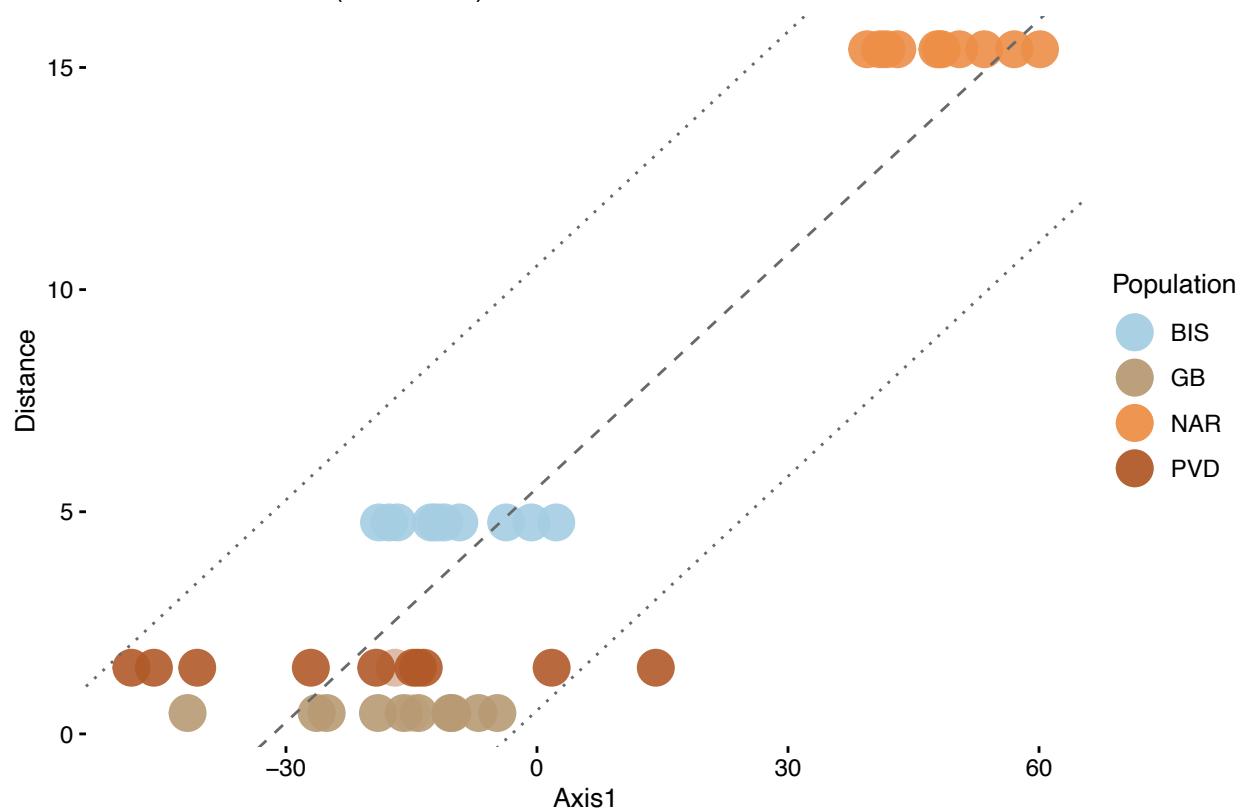
p = list(size=2)
plot(pviz, color = "Distance", draw.points = T, group.summary.labels = F, draw.pc.axes = T, geom.point.p
```

Axis1 vs. Axis2

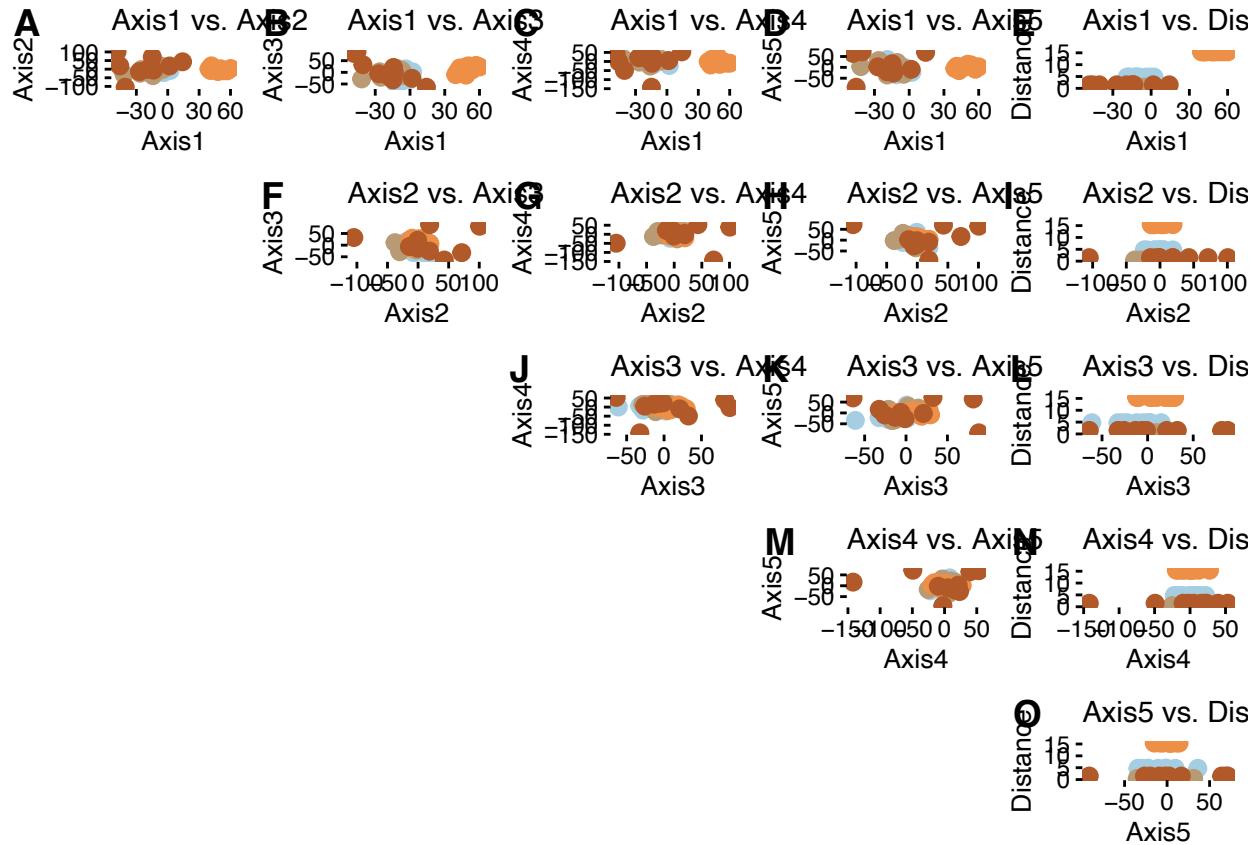


```
# Plot against distance
p = list(size=6)
plot(pviz, coords = c("Axis1","Distance"),
  show.legend = T,color = "Population",colors = col18,
  draw.points =T,group.summary.labels = F,geom.point.params= list(size = 6, alpha = 0.9),geom.point.
```

Axis1 vs. Distance ( $r^2 = 0.823$ )



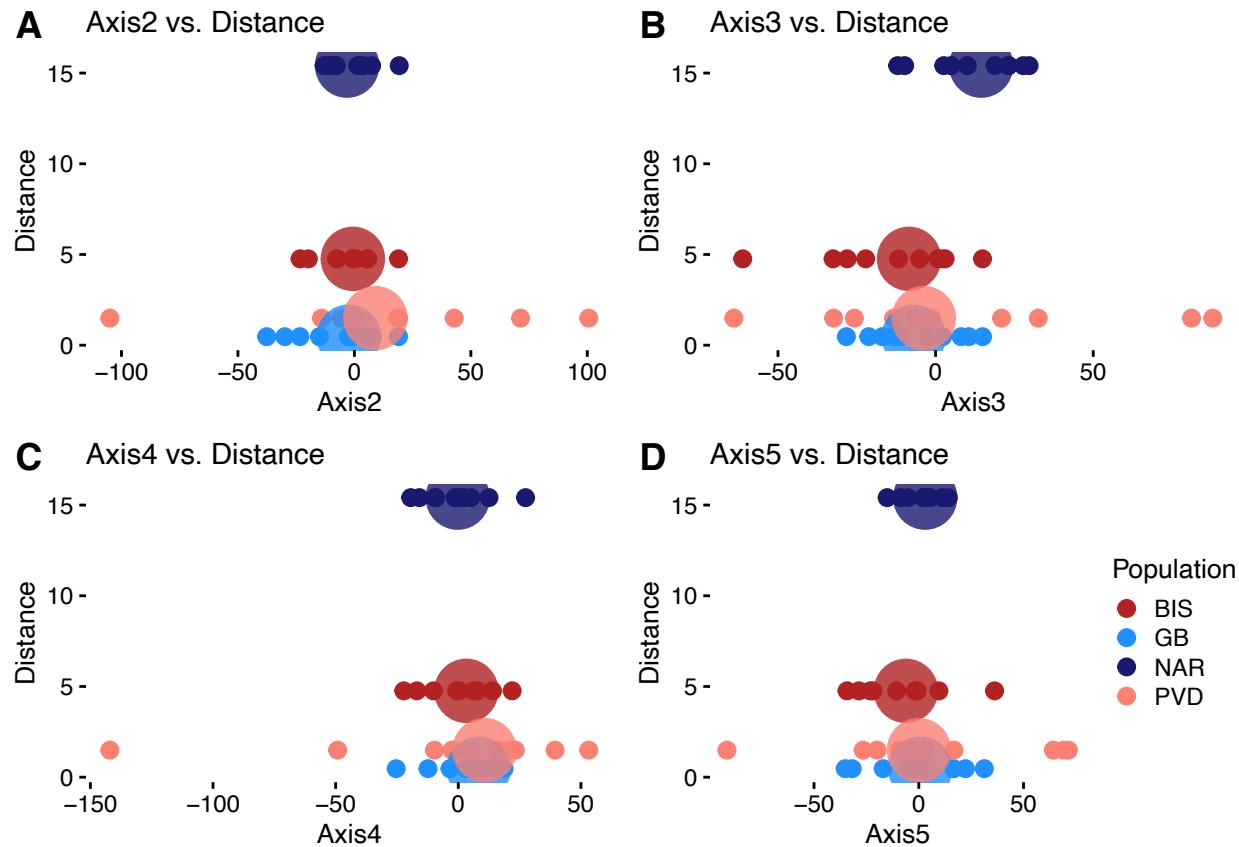
```
plot(pviz,coords = c("Axis1","Axis2","Axis3","Axis4","Axis5","Distance"),group = NULL,color= "Population")
```



```

p = list(size=6)
A2 <- plot(pviz,coords = c("Axis2","Distance"),group="Population",
            show.legend = F,color = "Population",colors = col6,draw.points = T,
            group.summary.labels = F,draw.linear.fit = F, coord_fixed=T)
A3 <- plot(pviz,coords = c("Axis3","Distance"),group="Population",
            show.legend = F,color = "Population",colors = col6,
            draw.points = T,group.summary.labels = F,draw.linear.fit = F)
A4 <- plot(pviz,coords = c("Axis4","Distance"),group="Population",
            show.legend = F,color = "Population",colors = col6,
            draw.points = T,group.summary.labels = F,draw.linear.fit = F)
A5 <- plot(pviz,coords = c("Axis5","Distance"),group="Population",
            show.legend = T,color = "Population",colors = col6,
            draw.points = T,group.summary.labels = F,draw.linear.fit = F)
plot_grid(A2,A3,A4,A5, labels = c('A','B','C','D'))

```



## Neutral Genic SNPs

---

### Making files

#### Make genind object

`neutralgene.recode.vcf` contains neutral SNPs from genic regions for populations PVD, GB, BIS, and NAR. Steps for generating this VCF file are located in `EecSeq_Cvirginica_NeutralLoci.md`.

Population NIN was removed from the VCF file following steps in `EecSeq_Cvirginica_OutlierDetection.md`.

`strata` contains population, environmental, and library information for each sample - can be accessed here.

```
my_vcf_gene <- read.vcfR("neutralgene.recode.vcf")
```

```
## Scanning file to determine attributes.
## File attributes:
##   meta lines: 63
##   header_line: 64
##   variant count: 66939
##   column count: 49
##
Meta line 63 read in.
```

```
## All meta lines processed.  
## gt matrix initialized.  
## Character matrix gt created.  
##   Character matrix gt rows: 66939  
##   Character matrix gt cols: 49  
##   skip: 0  
##   nrow: 66939  
##   row_num: 0  
##  
Processed variant 1000  
Processed variant 2000  
Processed variant 3000  
Processed variant 4000  
Processed variant 5000  
Processed variant 6000  
Processed variant 7000  
Processed variant 8000  
Processed variant 9000  
Processed variant 10000  
Processed variant 11000  
Processed variant 12000  
Processed variant 13000  
Processed variant 14000  
Processed variant 15000  
Processed variant 16000  
Processed variant 17000  
Processed variant 18000  
Processed variant 19000  
Processed variant 20000  
Processed variant 21000  
Processed variant 22000  
Processed variant 23000  
Processed variant 24000  
Processed variant 25000  
Processed variant 26000  
Processed variant 27000  
Processed variant 28000  
Processed variant 29000  
Processed variant 30000  
Processed variant 31000  
Processed variant 32000  
Processed variant 33000  
Processed variant 34000  
Processed variant 35000  
Processed variant 36000  
Processed variant 37000  
Processed variant 38000  
Processed variant 39000  
Processed variant 40000  
Processed variant 41000  
Processed variant 42000  
Processed variant 43000  
Processed variant 44000  
Processed variant 45000
```

```

Processed variant 46000
Processed variant 47000
Processed variant 48000
Processed variant 49000
Processed variant 50000
Processed variant 51000
Processed variant 52000
Processed variant 53000
Processed variant 54000
Processed variant 55000
Processed variant 56000
Processed variant 57000
Processed variant 58000
Processed variant 59000
Processed variant 60000
Processed variant 61000
Processed variant 62000
Processed variant 63000
Processed variant 64000
Processed variant 65000
Processed variant 66000
Processed variant: 66939
## All variants processed
strata <- read.table("strata", header=TRUE)

rad.filt_gene <- vcfR2genind(my_vcf_gene, strata = strata, pop = c(rep("BIS", 10),rep("GB", 10),rep("NAR", 10),rep("PVD", 10),rep("NAN", 10),rep("BRI", 10),rep("LAM", 10),rep("TUR", 10),rep("SOM", 10),rep("PER", 10)))

rad.filt_gene

## /// GENIND OBJECT ///////////
##
## // 40 individuals; 66,939 loci; 133,654 alleles; size: 61 Mb
##
## // Basic content
## @tab: 40 x 133654 matrix of allele counts
## @loc.n.all: number of alleles per locus (range: 1-2)
## @loc.fac: locus factor for the 133654 columns of @tab
## @all.names: list of allele names for each locus
## @ploidy: ploidy of each individual (range: 2-2)
## @type: codom
## @call: adegenet::df2genind(X = t(x), sep = sep, pop = ..2, strata = ..1)
##
## // Optional content
## @pop: population of each individual (group size range: 10-10)
## @strata: a data frame with 12 columns ( Individual, Population, Latitude, Longitude, Distance, SE, NAR, PVD, BIS, GB, TUR, LAM )
#Providing population names for plotting
pop_order <- c("BIS", "GB", "NAR", "PVD")

```

Read in the other info from .strata file and extract information such as locality, latitude, and longitude.

```

info <- as.data.frame(read.table("strata", header = T, sep = "\t", stringsAsFactors = F))
mystrats_gene <- as.data.frame(matrix(nrow = length(indNames(rad.filt_gene)), ncol=10))
colnames(mystrats_gene) <- c("Population", "Latitude", "Longitude", "Distance", "Temperature", "SE", "Salinity")

```

```

just.strats <- select(info,c("Population"))
stratated.filt_gene <- strata(rad.filt_gene, formula= Population, combine = TRUE,just.strats)
stratated.filt_gene@other <- select(info, Latitude,Longitude,Distance,SE,Temperature,Salinity,pH,Chlorophy
stratated.filt_gene

## /// GENIND OBJECT ///////////
##
## // 40 individuals; 66,939 loci; 133,654 alleles; size: 61 Mb
##
## // Basic content
## @tab: 40 x 133654 matrix of allele counts
## @loc.n.all: number of alleles per locus (range: 1-2)
## @loc.fac: locus factor for the 133654 columns of @tab
## @all.names: list of allele names for each locus
## @ploidy: ploidy of each individual (range: 2-2)
## @type: codom
## @call: adegenet::df2genind(X = t(x), sep = sep, pop = ..2, strata = ..1)
##
## // Optional content
## @pop: population of each individual (group size range: 10-10)
## @strata: a data frame with 1 columns ( Population )
## @other: a list containing: Latitude Longitude Distance SE Temperature Salinity pH Chlorophy

```

#### Repeat for quasi-independent set of SNPs

Steps for filtering are documented for OutFLANK in EecSeq\_Cvirginica\_OutlierDetection.md. However, I am repeating the steps here.

```
my_vcf_gene <- read.vcfR("neutralgene.recode.vcf")
```

```

## Scanning file to determine attributes.
## File attributes:
## meta lines: 63
## header_line: 64
## variant count: 66939
## column count: 49
##
Meta line 63 read in.
## All meta lines processed.
## gt matrix initialized.
## Character matrix gt created.
## Character matrix gt rows: 66939
## Character matrix gt cols: 49
## skip: 0
## nrows: 66939
## row_num: 0
##
Processed variant 1000
Processed variant 2000
Processed variant 3000
Processed variant 4000
Processed variant 5000
Processed variant 6000
Processed variant 7000
Processed variant 8000

```

Processed variant 9000  
Processed variant 10000  
Processed variant 11000  
Processed variant 12000  
Processed variant 13000  
Processed variant 14000  
Processed variant 15000  
Processed variant 16000  
Processed variant 17000  
Processed variant 18000  
Processed variant 19000  
Processed variant 20000  
Processed variant 21000  
Processed variant 22000  
Processed variant 23000  
Processed variant 24000  
Processed variant 25000  
Processed variant 26000  
Processed variant 27000  
Processed variant 28000  
Processed variant 29000  
Processed variant 30000  
Processed variant 31000  
Processed variant 32000  
Processed variant 33000  
Processed variant 34000  
Processed variant 35000  
Processed variant 36000  
Processed variant 37000  
Processed variant 38000  
Processed variant 39000  
Processed variant 40000  
Processed variant 41000  
Processed variant 42000  
Processed variant 43000  
Processed variant 44000  
Processed variant 45000  
Processed variant 46000  
Processed variant 47000  
Processed variant 48000  
Processed variant 49000  
Processed variant 50000  
Processed variant 51000  
Processed variant 52000  
Processed variant 53000  
Processed variant 54000  
Processed variant 55000  
Processed variant 56000  
Processed variant 57000  
Processed variant 58000  
Processed variant 59000  
Processed variant 60000  
Processed variant 61000  
Processed variant 62000

```

Processed variant 63000
Processed variant 64000
Processed variant 65000
Processed variant 66000
Processed variant: 66939
## All variants processed

geno_gene <- extract.gt(my_vcf_gene) # Character matrix containing the genotypes
position_gene <- getPOS(my_vcf_gene) # Positions in bp
chromosome_gene <- getCHROM(my_vcf_gene) # Chromosome information

G_gene <- matrix(NA, nrow = nrow(geno_gene), ncol = ncol(geno_gene))

G_gene[geno_gene %in% c("0/0", "0|0")] <- 0
G_gene[geno_gene %in% c("0/1", "1/0", "1|0", "0|1")] <- 1
G_gene[geno_gene %in% c("1/1", "1|1")] <- 2

# NA should be replaced with "9"
G_gene[is.na(G_gene)] <- 9

```

Chromosomes need to be of class integer for this to work.

```

# Visualizing chromosomes
chrom_unique_gene <- unique(chromosome_gene)
print(chrom_unique_gene)

## [1] "NC_035780.1" "NC_035781.1" "NC_035782.1" "NC_035783.1" "NC_035784.1"
## [6] "NC_035785.1" "NC_035786.1" "NC_035787.1" "NC_035788.1" "NC_035789.1"

# Removing "NC_" from the chromosome name so it can be converted to an integer
chrom_new_gene <- chromosome_gene %>% str_replace("NC_","")

# Converting the character string into an integer
chrom1_gene <- as.integer(chrom_new_gene)

```

chromosome and position need to be sorted for imputation to work.

```

chrom_sort_gene <- sort(chrom1_gene)
pos_sort_gene <- sort(position_gene)

```

*Note: This filtering program does not allow for missing genotype values.*

### Remove missing genotypes

```

# removing missing data
G_gene_miss <- matrix(NA, nrow = nrow(geno_gene), ncol = ncol(geno_gene))

G_gene_miss[geno_gene %in% c("0/0", "0|0")] <- 0
G_gene_miss[geno_gene %in% c("0/1", "1/0", "1|0", "0|1")] <- 1
G_gene_miss[geno_gene %in% c("1/1", "1|1")] <- 2

# removing missing data
G_gene_miss[na.omit(G_gene_miss)]

##      [1] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [37] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0

```

```

## [99217] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [99253] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [99289] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [99325] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [99361] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [99397] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [99433] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [99469] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [99505] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [99541] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [99577] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [99613] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [99649] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [99685] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [99721] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [99757] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [99793] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [99829] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [99865] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [99901] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [99937] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## [99973] 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
##   [ reached getOption("max.print") -- omitted 424072 entries ]

# Converting the genotype matrix to class FBM.code256 using the matrix where missing data was removed
G1_gene_miss <- add_code256(big_copy(t(G_gene_miss),type="raw"),code=bigsnpr:::CODE_012)

## Warning in replaceMat(x$address_rw, i, j, value): At least one value changed (nan -> 0)
##   while converting from R type 'double' to C type 'unsigned char (raw)'.

newpc_gene_miss <- snp_autoSVD(G1_gene_miss,infos.chr =chrom_sort_gene,infos.pos = pos_sort_gene)

##
## Phase of clumping (on MAF) at r^2 > 0.2.. keep 17139 SNPs.
## Discarding 5519 variants with MAC < 10.
##
## Iteration 1:
## Computing SVD..
## 0 outlier variant detected..
##
## Converged!

which_pruned_gene_miss <- attr(newpc_gene_miss, which="subset") # Indexes of remaining SNPs after pruning
length(which_pruned_gene_miss)

## [1] 11620

invisible(lapply(which_pruned_gene_miss, write, "pruned_data_gene.txt", append=TRUE))

```

In terminal

```

$ mawk '!/#/' neutralexon.recode.vcf | cut -f1,2 > totallociexon
$ NUM=(`cat totallociexon | wc -l`)
$ paste <(seq 1 $NUM) totallociexon > lociexon.plus.index
$ cat pruned_data_exon.txt | parallel "grep -w ^{} lociexon.plus.index" | cut -f2,3> pruned_data_exon.lci

$ head pruned_data_exon.loci.txt

```

```

output:
NC_035780.1 573785
NC_035780.1 573788
NC_035780.1 574229
NC_035780.1 574241
NC_035780.1 574307
NC_035780.1 574313
NC_035780.1 574568
NC_035780.1 574580
NC_035780.1 574616
NC_035780.1 620575

Create VCF file with just the pruned_data loci

$ vcftools --vcf neutralexon.recode.vcf --recode --recode-INFO-all --positions pruned_data_exon.loci.txt

output:
After filtering, kept 40 out of 40 Individuals
Outputting VCF file...
After filtering, kept 11377 out of a possible 50238 Sites
Run Time = 2.00 seconds

my_vcf_u_gene <- read.vcfR("pruned_data_gene.recode.vcf")

## Scanning file to determine attributes.
## File attributes:
##   meta lines: 63
##   header_line: 64
##   variant count: 11620
##   column count: 49
##
Meta line 63 read in.
## All meta lines processed.
## gt matrix initialized.
## Character matrix gt created.
##   Character matrix gt rows: 11620
##   Character matrix gt cols: 49
##   skip: 0
##   nrows: 11620
##   row_num: 0
##
Processed variant 1000
Processed variant 2000
Processed variant 3000
Processed variant 4000
Processed variant 5000
Processed variant 6000
Processed variant 7000
Processed variant 8000
Processed variant 9000
Processed variant 10000
Processed variant 11000
Processed variant: 11620
## All variants processed

```

```

strata <- read.table("strata", header=TRUE)

rad.u_gene <- vcfR2genind(my_vcf_u_gene, strata = strata, pop = c(rep("BIS", 10),rep("GB", 10),rep("NAR", 10))

rad.u_gene

## /// GENIND OBJECT ///////////
## 
## // 40 individuals; 11,620 loci; 23,240 alleles; size: 10.6 Mb
## 
## // Basic content
## @tab: 40 x 23240 matrix of allele counts
## @loc.n.all: number of alleles per locus (range: 2-2)
## @loc.fac: locus factor for the 23240 columns of @tab
## @all.names: list of allele names for each locus
## @ploidy: ploidy of each individual (range: 2-2)
## @type: codom
## @call: adegenet::df2genind(X = t(x), sep = sep, pop = ..2, strata = ..1)
## 
## // Optional content
## @pop: population of each individual (group size range: 10-10)
## @strata: a data frame with 12 columns ( Individual, Population, Latitude, Longitude, Distance, SE, 
stratted.u_gene <- strata(rad.u_gene, formula= Population, combine = TRUE,just.strats)
stratted.u_gene@other <- select(info, Latitude,Longitude, Distance, SE, Temperature,Salinity,pH,Chlorophy

stratted.u_gene

## /// GENIND OBJECT ///////////
## 
## // 40 individuals; 11,620 loci; 23,240 alleles; size: 10.6 Mb
## 
## // Basic content
## @tab: 40 x 23240 matrix of allele counts
## @loc.n.all: number of alleles per locus (range: 2-2)
## @loc.fac: locus factor for the 23240 columns of @tab
## @all.names: list of allele names for each locus
## @ploidy: ploidy of each individual (range: 2-2)
## @type: codom
## @call: adegenet::df2genind(X = t(x), sep = sep, pop = ..2, strata = ..1)
## 
## // Optional content
## @pop: population of each individual (group size range: 10-10)
## @strata: a data frame with 1 columns ( Population )
## @other: a list containing: Latitude Longitude Distance SE Temperature Salinity pH Chlorophy

Make hierfstat object

hf.filt_gene <- genind2hierfstat(rad.filt_gene, pop = c(rep("BIS", 10),rep("GB", 10),rep("NAR", 10), rep("PVD", 10))

hf.u_gene <- genind2hierfstat(rad.u_gene, pop = c(rep("BIS", 10),rep("GB", 10),rep("NAR", 10), rep("PVD", 10))
hf.u_gene <- hf.u_gene$hierfstat.no.imputation

```

## Estimating effective migration surfaces (EEMS)

### 1. datapath.diffs

Steps for generating `datapath.diffs` are completed in RStudio.

```
# V1 method to get .diffs matrix, preferred
bed2diffs_v1 <- function(Geno) {
  nIndiv <- nrow(Geno)
  nSites <- ncol(Geno)
  Diffs <- matrix(0, nIndiv, nIndiv)

  for (i in seq(nIndiv - 1)) {
    for (j in seq(i + 1, nIndiv)) {
      x <- Geno[i, ]
      y <- Geno[j, ]
      Diffs[i, j] <- mean((x - y)^2, na.rm = TRUE)
      Diffs[j, i] <- Diffs[i, j]
    }
  }
  Diffs
}

# V2 method to get .diffs matrix, only if V1 doesn't work
bed2diffs_v2 <- function(Geno) {
  nIndiv <- nrow(Geno)
  nSites <- ncol(Geno)
  Miss <- is.na(Geno)
  ## Impute NAs with the column means (= twice the allele frequencies)
  Mean <- matrix(colMeans(Geno, na.rm = TRUE), ## a row of means
                  nrow = nIndiv, ncol = nSites, byrow = TRUE) ## a matrix with nIndiv identical rows of means
  Mean[Miss == 0] <- 0 ## Set the means that correspond to observed genotypes to 0
  Geno[Miss == 1] <- 0 ## Set the missing genotypes to 0 (used to be NA)
  Geno <- Geno + Mean
  ## Compute similarities
  Sim <- Geno %*% t(Geno) / nSites
  SelfSim <- diag(Sim) ## self-similarities
  vector1s <- rep(1, nIndiv) ## vector of 1s
  ## This chunk generates a `diffs` matrix
  Diffs <- SelfSim %*% t(vector1s) + vector1s %*% t(SelfSim) - 2 * Sim
  Diffs
}

geno_gene <- stratted.filt_gene@tab

# Get rid of non-biallelic loci
multi.loci_gene <- names(which(stratted.filt_gene@loc.n.all != 2))
multi.cols_gene <- which(grep(paste0("^", multi.loci_gene, "\\\\.\\\\d+$", collapse = "|"), colnames(geno_
if (length(multi.cols_gene)) geno_gene <- geno_gene[, - multi.cols_gene]
nloci_gene <- dim(geno_gene)[2] / 2
dim(geno_gene)

## [1] 40 133430
stopifnot(identical(stratted.filt_gene@type, 'codom'))
```

```
# bed2diffs functions
diffs.v1_gene <- bed2diffs_v1(geno_gene)
diffs.v2_gene <- bed2diffs_v2(geno_gene)
# Round to 6 digits
diffs.v1_gene <- round(diffs.v1_gene, digits = 6)
diffs.v2_gene <- round(diffs.v2_gene, digits = 6)
```

Check that the dissimilarity matrix has one positive eigenvalue and nIndiv-1 negative eigenvalues, as required by a full-rank Euclidean distance matrix. If the V1 method does not make a Euclidean matrix, you must use V2.

```
tail(sort(round(eigen(diffs.v1_gene)$values, digits = 2)))

## [1] -0.49 -0.49 -0.48 -0.48 -0.48 21.15

tail(sort(round(eigen(diffs.v2_gene)$values, digits = 2)))

## [1] -0.48 -0.48 -0.47 -0.47 -0.46 20.83

# Set suffix for EEMS input files
suf_gene <- "neutrallexondata-filt"

# This saves the file to directory
write.table(diffs.v1_gene, paste(suf_gene, ".v1.diffs", sep = ""),
            col.names = FALSE, row.names = FALSE, quote = FALSE)
```

## 2. datapath.coord

datapath.coord are the sample coordinates, two coordinates per sample, one sample per line. The sampling locations should be given in the same order as the rows and columns of the dissimilarity matrix.

Steps for generating datapath.coord are completed in RStudio.

```
## Get gps coordinates from previously created info matrix
x0R.info <- dplyr::filter(info)
gps_matrix <- select(x0R.info, c("Longitude", "Latitude"))

#write .coord file
write.table(gps_matrix, paste(suf_gene, ".v1.coord", sep = ""),
            col.names = FALSE, row.names = FALSE, quote = FALSE)
```

## 3. datapath.outer

datapath.outer are the habitat coordinates, as a sequence of vertices that form a closed polygon. The habitat vertices should be listed counterclockwise and the first vertex should also be the last vertex, so that the outline is a closed ring. Otherwise, EEMS attempts to “correct” the polygon and prints a warning message.

datapath.outer is created manually in Excel, based on site coordinates gathered from Google Maps, copied into terminal using nano, and saved as the file with the appropriate extension.

\*\*runeems\_snps is then run in command-line following the steps documented in NB\_EEMS\_OutlierHap.md

Back in RStudio to plot runeems\_snps outputs

```
# Install rEEMSpots
library(rEEMSpots)

# Plotting EEMS after running runeems_snps
path = "./NB_EEMS_Neutral/"
dirs = c(paste0(path, "neutralgenedata-D200-chain1"), paste0(path, "neutralgenedata-D300-chain1"), paste0(path,
```

```

eems.plots(mcmcpath = c(paste0(path, "./neutralgenedata-D200-chain1"), paste0(path, "neutralgenedata-D300-
    longlat = T,add.grid=F,add.outline = T,add.demes = T,
    projection.in = "+proj=longlat +datum=WGS84",projection.out = "+proj=merc +datum=WGS84",
    add.map = T,add.abline = T, add.r.squared = T)

## Input projection: +proj=longlat +datum=WGS84
## Output projection: +proj=merc +datum=WGS84

## Loading rgdal (required by projection.in)
## Loading rworldmap (required by add.map)
## Loading rworldxtra (required by add.map)

## Using the default DarkOrange to Blue color scheme, with 'white' as the midpoint color.
## It combines two color schemes from the 'dichromat' package, which itself is based on
## a collection of color schemes for scientific data graphics:
## Light A and Bartlein PJ (2004). The End of the Rainbow? Color Schemes for Improved Data
## Graphics. EOS Transactions of the American Geophysical Union, 85(40), 385.
## See also http://geog.uoregon.edu/datagraphics/color\_scales.htm

## Using 'euclidean' distance to assign interpolation points to Voronoi tiles.

## Processing the following EEMS output directories :
## ./NB_EEMS_Neutral./neutralgenedata-D200-chain1./NB_EEMS_Neutral/neutralgenedata-D300-chain1./NB_EEMS_
## Plotting effective migration surface (posterior mean of m rates)
## ./NB_EEMS_Neutral./neutralgenedata-D200-chain1
## ./NB_EEMS_Neutral/neutralgenedata-D300-chain1
## ./NB_EEMS_Neutral/neutralgenedata-D600-chain1

## Using the default DarkOrange to Blue color scheme, with 'white' as the midpoint color.
## It combines two color schemes from the 'dichromat' package, which itself is based on
## a collection of color schemes for scientific data graphics:
## Light A and Bartlein PJ (2004). The End of the Rainbow? Color Schemes for Improved Data
## Graphics. EOS Transactions of the American Geophysical Union, 85(40), 385.
## See also http://geog.uoregon.edu/datagraphics/color\_scales.htm

## Plotting effective diversity surface (posterior mean of q rates)
## ./NB_EEMS_Neutral./neutralgenedata-D200-chain1
## ./NB_EEMS_Neutral/neutralgenedata-D300-chain1
## ./NB_EEMS_Neutral/neutralgenedata-D600-chain1

## Using the default DarkOrange to Blue color scheme, with 'white' as the midpoint color.
## It combines two color schemes from the 'dichromat' package, which itself is based on
## a collection of color schemes for scientific data graphics:
## Light A and Bartlein PJ (2004). The End of the Rainbow? Color Schemes for Improved Data
## Graphics. EOS Transactions of the American Geophysical Union, 85(40), 385.
## See also http://geog.uoregon.edu/datagraphics/color\_scales.htm

## Plotting posterior probability trace
## ./NB_EEMS_Neutral./neutralgenedata-D200-chain1
## ./NB_EEMS_Neutral/neutralgenedata-D300-chain1
## ./NB_EEMS_Neutral/neutralgenedata-D600-chain1

```

```

## Plotting average dissimilarities within and between demes
## ./NB_EEMS_Neutral./neutralgenedata-D200-chain1
## ./NB_EEMS_Neutral/neutralgenedata-D300-chain1
## Using 'euclidean' distance to assign interpolation points to Voronoi tiles.
##
##
##
## Using 'euclidean' distance to assign interpolation points to Voronoi tiles.
## EEMS results for at least two different population grids

```

## Pairwise Fst

```

fst.mat_gene <- pairwise.WCfst(hf.filt_gene)

gindF.fst.mat.triN_gene <- as.matrix(fst.mat_gene)
colnames(gindF.fst.mat.triN_gene) <- pop_order
rownames(gindF.fst.mat.triN_gene) <- pop_order

meltedN_gene <- melt(gindF.fst.mat.triN_gene, na.rm =TRUE)
round(gindF.fst.mat.triN_gene,4)

##          BIS      GB      NAR      PVD
## BIS      NA 0.0006 0.0091 -0.0003
## GB      0.0006      NA 0.0108  0.0000
## NAR     0.0091 0.0108      NA 0.0106
## PVD    -0.0003 0.0000 0.0106      NA

summary(meltedN_gene$value)

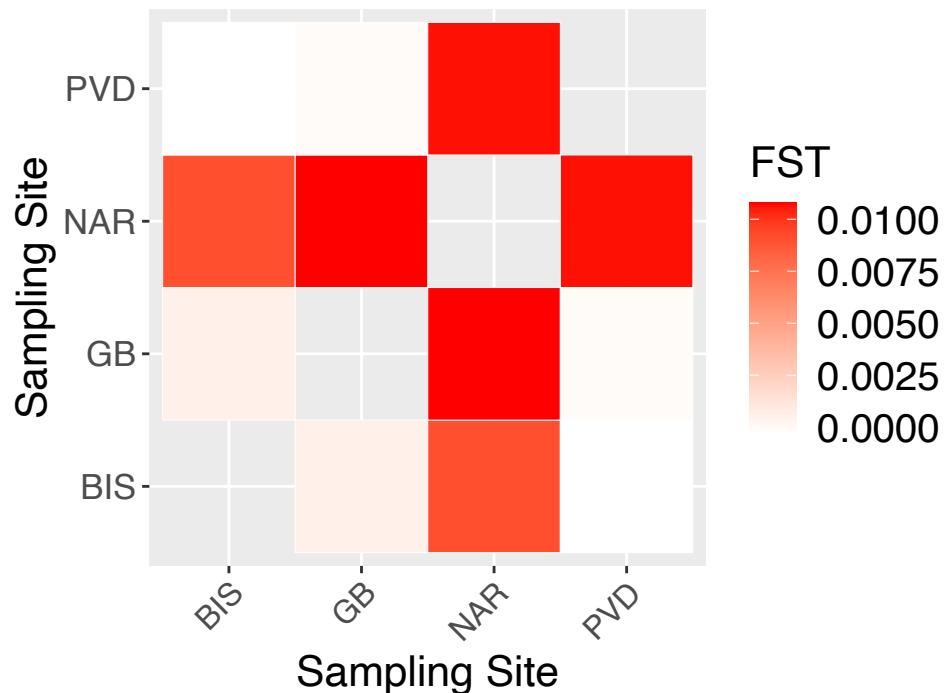
##      Min.   1st Qu.   Median   Mean   3rd Qu.   Max.
## -2.782e-04 -3.836e-05  4.828e-03  5.122e-03  1.063e-02  1.076e-02

#Plotting Pairwise fst
neutral <- ggplot(data = meltedN_gene, aes(Var2, Var1, fill = value)) + geom_tile(color = "white") +
  scale_fill_gradient(low = "white", high = "red", name="FST") +
  ggtitle(expression(atop("Pairwise FST, WC (1984) Neutral Gene SNPs", atop(italic("N = 40, L = 66,939"))))) +
  labs( x = "Sampling Site", y = "Sampling Site") +
  theme(axis.text.x = element_text(angle = 45, vjust = 1, size = 12, hjust = 1),axis.text.y = element_text(size = 12),
        axis.title = element_text(size = 16),legend.text = element_text(size = 15), legend.title = element_text(size = 17)) +
  theme(plot.title = element_text(size = 17)) +
  coord_fixed()
neutral

```

## Pairwise FST, WC (1984) Neutral Gene SNPs

$N = 40, L = 66,939$



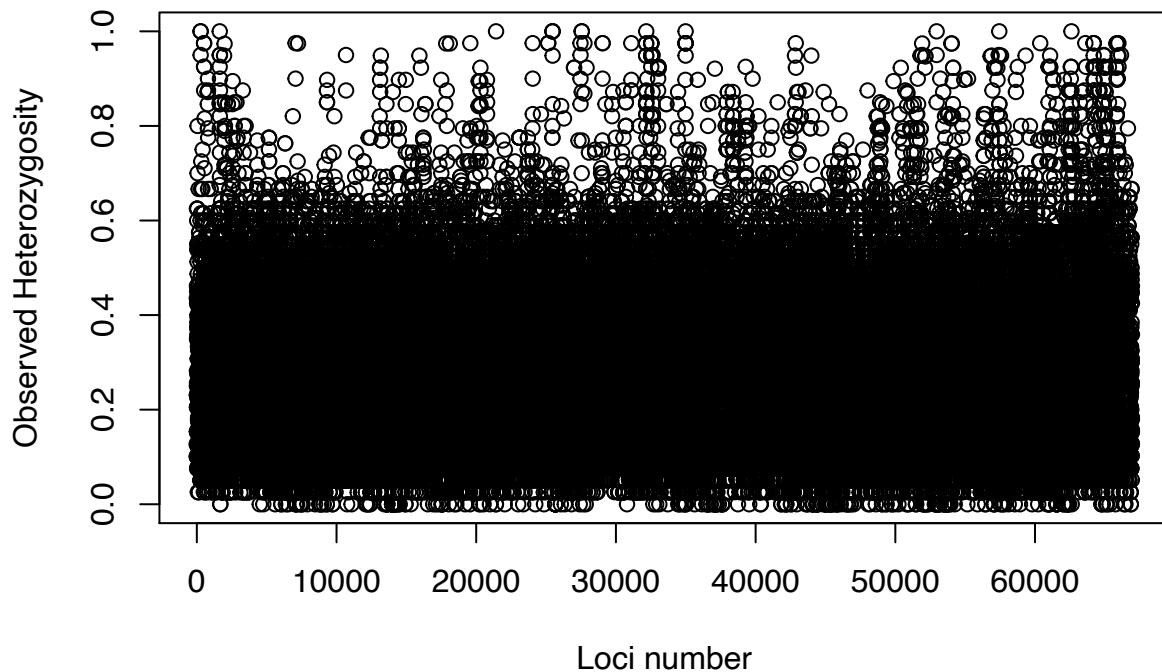
### Genetic diversity (observed and expected heterozygosity)

```
comb_gene <- summary(stratated.filt_gene)
names(comb_gene)

## [1] "n"          "n.by.pop"   "loc.n.all"  "pop.n.all"  "NA.perc"    "Hobs"
## [7] "Hexp"

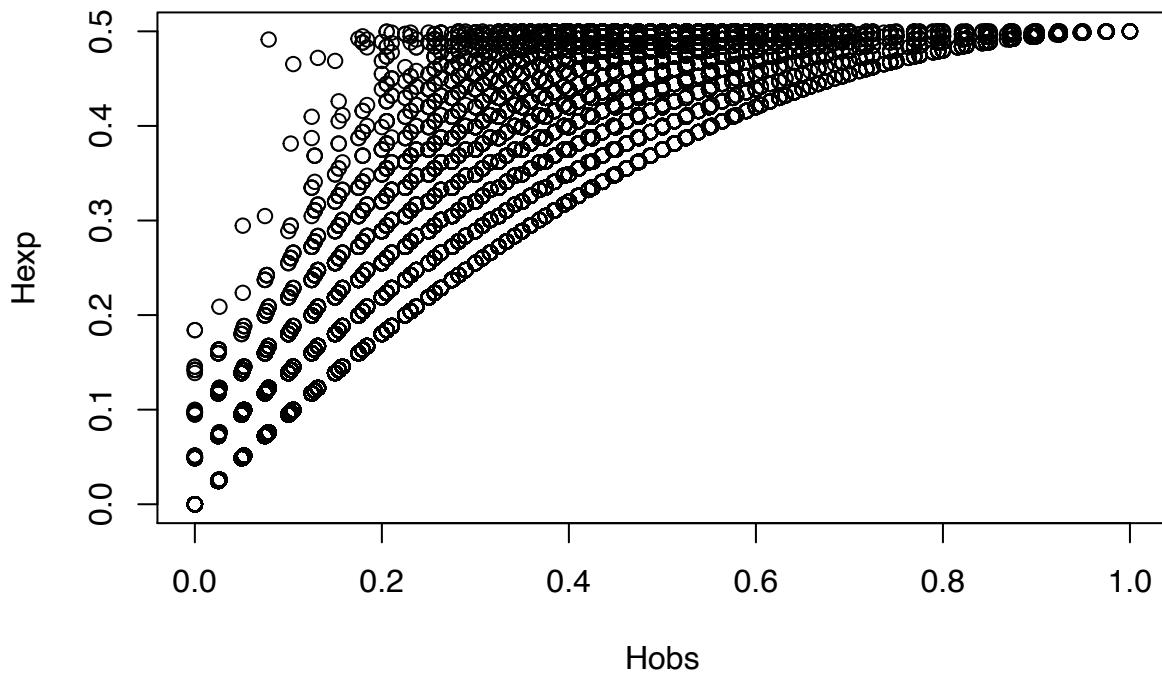
plot(comb_gene$Hobs, xlab="Loci number", ylab="Observed Heterozygosity",
     main="Observed heterozygosity per locus")
```

### Observed heterozygosity per locus



```
plot(comb_gene$Hobs,comb_gene$Hexp, xlab="Hobs", ylab="Hexp",
  main="Expected heterozygosity as a function of observed heterozygosity per locus")
```

## Expected heterozygosity as a function of observed heterozygosity per I



```
bartlett.test(list(comb_gene$Hexp, comb_gene$Hobs)) # a test : H0: Hexp = Hobs
```

```
##  
##  Bartlett test of homogeneity of variances  
##  
## data: list(comb_gene$Hexp, comb_gene$Hobs)  
## Bartlett's K-squared = 1440.3, df = 1, p-value < 2.2e-16
```

*Significant difference between Observed and expected heterozygosity.*

```
basicstat_gene <- basic.stats(hf.filt_gene, diploid = TRUE, digits = 3)  
  
as.data.frame(basicstat_gene$overall)
```

```
##      basicstat_gene$overall  
## Ho          0.277  
## Hs          0.272  
## Ht          0.274  
## Dst         0.001  
## Htp         0.274  
## Dstp        0.001  
## Fst         0.004  
## Fstp        0.005  
## Fis         -0.016  
## Dest        0.002
```

```
# get bootstrap confidence values for Fis  
boot_gene <- boot.ppfis(hf.filt_gene, nboot = 1000)
```

```

boot5_gene <- boot.ppfis(hf.filt_gene,nboot = 1000,quant = 0.5)

# combine all pop statistics
colnames(basicstat_gene$Ho) <- pop_order
Ho_gene <- colMeans(basicstat_gene$Ho,na.rm = T)
He_gene <- colMeans(basicstat_gene$Hs,na.rm = T)
Fis_gene <- boot5_gene$fis.ci$ll
y_gene <- cbind(pop_order, Ho_gene, He_gene, Fis_gene, boot_gene$fis.ci, latitude, longitude, distance, sewage)
y_gene

##      pop_order   Ho_gene   He_gene Fis_gene      ll      hl latitude longitude
## BIS      BIS 0.2775137 0.2730459 -0.0164 -0.0194 -0.0136  41.545  -71.431
## GB       GB 0.2816576 0.2741167 -0.0276 -0.0305 -0.0244  41.654  -71.445
## NAR      NAR 0.2763264 0.2701158 -0.0230 -0.0264 -0.0200  41.505  -71.453
## PVD      PVD 0.2718104 0.2726113  0.0028 -0.0002  0.0060  41.816  -71.391
##      distance sewage temperature salinity pH Chlor_a DO
## BIS      4.76    8.82        23     30 7.9    4.9 8.2
## GB       0.47   14.60        24     28 7.4   18.8 5.7
## NAR     15.41    2.03        25     18 7.6    4.6 7.0
## PVD      1.49   59.86        23     25 7.4    8.1 4.9

summary(He_gene)

##      Min. 1st Qu. Median      Mean 3rd Qu.      Max.
##  0.2701 0.2720 0.2728 0.2725 0.2733 0.2741

summary(Fis_gene)

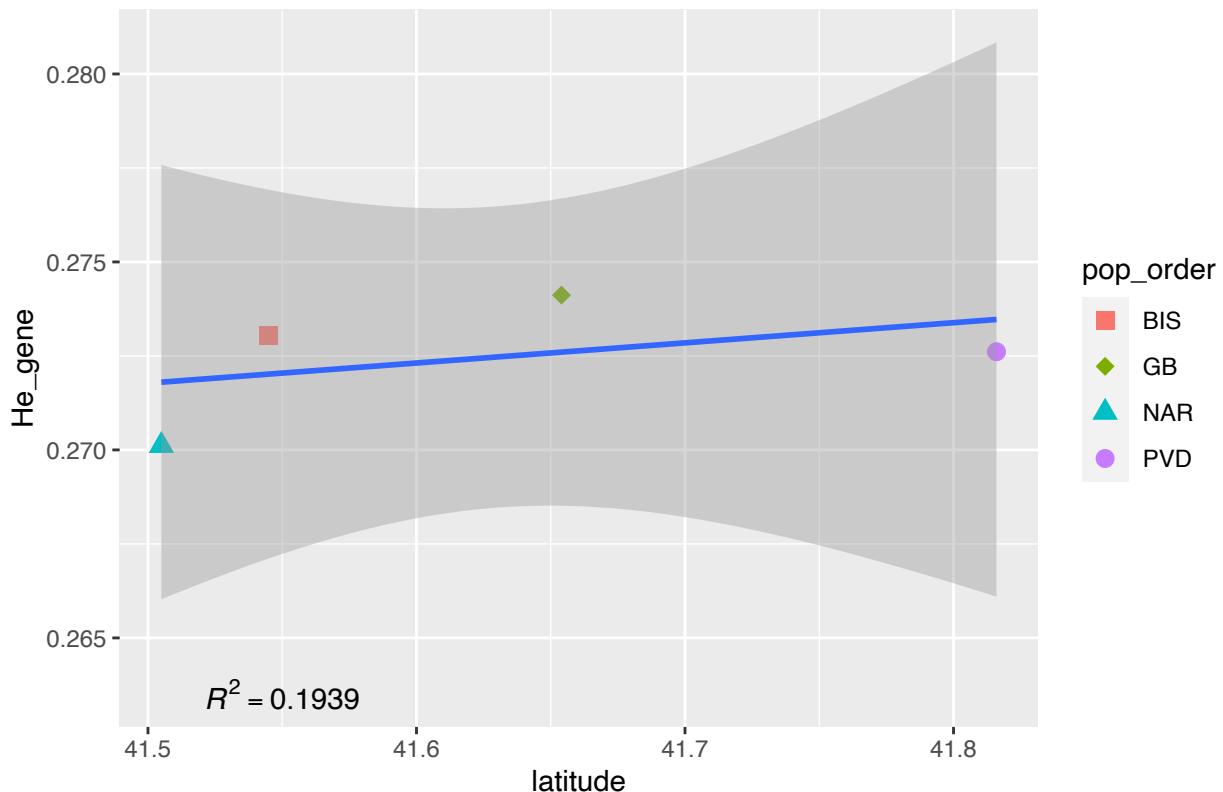
##      Min. 1st Qu. Median      Mean 3rd Qu.      Max.
## -0.02760 -0.02415 -0.01970 -0.01605 -0.01160  0.00280

# Plot He vs Latitude
R2_gene = round(summary(lm(y_gene$He_gene ~ y_gene$latitude))$r.squared, 4)
ggplot(y_gene, aes(x = latitude, y = He_gene)) + geom_point(aes(shape=pop_order, color=pop_order), size = scale_shape_manual(values = c(15,18,17,16))) +
geom_smooth(method=lm) +
ggttitle("Expected heterozygosity vs Latitude, Neutral Gene SNPs") +
annotate(geom = "text", label=paste("italic(R^2)==",R2_gene), x=41.55, y=0.2635, parse=T) +
scale_x_continuous()

## `geom_smooth()` using formula 'y ~ x'

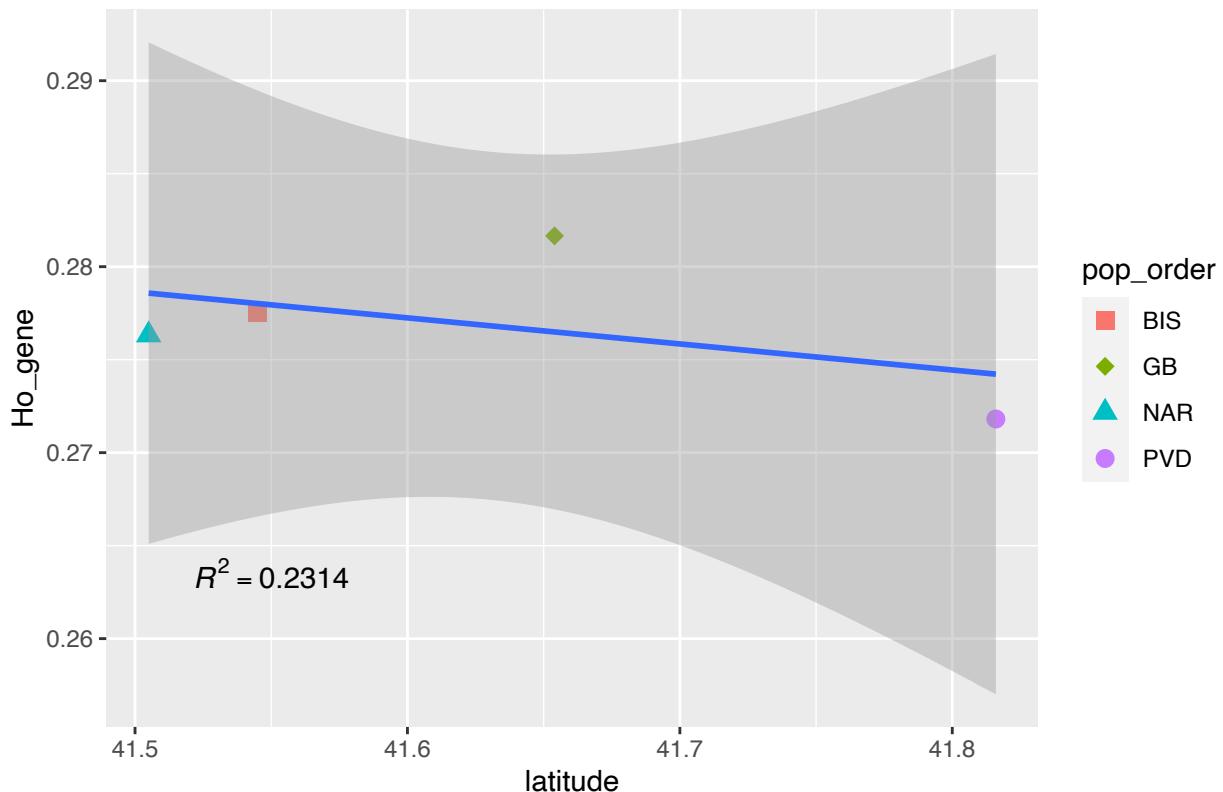
```

### Expected heterozygosity vs Latitude, Neutral Gene SNPs



```
# Plot Ho vs Latitude
R2_gene = round(summary(lm(y_gene$Ho_gene ~ y_gene$latitude))$r.squared, 4)
ggplot(y_gene, aes(x = latitude, y = Ho_gene)) + geom_point(aes(shape=pop_order, color=pop_order), size = scale_shape_manual(values = c(15,18,17,16))) +
  geom_smooth(method=lm) +
  ggttitle("Observed heterozygosity vs Latitude, Neutral Gene SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R2_gene), x=41.55, y=0.2635, parse=T) +
  scale_x_continuous()
## `geom_smooth()` using formula 'y ~ x'
```

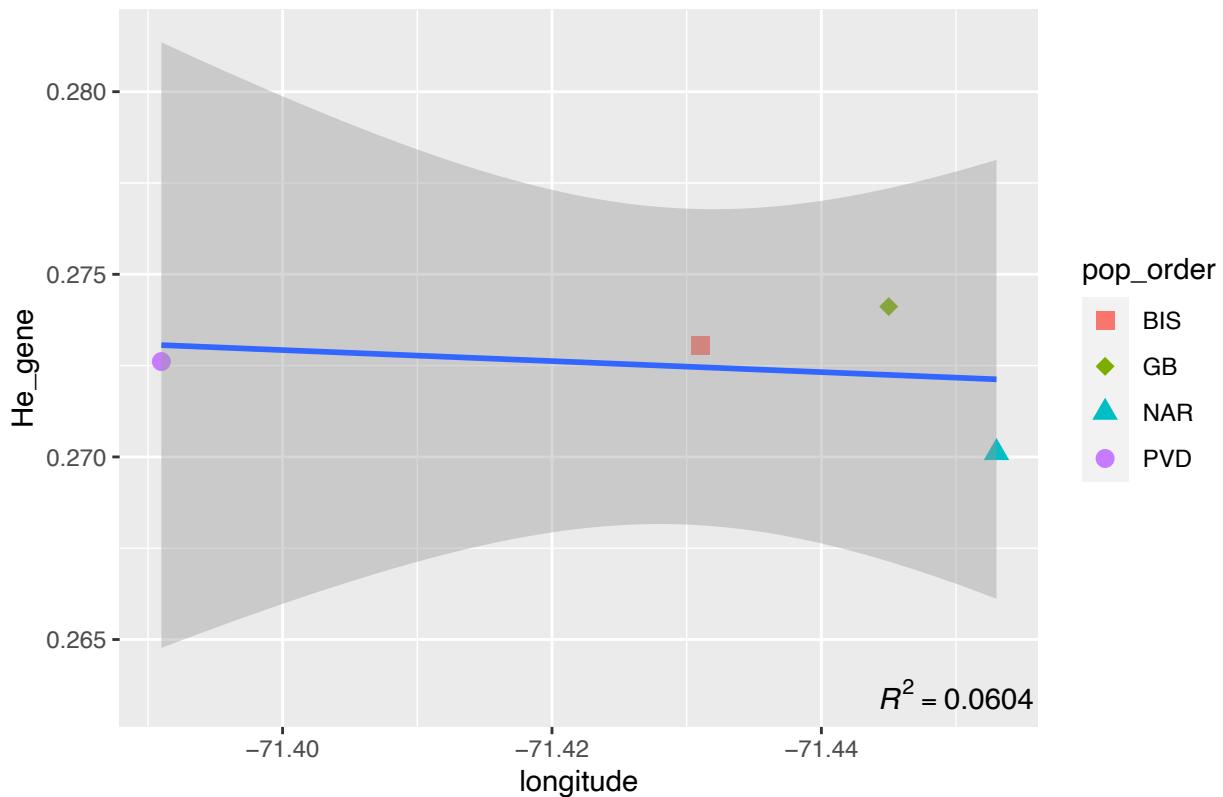
### Observed heterozygosity vs Latitude, Neutral Gene SNPs



```
#Plot He vs Longitude
R3_gene = round(summary(lm(y_gene$He_gene ~ y_gene$longitude))$r.squared, 4)
ggplot(y_gene, aes(x = longitude, y = He_gene)) + geom_point(aes(shape=pop_order, color=pop_order), size=10) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggttitle("Expected heterozygosity vs Longitude, Neutral Gene SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R3_gene), x=-71.45, y=0.2635, parse=T) +
  scale_x_reverse()

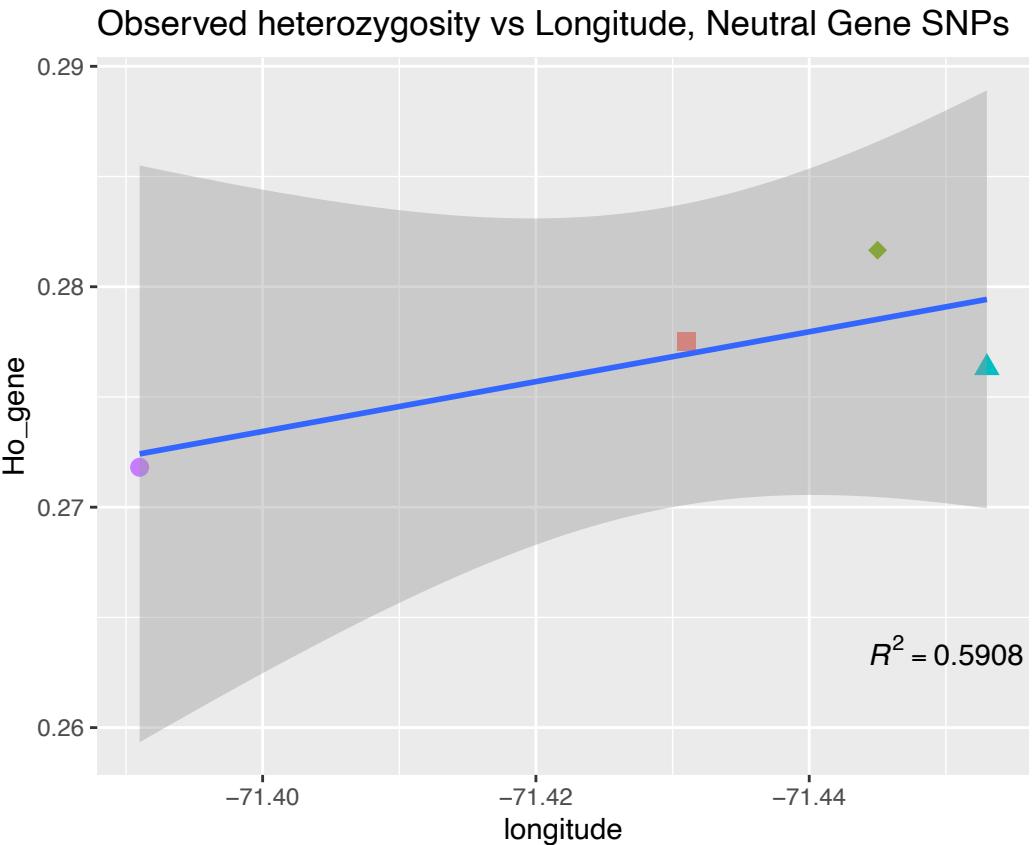
## `geom_smooth()` using formula 'y ~ x'
```

### Expected heterozygosity vs Longitude, Neutral Gene SNPs



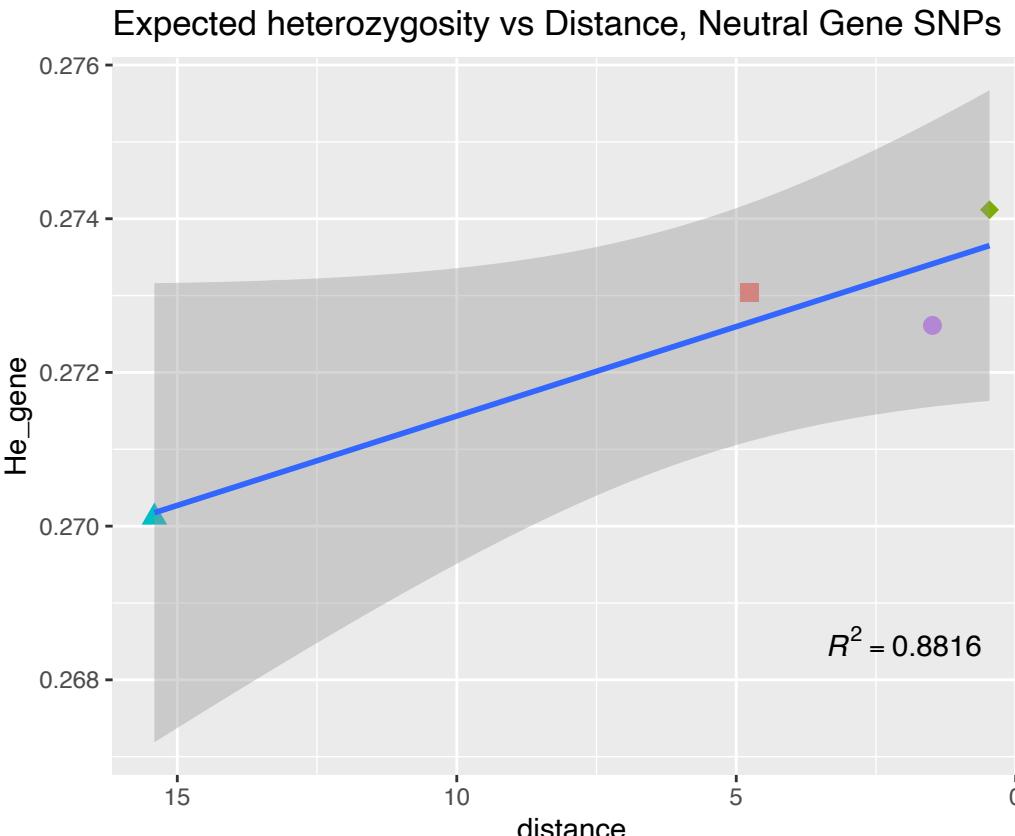
```
#Plot Ho vs Longitude
R3_gene = round(summary(lm(y_gene$Ho_gene ~ y_gene$longitude))$r.squared, 4)
ggplot(y_gene, aes(x = longitude, y = Ho_gene)) + geom_point(aes(shape=pop_order, color=pop_order), size=10) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  gtitle("Observed heterozygosity vs Longitude, Neutral Gene SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R3_gene), x=-71.45, y=0.2635, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```



```
#Plot He vs Distance from sewage outflow
R4_gene = round(summary(lm(y_gene$He_gene ~ y_gene$distance))$r.squared, 4)
ggplot(y_gene, aes(x = distance, y = He_gene)) + geom_point(aes(shape=pop_order, color=pop_order), size = scale_shape_manual(values = c(15,18,17,16))) +
  geom_smooth(method=lm) +
  ggtitle("Expected heterozygosity vs Distance, Neutral Gene SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R4_gene), x=2, y=0.2685, parse=T) +
  scale_x_reverse()

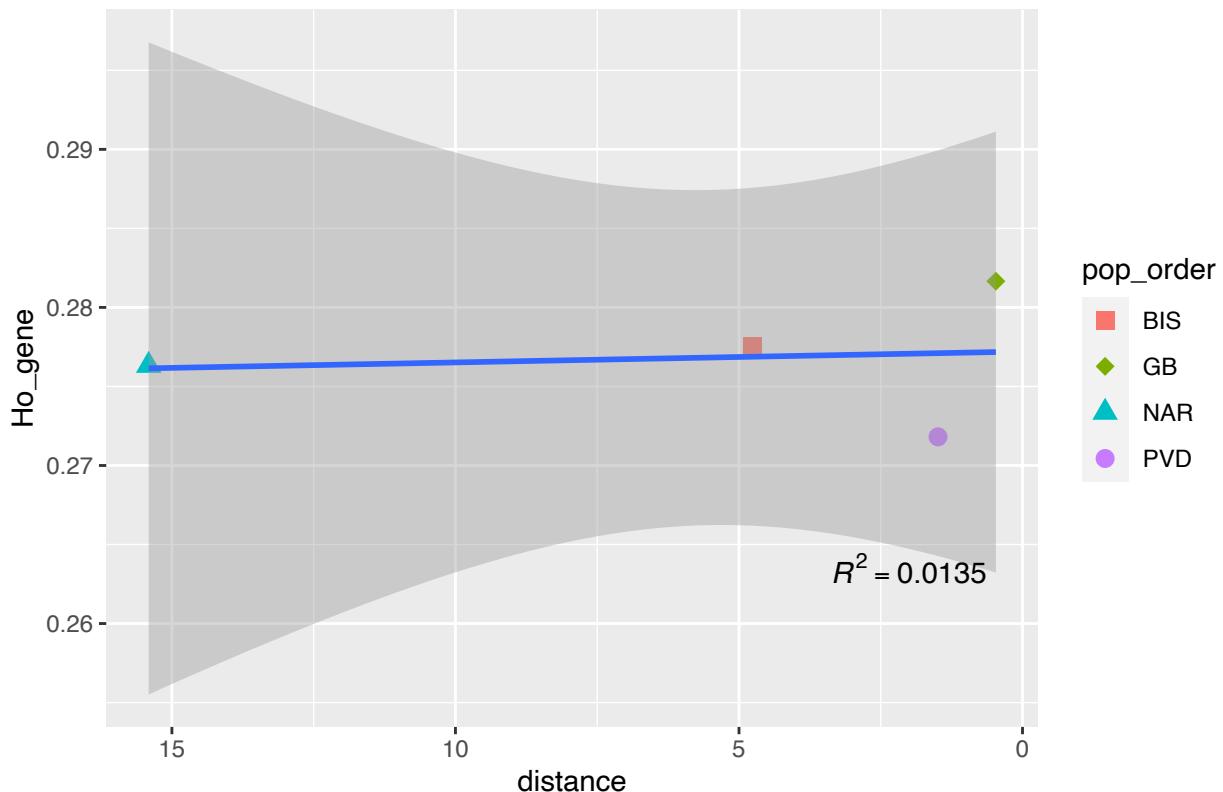
## `geom_smooth()` using formula 'y ~ x'
```



```
#Plot Ho vs Distance from sewage outflow
R4_gene = round(summary(lm(y_gene$Ho_gene ~ y_gene$distance))$r.squared, 4)
ggplot(y_gene, aes(x = distance, y = Ho_gene)) + geom_point(aes(shape=pop_order, color=pop_order), size=scale_shape_manual(values = c(15,18,17,16))) +
  geom_smooth(method=lm) +
  ggttitle("Observed heterozygosity vs Distance, Neutral Gene SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R4_gene), x=2, y=0.2635, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

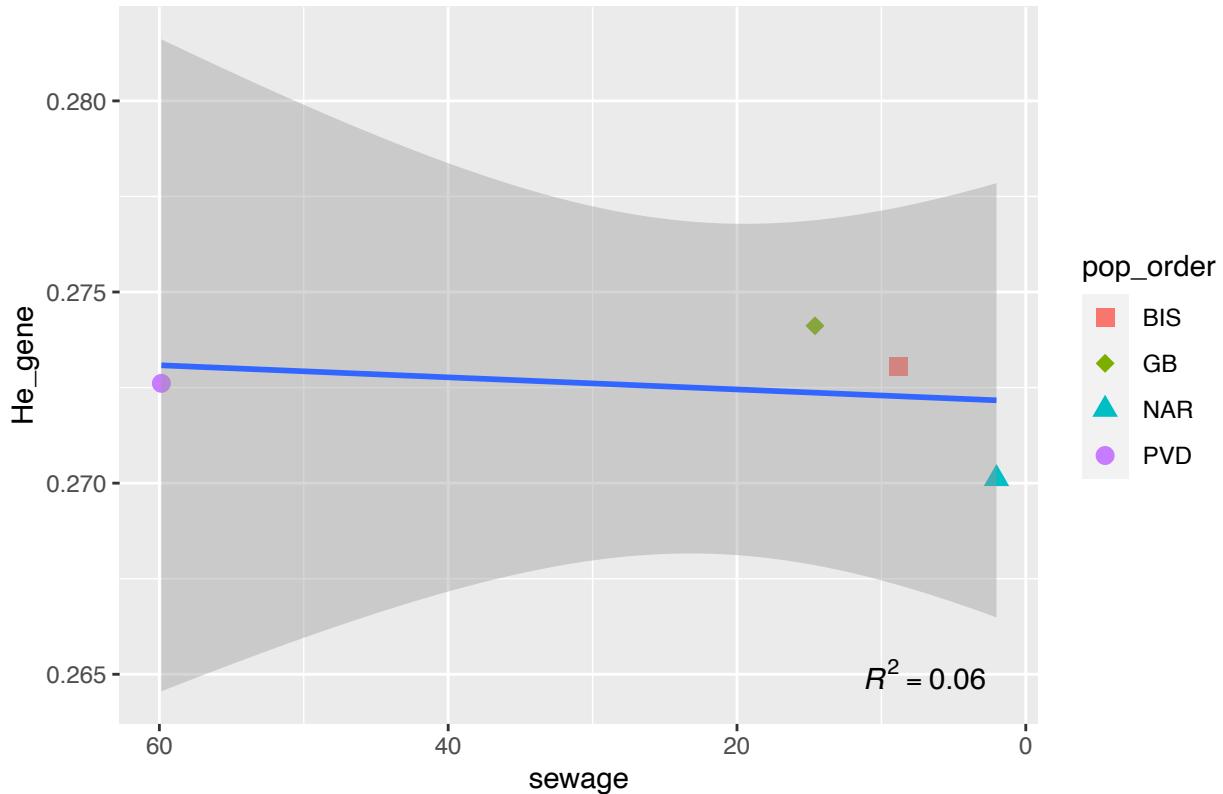
### Observed heterozygosity vs Distance, Neutral Gene SNPs



```
#Plot He vs Sewage effluent
R_gene = round(summary(lm(y_gene$He_gene ~ y_gene$sewage))$r.squared, 4)
ggplot(y_gene, aes(x = sewage, y = He_gene)) + geom_point(aes(shape=pop_order, color=pop_order), size =
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggttitle("Expected heterozygosity vs Sewage Effluent, Neutral Gene SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R_gene), x=7, y=0.265, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

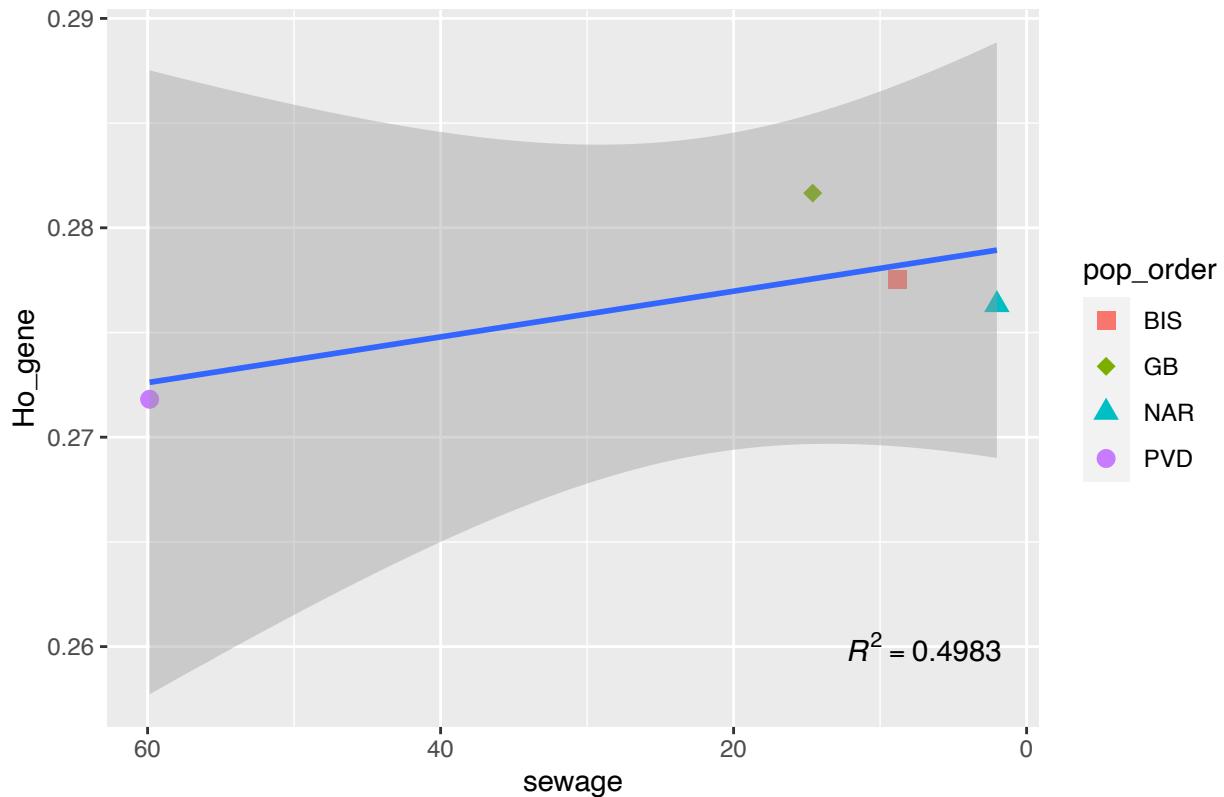
## Expected heterozygosity vs Sewage Effluent, Neutral Gene SNPs



```
#Plot Ho vs Sewage effluent
R_gene = round(summary(lm(y_gene$Ho_gene ~ y_gene$sewage))$r.squared, 4)
ggplot(y_gene, aes(x = sewage, y = Ho_gene)) + geom_point(aes(shape=pop_order, color=pop_order), size =
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggttitle("Observed heterozygosity vs Sewage Effluent, Neutral Gene SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R_gene), x=7, y=0.26, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

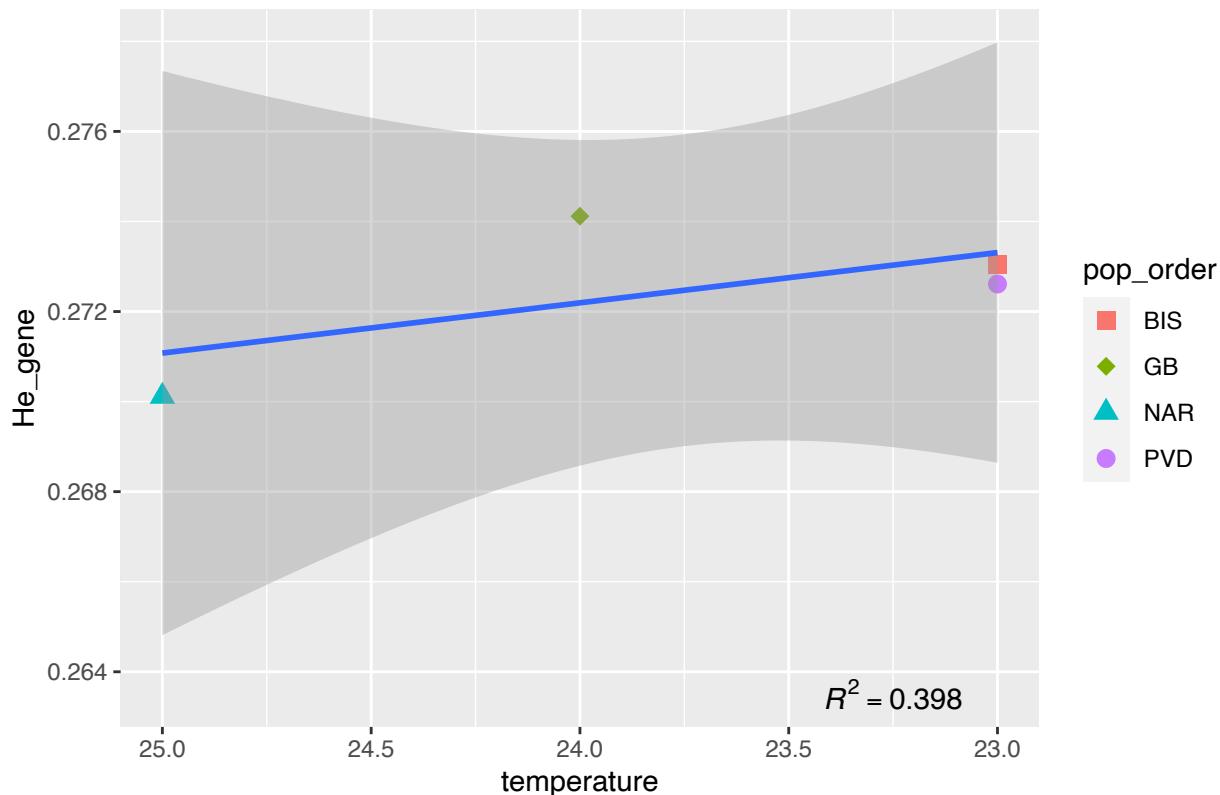
### Observed heterozygosity vs Sewage Effluent, Neutral Gene SNPs



```
#Plot He vs Temperature
R5_gene = round(summary(lm(y_gene$He_gene ~ y_gene$temperature))$r.squared, 4)
ggplot(y_gene, aes(x = temperature, y = He_gene)) + geom_point(aes(shape=pop_order, color=pop_order), size=5) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggtitle("Expected heterozygosity vs Temperature, Neutral Gene SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R5_gene), x=23.25, y=0.2635, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

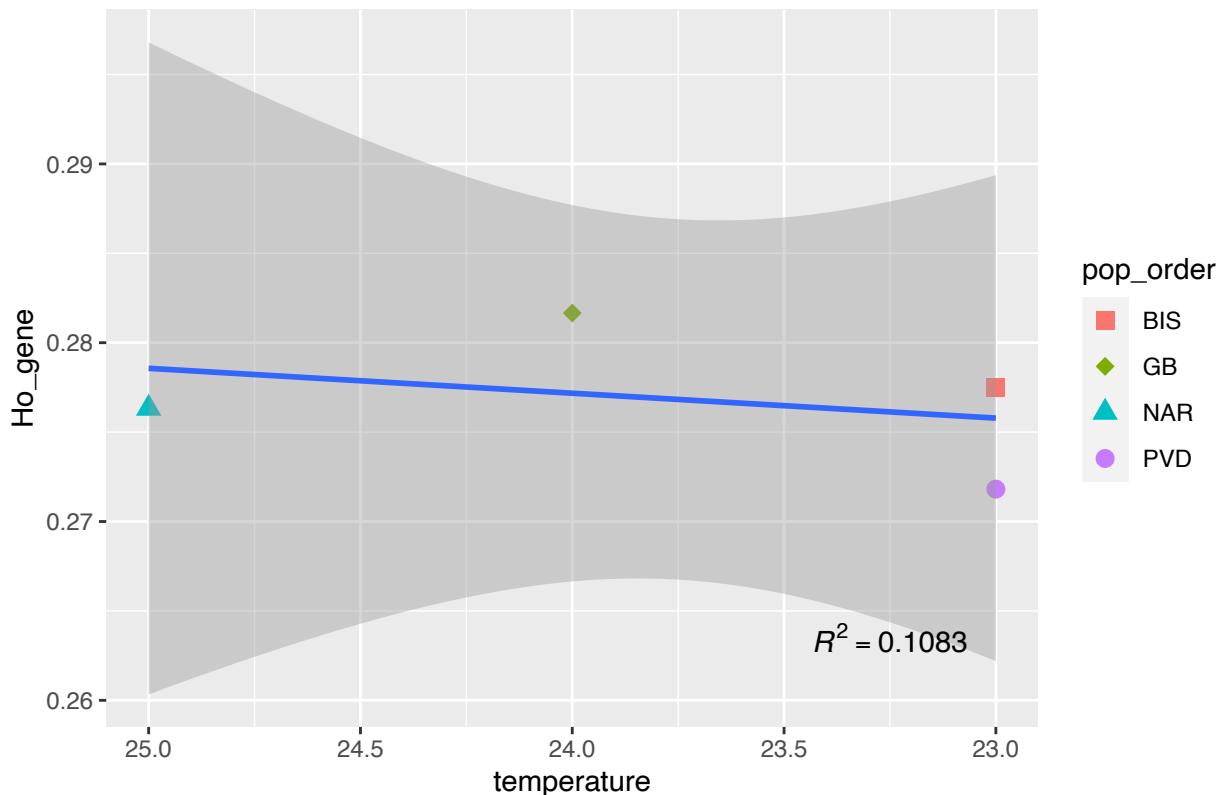
## Expected heterozygosity vs Temperature, Neutral Gene SNPs



```
#Plot Ho vs Temperature
R5_gene = round(summary(lm(y_gene$Ho_gene ~ y_gene$temperature))$r.squared, 4)
ggplot(y_gene, aes(x = temperature, y = Ho_gene)) + geom_point(aes(shape=pop_order, color=pop_order), s
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  gtitle("Observed heterozygosity vs Temperature, Neutral Gene SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R5_gene), x=23.25, y=0.2635, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

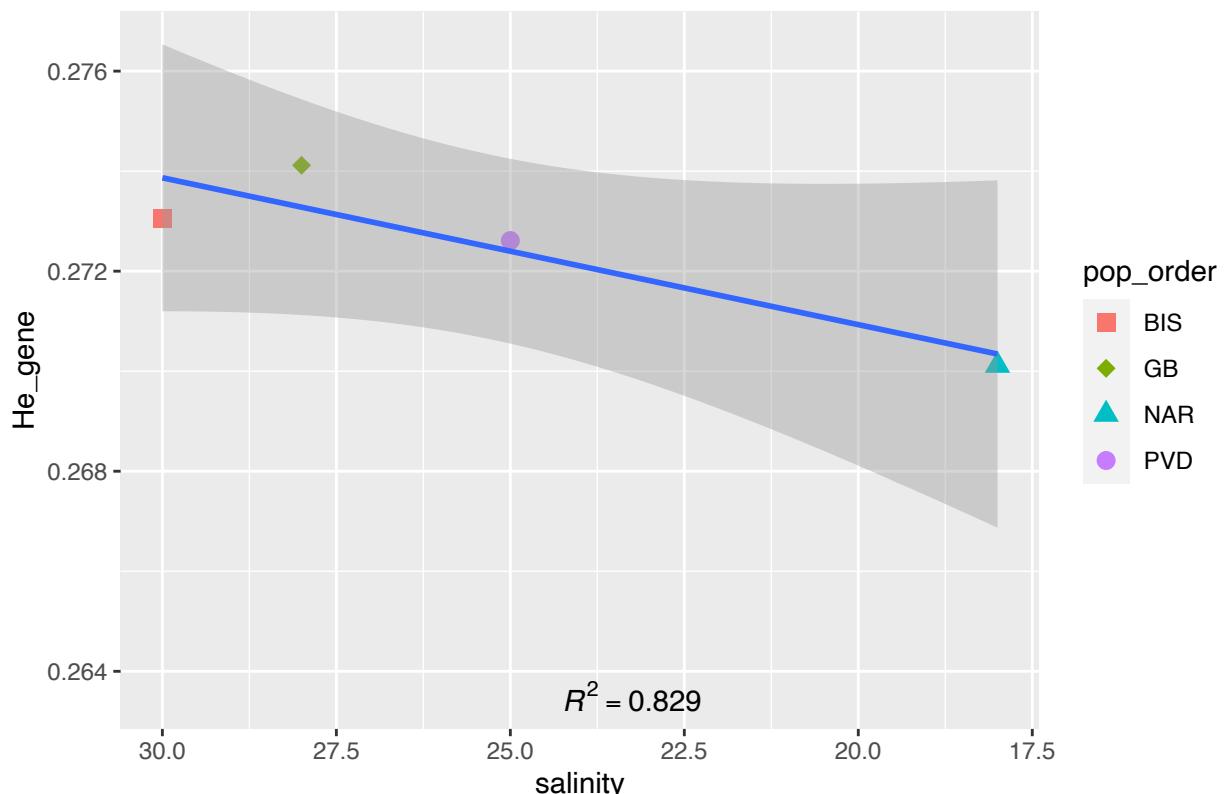
## Observed heterozygosity vs Temperature, Neutral Gene SNPs



```
#Plot He vs Salinity
R6_gene = round(summary(lm(y_gene$He_gene ~ y_gene$salinity))$r.squared, 4)
ggplot(y_gene, aes(x = salinity, y = He_gene)) + geom_point(aes(shape=pop_order, color=pop_order), size = 15) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggttitle("Expected heterozygosity vs Salinity, Neutral Gene SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R6_gene), x=23.25, y=0.2635, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

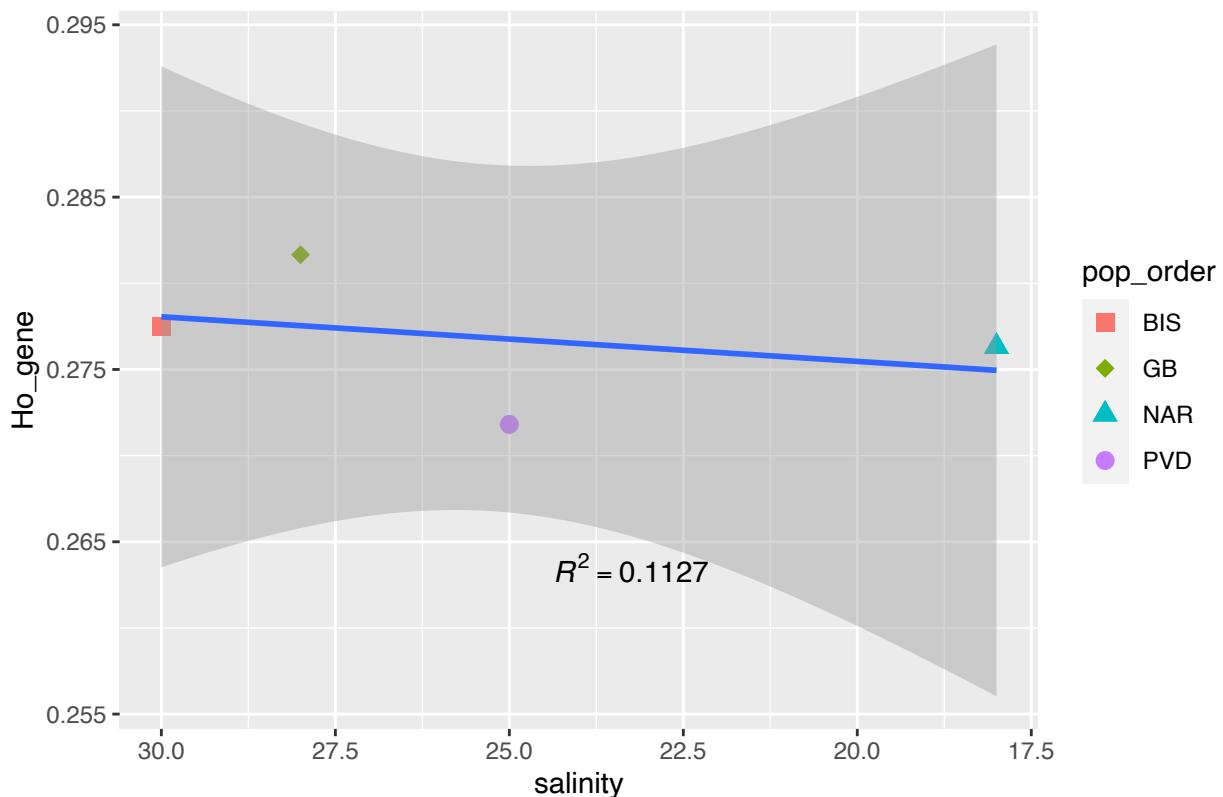
### Expected heterozygosity vs Salinity, Neutral Gene SNPs



```
#Plot Ho vs Salinity
R6_gene = round(summary(lm(y_gene$Ho_gene ~ y_gene$salinity))$r.squared, 4)
ggplot(y_gene, aes(x = salinity, y = Ho_gene)) + geom_point(aes(shape=pop_order, color=pop_order), size = scale_shape_manual(values = c(15,18,17,16))) +
  geom_smooth(method=lm) +
  ggttitle("Observed heterozygosity vs Salinity, Neutral Gene SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R6_gene), x=23.25, y=0.2635, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

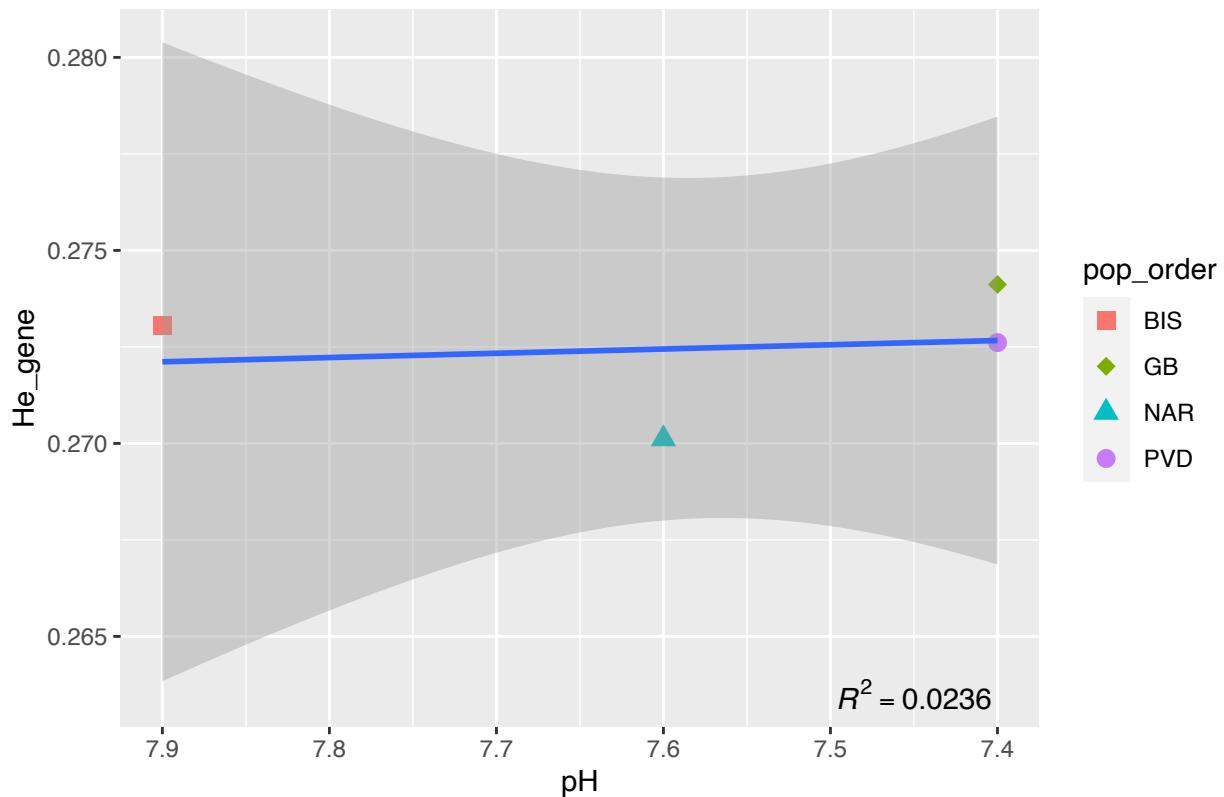
### Observed heterozygosity vs Salinity, Neutral Gene SNPs



```
#Plot He vs pH
R7_gene = round(summary(lm(y_gene$He_gene ~ y_gene$pH))$r.squared, 4)
ggplot(y_gene, aes(x = pH, y = He_gene)) + geom_point(aes(shape=pop_order, color=pop_order), size = 3) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggttitle("Expected heterozygosity vs pH, Neutral Gene SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R7_gene), x=7.45, y=0.2635, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

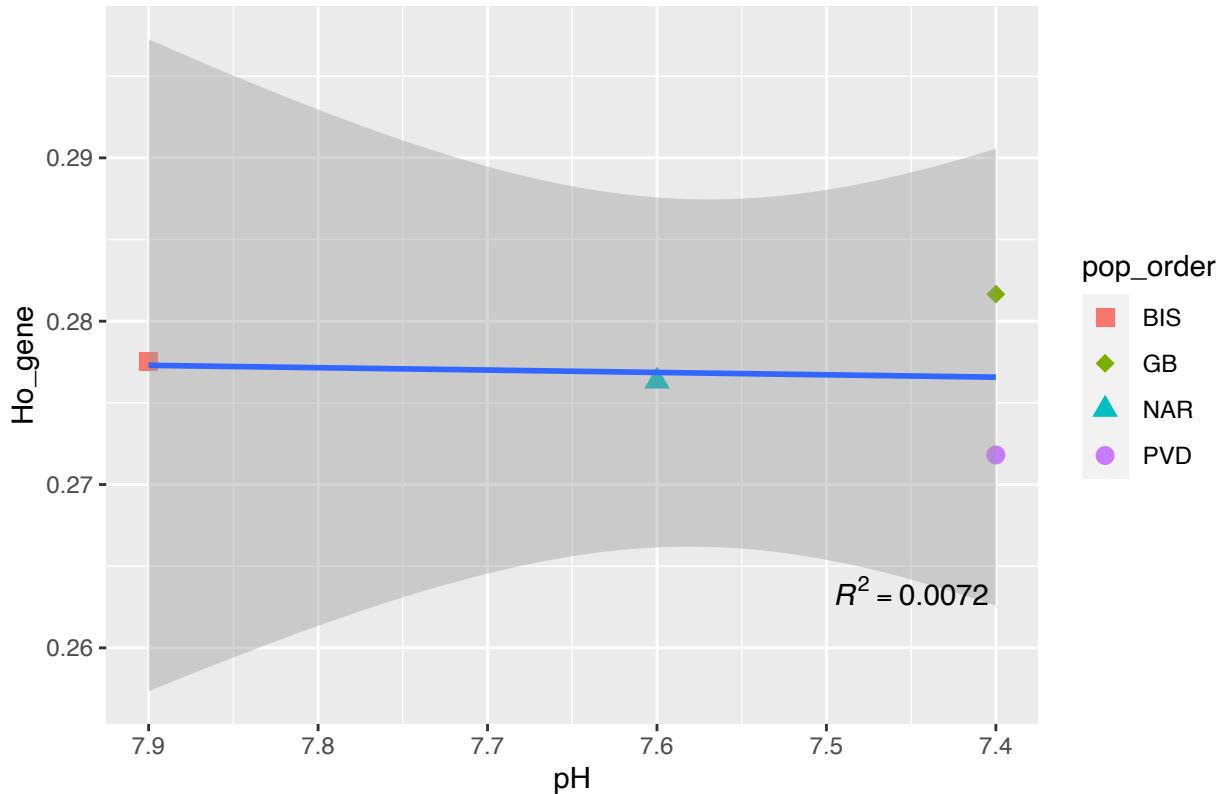
### Expected heterozygosity vs pH, Neutral Gene SNPs



```
#Plot Ho vs pH
R7_gene = round(summary(lm(y_gene$Ho_gene ~ y_gene$pH))$r.squared, 4)
ggplot(y_gene, aes(x = pH, y = Ho_gene)) + geom_point(aes(shape=pop_order, color=pop_order), size = 3) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggtitle("Observed heterozygosity vs pH, Neutral Gene SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R7_gene), x=7.45, y=0.2635, parse=T) +
  scale_x_reverse()

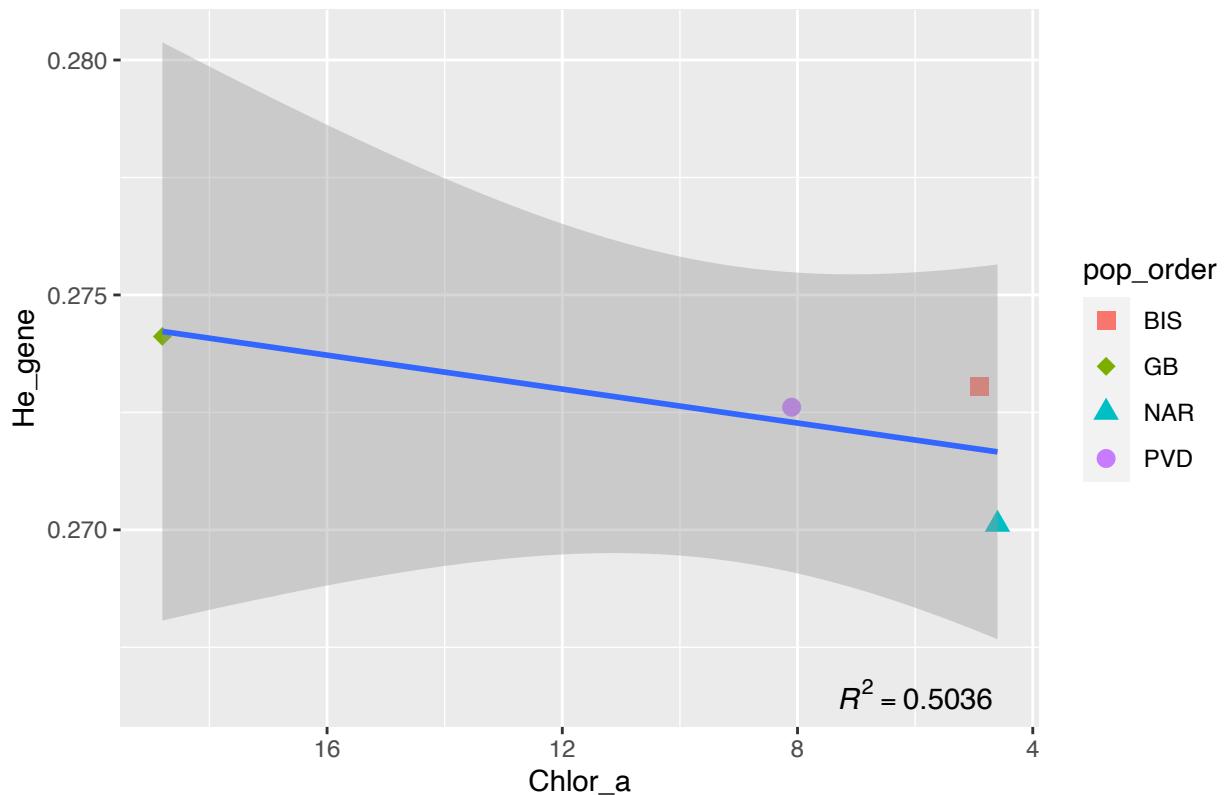
## `geom_smooth()` using formula 'y ~ x'
```

### Observed heterozygosity vs pH, Neutral Gene SNPs



```
#Plot He vs Chlorophyll a
R8_gene = round(summary(lm(y_gene$He_gene ~ y_gene$Chlor_a))$r.squared, 4)
ggplot(y_gene, aes(x = Chlor_a, y = He_gene)) + geom_point(aes(shape=pop_order, color=pop_order), size =
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggttitle("Expected heterozygosity vs Chlorophyll a, Neutral Gene SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R8_gene), x=6, y=0.2665, parse=T) +
  scale_x_reverse()
## `geom_smooth()` using formula 'y ~ x'
```

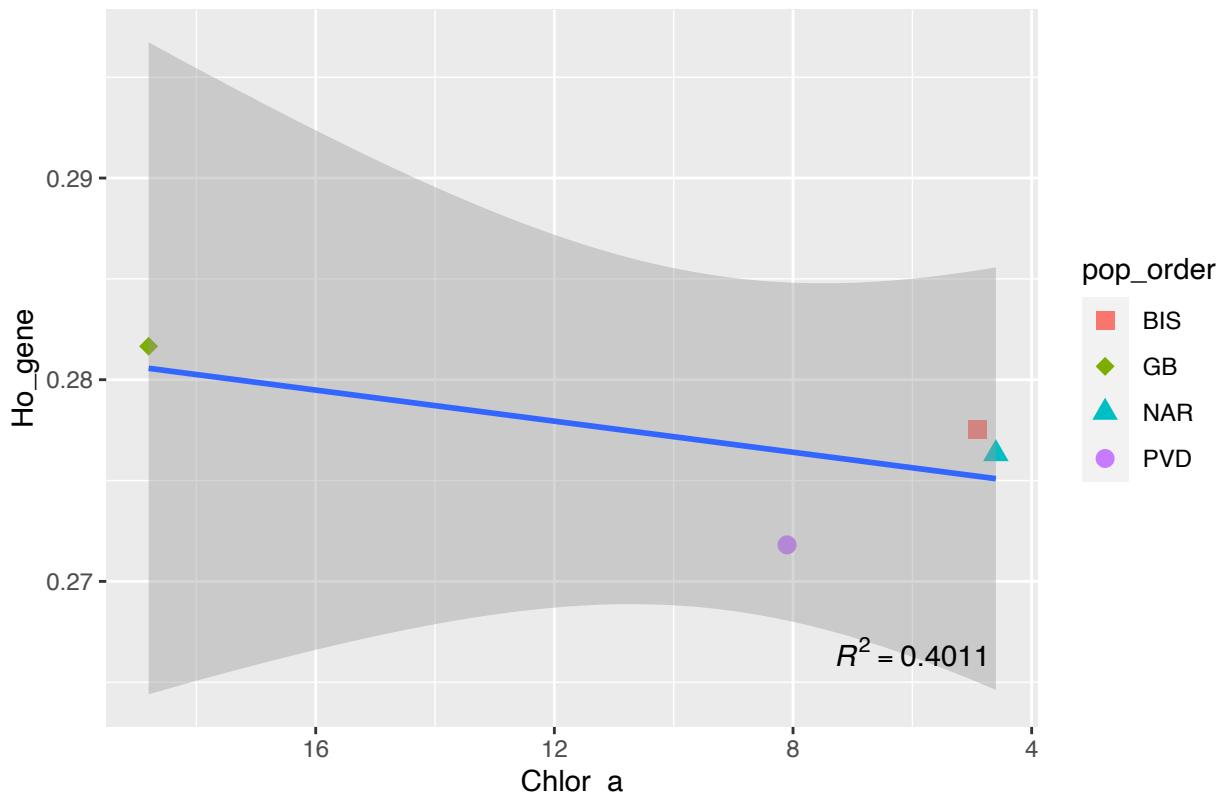
### Expected heterozygosity vs Chlorophyll a, Neutral Gene SNPs



```
#Plot Ho vs Chlorophyll a
R8_gene = round(summary(lm(y_gene$Ho_gene ~ y_gene$Chlor_a))$r.squared, 4)
ggplot(y_gene, aes(x = Chlor_a, y = Ho_gene)) + geom_point(aes(shape=pop_order, color=pop_order), size =
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggttitle("Observed heterozygosity vs Chlorophyll a, Neutral Gene SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R8_gene), x=6, y=0.2665, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

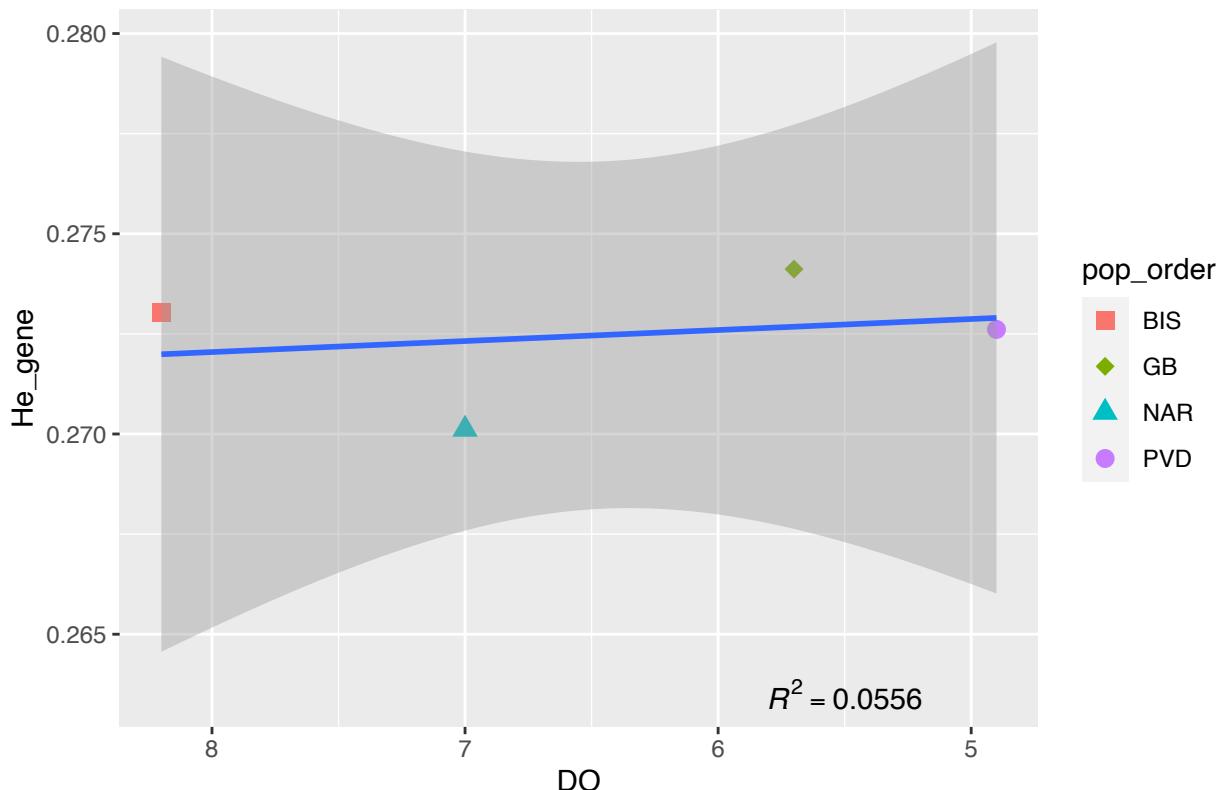
## Observed heterozygosity vs Chlorophyll a, Neutral Gene SNPs



```
#Plot He vs Dissolved Oxygen
R9_gene = round(summary(lm(y_gene$He_gene ~ y_gene$DO))$r.squared, 4)
ggplot(y_gene, aes(x = DO, y = He_gene)) + geom_point(aes(shape=pop_order, color=pop_order), size = 3) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggttitle("Expected heterozygosity vs Dissolved Oxygen, Neutral Gene SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R9_gene), x=5.5, y=0.2635, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

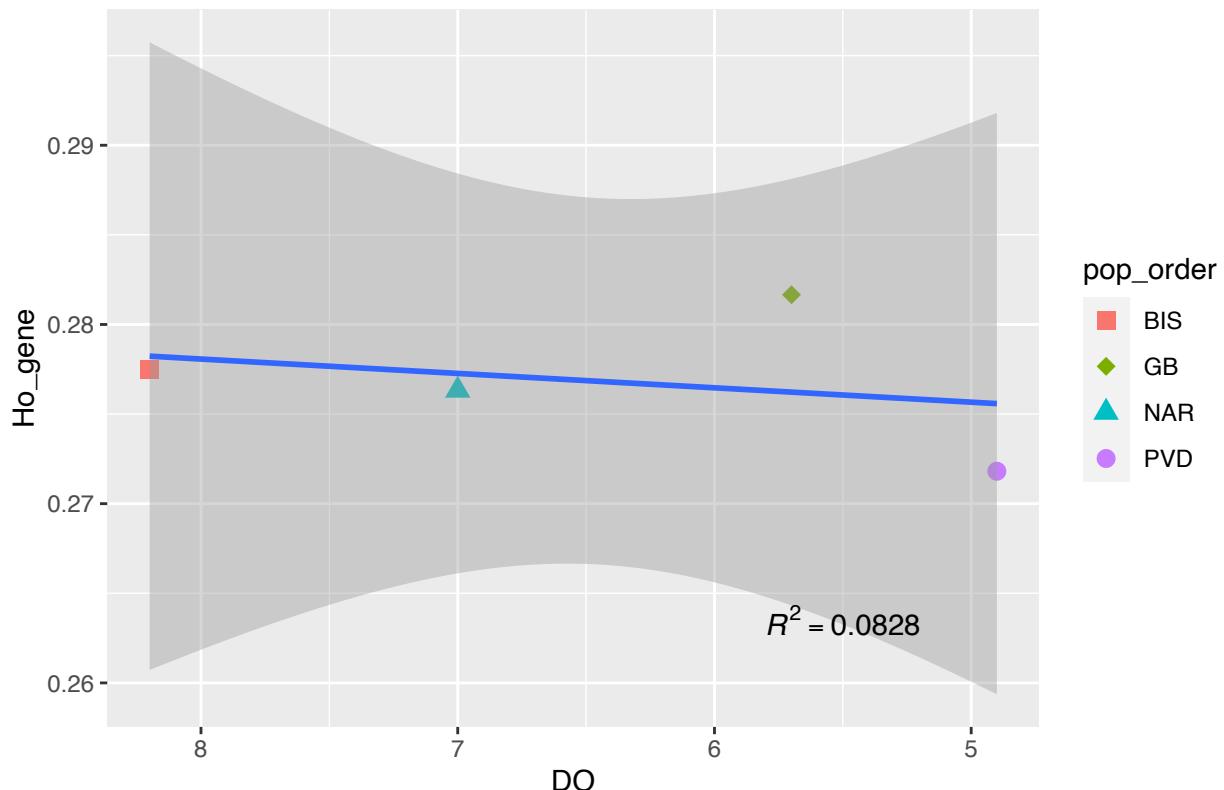
### Expected heterozygosity vs Dissolved Oxygen, Neutral Gene SNPs



```
#Plot Ho vs Dissolved Oxygen
R9_gene = round(summary(lm(y_gene$Ho_gene ~ y_gene$DO))$r.squared, 4)
ggplot(y_gene, aes(x = DO, y = Ho_gene)) + geom_point(aes(shape=pop_order, color=pop_order), size = 3) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggttitle("Observed heterozygosity vs Dissolved Oxygen, Neutral Gene SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R9_gene), x=5.5, y=0.2635, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

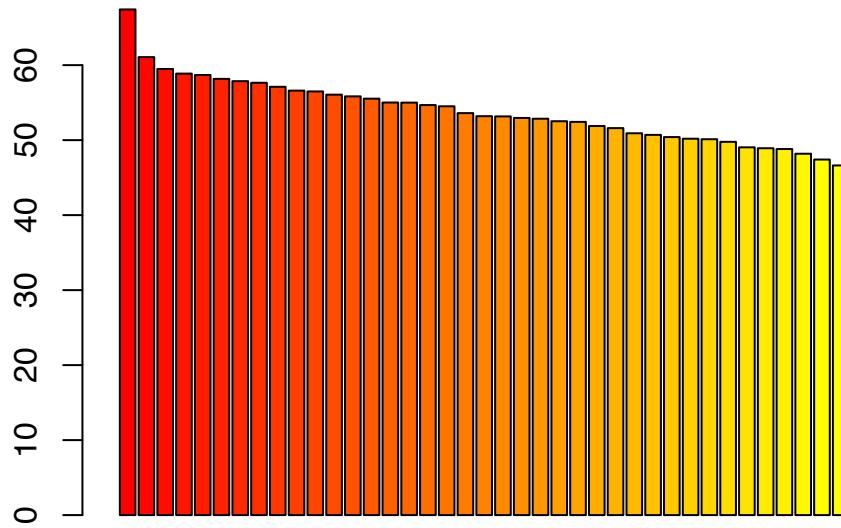
## Observed heterozygosity vs Dissolved Oxygen, Neutral Gene SNPs



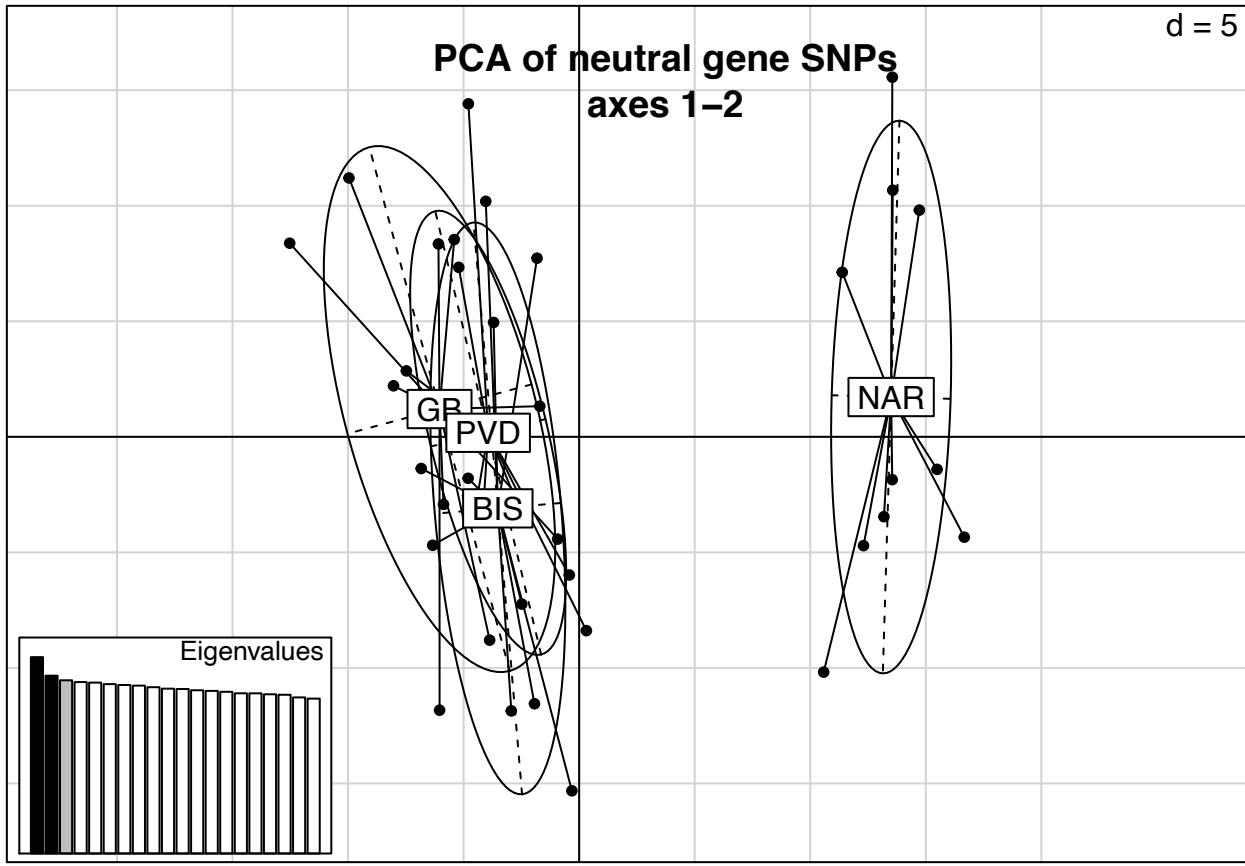
## PCA

```
X_gene <- tab(stratated.u_gene, freq = TRUE, NA.method = "mean")
pca1_gene <- dudi.pca(X_gene, scale = FALSE, scannf = FALSE, nf = 3)
barplot(pca1_gene$eig[1:50], main = "PCA eigenvalues", col = heat.colors(50))
```

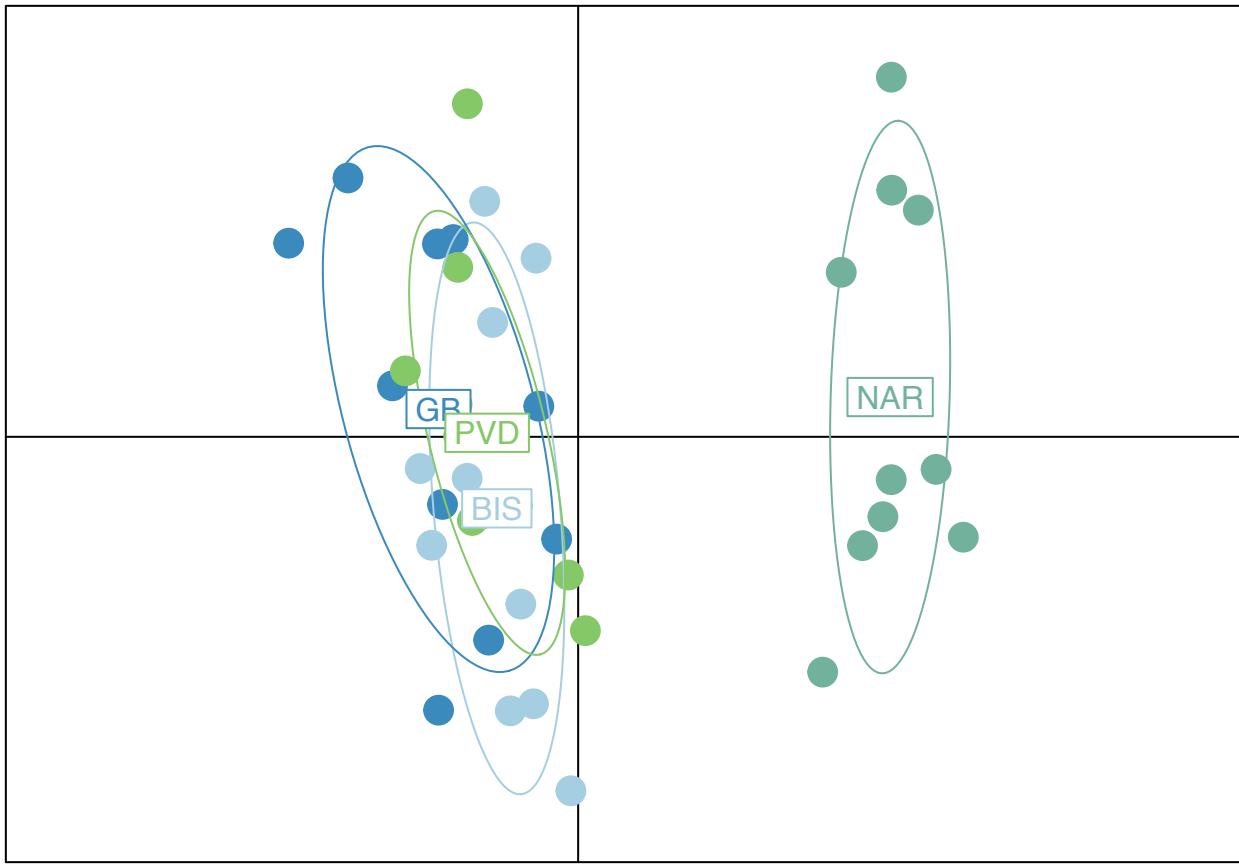
## PCA eigenvalues



```
s.class(pca1_gene$li, pop(stratated.filt_gene))
title("PCA of neutral gene SNPs\naxes 1-2")
add.scatter.eig(pca1_gene$eig[1:20], 3,1,2)
```



```
col <- funky(15)
s.class(pca1_gene$li, pop(stratated.filt_gene), xax=1, yax=2, col=col, axesell=FALSE, cstarr=0, cpoint=3, g
```



## Discriminant Analysis of Principal Components (DAPC)

The first line of the code block below is interactive. In the console, it will tell you:

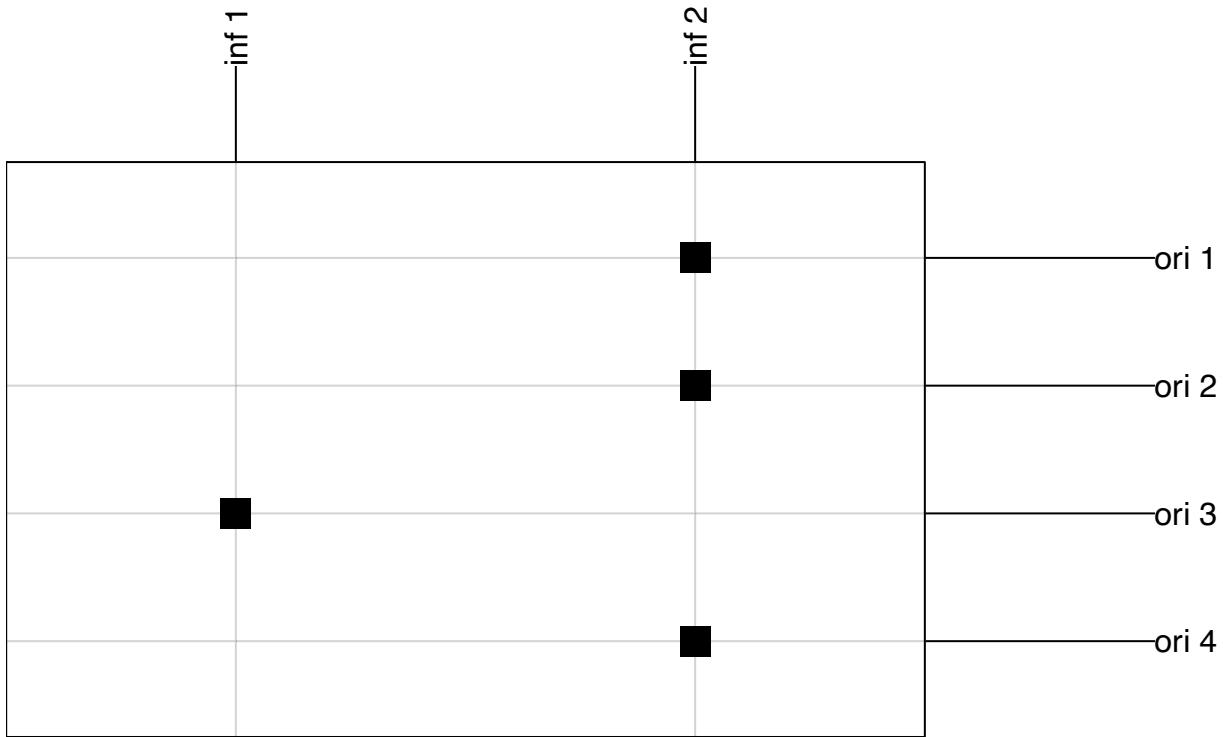
1. Choose the number PCs to retain ( $\geq 1$ ): I picked 1 2. Choose the number discriminant functions to retain ( $\geq 1$ ): I picked 2

```
grp_gene <- find.clusters(stratated.u_gene, n.pca = 1, choose.n.clust = FALSE)
table(pop(stratated.u_gene), grp_gene$grp)
```

```
##
```

```
##      1 2
## BIS  0 10
## GB   0 10
## NAR  10 0
## PVD  0 10
```

```
table.value(table(pop(stratated.u_gene), grp_gene$grp), col.lab=paste("inf", 1:2), row.lab=paste("ori",
```

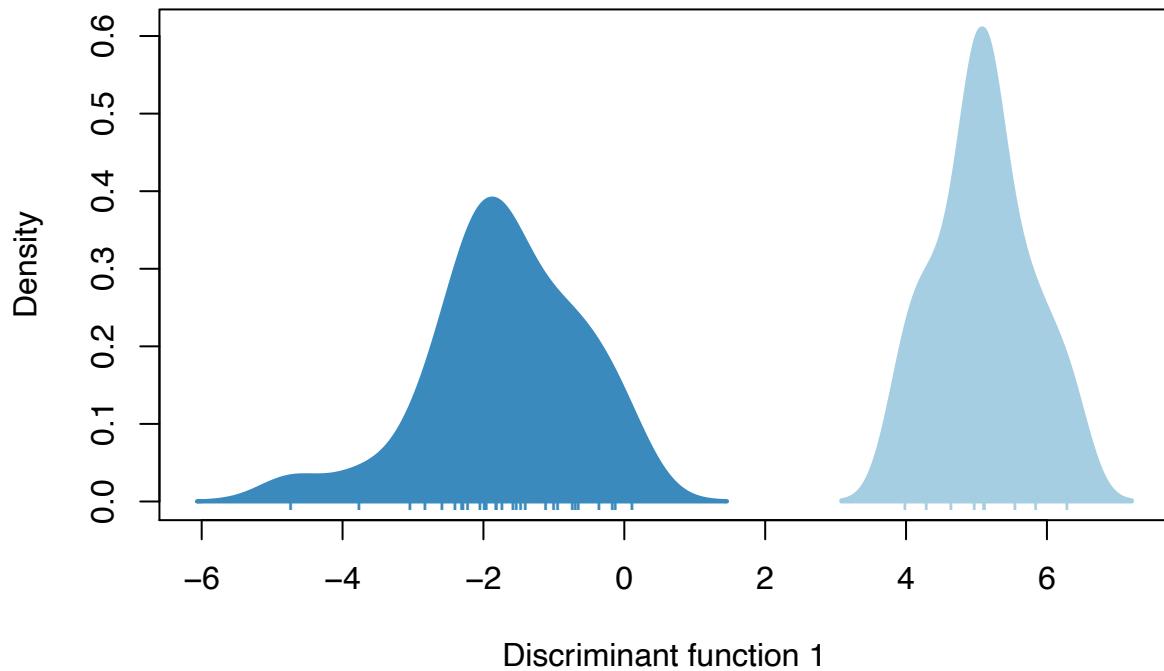


**1 • 3 ■ 5 ■ 7 ■ 9 ■**

Again, the first line of the code block below is interactive. In the console, it will tell you:

1. Choose the number PCs to retain ( $\geq 1$ ): I picked 1 2. Choose the number discriminant functions to retain ( $\geq 1$ ): I picked 1

```
dapc1_gene <- dapc(stratated.u_gene, grp_gene$grp, n.pca = 1, n.da= 1)
scatter(dapc1_gene,col=col,bg="white", solid=1)
```

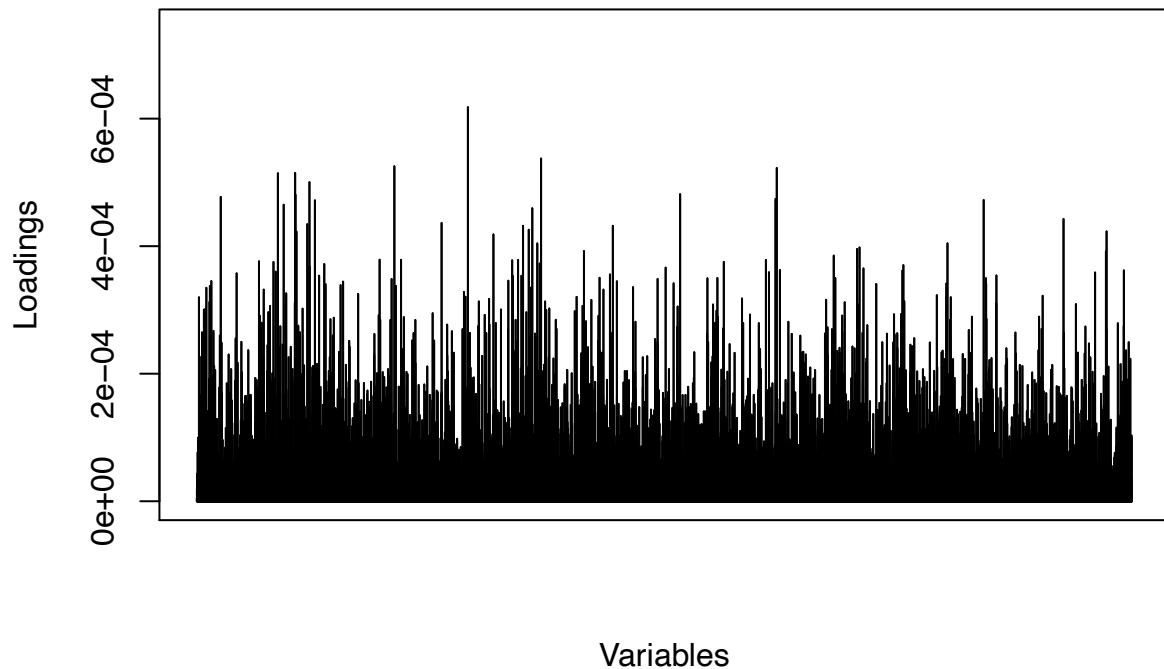


Again, the first line of the code block below is interactive. In the console, it will tell you:

1. Choose the number PCs to retain ( $\geq 1$ ): I picked 1 2. Choose the number discriminant functions to retain ( $\geq 1$ ): I picked 1

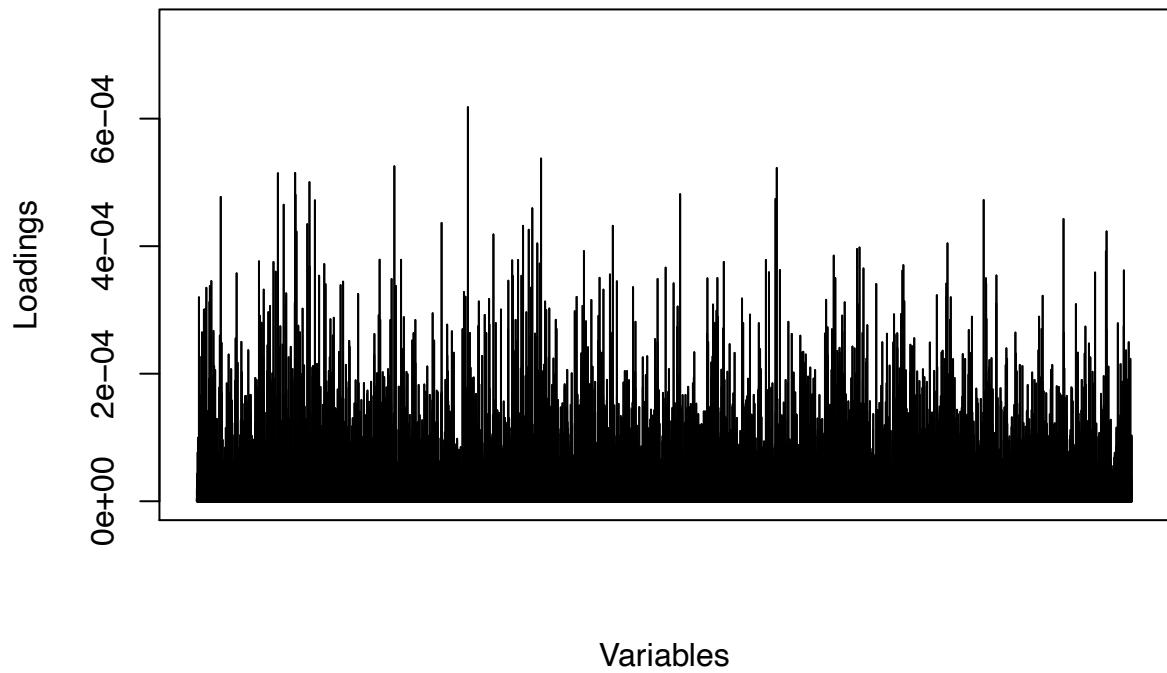
```
contrib_gene <- loadingplot(dapc1_gene$var.contr, axis=1, thres=.01, lab.jitter=1)
```

## Loading plot

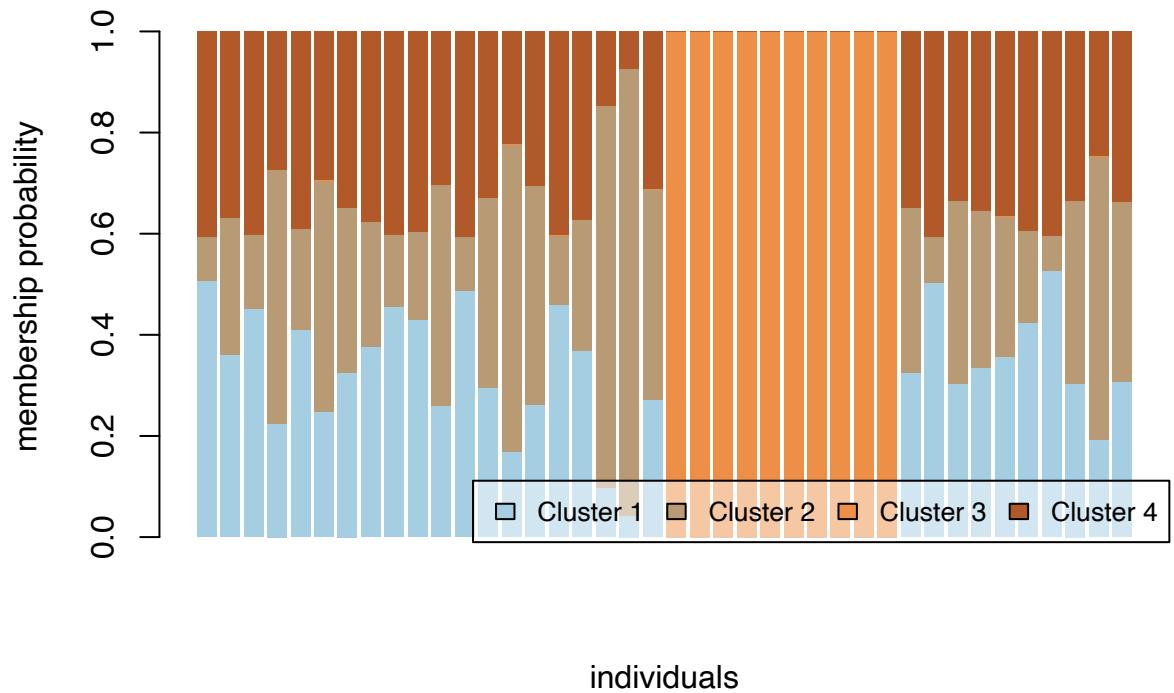


```
contrib_gene  
## NULL  
setPop(rad.filt_gene) <- ~Library  
dapc1_gene <- dapc(stratified.u_gene, pop=stratified.u_gene), n.pca = 1, n.da = 1)  
contrib_gene <- loadingplot(dapc1_gene$var.contr, axis=1, thres=.05, lab.jitter=1)
```

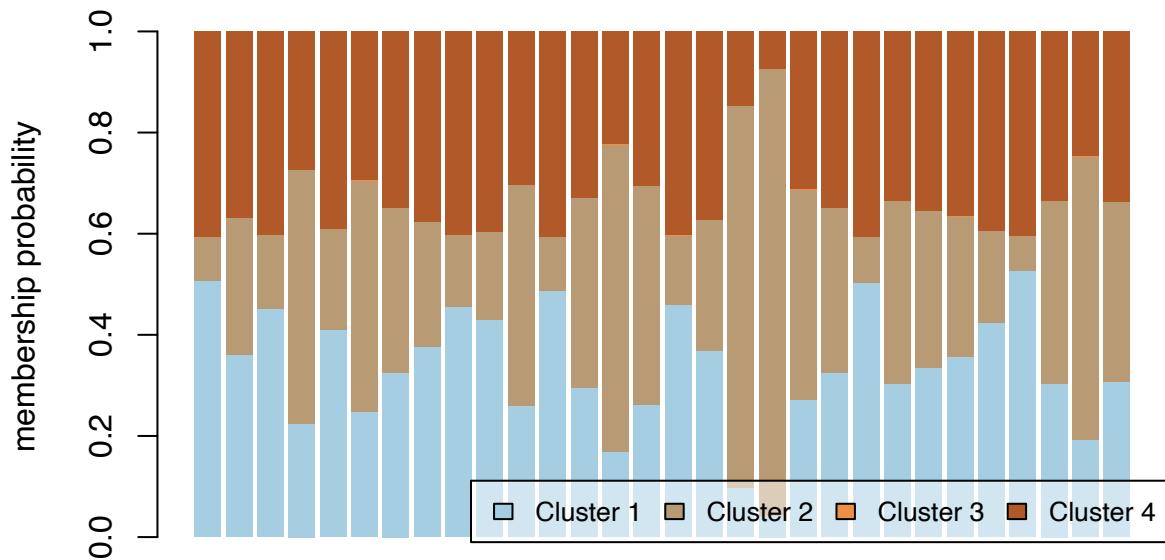
## Loading plot



```
#Structure Like
compoplot(dapc1_gene, posi="bottomright",txt.leg=paste("Cluster", 1:4), lab="", ncol=1, xlab="individual")
```



```
temp_gene <- which(apply(dapc1_gene$posterior, 1, function(e) all(e<0.9)))
compoplot(dapc1_gene, subset=temp_gene, posi="bottomright", txt.leg=paste("Cluster", 1:4), ncol=2)
```



## PCAviz

```

NA.afDraw<- function(ind){
  ind.mat <- ind@tab
  new.mat <- ind.mat
  af = colSums(ind.mat[,seq(1,ncol(ind.mat)-1,2)],na.rm = TRUE) /
    (2*apply(ind.mat[,seq(1,ncol(ind.mat)-1,2)],2,function(x) sum(!is.na(x))))
  af.Draw <- function(geno, af){
    new <- function(geno,af){
      if(is.na(geno)){
        newA = rbinom(1,2,af)
      } else {newA <- geno}
      return(newA)
    }
    new.row <- mapply(geno,af,FUN = new)
    return(new.row)
  }

  new.mat[,seq(1,ncol(ind.mat)-1,2)] <- t(apply(ind.mat[,seq(1,ncol(ind.mat)-1,2)],1,af.Draw,af))
  new.mat[,seq(2,ncol(ind.mat),2)] <- 2-new.mat[,seq(1,ncol(ind.mat)-1,2)]
  new.ind <- ind
  new.ind@tab <- new.mat
  return(new.ind)
}

```

```

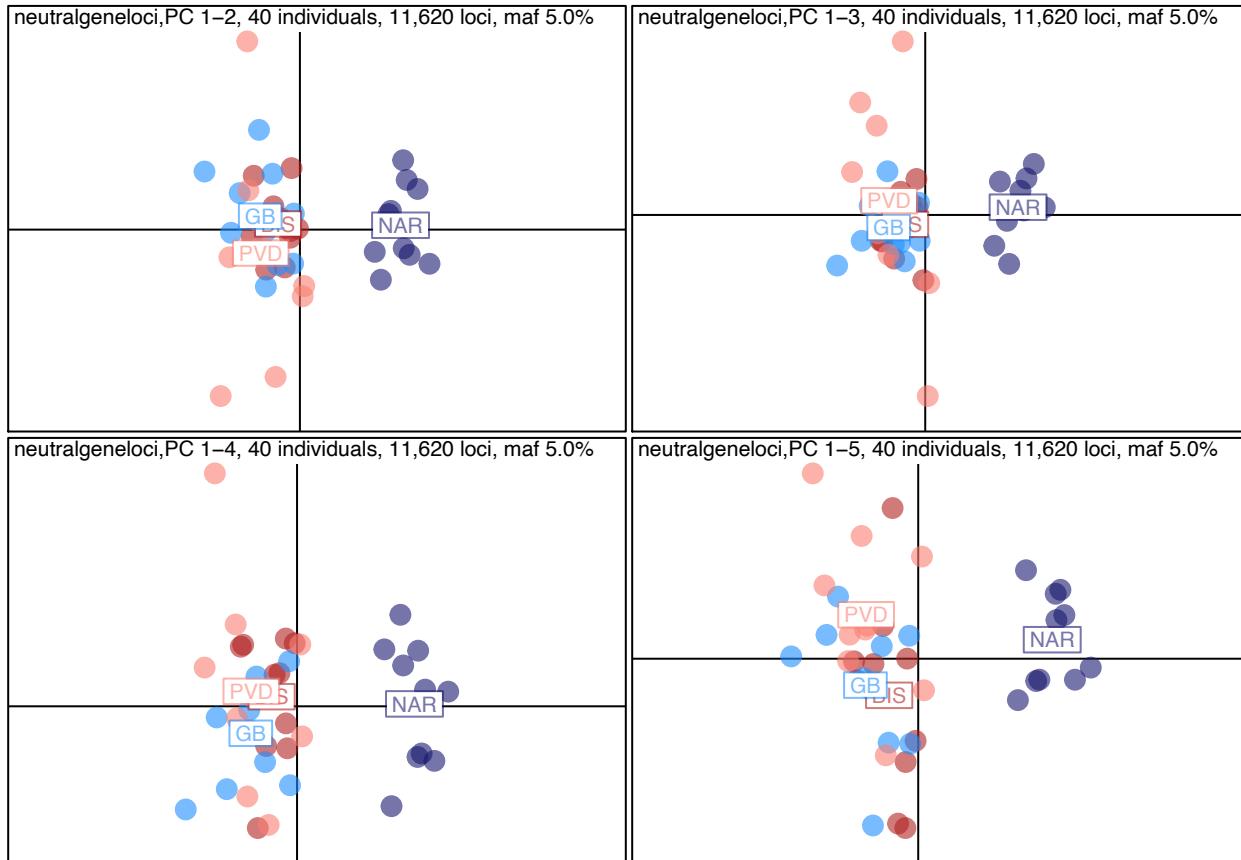
u.na_gene <- NA.afDraw(stratified.u_gene)

pca_gene <- dudi.pca(u.na_gene, center=TRUE, scale=TRUE, scannf = F, nf = 30)

col18 <- funky(length(unique(u.na_gene$strata$Population)))
#Colors that match the neutral Structure results
col6 <- c("firebrick", "dodgerblue", "midnightblue", "salmon")
par(mfrow=c(2,2))

s.class(pca_gene$li, strata(u.na_gene)$Population, xax=1, yax=2,
        sub = "neutralgeneloci, PC 1-2, 40 individuals, 11,620 loci, maf 5.0%",
        possub = "topleft", col=transp(col6,.6), axesell=FALSE,
        cstarc=0, cpoint=3, grid=FALSE, cellipse = 0)
s.class(pca_gene$li, strata(u.na_gene)$Population, xax=1, yax=3,
        sub = "neutralgeneloci, PC 1-3, 40 individuals, 11,620 loci, maf 5.0%",
        possub = "topleft", col=transp(col6,.6), axesell=FALSE,
        cstarc=0, cpoint=3, grid=FALSE, cellipse = 0)
s.class(pca_gene$li, strata(u.na_gene)$Population, xax=1, yax=4,
        sub = "neutralgeneloci, PC 1-4, 40 individuals, 11,620 loci, maf 5.0%",
        possub = "topleft", col=transp(col6,.6), axesell=FALSE,
        cstarc=0, cpoint=3, grid=FALSE, cellipse = 0)
s.class(pca_gene$li, strata(u.na_gene)$Population, xax=1, yax=5,
        sub = "neutralgeneloci, PC 1-5, 40 individuals, 11,620 loci, maf 5.0%",
        possub = "topleft", col=transp(col6,.6), axesell=FALSE,
        cstarc=0, cpoint=3, grid=FALSE, cellipse = 0)

```



# Neutral Intergenic SNPs

---

## Making files

### Make genind object

`neutralintergenic.recode.vcf` contains neutral SNPs from intergenic regions for populations PVD, GB, BIS, and NAR. Steps for generating this VCF file are located in `EecSeq_Cvirginica_NeutralLoci.md`.

Population NIN was removed from the VCF file following steps in `EecSeq_Cvirginica_OutlierDetection.md`.

`strata` contains population, environmental, and library information for each sample - can be accessed here.

```
my_vcf_inter <- read.vcfR("neutralintergenic.recode.vcf")
```

```
## Scanning file to determine attributes.
```

```
## File attributes:
```

```
##   meta lines: 63
```

```
##   header_line: 64
```

```
##   variant count: 7581
```

```
##   column count: 49
```

```
##
```

```
Meta line 63 read in.
```

```
## All meta lines processed.
```

```
## gt matrix initialized.
```

```
## Character matrix gt created.
```

```
##   Character matrix gt rows: 7581
```

```
##   Character matrix gt cols: 49
```

```
##   skip: 0
```

```
##   nrows: 7581
```

```
##   row_num: 0
```

```
##
```

```
Processed variant 1000
```

```
Processed variant 2000
```

```
Processed variant 3000
```

```
Processed variant 4000
```

```
Processed variant 5000
```

```
Processed variant 6000
```

```
Processed variant 7000
```

```
Processed variant: 7581
```

```
## All variants processed
```

```
strata <- read.table("strata", header=TRUE)
```

```
rad.filt_inter <- vcfR2genind(my_vcf_inter, strata = strata, pop = c(rep("BIS", 10),rep("GB", 10),rep("NAR", 10),rep("PVD", 10)))
```

```
rad.filt_inter
```

```
## /// GENIND OBJECT //////////
```

```
##
```

```
##   // 40 individuals; 7,581 loci; 15,140 alleles; size: 6.9 Mb
```

```
##
```

```
##   // Basic content
```

```
##   @tab: 40 x 15140 matrix of allele counts
```

```

##  @loc.n.all: number of alleles per locus (range: 1-2)
##  @loc.fac: locus factor for the 15140 columns of @tab
##  @call.names: list of allele names for each locus
##  @ploidy: ploidy of each individual (range: 2-2)
##  @type: codom
##  @call: adegenet::df2genind(X = t(x), sep = sep, pop = ..2, strata = ..1)
##
## // Optional content
##  @pop: population of each individual (group size range: 10-10)
##  @strata: a data frame with 12 columns ( Individual, Population, Latitude, Longitude, Distance, SE
#Providing population names for plotting
pop_order <- c("BIS", "GB", "NAR", "PVD")

```

Read in the other info from .strata file and extract information such as locality, latitude, and longitude.

```

info <- as.data.frame(read.table("strata", header = T, sep = "\t", stringsAsFactors = F))
mystrats_inter <- as.data.frame(matrix(nrow = length(indNames(rad.filt_inter)), ncol=10))
colnames(mystrats_inter) <- c("Population", "Latitude", "Longitude", "Distance", "Temperature", "SE", "Salinity", "just.strats" <- select(info, c("Population"))
stratified.filt_inter <- strata(rad.filt_inter, formula= Population, combine = TRUE, just.strats)
stratified.filt_inter@other <- select(info, Latitude, Longitude, Distance, SE, Temperature, Salinity, pH, Chlorophyll)
stratified.filt_inter

```

```

## /// GENIND OBJECT ///////////
##
## // 40 individuals; 7,581 loci; 15,140 alleles; size: 6.9 Mb
##
## // Basic content
##  @tab: 40 x 15140 matrix of allele counts
##  @loc.n.all: number of alleles per locus (range: 1-2)
##  @loc.fac: locus factor for the 15140 columns of @tab
##  @call.names: list of allele names for each locus
##  @ploidy: ploidy of each individual (range: 2-2)
##  @type: codom
##  @call: adegenet::df2genind(X = t(x), sep = sep, pop = ..2, strata = ..1)
##
## // Optional content
##  @pop: population of each individual (group size range: 10-10)
##  @strata: a data frame with 1 columns ( Population )
##  @other: a list containing: Latitude Longitude Distance SE Temperature Salinity pH Chlorophyll

```

Repeat for quasi-independent set of SNPs

Steps for filtering are documented for OutFLANK in EecSeq\_Cvirginica\_OutlierDetection.md. However, I am repeating the steps here.

```
my_vcf_inter <- read.vcfR("neutralintergenic.recode.vcf")
```

```

## Scanning file to determine attributes.
## File attributes:
##  meta lines: 63
##  header_line: 64
##  variant count: 7581
##  column count: 49
## 
```

```

Meta line 63 read in.
## All meta lines processed.
## gt matrix initialized.
## Character matrix gt created.
##   Character matrix gt rows: 7581
##   Character matrix gt cols: 49
##   skip: 0
##   nrow: 7581
##   row_num: 0
##
Processed variant 1000
Processed variant 2000
Processed variant 3000
Processed variant 4000
Processed variant 5000
Processed variant 6000
Processed variant 7000
Processed variant: 7581
## All variants processed

geno_inter <- extract.gt(my_vcf_inter) # Character matrix containing the genotypes
position_inter <- getPOS(my_vcf_inter) # Positions in bp
chromosome_inter <- getCHROM(my_vcf_inter) # Chromosome information

G_inter <- matrix(NA, nrow = nrow(geno_inter), ncol = ncol(geno_inter))

G_inter[geno_inter %in% c("0/0", "0|0")] <- 0
G_inter[geno_inter %in% c("0/1", "1/0", "1|0", "0|1")] <- 1
G_inter[geno_inter %in% c("1/1", "1|1")] <- 2

# NA should be replaced with "9"
G_inter[is.na(G_inter)] <- 9

```

Chromosomes need to be of class integer for this to work.

```

# Visualizing chromosomes
chrom_unique_inter <- unique(chromosome_inter)
print(chrom_unique_inter)

## [1] "NC_035780.1" "NC_035781.1" "NC_035782.1" "NC_035783.1" "NC_035784.1"
## [6] "NC_035785.1" "NC_035786.1" "NC_035787.1" "NC_035788.1" "NC_035789.1"

# Removing "NC_" from the chromosome name so it can be converted to an integer
chrom_new_inter <- chromosome_inter %>% str_replace("NC_","")

# Converting the character string into an integer
chrom1_inter <- as.integer(chrom_new_inter)

```

chromosome and position need to be sorted for imputation to work.

```

chrom_sort_inter <- sort(chrom1_inter)
pos_sort_inter <- sort(position_inter)

```

*Note: This filtering program does not allow for missing genotype values.*

## Remove missing genotypes

```
# removing missing data
G_inter_miss <- matrix(NA, nrow = nrow(geno_inter), ncol = ncol(geno_inter))

G_inter_miss[geno_inter %in% c("0/0", "0|0")] <- 0
G_inter_miss[geno_inter %in% c("0/1", "1/0", "1|0", "0|1")] <- 1
G_inter_miss[geno_inter %in% c("1/1", "1|1")] <- 2

# removing missing data
G_inter_miss[na.omit(G_inter_miss)]
```

```
##      [1] 1 1 1 0 0 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 0 1 1 1 1 1 1 1 1 1 0 0 0 0 1 1 0 0
##      [37] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 0 0 0 1 1 1 1 1 0 1 0 0
##      [73] 1 0 1 0 1 1 1 1 1 1 1 0 0 0 0 1 1 1 1 1 1 1 1 1 1 1 0 0 0 1 1 1 1 0 0 0 0 1 1 1 1 0 0 0 0 1 1 1
##     [109] 0 1 1 1 1 1 1 1 1 0 1 0 0 1 1 1 1 1 1 1 1 1 1 0 1 0 1 1 1 1 1 1 0 1 1 1 1 0 1 1 1 1 0 0 1 1 1
##     [145] 1 1 1 1 1 1 0 1 0 0 1 0 1 1 1 1 1 1 1 0 1 0 1 1 1 1 0 1 0 1 1 1 1 0 1 0 1 1 1 0 1 0 0 0 0 1
##     [181] 0 0 0 1 1 1 1 0 0 1 1 1 1 0 1 0 1 1 0 1 1 0 0 0 0 0 0 1 1 1 1 0 0 1 1 1 1 0 0 1 1
##     [217] 1 1 1 1 1 1 1 1 1 0 0 1 1 1 0 0 0 0 1 0 0 1 0 0 0 1 1 0 0 1 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0
##     [253] 0 0 0 0 0 1 1 1 0 0 1 0 1 1 1 1 1 1 0 0 1 1 1 0 1 1 1 1 0 1 0 1 1 1 1 0 1 0 1 1 1 1 1 1 1
##     [289] 1 1 1 1 1 0 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##     [325] 1 1 1 1 1 1 1 1 0 0 0 0 1 1 1 1 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 0 1 1 1 1 1 1 1 1
##     [361] 1 1 0 1 1 1 0 0 0 0 0 1 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 1 1 1 1 1 1 1 1 1
##     [397] 1 1 1 1 0 0 0 1 0 0 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 1 1 1 1 1 0 1 1 1
##     [433] 1 0 1 0 0 0 0 0 1 1 0 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0
##     [469] 0 1 1 0 1 0 1 1 1 0 0 1 1 1 0 1 1 0 0 1 1 0 1 1 0 1 1 0 1 1 0 0 1 1 1 1 1 1 1 1 1 1 1 1 0
##     [505] 1 1 1 0 1 1 1 0 1 1 0 1 0 0 0 0 0 0 0 0 0 1 1 1 1 0 1 0 1 0 1 0 1 0 1 1 1 1 1 1 1 1 1 1 1
##     [541] 1 0 0 0 0 0 1 0 0 0 0 0 0 1 1 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 0 1 1 1 1 1 1 1 1
##     [577] 1 1 1 1 1 0 1 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 1 1 1 1 1 1 1 1 1
##     [613] 1 1 1 1 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 1 1 1 1 1 1 1 0 1 1
##     [649] 1 1 1 1 1 1 1 1 1 0 0 0 1 1 0 1 1 0 0 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 1 1 1 1 1 1 1 1 1
##     [685] 0 0 1 1 1 1 1 1 1 1 1 1 1 1 0 0 1 1 0 0 1 1 1 0 1 1 1 1 1 1 1 1 1 1 0 1 0 0 1 1 0 0 1 1 0
##     [721] 0 0 0 0 0 0 0 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 1 1 1 1 1 1 0 0 0 1
##     [757] 0 1 1 1 1 1 1 1 1 1 1 1 1 0 1 1 0 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 0 1 1 1 1 0 1 1 0 1 0 1
##     [793] 0 0 1 1 1 1 1 0 0 0 1 1 0 0 0 1 0 1 1 0 0 0 0 0 1 1 0 0 1 0 0 0 0 0 1 1 0 0 1 0 0 0 0 0 1 1
##     [829] 1 0 1 0 0 0 1 1 1 0 0 0 1 0 0 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 1 0 1 1 1 1 1 1 1 1
##     [865] 1 1 1 1 1 1 1 0 0 1 1 0 0 0 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 0 1 0 0 0 0 0 0 1 1 1
##     [901] 1 1 1 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 1 1 0 1 1 1 0 1 1 1 0
##     [937] 0 1 1 1 1 1 1 1 1 1 1 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 0 1 0 1 1 1 1 1 1 1 1 0 1
##     [973] 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##     [1009] 1 1 1 1 1 1 1 1 1 1 1 1 1 0 0 0 0 0 1 0 0 1 0 0 0 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##     [1045] 1 0 1 1 0 1 1 0 0 1 1 1 1 0 1 1 0 0 1 0 1 1 1 1 1 0 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1
##     [1081] 1 0 1 1 0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 0 0 0
##     [1117] 1 1 1 0 0 1 1 1 0 1 1 0 0 0 0 1 1 1 1 0 1 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##     [1153] 1 0 0 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##     [1189] 0 1 1 1 1 0 1 0 0 1 0 1 1 1 1 1 1 1 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0
##     [1225] 0 1 1 1 0 0 1 0 0 1 0 0 0 0 0 0 0 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 0 0 0 1 1 0 0
##     [1261] 0 0 1 1 1 1 0 1 0 0 1 1 0 1 1 0 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##     [1297] 1 1 1 0 0 1 1 1 1 0 1 1 1 1 1 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 1 1 1
##     [1333] 0 0 1 1 1 1 1 1 0 0 1 0 1 1 1 1 1 1 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0
##     [1369] 0 1 1 1 0 1 1 0 1 0 0 1 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 0 0 1
##     [1405] 1 1 1 1 0 1 1 1 0 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 0
##     [1441] 0 1 0 0 1 1 1 0 1 1 1 1 1 0 1 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 0
##     [1477] 0 0 1 1 1 1 1 0 0 1 0 1 1 1 1 1 1 1 0 0 0 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0
```

```

## [52057] 0 0 0 0 0 0 1 1 1 1 1 1 0 0 1 0 1 1 1 1 0 0 0 0 0 1 1 1 1 1 1 1
## [52093] 1 1 0 1 1 0 0 0 0 1 1 1 0 1 1 1 0 0 1 1 1 1 1 0 0 1 1 0 1 0 1 0 0 0 0 0 0
## [52129] 1 1 1 1 1 1 1 1 1 0 0 0 1 1 1 1 1 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 0 0 0
## [52165] 1 1 1 1 1 1 0 1 1 1 1 1 1 1 1 1 1 0 0 0 0 0 0 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1
## [52201] 1 1 1 1 1 1 1 0 1 1 1 1 1 1 1 1 1 0 0 1 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
## [52237] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
## [52273] 1 1 1 1 1 1 0 1 1 1 1 0 0 0 0 1 1 1 1 0 0 0 1 1 1 1 0 1 1 1 1 1 1 1 1 1 1 1
## [52309] 1 1 1 1 1 1 1 1 1 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 0
## [52345] 1 1 1 0 1 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 1 1 1 0 0 0 1 0 1 0 0
## [52381] 0 0 0 0 0 0 0 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 0 0 0 1 0 0 1 1 1 1
## [52417] 0 1 1 1 0 1 1 0 1 1 1 1 1 0 1 0 1 1 1 1 1 1 1 1 1 0 0 1 1 1 0 1 1 1 1 1 1
## [52453] 1 1 1 0 0 0 0 0 0 0 1 1 1 1 1 1 1 0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 1 1 0
## [52489] 0 1 1 1 1 0 0 1 0 1 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 0 1 0 0 0 0 1 0 0 0
## [52525] 0 0 0 1 0 0 0 0 0 0 1 1 1 1 1 1 1 0 0 1 1 1 1 1 0 1 1 1 1 1 1 0 1 1 1 0 1
## [52561] 1 1 1 0 0 0 1 0 1 1 1 1 0 0 0 1 1 1 1 0 0 0 0 1 0 1 0 0 1 0 1 1 1 0 1 1
## [52597] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 0 1 1 1 1 1 0 0 0 0
## [52633] 1 0 0 1 1 0 0 0 0 0 0 1 1 1 1 1 0 1 1 1 1 1 1 1 1 0 1 0 1 1 1 1 1 1 1 1 0
## [52669] 0 1 1 0 0 0 0 0 0 0 1 0 0 1 1 1 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
## [52705] 1 1 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 1 1 0
## [52741] 1 1 1 0 1 1 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 1 1 0
## [52777] 0 1 1 1 1 1 1 1 0 1 1 1 0 1 1 1 1 1 1 1 1 1 0 0 1 1 1 1 1 1 0 1 1 1 1 1 1
## [52813] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 1 0 1 0
## [52849] 1 0 0 0 0 1 1 1 1 1 1 1 1 0 1 1 0 1 1 1 0 0 1 1 1 1 1 1 1 1 1 0 0 1 1 1 1 1
## [52885] 1 1 1 1 1 1 1 0 1 1 1 1 1 0 0 0 0 1 1 1 1 1 0 1 1 1 1 1 0 0 0 1 0 1 1 1 0 1
## [52921] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 1 1 1 0 1 1 1 1 1 1 1 1
## [52957] 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1
## [52993] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1
## [53029] 1 1 1 1 1 1 1 1 1 1 1 1 0 0 0 0 1 1 0 1 1 1 1 1 1 0 0 0 0 1 0 1 1 1 1 1 1 1
## [53065] 0 0 1 0 1 0 1 1 0 1 1 0 0 0 0 1 1 0 1 1 0 0 1 1 1 1 1 1 1 1 0 0 1 1 1 0
## [53101] 1 1 0 0 1 1 1 1 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 0 0 0 0 1 1 1 0 1 1 1 0 0 1
## [53137] 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1
## [53173] 1 1 1 1 1 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 0 0 1 0 1 1 1 1 1 0 1 1 1 1
## [53209] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 1 0 0 0 0
## [53245] 0 1 1 1 1 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 1 1 1 1 0
## [53281] 1

# Converting the genotype matrix to class FBM.code256 using the matrix where missing data was removed
G1_inter_miss <- add_code256(big_copy(t(G_inter_miss), type="raw"), code=bigsnpr:::CODE_012)

## Warning in replaceMat(x$address_rw, i, j, value): At least one value changed (nan -> 0)
##   while converting from R type 'double' to C type 'unsigned char (raw)'.

newpc_inter_miss <-.snp_autoSVD(G1_inter_miss, infos.chr = chrom_sort_inter, infos.pos = pos_sort_inter, r

## 
## Phase of clumping (on MAF) at r^2 > 0.2.. keep 3422 SNPs.
## Discarding 1403 variants with MAC < 10.
## 
## Iteration 1:
## Computing SVD..
## 0 outlier variant detected..
## 
## Converged!

# maximum roll.size I could use is 35
which_pruned_inter_miss <- attr(newpc_inter_miss, which="subset") # Indexes of remaining SNPs after pru

```

```

length(which_pruned_inter_miss)

## [1] 2019

invisible(lapply(which_pruned_inter_miss, write, "pruned_data_inter.txt", append=TRUE))

```

### In terminal

```

$ mawk '!/#/' neutralintergenic.recode.vcf | cut -f1,2 > totallociinter
$ NUM=(`cat totallociinter | wc -l`)
$ paste <(seq 1 $NUM) totallociinter > lociinter.plus.index
$ cat pruned_data_inter.txt | parallel "grep -w ^{} lociinter.plus.index" | cut -f2,3> pruned_data_inter

```

```
$ head pruned_data_inter.loci.txt
```

```

output:
NC_035780.1 623558
NC_035780.1 623903
NC_035780.1 623906
NC_035780.1 1490085
NC_035780.1 1490100
NC_035780.1 1695081
NC_035780.1 1695087
NC_035780.1 1756507
NC_035780.1 1882142
NC_035780.1 2425975

```

Create VCF file with just the pruned\_data loci

```
$ vcftools --vcf neutralintergenic.recode.vcf --recode --recode-INFO-all --positions pruned_data_inter.loci.txt
```

```

output:
After filtering, kept 40 out of 40 Individuals
Outputting VCF file...
After filtering, kept 2019 out of a possible 7581 Sites
Run Time = 1.00 seconds

```

```
my_vcf_u_inter <- read.vcfR("pruned_data_inter.recode.vcf")
```

```

## Scanning file to determine attributes.
## File attributes:
##   meta lines: 63
##   header_line: 64
##   variant count: 2019
##   column count: 49
##
Meta line 63 read in.
## All meta lines processed.
## gt matrix initialized.
## Character matrix gt created.
##   Character matrix gt rows: 2019
##   Character matrix gt cols: 49
##   skip: 0
##   nrows: 2019
##   row_num: 0
##
Processed variant 1000

```

```

Processed variant 2000
Processed variant: 2019
## All variants processed
strata <- read.table("strata", header=TRUE)

rad.u_inter <- vcfR2genind(my_vcf_u_inter, strata = strata, pop = c(rep("BIS", 10),rep("GB", 10),rep("N", 10),rep("P", 10),rep("S", 10),rep("T", 10),rep("U", 10),rep("V", 10),rep("W", 10),rep("X", 10),rep("Y", 10)),geno=2,checkgeno=TRUE)

rad.u_inter

## /// GENIND OBJECT ///////////
##
## // 40 individuals; 2,019 loci; 4,038 alleles; size: 1.9 Mb
##
## // Basic content
## @tab: 40 x 4038 matrix of allele counts
## @loc.n.all: number of alleles per locus (range: 2-2)
## @loc.fac: locus factor for the 4038 columns of @tab
## @all.names: list of allele names for each locus
## @ploidy: ploidy of each individual (range: 2-2)
## @type: codom
## @call: adegenet::df2genind(X = t(x), sep = sep, pop = ..2, strata = ..1)
##
## // Optional content
## @pop: population of each individual (group size range: 10-10)
## @strata: a data frame with 12 columns ( Individual, Population, Latitude, Longitude, Distance, SE, Chlorophyll, Salinity, Temperature, pH, Chlorophyll, Salinity )
stratated.u_inter <- strata(rad.u_inter, formula= Population, combine = TRUE,just.strats)
stratated.u_inter@other <- select(info, Latitude,Longitude, Distance, SE, Temperature,Salinity,pH,Chlorophyll,Salinity)

stratated.u_inter

## /// GENIND OBJECT ///////////
##
## // 40 individuals; 2,019 loci; 4,038 alleles; size: 1.9 Mb
##
## // Basic content
## @tab: 40 x 4038 matrix of allele counts
## @loc.n.all: number of alleles per locus (range: 2-2)
## @loc.fac: locus factor for the 4038 columns of @tab
## @all.names: list of allele names for each locus
## @ploidy: ploidy of each individual (range: 2-2)
## @type: codom
## @call: adegenet::df2genind(X = t(x), sep = sep, pop = ..2, strata = ..1)
##
## // Optional content
## @pop: population of each individual (group size range: 10-10)
## @strata: a data frame with 1 columns ( Population )
## @other: a list containing: Latitude Longitude Distance SE Temperature Salinity pH Chlorophyll Salinity
Make hierfstat object

hf.filt_inter <- genind2hierfstat(rad.filt_inter, pop = c(rep("BIS", 10),rep("GB", 10),rep("NAR", 10), rep("P", 10),rep("S", 10),rep("T", 10),rep("U", 10),rep("V", 10),rep("W", 10),rep("X", 10),rep("Y", 10)),geno=2,checkgeno=TRUE)

hf.u_inter <- genind2hierfstat(rad.u_inter, pop = c(rep("BIS", 10),rep("GB", 10),rep("NAR", 10), rep("P", 10),rep("S", 10),rep("T", 10),rep("U", 10),rep("V", 10),rep("W", 10),rep("X", 10),rep("Y", 10)),geno=2,checkgeno=TRUE)

hf.u_inter <- hf.u_inter$hierfstat.no.imputation

```

## Estimating effective migration surfaces (EEMS)

### 1. datapath.diffs

Steps for generating `datapath.diffs` are completed in RStudio.

```
# V1 method to get .diffs matrix, preferred
bed2diffs_v1 <- function(Geno) {
  nIndiv <- nrow(Geno)
  nSites <- ncol(Geno)
  Diffs <- matrix(0, nIndiv, nIndiv)

  for (i in seq(nIndiv - 1)) {
    for (j in seq(i + 1, nIndiv)) {
      x <- Geno[i, ]
      y <- Geno[j, ]
      Diffs[i, j] <- mean((x - y)^2, na.rm = TRUE)
      Diffs[j, i] <- Diffs[i, j]
    }
  }
  Diffs
}

# V2 method to get .diffs matrix, only if V1 doesn't work
bed2diffs_v2 <- function(Geno) {
  nIndiv <- nrow(Geno)
  nSites <- ncol(Geno)
  Miss <- is.na(Geno)
  ## Impute NAs with the column means (= twice the allele frequencies)
  Mean <- matrix(colMeans(Geno, na.rm = TRUE), ## a row of means
                  nrow = nIndiv, ncol = nSites, byrow = TRUE) ## a matrix with nIndiv identical rows of means
  Mean[Miss == 0] <- 0 ## Set the means that correspond to observed genotypes to 0
  Geno[Miss == 1] <- 0 ## Set the missing genotypes to 0 (used to be NA)
  Geno <- Geno + Mean
  ## Compute similarities
  Sim <- Geno %*% t(Geno) / nSites
  SelfSim <- diag(Sim) ## self-similarities
  vector1s <- rep(1, nIndiv) ## vector of 1s
  ## This chunk generates a `diffs` matrix
  Diffs <- SelfSim %*% t(vector1s) + vector1s %*% t(SelfSim) - 2 * Sim
  Diffs
}

geno_inter <- stratted.filt_inter@tab

# Get rid of non-biallelic loci
multi.loci_inter <- names(which(stratted.filt_inter@loc.n.all != 2))
multi.cols_inter <- which(grep0("^\.", multi.loci_inter, "\\\.\d+$", collapse = "|"), colnames(geno_inter))
if (length(multi.cols_inter)) geno_inter <- geno_inter[, - multi.cols_inter]
nloci_inter <- dim(geno_inter)[2] / 2
dim(geno_inter)

## [1] 40 15118
stopifnot(identical(stratted.filt_inter@type, 'codom'))
```

```
# bed2diffs functions
diffs.v1_inter <- bed2diffs_v1(geno_inter)
diffs.v2_inter <- bed2diffs_v2(geno_inter)
# Round to 6 digits
diffs.v1_inter <- round(diffs.v1_inter, digits = 6)
diffs.v2_inter <- round(diffs.v2_inter, digits = 6)
```

Check that the dissimilarity matrix has one positive eigenvalue and nIndiv-1 negative eigenvalues, as required by a full-rank Euclidean distance matrix. If the V1 method does not make a Euclidean matrix, you must use V2.

```
tail(sort(round(eigen(diffs.v1_inter)$values, digits = 2)))

## [1] -0.40 -0.40 -0.39 -0.38 -0.37 20.20

tail(sort(round(eigen(diffs.v2_inter)$values, digits = 2)))

## [1] -0.40 -0.39 -0.39 -0.38 -0.37 19.83

# Set suffix for EEMS input files
suf_inter <- "neutralinterdata-filt"

# This saves the file to directory
write.table(diffs.v1_inter, paste(suf_inter,".v1.diffs",sep=""),
            col.names = FALSE, row.names = FALSE, quote = FALSE)
```

## 2. datapath.coord

datapath.coord are the sample coordinates, two coordinates per sample, one sample per line. The sampling locations should be given in the same order as the rows and columns of the dissimilarity matrix.

Steps for generating datapath.coord are completed in RStudio.

```
## Get gps coordinates from previously created info matrix
x0R.info <- dplyr::filter(info)
gps_matrix <- select(x0R.info,c("Longitude","Latitude"))

#write .coord file
write.table(gps_matrix, paste(suf_inter,".v1.coord",sep=""),col.names = FALSE, row.names = FALSE,quote = FALSE)
```

## 3. datapath.outer

datapath.outer are the habitat coordinates, as a sequence of vertices that form a closed polygon. The habitat vertices should be listed counterclockwise and the first vertex should also be the last vertex, so that the outline is a closed ring. Otherwise, EEMS attempts to “correct” the polygon and prints a warning message.

datapath.outer is created manually in Excel,based on site coordinates gathered from Google Maps, copied into terminal using nano, and saved as the file with the appropriate extension.

\*\*runeems\_snps is then run in command-line following the steps documented in NB\_EEMS\_OutlierHap.md

Back in RStudio to plot runeems\_snps outputs

```
# Install rEEMSpots
library(rEEMSpots)

# Plotting EEMS after running runeems_snps
path = "./NB_EEMS_Neutral/"
dirs = c(paste0(path,"neutralinterdata-D200-chain1"), paste0(path,"neutralinterdata-D300-chain1"), paste0(path,"neutralinterdata-D400-chain1"))
```

```

eems.plots(mcmcpath = c(paste0(path, "./neutralinterdata-D200-chain1"), paste0(path,"neutralinterdata-D300-chain1"),
                        longlat = T,add.grid=F,add.outline = T,add.demes = T,
                        projection.in = "+proj=longlat +datum=WGS84",projection.out = "+proj=merc +datum=WGS84",
                        add.map = T,add.abline = T, add.r.squared = T)

## Input projection: +proj=longlat +datum=WGS84
## Output projection: +proj=merc +datum=WGS84

## Loading rgdal (required by projection.in)
## Loading rworldmap (required by add.map)
## Loading rworldxtra (required by add.map)

## Using the default DarkOrange to Blue color scheme, with 'white' as the midpoint color.
## It combines two color schemes from the 'dichromat' package, which itself is based on
## a collection of color schemes for scientific data graphics:
## Light A and Bartlein PJ (2004). The End of the Rainbow? Color Schemes for Improved Data
## Graphics. EOS Transactions of the American Geophysical Union, 85(40), 385.
## See also http://geog.uoregon.edu/datagraphics/color\_scales.htm

## Using 'euclidean' distance to assign interpolation points to Voronoi tiles.

## Processing the following EEMS output directories :
## ./NB_EEMS_Neutral./neutralinterdata-D200-chain1./NB_EEMS_Neutral/neutralinterdata-D300-chain1./NB_EEMS_Neutral/neutralinterdata-D600-chain1

## Plotting effective migration surface (posterior mean of m rates)
## ./NB_EEMS_Neutral./neutralinterdata-D200-chain1
## ./NB_EEMS_Neutral/neutralinterdata-D300-chain1
## ./NB_EEMS_Neutral/neutralinterdata-D600-chain1

## Using the default DarkOrange to Blue color scheme, with 'white' as the midpoint color.
## It combines two color schemes from the 'dichromat' package, which itself is based on
## a collection of color schemes for scientific data graphics:
## Light A and Bartlein PJ (2004). The End of the Rainbow? Color Schemes for Improved Data
## Graphics. EOS Transactions of the American Geophysical Union, 85(40), 385.
## See also http://geog.uoregon.edu/datagraphics/color\_scales.htm

## Plotting effective diversity surface (posterior mean of q rates)
## ./NB_EEMS_Neutral./neutralinterdata-D200-chain1
## ./NB_EEMS_Neutral/neutralinterdata-D300-chain1
## ./NB_EEMS_Neutral/neutralinterdata-D600-chain1

## Using the default DarkOrange to Blue color scheme, with 'white' as the midpoint color.
## It combines two color schemes from the 'dichromat' package, which itself is based on
## a collection of color schemes for scientific data graphics:
## Light A and Bartlein PJ (2004). The End of the Rainbow? Color Schemes for Improved Data
## Graphics. EOS Transactions of the American Geophysical Union, 85(40), 385.
## See also http://geog.uoregon.edu/datagraphics/color\_scales.htm

## Plotting posterior probability trace
## ./NB_EEMS_Neutral./neutralinterdata-D200-chain1
## ./NB_EEMS_Neutral/neutralinterdata-D300-chain1
## ./NB_EEMS_Neutral/neutralinterdata-D600-chain1

```

```

## Plotting average dissimilarities within and between demes
## ./NB_EEMS_Neutral./neutralinterdata-D200-chain1
## ./NB_EEMS_Neutral/neutralinterdata-D300-chain1
## Using 'euclidean' distance to assign interpolation points to Voronoi tiles.
##
##
## Using 'euclidean' distance to assign interpolation points to Voronoi tiles.
## EEMS results for at least two different population grids

```

## Pairwise Fst

```

fst.mat_inter <- pairwise.WCfst(hf.filt_inter)

gindF.fst.mat.triN_inter <- as.matrix(fst.mat_inter)
colnames(gindF.fst.mat.triN_inter) <- pop_order
rownames(gindF.fst.mat.triN_inter) <- pop_order

meltedN_inter <- melt(gindF.fst.mat.triN_inter, na.rm =TRUE)
round(gindF.fst.mat.triN_inter,4)

##          BIS      GB      NAR      PVD
## BIS      NA -0.0001  0.0064 -0.0016
## GB   -0.0001      NA  0.0079 -0.0004
## NAR   0.0064  0.0079      NA  0.0066
## PVD  -0.0016 -0.0004  0.0066      NA

summary(meltedN_inter$value)

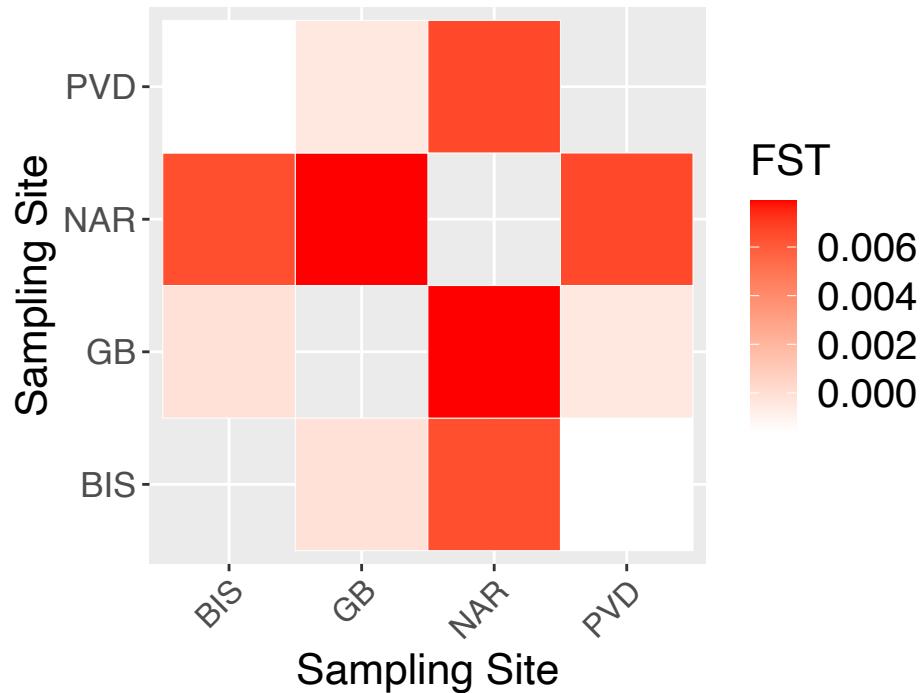
##      Min.    1st Qu.    Median    Mean    3rd Qu.    Max.
## -0.0015667 -0.0004375  0.0031588  0.0031339  0.0066106  0.0078794

#Plotting Pairwise fst
neutral_inter <- ggplot(data = meltedN_inter, aes(Var2, Var1, fill = value)) + geom_tile(color = "white")
  scale_fill_gradient(low = "white", high = "red", name="FST") +
  ggtitle(expression(atop("Pairwise FST, WC (1984) Neutral Intergenic SNPs", atop(italic("N = 40, L = 700000000"))))) +
  labs( x = "Sampling Site", y = "Sampling Site") +
  theme(axis.text.x = element_text(angle = 45, vjust = 1, size = 12, hjust = 1),axis.text.y = element_text(size = 12),
        theme(axis.title = element_text(size = 16),legend.text = element_text(size = 15), legend.title = element_text(size = 15)),
        theme(plot.title = element_text(size = 17))) +
  coord_fixed()
neutral_inter

```

## Pairwise FST, WC (1984) Neutral Intergenic SN

$N = 40, L = 7,581$



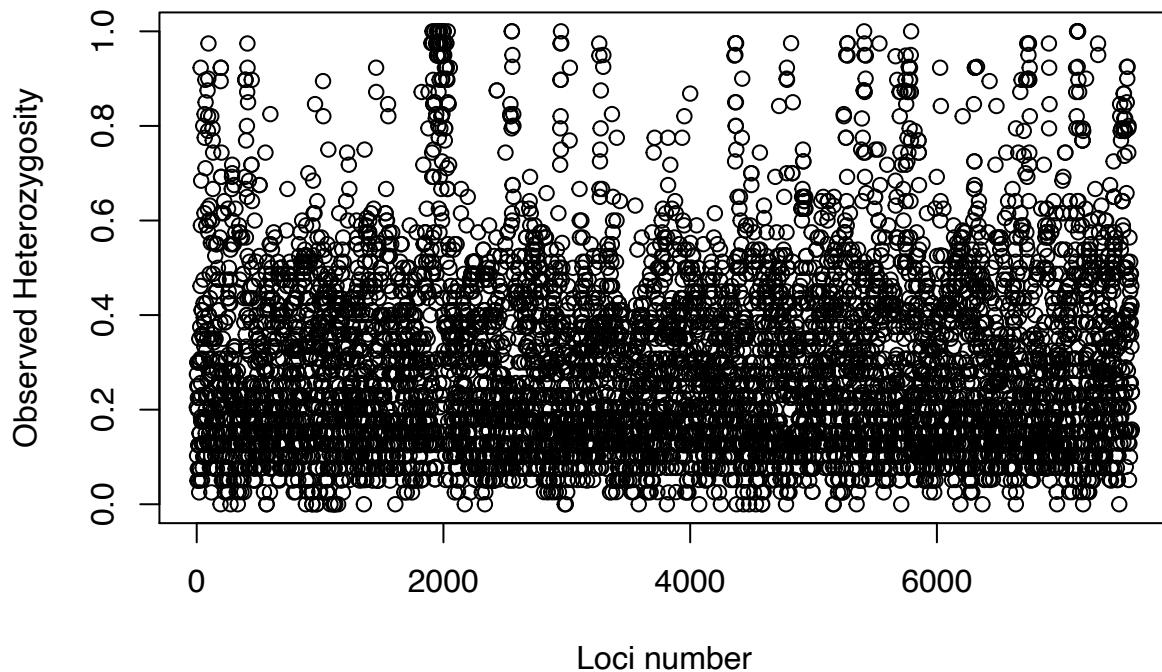
### Genetic diversity (observed and expected heterozygosity)

```
comb_inter <- summary(stratated.filt_inter)
names(comb_inter)

## [1] "n"          "n.by.pop"   "loc.n.all"  "pop.n.all"  "NA.perc"    "Hobs"
## [7] "Hexp"

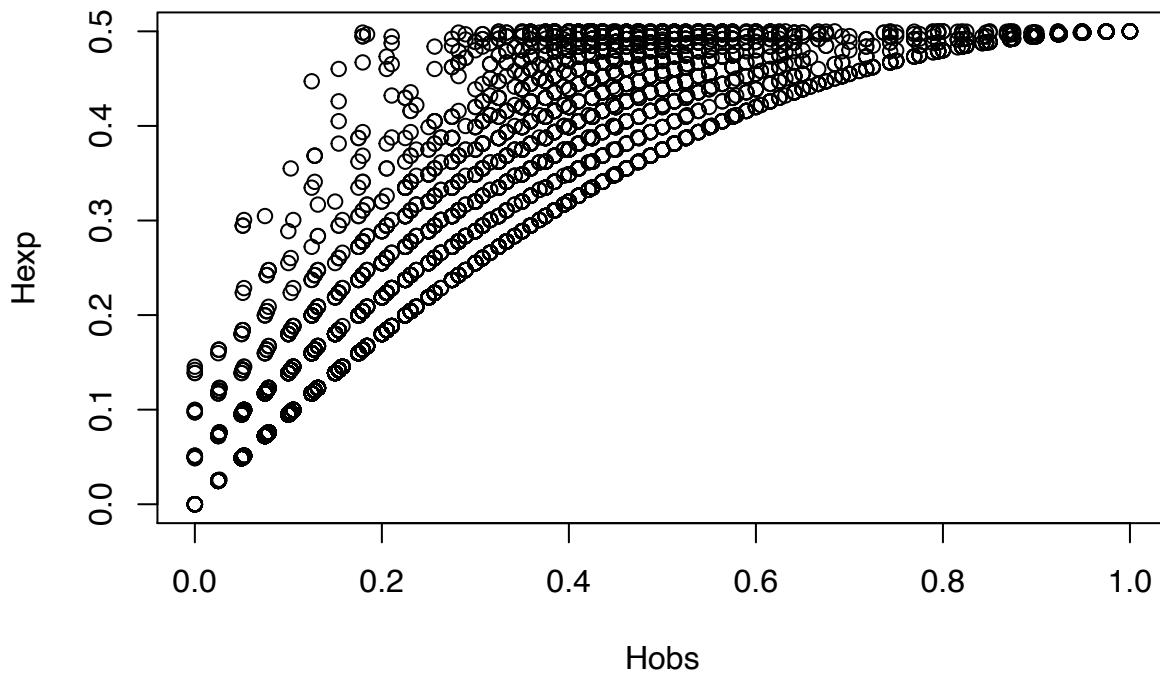
plot(comb_inter$Hobs, xlab="Loci number", ylab="Observed Heterozygosity",
     main="Observed heterozygosity per locus")
```

## Observed heterozygosity per locus



```
plot(comb_inter$Hobs,comb_inter$Hexp, xlab="Hobs", ylab="Hexp",
 main="Expected heterozygosity as a function of observed heterozygosity per locus")
```

## Expected heterozygosity as a function of observed heterozygosity per I



```
bartlett.test(list(comb_inter$Hexp, comb_inter$Hobs)) # a test : H0: Hexp = Hobs

##
##  Bartlett test of homogeneity of variances
##
## data: list(comb_inter$Hexp, comb_inter$Hobs)
## Bartlett's K-squared = 767.41, df = 1, p-value < 2.2e-16

Significant difference between Observed and expected heterozygosity.

basicstat_inter <- basic.stats(hf.filt_inter, diploid = TRUE, digits = 3)

as.data.frame(basicstat_inter$overall)

##      basicstat_inter$overall
##  Ho          0.294
##  Hs          0.275
##  Ht          0.276
##  Dst         0.001
##  Htp         0.276
##  Dstp        0.001
##  Fst         0.002
##  Fstp        0.003
##  Fis         -0.068
##  Dest        0.001

# get bootstrap confidence values for Fis
boot_inter <- boot.ppfis(hf.filt_inter, nboot = 1000)
```

```

boot5_inter <- boot.ppfis(hf.filt_inter,nboot = 1000,quant = 0.5)

# combine all pop statistics
colnames(basicstat_inter$Ho) <- pop_order
Ho_inter <- colMeans(basicstat_inter$Ho,na.rm = T)
He_inter <- colMeans(basicstat_inter$Hs,na.rm = T)
Fis_inter <- boot5_inter$fis.ci$ll
y_inter <- cbind(pop_order, Ho_inter, He_inter, Fis_inter, boot_inter$fis.ci, latitude, longitude, distance,
y_inter

##      pop_order  Ho_inter  He_inter  Fis_inter      ll      hl latitude longitude
## BIS      BIS 0.2947197 0.2734420 -0.0779 -0.0886 -0.0667  41.545 -71.431
## GB       GB 0.2967779 0.2780929 -0.0671 -0.0771 -0.0566  41.654 -71.445
## NAR      NAR 0.2927419 0.2725640 -0.0742 -0.0848 -0.0623  41.505 -71.453
## PVD      PVD 0.2908368 0.2764233 -0.0520 -0.0622 -0.0423  41.816 -71.391
##      distance sewage temperature salinity   pH Chlor_a    DO
## BIS      4.76   8.82        23     30 7.9    4.9 8.2
## GB       0.47  14.60        24     28 7.4   18.8 5.7
## NAR     15.41   2.03        25     18 7.6    4.6 7.0
## PVD      1.49  59.86        23     25 7.4    8.1 4.9

summary(He_inter)

##      Min. 1st Qu. Median Mean 3rd Qu. Max.
## 0.2726 0.2732 0.2749 0.2751 0.2768 0.2781

summary(Fis_inter)

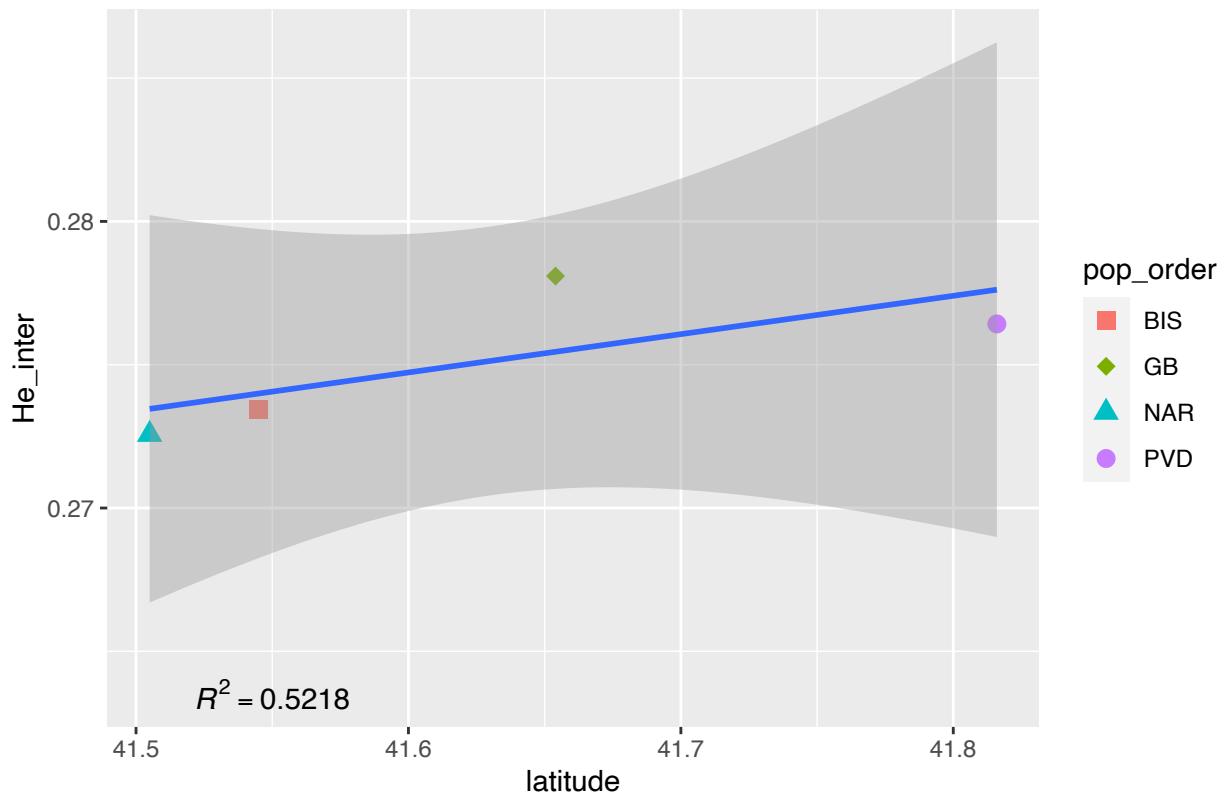
##      Min. 1st Qu. Median Mean 3rd Qu. Max.
## -0.07790 -0.07512 -0.07065 -0.06780 -0.06333 -0.05200

# Plot He vs Latitude
R2_inter = round(summary(lm(y_inter$He_inter ~ y_inter$latitude))$r.squared, 4)
ggplot(y_inter, aes(x = latitude, y = He_inter)) + geom_point(aes(shape=pop_order, color=pop_order), size=10) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggtitle("Expected heterozygosity vs Latitude, Neutral Intergenic SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R2_inter), x=41.55, y=0.2635, parse=T) +
  scale_x_continuous()

## `geom_smooth()` using formula 'y ~ x'

```

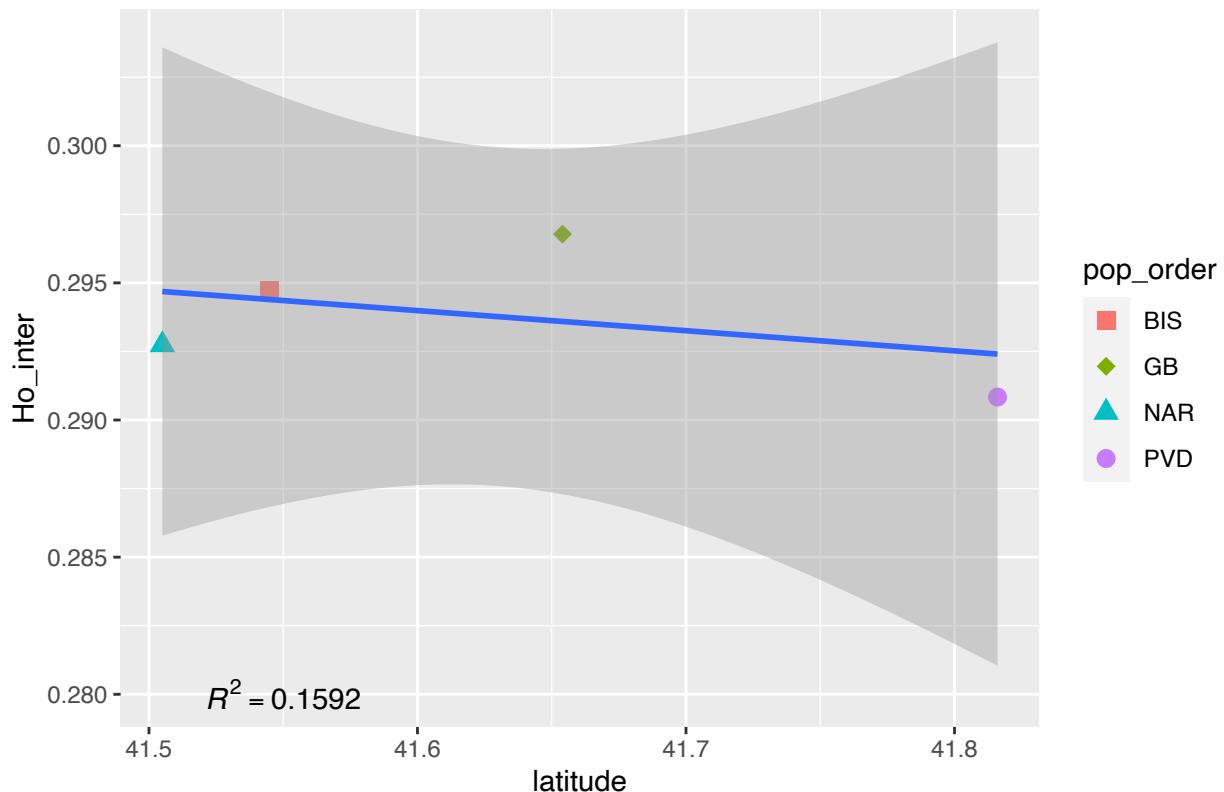
## Expected heterozygosity vs Latitude, Neutral Intergenic SNPs



```
# Plot Ho vs Latitude
R2_inter = round(summary(lm(y_inter$Ho_inter ~ y_inter$latitude))$r.squared, 4)
ggplot(y_inter, aes(x = latitude, y = Ho_inter)) + geom_point(aes(shape=pop_order, color=pop_order), size=10) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggttitle("Observed heterozygosity vs Latitude, Neutral Intergenic SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R2_inter), x=41.55, y=0.28, parse=T) +
  scale_x_continuous()
```

## `geom\_smooth()` using formula 'y ~ x'

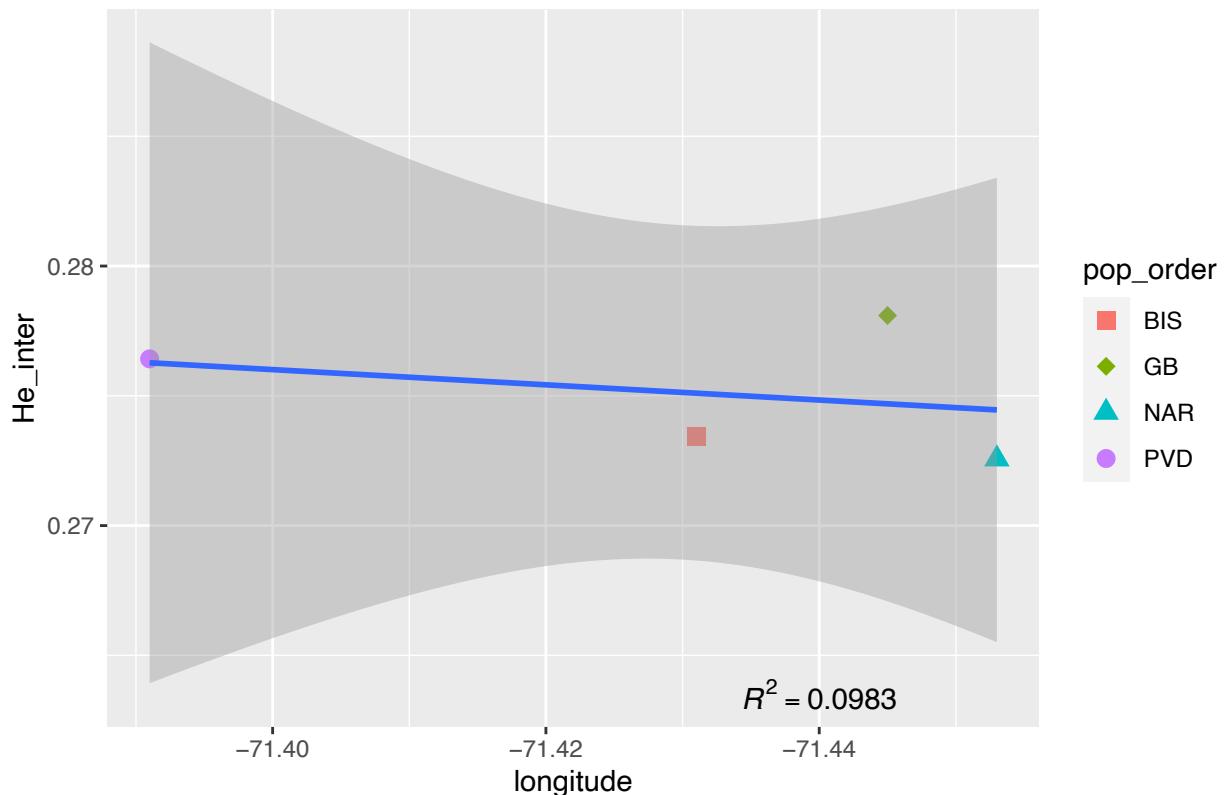
### Observed heterozygosity vs Latitude, Neutral Intergenic SNPs



```
#Plot He vs Longitude
R3_inter = round(summary(lm(y_inter$He_inter ~ y_inter$longitude))$r.squared, 4)
ggplot(y_inter, aes(x = longitude, y = He_inter)) + geom_point(aes(shape=pop_order, color=pop_order), s
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggttitle("Expected heterozygosity vs Longitude, Neutral Intergenic SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R3_inter), x=-71.44, y=0.2635, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

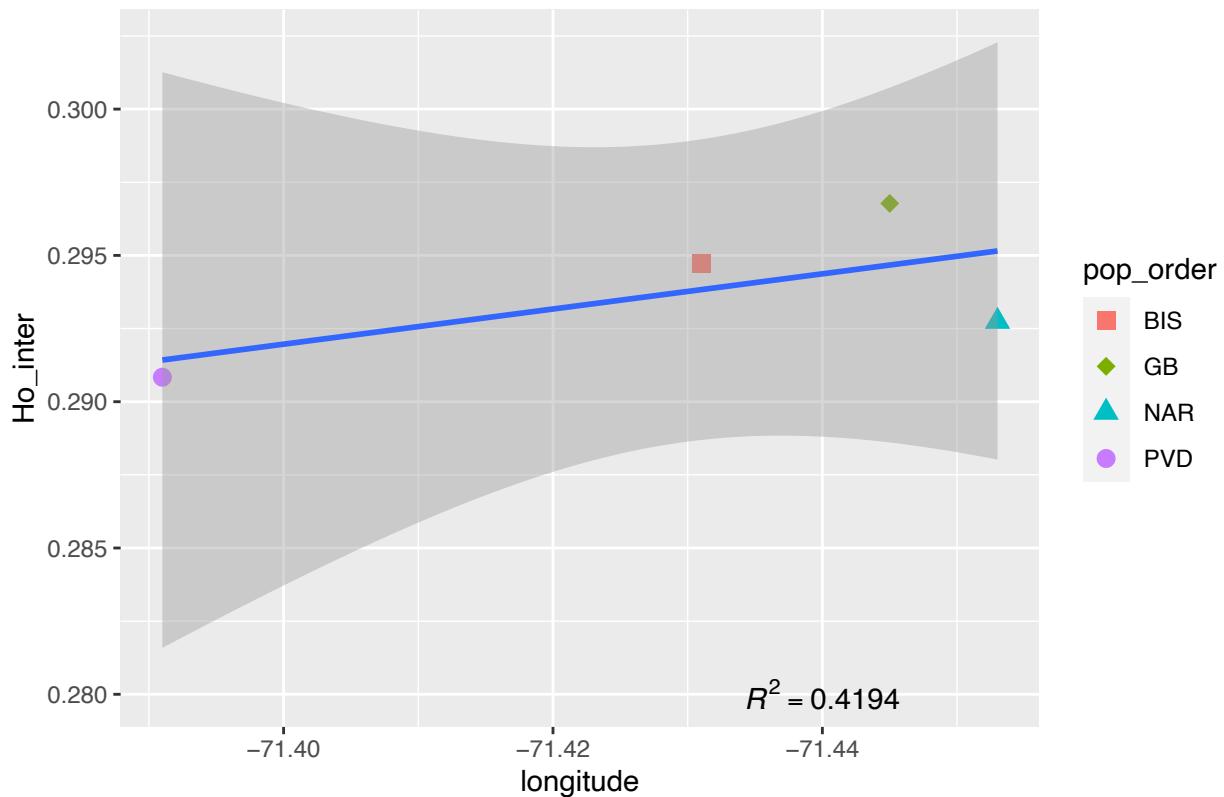
## Expected heterozygosity vs Longitude, Neutral Intergenic SNPs



```
#Plot Ho vs Longitude
R3_inter = round(summary(lm(y_inter$Ho_inter ~ y_inter$longitude))$r.squared, 4)
ggplot(y_inter, aes(x = longitude, y = Ho_inter)) + geom_point(aes(shape=pop_order, color=pop_order), s
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  gtitle("Observed heterozygosity vs Longitude, Neutral Intergenic SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R3_inter), x=-71.44, y=0.28, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

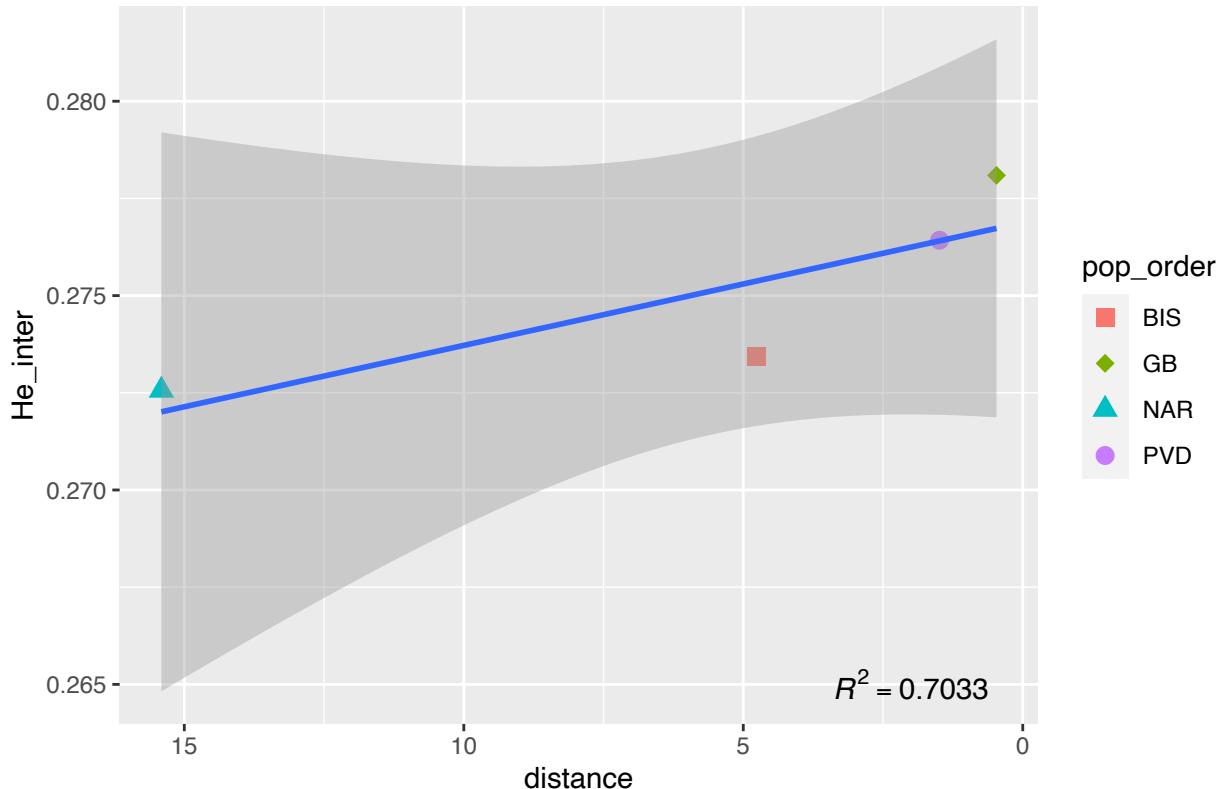
## Observed heterozygosity vs Longitude, Neutral Intergenic SNPs



```
#Plot He vs Distance from sewage outflow
R4_inter = round(summary(lm(y_inter$He_inter ~ y_inter$distance))$r.squared, 4)
ggplot(y_inter, aes(x = distance, y = He_inter)) + geom_point(aes(shape=pop_order, color=pop_order), size=10) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggttitle("Expected heterozygosity vs Distance, Neutral Intergenic SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R4_inter), x=2, y=0.265, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

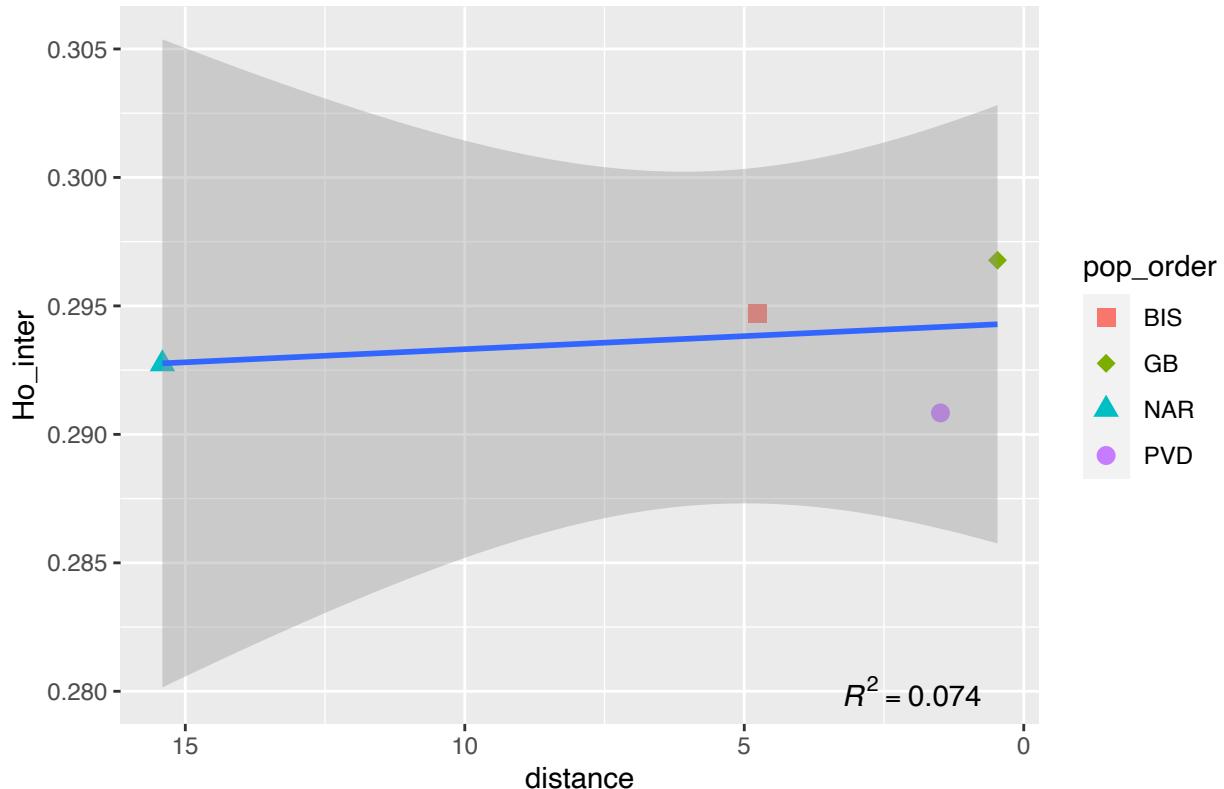
## Expected heterozygosity vs Distance, Neutral Intergenic SNPs



```
#Plot Ho vs Distance from sewage outflow
R4_inter = round(summary(lm(y_inter$Ho_inter ~ y_inter$distance))$r.squared, 4)
ggplot(y_inter, aes(x = distance, y = Ho_inter)) + geom_point(aes(shape=pop_order, color=pop_order), size=10) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggttitle("Observed heterozygosity vs Distance, Neutral Intergenic SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R4_inter), x=2, y=0.28, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

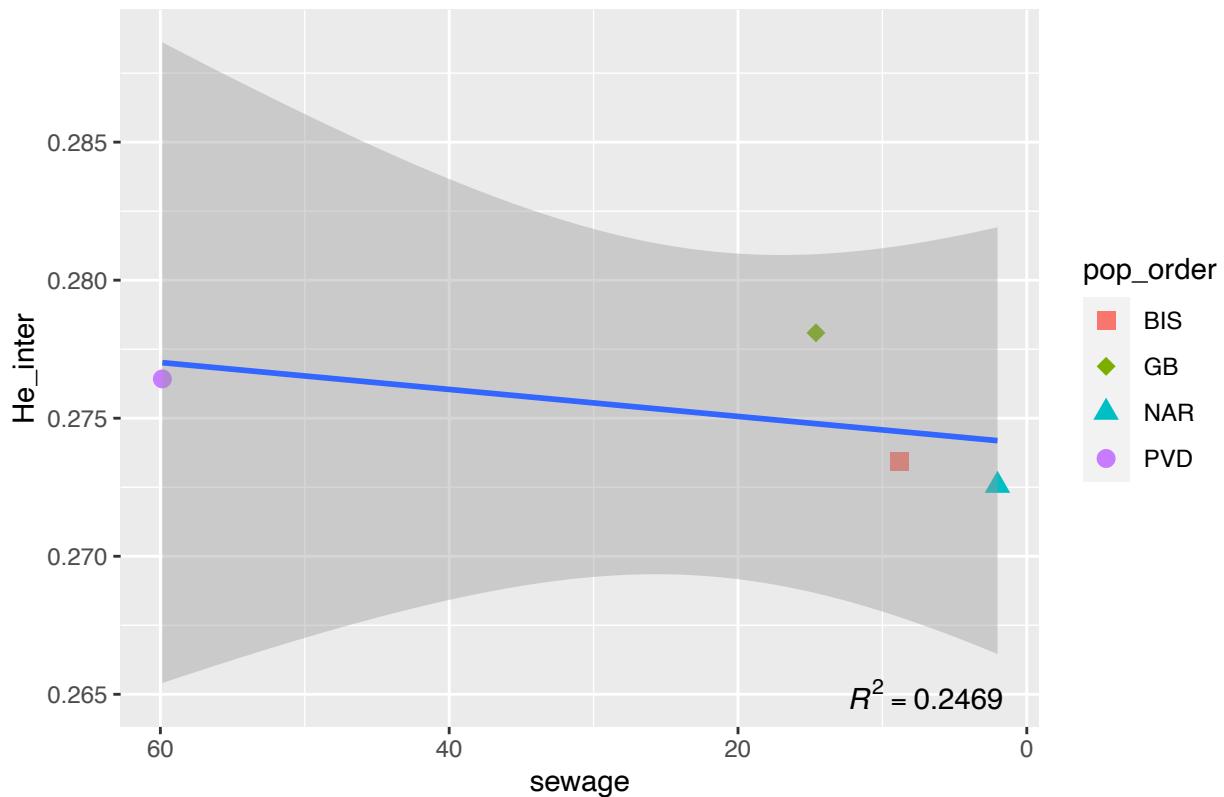
### Observed heterozygosity vs Distance, Neutral Intergenic SNPs



```
#Plot He vs Sewage effluent
R_inter = round(summary(lm(y_inter$He_inter ~ y_inter$sewage))$r.squared, 4)
ggplot(y_inter, aes(x = sewage, y = He_inter)) + geom_point(aes(shape=pop_order, color=pop_order), size=3) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggttitle("Expected heterozygosity vs Sewage Effluent, Neutral Intergenic SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R_inter), x=7, y=0.265, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

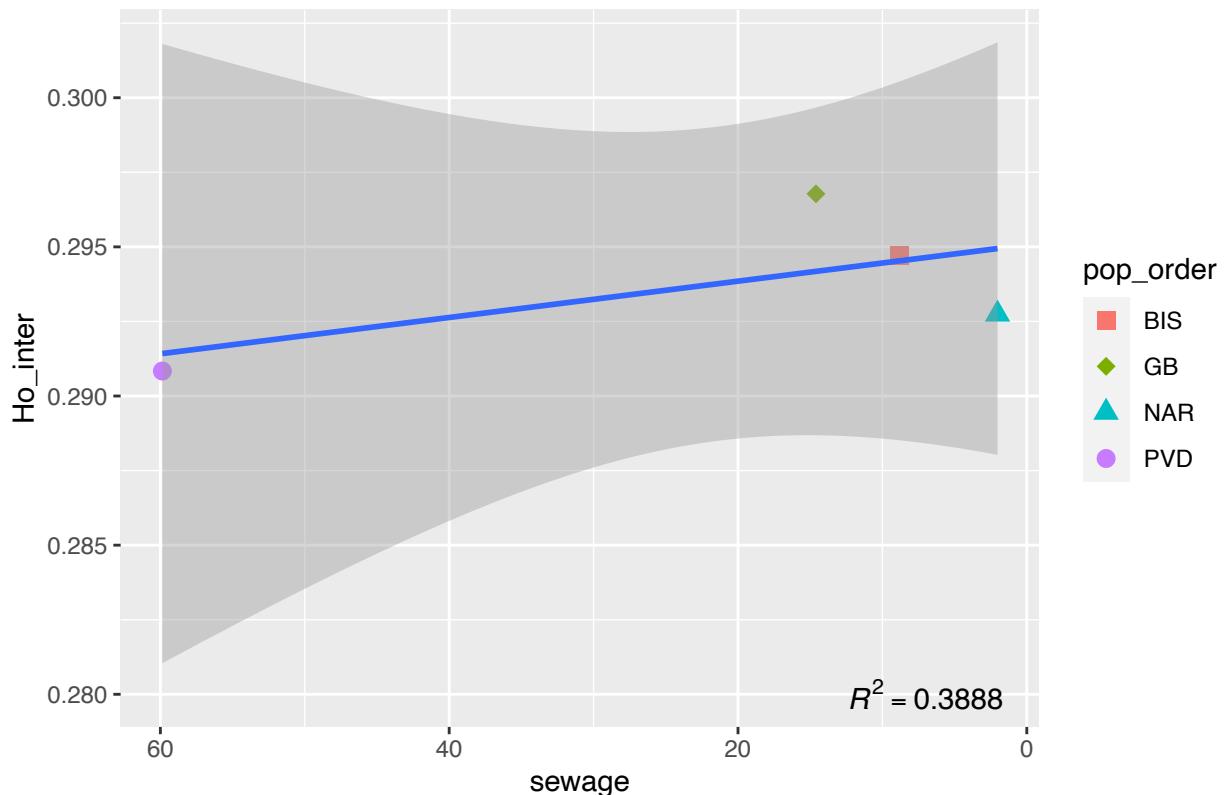
## Expected heterozygosity vs Sewage Effluent, Neutral Intergenic SNPs



```
#Plot Ho vs Sewage effluent
R_inter = round(summary(lm(y_inter$Ho_inter ~ y_inter$sewage))$r.squared, 4)
ggplot(y_inter, aes(x = sewage, y = Ho_inter)) + geom_point(aes(shape=pop_order, color=pop_order), size=5) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggttitle("Observed heterozygosity vs Sewage Effluent, Neutral Intergenic SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R_inter), x=7, y=0.28, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

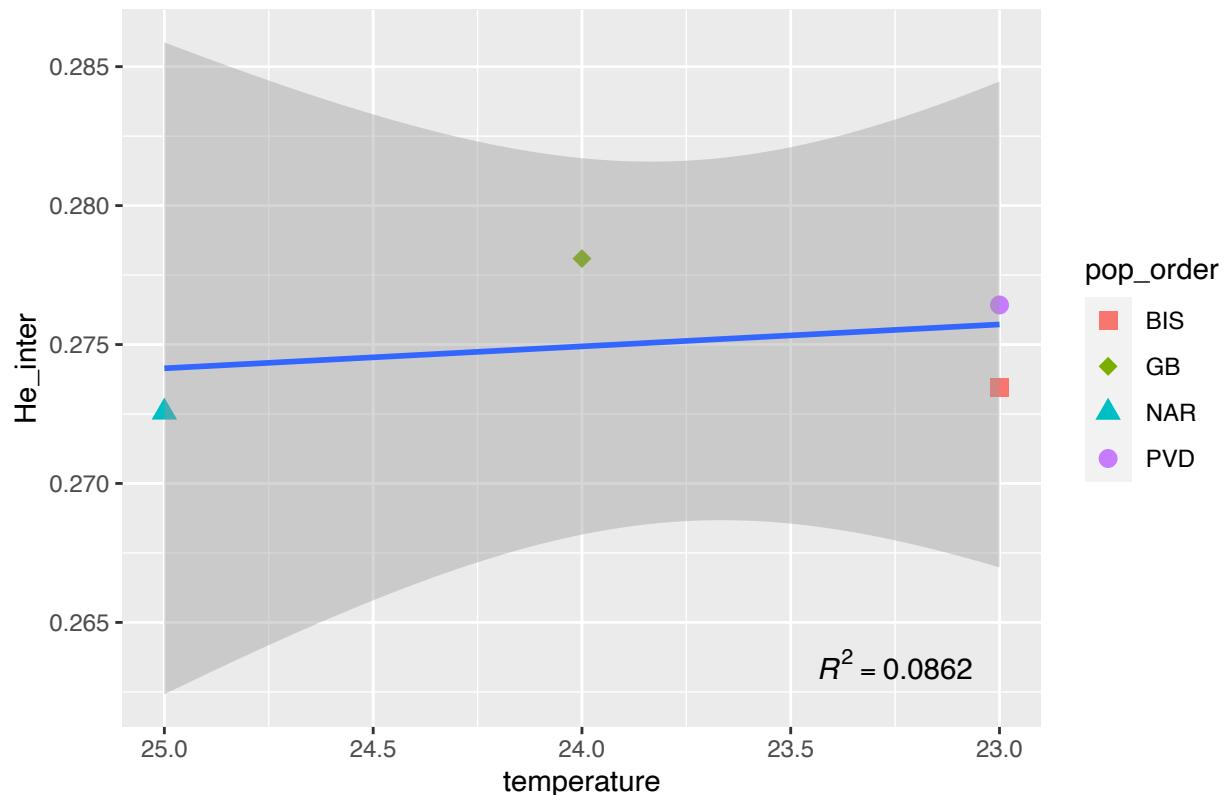
## Observed heterozygosity vs Sewage Effluent, Neutral Intergenic SNPs



```
#Plot He vs Temperature
R5_inter = round(summary(lm(y_inter$He_inter ~ y_inter$temperature))$r.squared, 4)
ggplot(y_inter, aes(x = temperature, y = He_inter)) + geom_point(aes(shape=pop_order, color=pop_order),
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggttitle("Expected heterozygosity vs Temperature, Neutral Intergenic SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R5_inter), x=23.25, y=0.2635, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

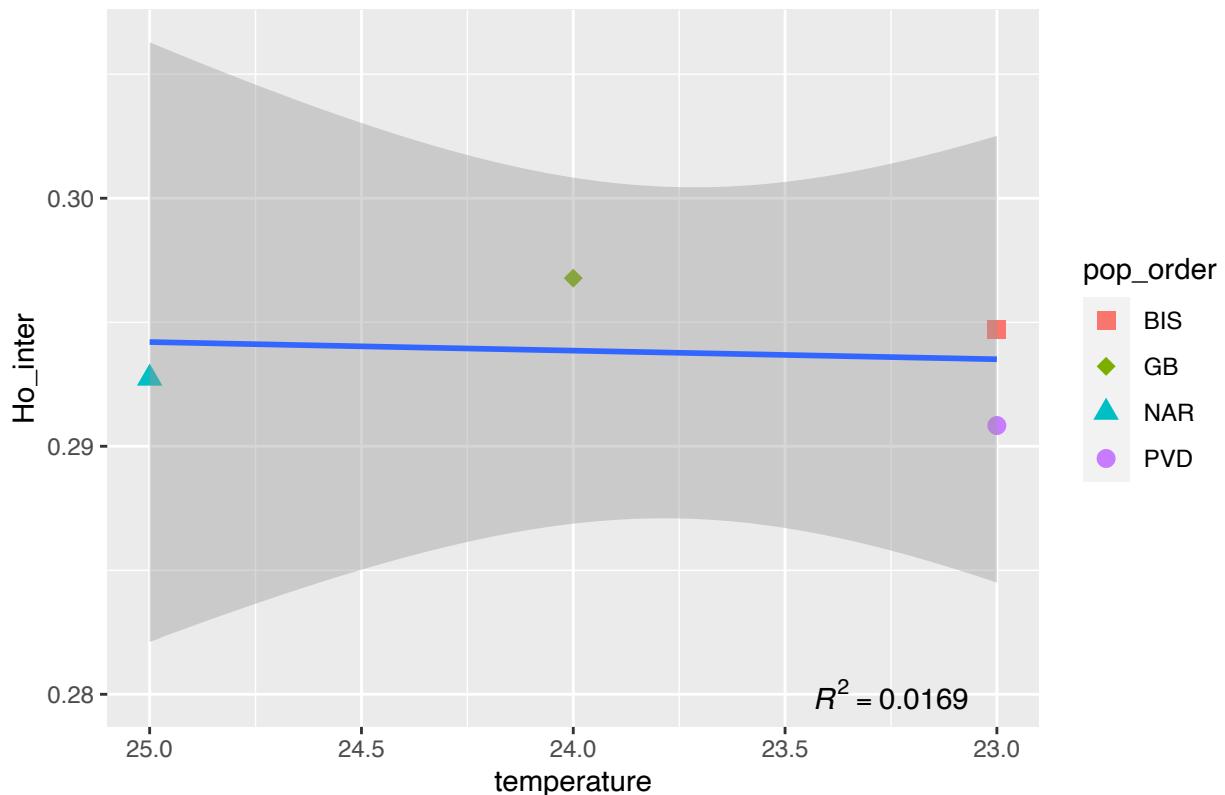
## Expected heterozygosity vs Temperature, Neutral Intergenic SNPs



```
#Plot Ho vs Temperature
R5_inter = round(summary(lm(y_inter$Ho_inter ~ y_inter$temperature))$r.squared, 4)
ggplot(y_inter, aes(x = temperature, y = Ho_inter)) + geom_point(aes(shape=pop_order, color=pop_order),
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  gtitle("Observed heterozygosity vs Temperature, Neutral Intergenic SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R5_inter), x=23.25, y=0.28, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

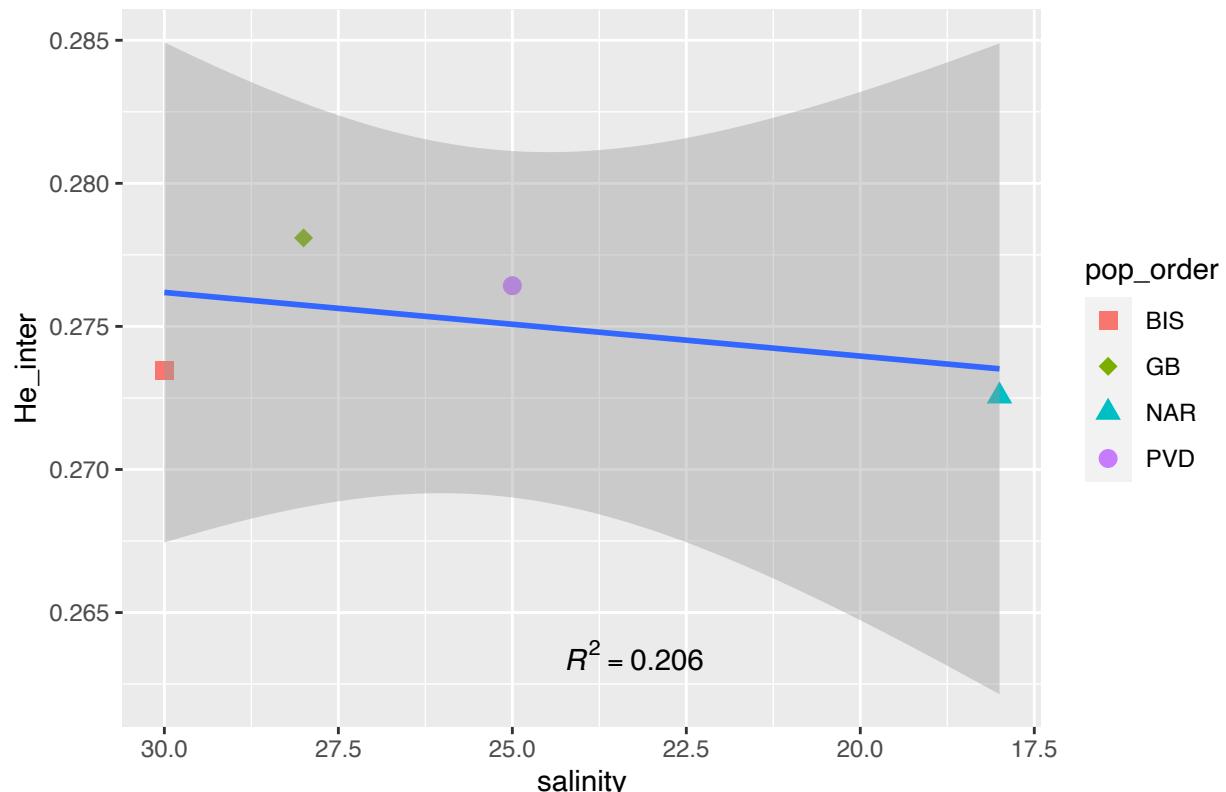
## Observed heterozygosity vs Temperature, Neutral Intergenic SNPs



```
#Plot He vs Salinity
R6_inter = round(summary(lm(y_inter$He_inter ~ y_inter$salinity))$r.squared, 4)
ggplot(y_inter, aes(x = salinity, y = He_inter)) + geom_point(aes(shape=pop_order, color=pop_order), size=10) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggttitle("Expected heterozygosity vs Salinity, Neutral Intergenic SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R6_inter), x=23.25, y=0.2635, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

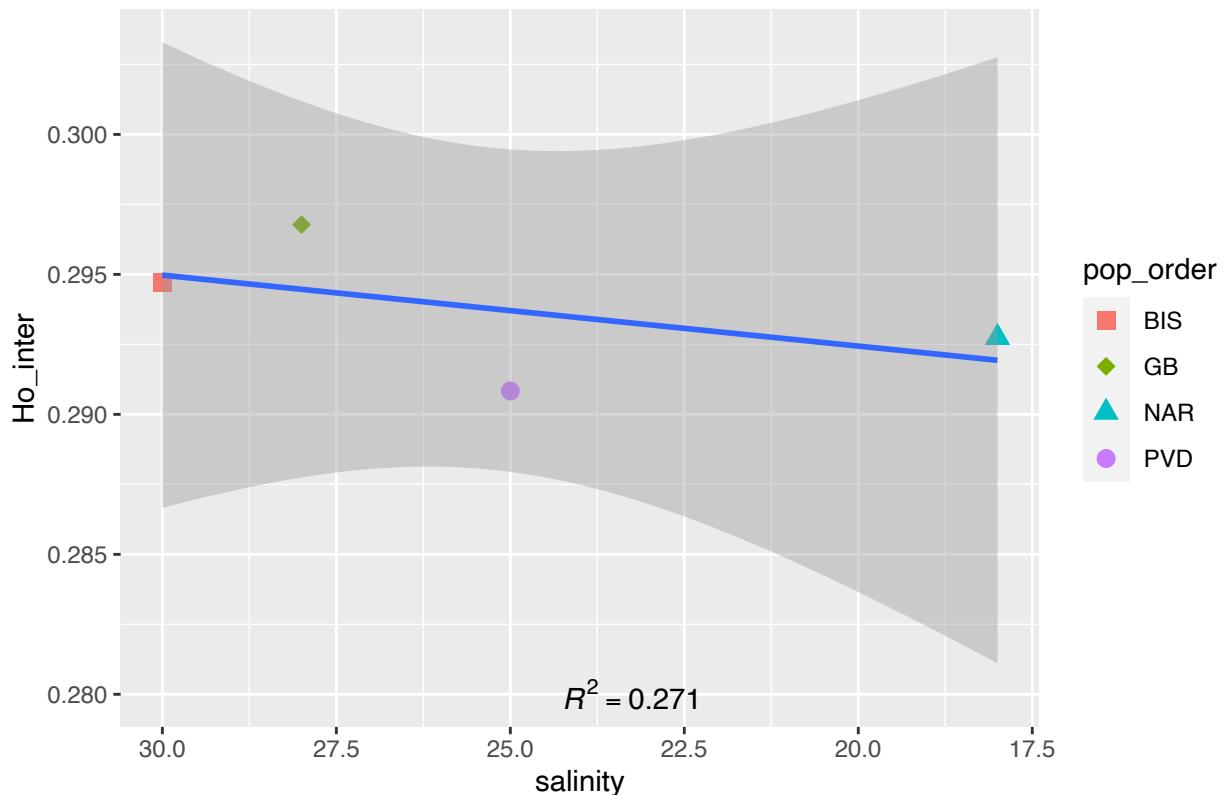
### Expected heterozygosity vs Salinity, Neutral Intergenic SNPs



```
#Plot Ho vs Salinity
R6_inter = round(summary(lm(y_inter$Ho_inter ~ y_inter$salinity))$r.squared, 4)
ggplot(y_inter, aes(x = salinity, y = Ho_inter)) + geom_point(aes(shape=pop_order, color=pop_order), size=10) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggttitle("Observed heterozygosity vs Salinity, Neutral Intergenic SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R6_inter), x=23.25, y=0.28, parse=T) +
  scale_x_reverse()

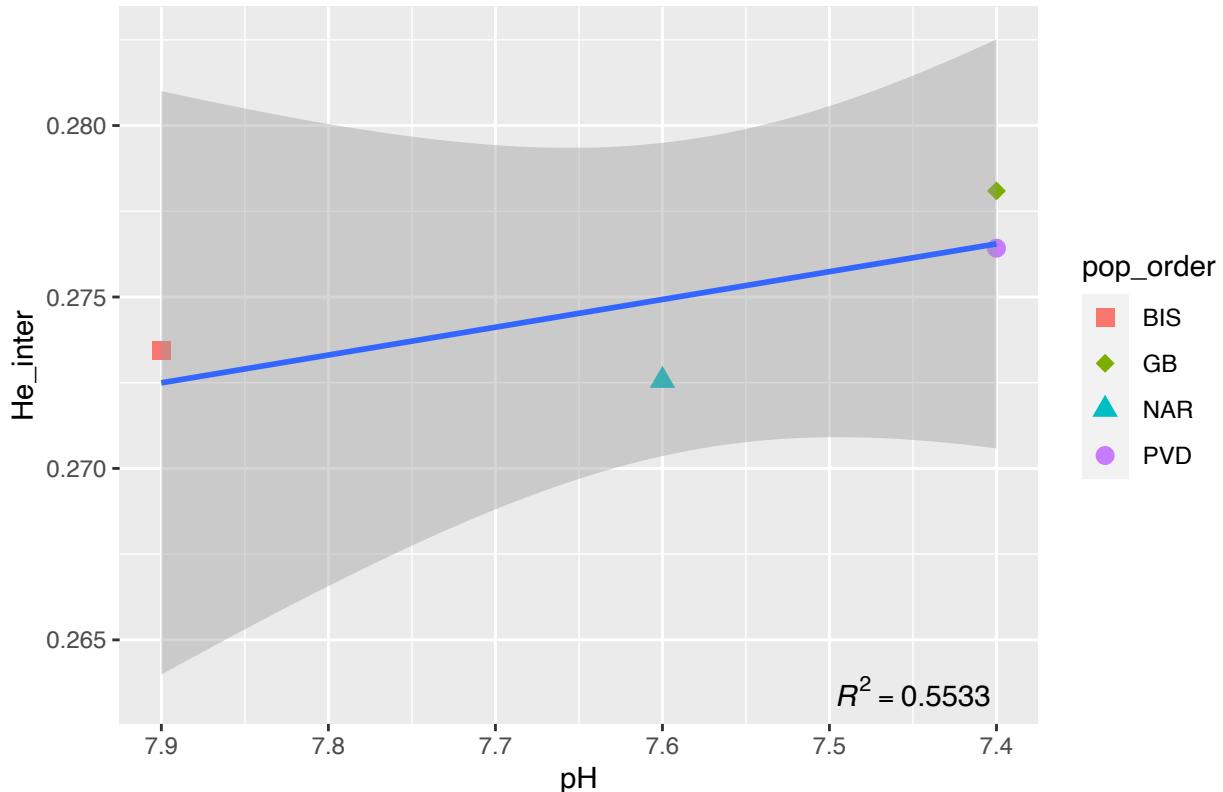
## `geom_smooth()` using formula 'y ~ x'
```

## Observed heterozygosity vs Salinity, Neutral Intergenic SNPs



```
#Plot He vs pH
R7_inter = round(summary(lm(y_inter$He_inter ~ y_inter$pH))$r.squared, 4)
ggplot(y_inter, aes(x = pH, y = He_inter)) + geom_point(aes(shape=pop_order, color=pop_order), size = 3)
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggtitle("Expected heterozygosity vs pH, Neutral Intergenic SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R7_inter), x=7.45, y=0.2635, parse=T) +
  scale_x_reverse()
## `geom_smooth()` using formula 'y ~ x'
```

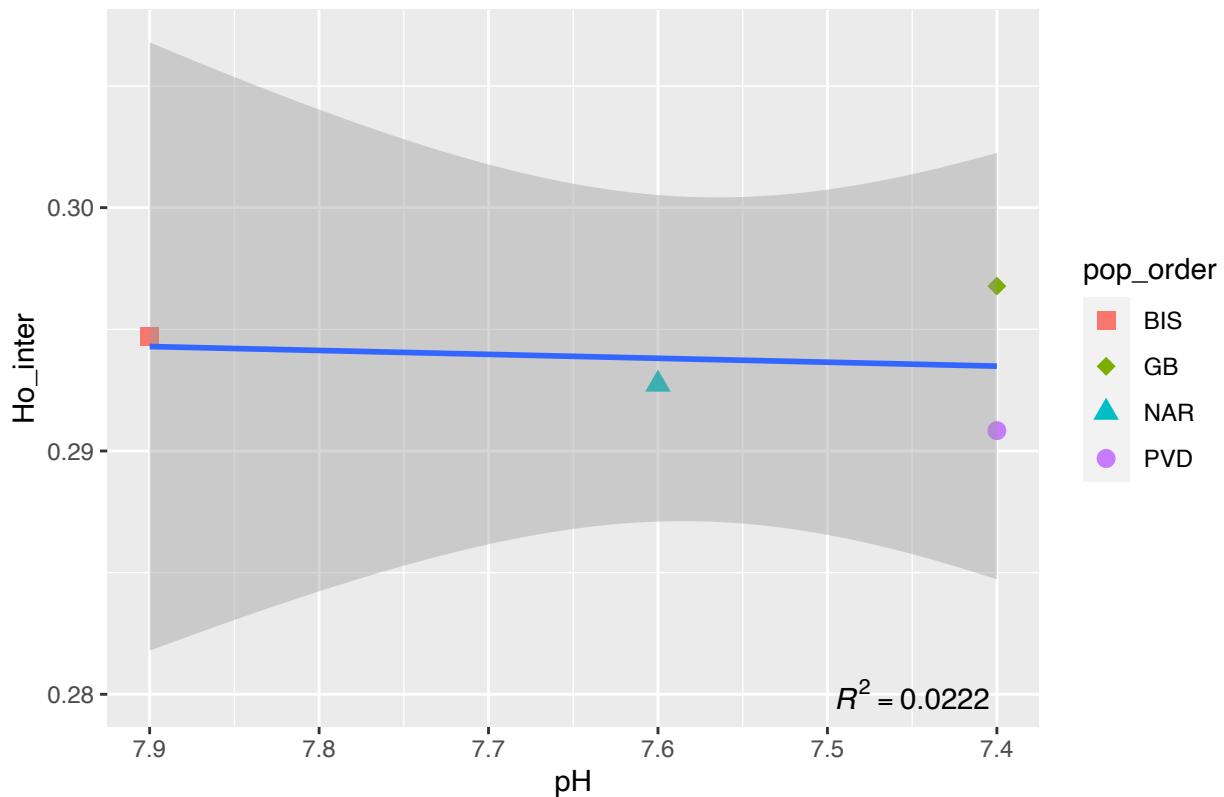
## Expected heterozygosity vs pH, Neutral Intergenic SNPs



```
#Plot Ho vs pH
R7_inter = round(summary(lm(y_inter$Ho_inter ~ y_inter$pH))$r.squared, 4)
ggplot(y_inter, aes(x = pH, y = Ho_inter)) + geom_point(aes(shape=pop_order, color=pop_order), size = 3)
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggtitle("Observed heterozygosity vs pH, Neutral Intergenic SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R7_inter), x=7.45, y=0.28, parse=T) +
  scale_x_reverse()

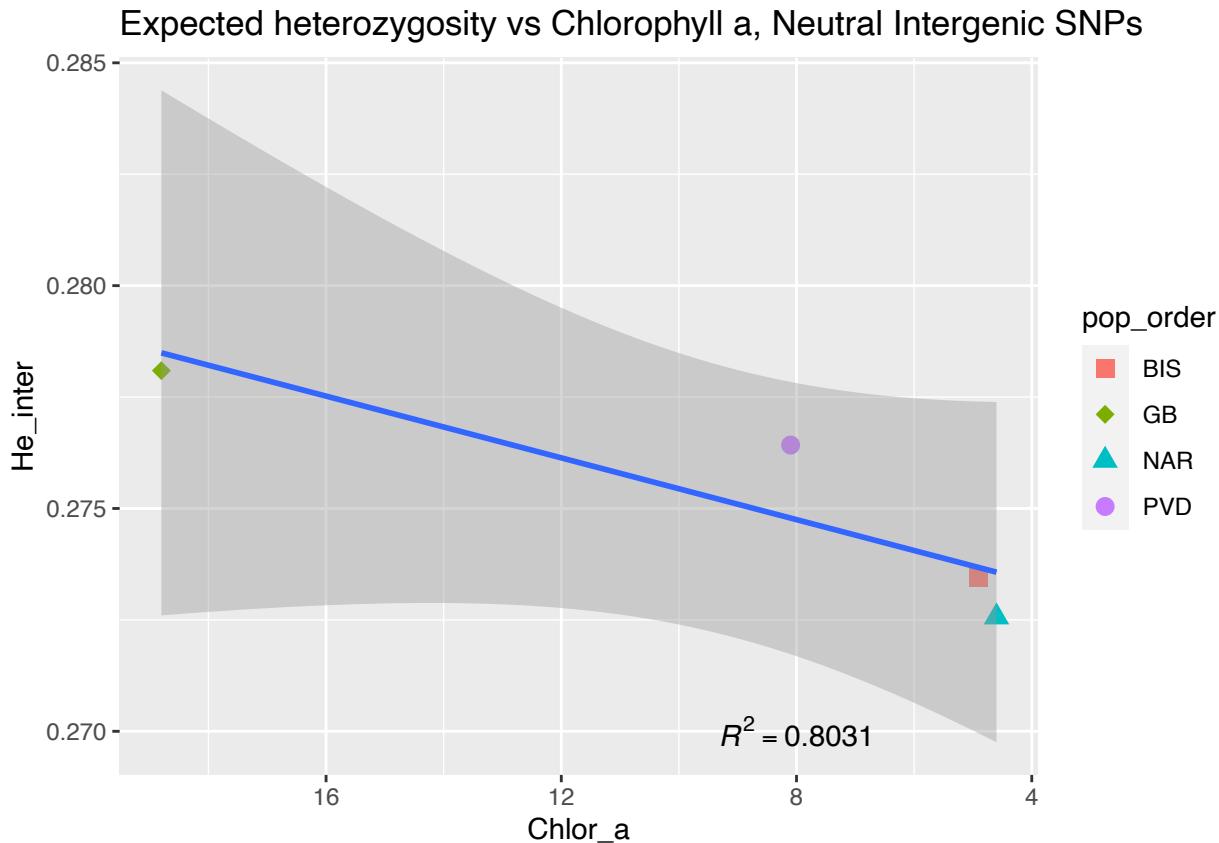
## `geom_smooth()` using formula 'y ~ x'
```

## Observed heterozygosity vs pH, Neutral Intergenic SNPs



```
#Plot He vs Chlorophyll a
R8_inter = round(summary(lm(y_inter$He_inter ~ y_inter$Chlor_a))$r.squared, 4)
ggplot(y_inter, aes(x = Chlor_a, y = He_inter)) + geom_point(aes(shape=pop_order, color=pop_order), size=10) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggttitle("Expected heterozygosity vs Chlorophyll a, Neutral Intergenic SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R8_inter), x=8, y=0.27, parse=T) +
  scale_x_reverse()

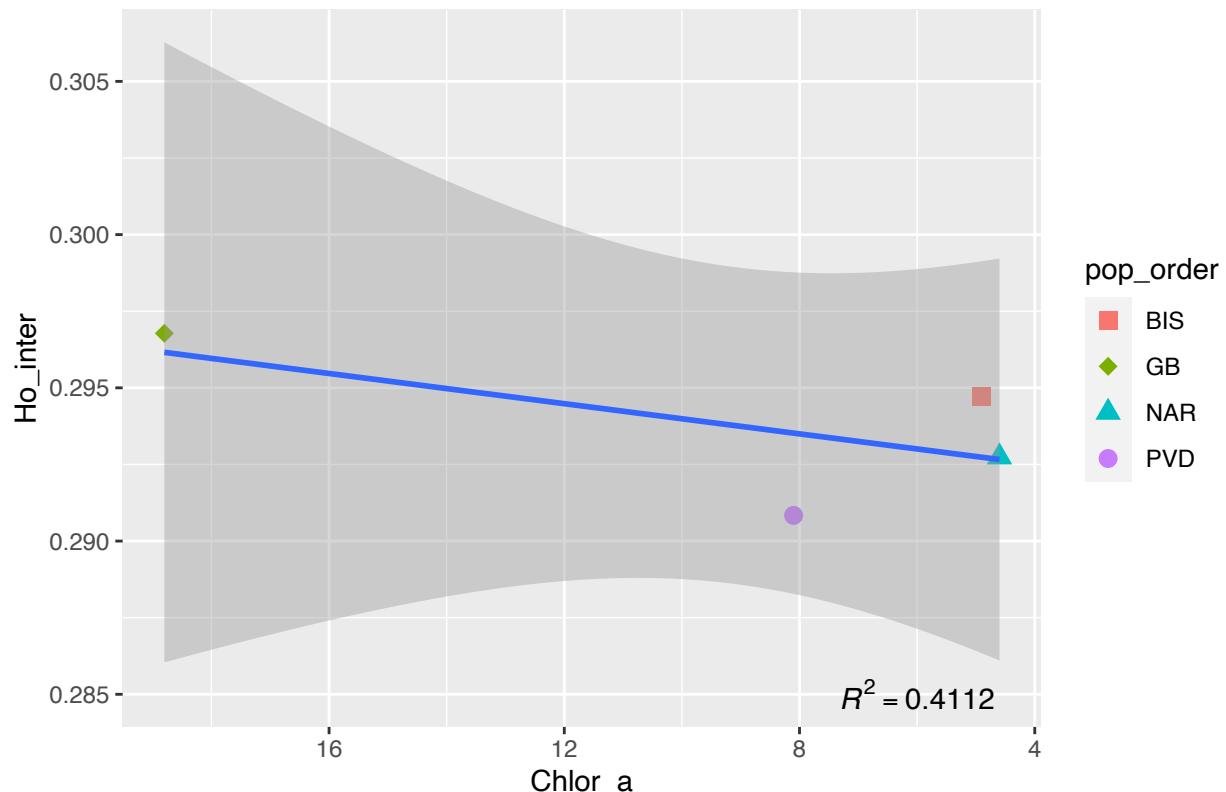
## `geom_smooth()` using formula 'y ~ x'
```



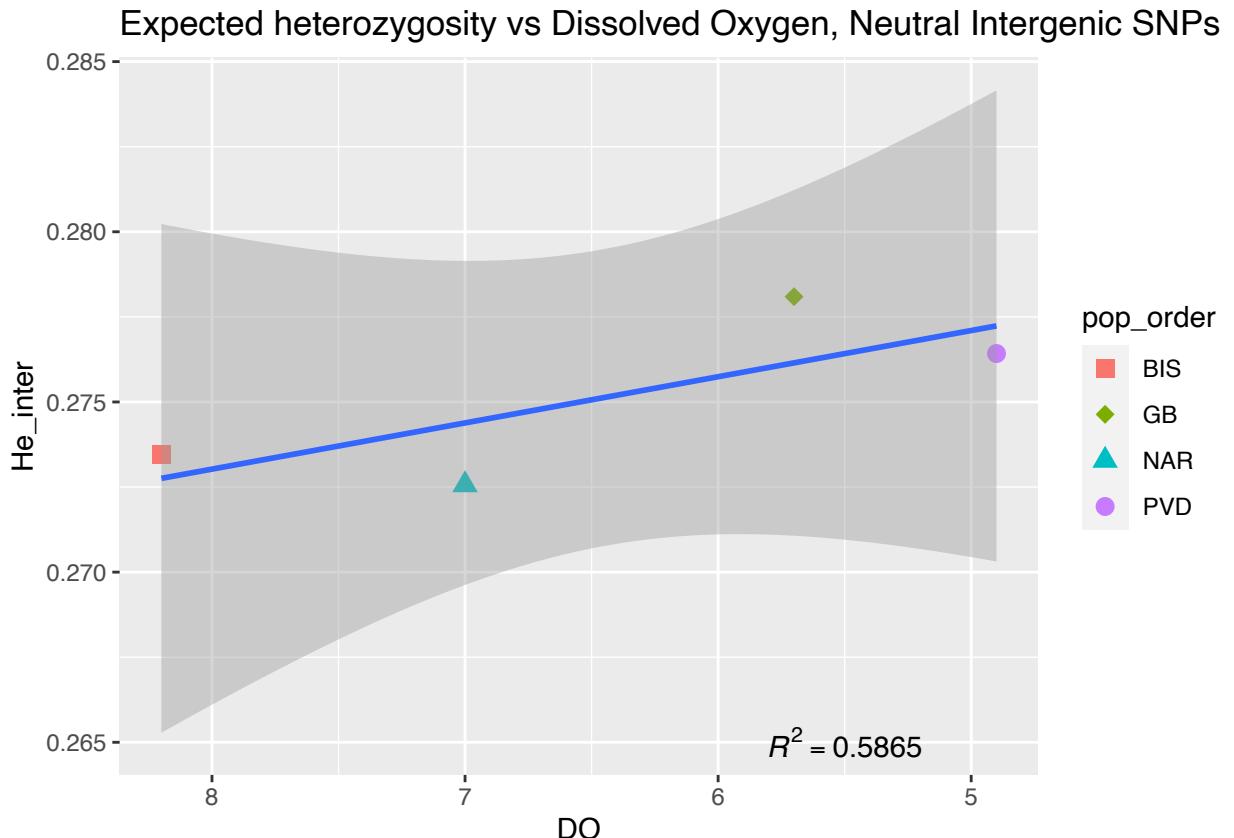
```
#Plot Ho vs Chlorophyll a
R8_inter = round(summary(lm(y_inter$Ho_inter ~ y_inter$Chlor_a))$r.squared, 4)
ggplot(y_inter, aes(x = Chlor_a, y = Ho_inter)) + geom_point(aes(shape=pop_order, color=pop_order), size=10) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggttitle("Observed heterozygosity vs Chlorophyll a, Neutral Intergenic SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R8_inter), x=6, y=0.285, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

### Observed heterozygosity vs Chlorophyll a, Neutral Intergenic SNPs



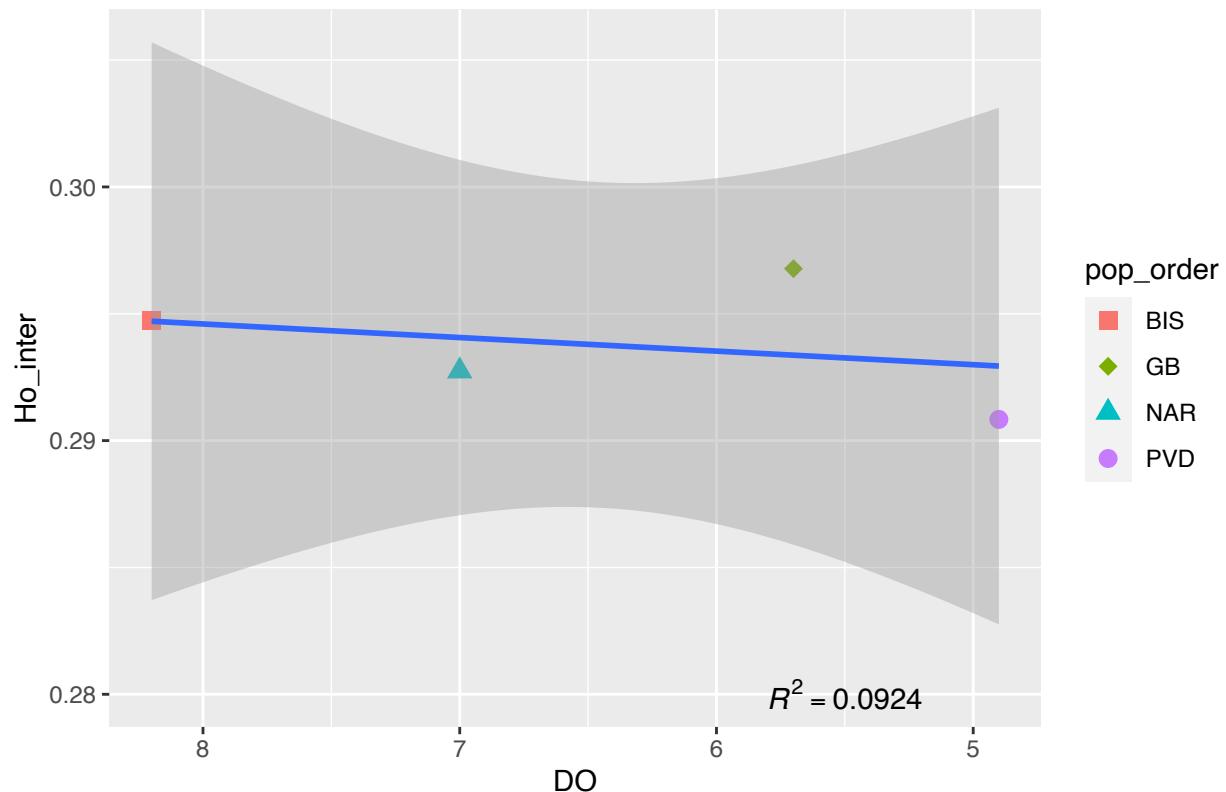
```
#Plot He vs Dissolved Oxygen
R9_inter = round(summary(lm(y_inter$He_inter ~ y_inter$D0))$r.squared, 4)
ggplot(y_inter, aes(x = D0, y = He_inter)) + geom_point(aes(shape=pop_order, color=pop_order), size = 3)
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggttitle("Expected heterozygosity vs Dissolved Oxygen, Neutral Intergenic SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R9_inter), x=5.5, y=0.265, parse=T) +
  scale_x_reverse()
## `geom_smooth()` using formula 'y ~ x'
```



```
#Plot Ho vs Dissolved Oxygen
R9_inter = round(summary(lm(y_inter$Ho_inter ~ y_inter$DO))$r.squared, 4)
ggplot(y_inter, aes(x = DO, y = Ho_inter)) + geom_point(aes(shape=pop_order, color=pop_order), size = 3)
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggttitle("Observed heterozygosity vs Dissolved Oxygen, Neutral Intergenic SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R9_inter), x=5.5, y=0.28, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

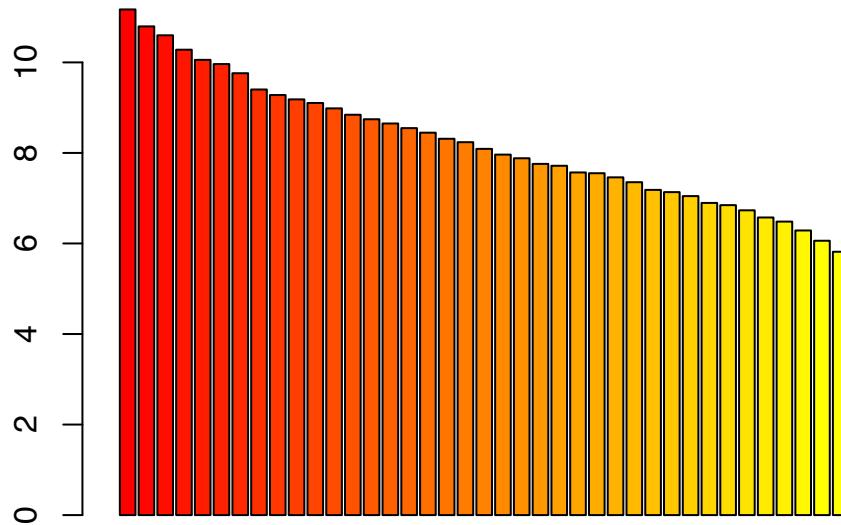
### Observed heterozygosity vs Dissolved Oxygen, Neutral Intergenic SNPs



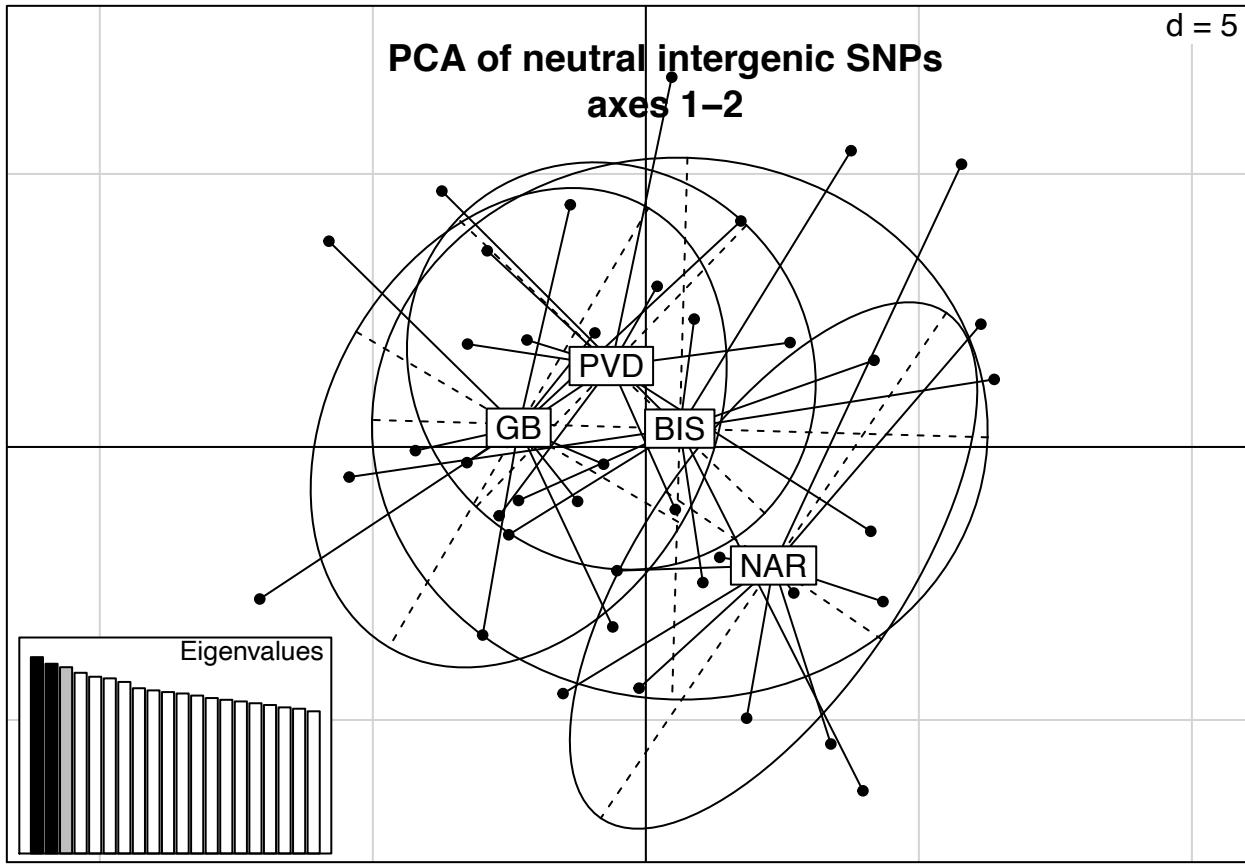
### PCA

```
X_inter <- tab(stratified.u_inter, freq = TRUE, NA.method = "mean")
pca1_inter <- dudi.pca(X_inter, scale = FALSE, scannf = FALSE, nf = 3)
barplot(pca1_inter$eig[1:50], main = "PCA eigenvalues", col = heat.colors(50))
```

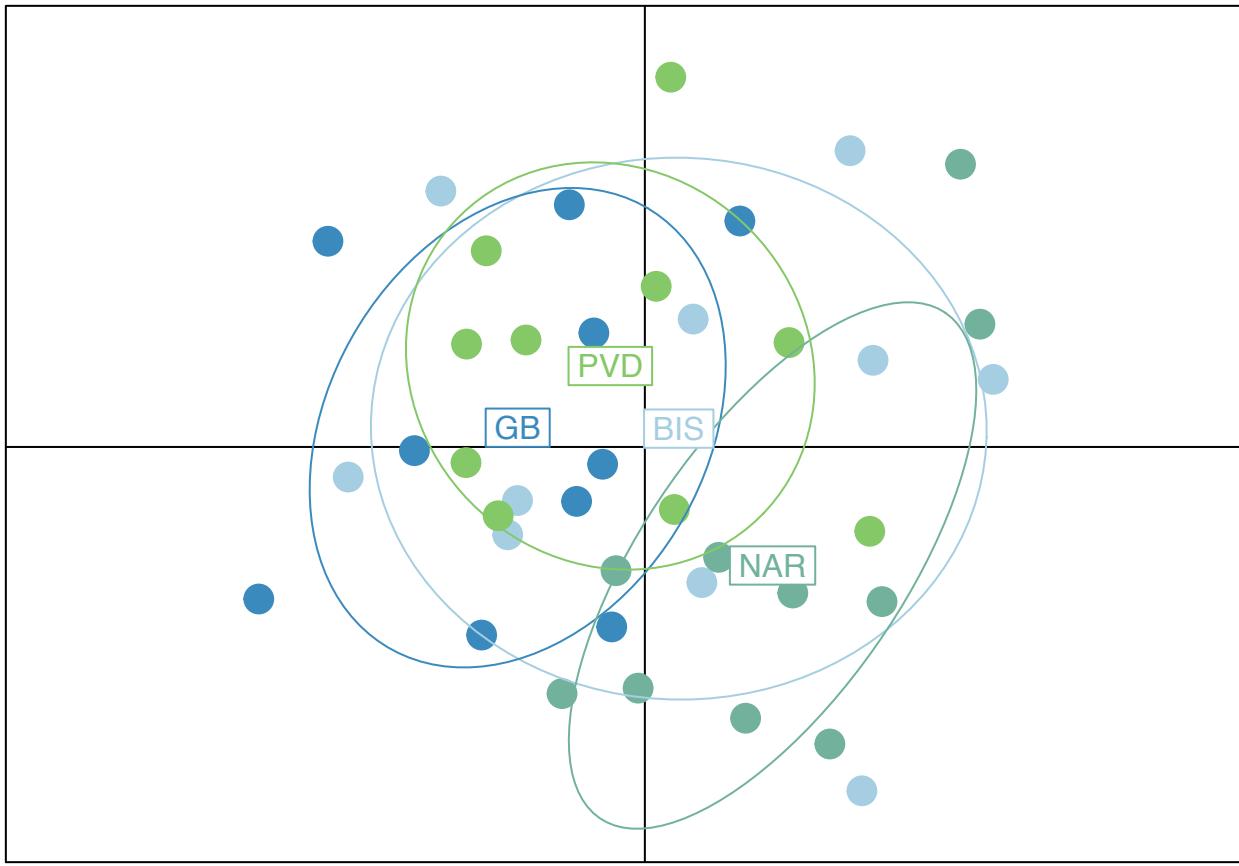
## PCA eigenvalues



```
s.class(pca1_inter$li, pop(stratated.filt_inter))
title("PCA of neutral intergenic SNPs\naxes 1-2")
add.scatter.eig(pca1_inter$eig[1:20], 3,1,2)
```



```
col <- funky(15)
s.class(pca1_inter$li, pop(stratated.filt_inter), xax=1, yax=2, col=col, axesell=FALSE, cstar=0, cpoint=3,
```



## Discriminant Analysis of Principal Components (DAPC)

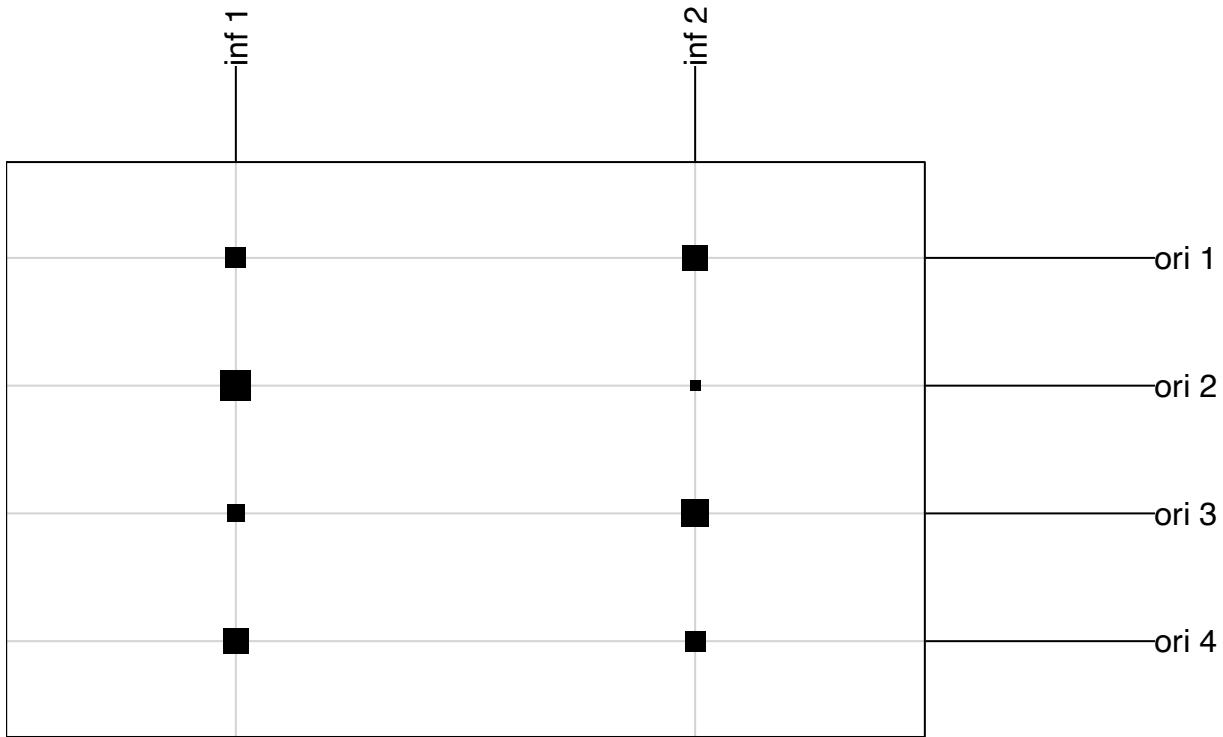
The first line of the code block below is interactive. In the console, it will tell you:

1. Choose the number PCs to retain ( $\geq 1$ ): I picked 1 2. Choose the number discriminant functions to retain ( $\geq 1$ ): I picked 2

```
grp_inter <- find.clusters(stratated.u_inter, n.pca = 1, choose.n.clust = FALSE)
table(pop(stratated.u_inter), grp_inter$grp)
```

```
##          1 2
##  BIS 4 6
##  GB  9 1
##  NAR 3 7
##  PVD 6 4

table.value(table(pop(stratated.u_inter), grp_inter$grp), col.lab=paste("inf", 1:2), row.lab=paste("ori"
```

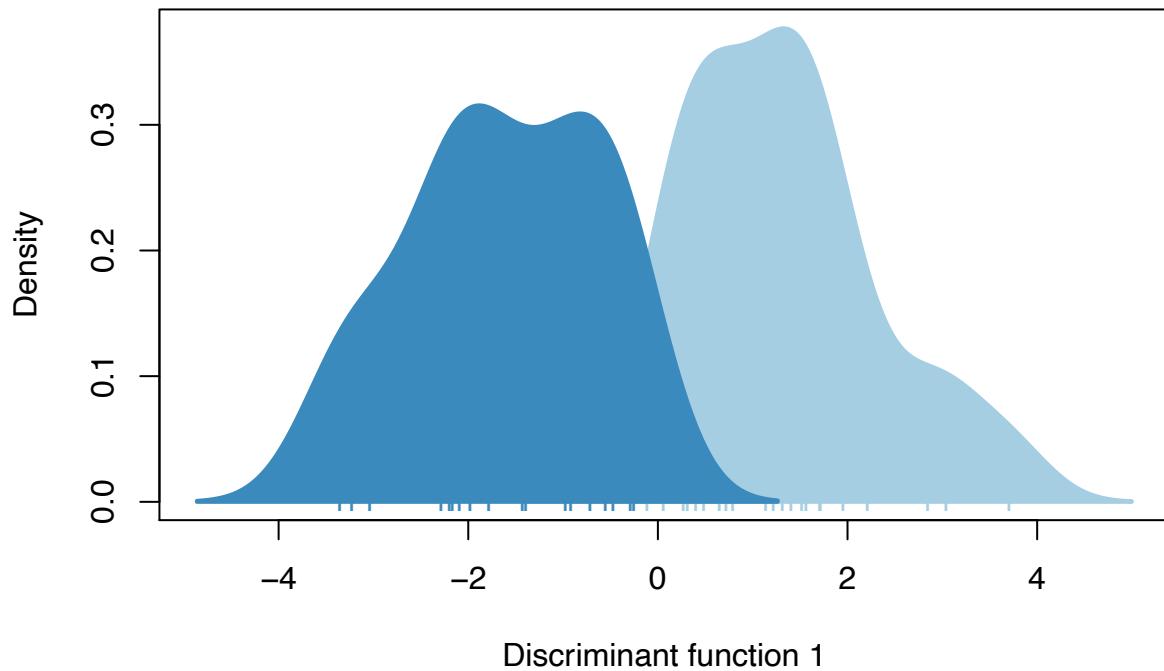


**1 • 3 ■ 5 ■ 7 ■ 9 ■**

Again, the first line of the code block below is interactive. In the console, it will tell you:

1. Choose the number PCs to retain ( $\geq 1$ ): I picked 1
2. Choose the number discriminant functions to retain ( $\geq 1$ ): I picked 1

```
dapc1_inter <- dapc(stratated.u_inter, grp_inter$grp, n.pca = 1, n.da= 1)
scatter(dapc1_inter,col=col,bg="white", solid=1)
```

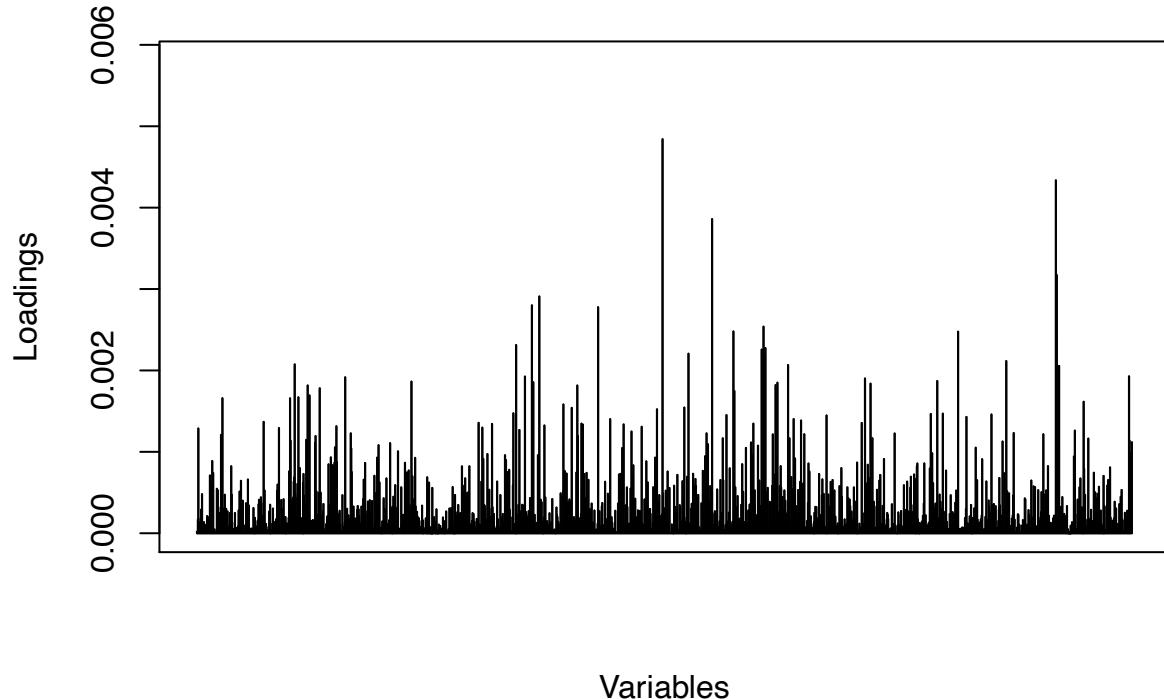


Again, the first line of the code block below is interactive. In the console, it will tell you:

1. Choose the number PCs to retain ( $\geq 1$ ): I picked 1 2. Choose the number discriminant functions to retain ( $\geq 1$ ): I picked 1

```
contrib_inter <- loadingplot(dapc1_inter$var.contr, axis=1, thres=.01, lab.jitter=1)
```

## Loading plot

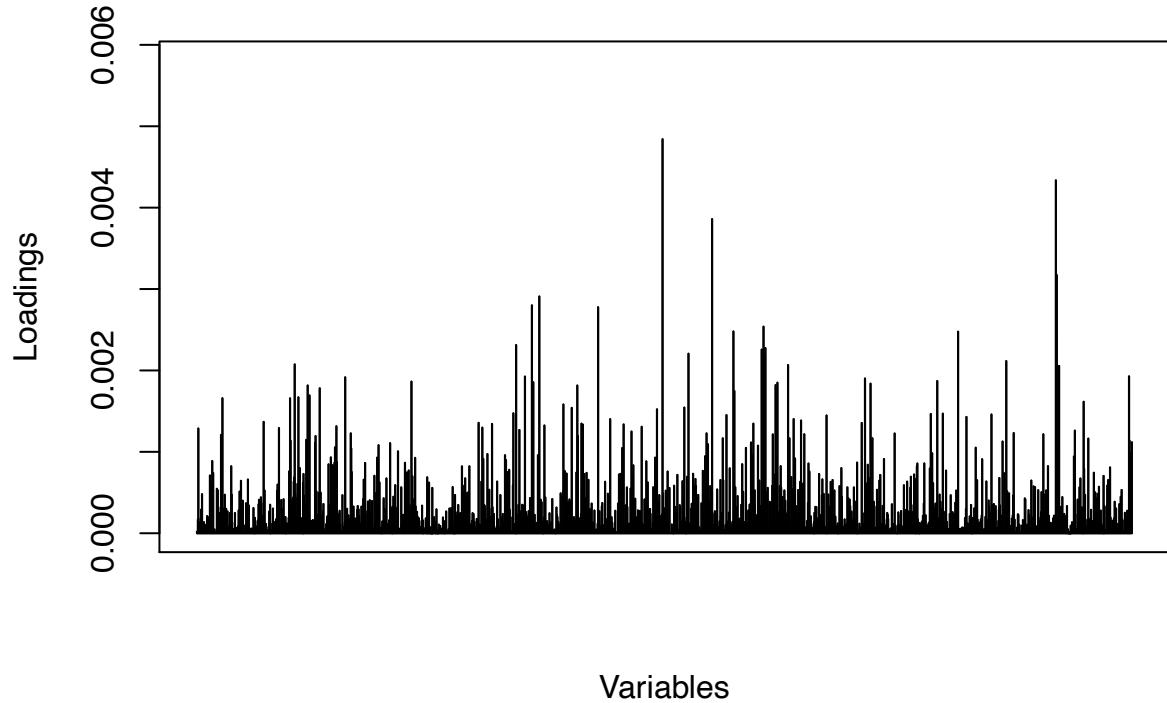


## Variables

```
contrib_inter
```

```
## NULL
setPop(rad.filt_inter) <- ~Library
dapc1_inter <- dapc(stratated.u_inter, pop=stratated.u_inter), n.pca = 1, n.da = 1)
contrib_inter <- loadingplot(dapc1_inter$var.contr, axis=1, thres=.05, lab.jitter=1)
```

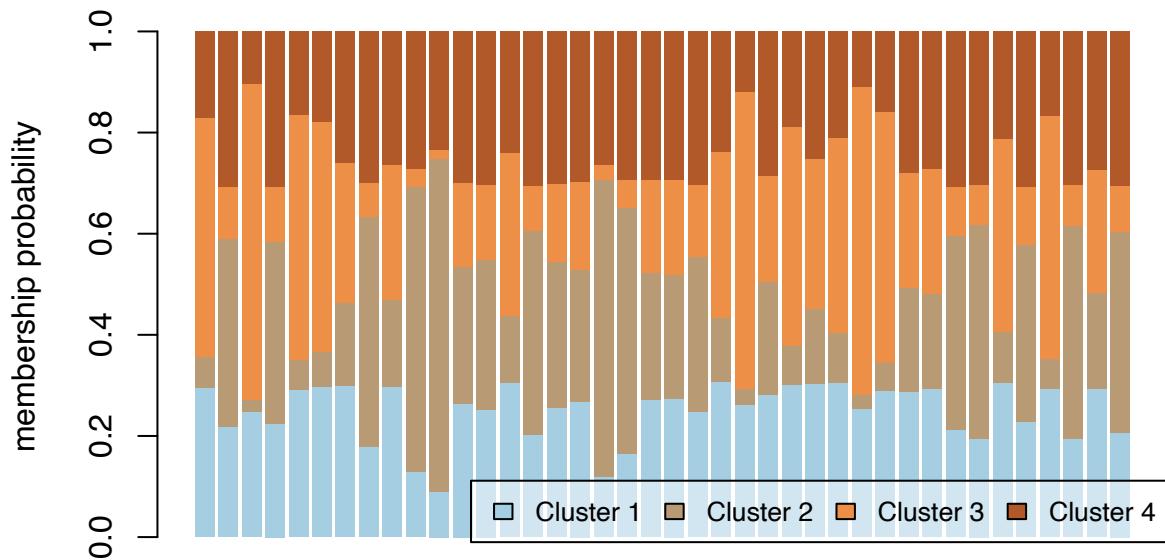
## Loading plot



```
#Structure Like
compoplot(dapc1_inter, posi="bottomright", txt.leg=paste("Cluster", 1:4), lab="", ncol=1, xlab="individual")
```



```
temp_inter <- which(apply(dapc1_inter$posterior, 1, function(e) all(e<0.9)))
compoplot(dapc1_inter, subset=temp_inter, posi="bottomright", txt.leg=paste("Cluster", 1:4), ncol=2)
```



## PCAviz

```

NA.afDraw<- function(ind){
  ind.mat <- ind@tab
  new.mat <- ind.mat
  af = colSums(ind.mat[,seq(1,ncol(ind.mat)-1,2)],na.rm = TRUE) /
    (2*apply(ind.mat[,seq(1,ncol(ind.mat)-1,2)],2,function(x) sum(!is.na(x))))
  af.Draw <- function(geno, af){
    new <- function(geno,af){
      if(is.na(geno)){
        newA = rbinom(1,2,af)
      } else {newA <- geno}
      return(newA)
    }
    new.row <- mapply(geno,af,FUN = new)
    return(new.row)
  }

  new.mat[,seq(1,ncol(ind.mat)-1,2)] <- t(apply(ind.mat[,seq(1,ncol(ind.mat)-1,2)],1,af.Draw,af))
  new.mat[,seq(2,ncol(ind.mat),2)] <- 2-new.mat[,seq(1,ncol(ind.mat)-1,2)]
  new.ind <- ind
  new.ind@tab <- new.mat
  return(new.ind)
}

```

```

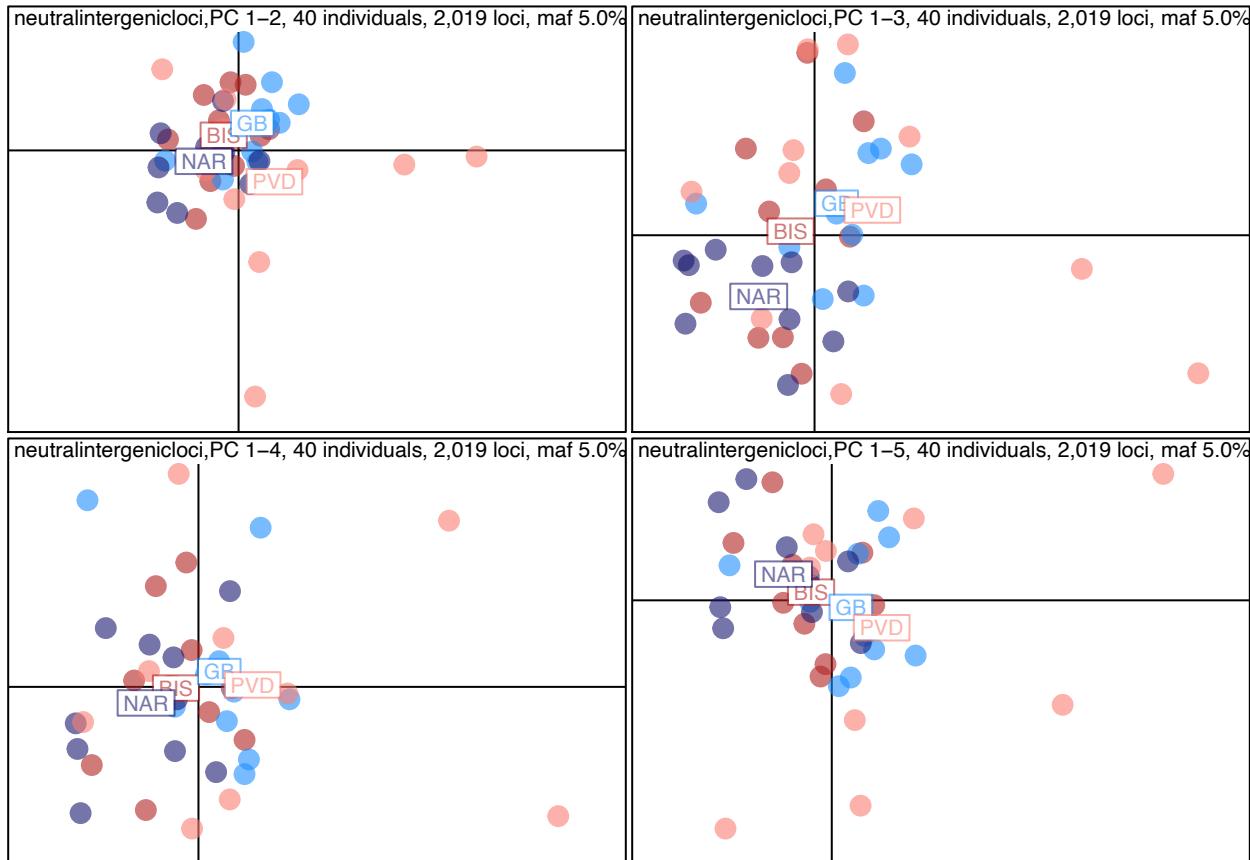
u.na_inter <- NA.afDraw(stratified.u_inter)

pca_inter <- dudi.pca(u.na_inter,center=TRUE,scale=TRUE,scannf = F, nf = 30)

col18 <- funky(length(unique(u.na_inter$strata$Population)))
#Colors that match the neutral Structure results
col6 <- c("firebrick","dodgerblue","midnightblue","salmon")
par(mfrow=c(2,2))

s.class(pca_inter$li, strata(u.na_inter)$Population,xax=1,yax=2,
         sub = "neutralintergenicloci,PC 1-2, 40 individuals, 2,019 loci, maf 5.0%",
         possub = "topleft",col=transp(col6,.6),axesell=FALSE,
         cstarc=0, cpoint=3, grid=FALSE, cellipse = 0)
s.class(pca_inter$li, strata(u.na_inter)$Population,xax=1,yax=3,
         sub = "neutralintergenicloci,PC 1-3, 40 individuals, 2,019 loci, maf 5.0%",
         possub = "topleft",col=transp(col6,.6),axesell=FALSE,
         cstarc=0, cpoint=3, grid=FALSE, cellipse = 0)
s.class(pca_inter$li, strata(u.na_inter)$Population,xax=1,yax=4,
         sub = "neutralintergenicloci,PC 1-4, 40 individuals, 2,019 loci, maf 5.0%",
         possub = "topleft",col=transp(col6,.6),axesell=FALSE,
         cstarc=0, cpoint=3, grid=FALSE, cellipse = 0)
s.class(pca_inter$li, strata(u.na_inter)$Population,xax=1,yax=5,
         sub = "neutralintergenicloci,PC 1-5, 40 individuals, 2,019 loci, maf 5.0%",
         possub = "topleft",col=transp(col6,.6),axesell=FALSE,
         cstarc=0, cpoint=3, grid=FALSE, cellipse = 0)

```



## Neutral Intron SNPs

---

### Making files

#### Make genind object

`neutralintron.recode.vcf` contains neutral SNPs from intron regions for populations PVD, GB, BIS, and NAR. Steps for generating this VCF file are located in `EecSeq_Cvirginica_NeutralLoci.md`.

Population NIN was removed from the VCF file following steps in `EecSeq_Cvirginica_OutlierDetection.md`.

`strata` contains population, environmental, and library information for each sample - can be accessed here.

```
my_vcf_intron <- read.vcfR("neutralintron.recode.vcf")

## Scanning file to determine attributes.
## File attributes:
##   meta lines: 63
##   header_line: 64
##   variant count: 16701
##   column count: 49
##
Meta line 63 read in.
## All meta lines processed.
## gt matrix initialized.
## Character matrix gt created.
##   Character matrix gt rows: 16701
##   Character matrix gt cols: 49
##   skip: 0
##   nrows: 16701
##   row_num: 0
##
Processed variant 1000
Processed variant 2000
Processed variant 3000
Processed variant 4000
Processed variant 5000
Processed variant 6000
Processed variant 7000
Processed variant 8000
Processed variant 9000
Processed variant 10000
Processed variant 11000
Processed variant 12000
Processed variant 13000
Processed variant 14000
Processed variant 15000
Processed variant 16000
Processed variant: 16701
## All variants processed

strata <- read.table("strata", header=TRUE)
```

```

rad.filt_introns <- vcfR2genind(my_vcf_introns, strata = strata, pop = c(rep("BIS", 10), rep("GB", 10), rep("NAR", 10), rep("PWD", 10)))
rad.filt_introns

## /// GENIND OBJECT ///////////
##
## // 40 individuals; 16,701 loci; 33,351 alleles; size: 15.2 Mb
##
## // Basic content
## @tab: 40 x 33351 matrix of allele counts
## @loc.n.all: number of alleles per locus (range: 1-2)
## @loc.fac: locus factor for the 33351 columns of @tab
## @all.names: list of allele names for each locus
## @ploidy: ploidy of each individual (range: 2-2)
## @type: codom
## @call: adegenet::df2genind(X = t(x), sep = sep, pop = .2, strata = .1)
##
## // Optional content
## @pop: population of each individual (group size range: 10-10)
## @strata: a data frame with 12 columns ( Individual, Population, Latitude, Longitude, Distance, SE
#Providing population names for plotting
pop_order <- c("BIS", "GB", "NAR", "PWD")

```

Read in the other info from .strata file and extract information such as locality, latitude, and longitude.

```

info <- as.data.frame(read.table("strata", header = T, sep = "\t", stringsAsFactors = F))
mystrats_introns <- as.data.frame(matrix(nrow = length(indNames(rad.filt_introns)), ncol=10))
colnames(mystrats_introns) <- c("Population", "Latitude", "Longitude", "Distance", "Temperature", "SE", "Salinity", "pH", "Chlorophy")
just.strats <- select(info, c("Population"))
stratified.filt_introns <- strata(rad.filt_introns, formula = Population, combine = TRUE, just.strats)
stratified.filt_introns@other <- select(info, Latitude, Longitude, Distance, SE, Temperature, Salinity, pH, Chlorophy)

stratified.filt_introns

## /// GENIND OBJECT ///////////
##
## // 40 individuals; 16,701 loci; 33,351 alleles; size: 15.2 Mb
##
## // Basic content
## @tab: 40 x 33351 matrix of allele counts
## @loc.n.all: number of alleles per locus (range: 1-2)
## @loc.fac: locus factor for the 33351 columns of @tab
## @all.names: list of allele names for each locus
## @ploidy: ploidy of each individual (range: 2-2)
## @type: codom
## @call: adegenet::df2genind(X = t(x), sep = sep, pop = .2, strata = .1)
##
## // Optional content
## @pop: population of each individual (group size range: 10-10)
## @strata: a data frame with 1 columns ( Population )
## @other: a list containing: Latitude Longitude Distance SE Temperature Salinity pH Chlorophy

```

#### Repeat for quasi-independent set of SNPs

Steps for filtering are documented for OutFLANK in EecSeq\_Cvirginica\_OutlierDetection.md. However, I

am repeating the steps here.

```
my_vcf_intron <- read.vcfR("neutralintron.recode.vcf")

## Scanning file to determine attributes.
## File attributes:
##   meta lines: 63
##   header_line: 64
##   variant count: 16701
##   column count: 49
##
## Meta line 63 read in.
## All meta lines processed.
## gt matrix initialized.
## Character matrix gt created.
##   Character matrix gt rows: 16701
##   Character matrix gt cols: 49
##   skip: 0
##   nrows: 16701
##   row_num: 0
##
## Processed variant 1000
## Processed variant 2000
## Processed variant 3000
## Processed variant 4000
## Processed variant 5000
## Processed variant 6000
## Processed variant 7000
## Processed variant 8000
## Processed variant 9000
## Processed variant 10000
## Processed variant 11000
## Processed variant 12000
## Processed variant 13000
## Processed variant 14000
## Processed variant 15000
## Processed variant 16000
## Processed variant: 16701
## All variants processed

geno_intron <- extract.gt(my_vcf_intron) # Character matrix containing the genotypes
position_intron <- getPOS(my_vcf_intron) # Positions in bp
chromosome_intron <- getCHROM(my_vcf_intron) # Chromosome information

G_intron <- matrix(NA, nrow = nrow(geno_intron), ncol = ncol(geno_intron))

G_intron[geno_intron %in% c("0/0", "0|0")] <- 0
G_intron[geno_intron %in% c("0/1", "1/0", "1|0", "0|1")] <- 1
G_intron[geno_intron %in% c("1/1", "1|1")] <- 2

# NA should be replaced with "9"
G_intron[is.na(G_intron)] <- 9
```

Chromosomes need to be of class integer for this to work.

```

# Visualizing chromosomes
chrom_unique_intron <- unique(chromosome_intron)
print(chrom_unique_intron)

## [1] "NC_035780.1" "NC_035781.1" "NC_035782.1" "NC_035783.1" "NC_035784.1"
## [6] "NC_035785.1" "NC_035786.1" "NC_035787.1" "NC_035788.1" "NC_035789.1"

# Removing "NC_" from the chromosome name so it can be converted to an integer
chrom_new_intron <- chromosome_intron %>% str_replace("NC_","")

# Converting the character string into an integer
chrom1_intron <- as.integer(chrom_new_intron)

```

chromosome and position need to be sorted for imputation to work.

```

chrom_sort_intron <- sort(chrom1_intron)
pos_sort_intron <- sort(position_intron)

```

*Note: This filtering program does not allow for missing genotype values.*

### Remove missing genotypes

```

# removing missing data
G_intron_miss <- matrix(NA, nrow = nrow(geno_intron), ncol = ncol(geno_intron))

G_intron_miss[geno_intron %in% c("0/0", "0|0")] <- 0
G_intron_miss[geno_intron %in% c("0/1", "1/0", "1|0", "0|1")] <- 1
G_intron_miss[geno_intron %in% c("1/1", "1|1")] <- 2

# removing missing data
G_intron_miss[na.omit(G_intron_miss)]

##      [1] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1
##     [37] 0 0 0 0 1 1 1 1 1 1 1 1 0 0 0 0 0 0 0 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 0 1 1 1 1 1 1 1 1 1 1
##     [73] 1 1 1 0 1 1 0 1 1 1 1 0 0 1 1 1 1 1 1 1 0 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 0
##    [109] 0 0 0 0 1 0 0 0 1 1 1 0 0 0 1 1 1 1 1 1 0 0 1 0 1 1 1 1 1 0 1 0 1 1 1 1 1 0 1 0 1 0 1 0 0
##    [145] 0 0 0 0 1 0 0 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 0 1 0 1 1 1 1 1 0 1 1 1 1 0 1 1
##    [181] 0 1 1 0 0 0 0 0 0 0 0 0 1 0 1 0 1 1 1 0 1 0 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0
##    [217] 0 1 1 1 1 1 1 1 0 0 0 1 0 0 0 1 1 1 1 1 1 0 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##    [253] 1 1 1 1 1 0 1 0 0 1 0 0 1 1 0 1 0 1 1 1 0 0 0 1 0 0 0 1 1 1 1 1 1 0 0 1 0 1 0 1 1 1
##    [289] 1 0 1 1 0 1 1 0 0 1 1 0 0 1 1 1 0 0 1 1 0 0 1 1 1 1 1 1 0 0 1 1 1 1 1 1 1 1 1 1 1 1
##    [325] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 0 0 1 1 1 1 0 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##    [361] 1 1 1 1 1 1 1 1 1 1 1 1 0 0 1 1 1 0 1 1 1 0 0 0 0 1 0 0 1 1 1 0 1 0 1 0 0 1 0 0 1
##    [397] 1 1 0 0 0 1 1 1 0 0 1 0 0 0 0 0 1 1 1 1 0 0 1 1 0 1 1 1 1 1 1 0 0 1 1 1 1 1 1 1 1 1
##    [433] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 0 0 0 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##    [469] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 1 1 1 0 1 1 1 1 1 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0
##    [505] 0 1 0 0 0 0 0 0 1 1 0 0 0 1 1 1 1 0 1 1 1 1 1 1 1 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##    [541] 0 1 1 1 1 0 0 0 1 1 0 0 0 0 0 0 0 1 1 1 0 0 0 0 0 1 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1
##    [577] 1 0 1 1 1 1 1 0 0 1 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1
##    [613] 1 1 1 1 1 0 1 1 1 0 1 1 1 1 1 0 1 1 0 0 0 1 1 1 1 1 1 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1
##    [649] 1 1 1 1 1 1 1 1 1 1 0 1 0 0 1 1 1 1 1 1 1 0 1 1 0 1 0 1 1 0 0 1 0 1 0 1 0 1 0 1 0 1
##    [685] 1 1 1 1 1 1 1 1 1 1 0 1 1 0 1 1 0 1 1 0 1 1 0 0 1 1 1 1 1 1 1 0 1 1 1 1 1 1 1 1 1 1
##    [721] 0 0 0 0 0 0 1 0 0 0 0 0 1 1 1 0 0 1 1 1 1 1 1 1 0 1 1 0 1 0 1 1 1 1 1 1 1 1 1 1 1 1
##    [757] 0 1 1 1 1 1 0 1 0 0 1 1 1 1 1 1 1 1 1 0 0 0 0 1 0 0 1 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1

```

```

## [99937] 1 1 0 1 1 1 1 1 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 1 1 1 0 0
## [99973] 0 1 1 1 0 1 1 1 0 1 0 1 1 1 1 1 1 0 1 1 1 1 1 1 1 0 1
## [ reached getOption("max.print") -- omitted 17772 entries ]
# Converting the genotype matrix to class FBM.code256 using the matrix where missing data was removed
G1_intron_miss <- add_code256(big_copy(t(G_intron_miss),type="raw"),code=bigsnpr:::CODE_012)

## Warning in replaceMat(x$address_rw, i, j, value): At least one value changed (nan -> 0)
## while converting from R type 'double' to C type 'unsigned char (raw)'.

newpc_intron_miss <- snp_autoSVD(G1_intron_miss,infos.chr =chrom_sort_intron,infos.pos = pos_sort_intron

##
## Phase of clumping (on MAF) at r^2 > 0.2.. keep 6212 SNPs.
## Discarding 2446 variants with MAC < 10.
##
## Iteration 1:
## Computing SVD..
## 0 outlier variant detected..
##
## Converged!

# maximum roll.size I could use is 32
which_pruned_intron_miss <- attr(newpc_intron_miss, which="subset") # Indexes of remaining SNPs after pruning
length(which_pruned_intron_miss)

## [1] 3766

invisible(lapply(which_pruned_intron_miss, write, "pruned_data_intron.txt", append=TRUE))

```

### In terminal

```

$ mawk '!/#/' neutralintron.recode.vcf | cut -f1,2 > totallociintron
$ NUM=(`cat totallociintron | wc -l`)
$ paste <(seq 1 $NUM) totallociintron > lociintron.plus.index
$ cat pruned_data_intron.txt | parallel "grep -w {} lociintron.plus.index" | cut -f2,3> pruned_data_intron

$ head pruned_data_intron.loci.txt

output:
NC_035780.1 574321
NC_035780.1 574444
NC_035780.1 574657
NC_035780.1 612658
NC_035780.1 618706
NC_035780.1 722826
NC_035780.1 771520
NC_035780.1 771800
NC_035780.1 771812
NC_035780.1 780771

Create VCF file with just the pruned_data loci

$ vcftools --vcf neutralintron.recode.vcf --recode --recode-INFO-all --positions pruned_data_intron.loci.txt

output:
After filtering, kept 40 out of 40 Individuals
Outputting VCF file...

```

```

After filtering, kept 3766 out of a possible 16701 Sites
Run Time = 1.00 seconds
my_vcf_u_introns <- read.vcfR("pruned_data_introns.recode.vcf")

## Scanning file to determine attributes.
## File attributes:
##   meta lines: 63
##   header_line: 64
##   variant count: 3766
##   column count: 49
##
Meta line 63 read in.
## All meta lines processed.
## gt matrix initialized.
## Character matrix gt created.
##   Character matrix gt rows: 3766
##   Character matrix gt cols: 49
##   skip: 0
##   nrows: 3766
##   row_num: 0
##
Processed variant 1000
Processed variant 2000
Processed variant 3000
Processed variant: 3766
## All variants processed
strata <- read.table("strata", header=TRUE)

rad.u_introns <- vcfR2genind(my_vcf_u_introns, strata = strata, pop = c(rep("BIS", 10),rep("GB", 10),rep("CHL", 10),rep("SE", 10),rep("T", 10),rep("L", 10),rep("D", 10),rep("P", 10),rep("O", 10),rep("S", 10),rep("C", 10)), n.individuals = 40, n.alleles = 7532, size = 3.5, pop.names = strata$Population, locus.factor = 1, call.function = adegenet::df2genind, sep = "\t", strata.name = "strata", just.strats = TRUE)
rad.u_introns

## /// GENIND OBJECT ///////////
##
## // 40 individuals; 3,766 loci; 7,532 alleles; size: 3.5 Mb
##
## // Basic content
##   @tab: 40 x 7532 matrix of allele counts
##   @loc.n.all: number of alleles per locus (range: 2-2)
##   @loc.fac: locus factor for the 7532 columns of @tab
##   @call.names: list of allele names for each locus
##   @ploidy: ploidy of each individual (range: 2-2)
##   @type: codom
##   @call: adegenet::df2genind(X = t(x), sep = sep, pop = ..2, strata = ..1)
##
## // Optional content
##   @pop: population of each individual (group size range: 10-10)
##   @strata: a data frame with 12 columns ( Individual, Population, Latitude, Longitude, Distance, SE, Temperature, Salinity, pH, Chlorophyll, Depth )
stratified.u_introns <- strata(rad.u_introns, formula= Population, combine = TRUE, just.strats)
stratified.u_introns@other <- select(info, Latitude, Longitude, Distance, SE, Temperature, Salinity, pH, Chlorophyll, Depth)

stratified.u_introns

## /// GENIND OBJECT ///////////

```

```

## // 40 individuals; 3,766 loci; 7,532 alleles; size: 3.4 Mb
##
## // Basic content
##   @tab: 40 x 7532 matrix of allele counts
##   @loc.n.all: number of alleles per locus (range: 2-2)
##   @loc.fac: locus factor for the 7532 columns of @tab
##   @all.names: list of allele names for each locus
##   @ploidy: ploidy of each individual (range: 2-2)
##   @type: codom
##   @call: adegenet::df2genind(X = t(x), sep = sep, pop = ..2, strata = ..1)
##
## // Optional content
##   @pop: population of each individual (group size range: 10-10)
##   @strata: a data frame with 1 columns ( Population )
##   @other: a list containing: Latitude Longitude Distance SE Temperature Salinity pH Chlorophy

Make hierfstat object

hf.filt_introns <- genind2hierfstat(rad.filt_introns, pop = c(rep("BIS", 10), rep("GB", 10), rep("NAR", 10))

hf.u_introns <- genind2hierfstat(rad.u_introns, pop = c(rep("BIS", 10), rep("GB", 10), rep("NAR", 10), rep(
hf.u_introns <- hf.u_introns$hierfstat.no.imputation

```

## Estimating effective migration surfaces (EEMS)

### 1. datapath.diffs

Steps for generating `datapath.diffs` are completed in RStudio.

```

# V1 method to get .diffs matrix, preferred
bed2diffs_v1 <- function(Geno) {
  nIndiv <- nrow(Geno)
  nSites <- ncol(Geno)
  Diffs <- matrix(0, nIndiv, nIndiv)

  for (i in seq(nIndiv - 1)) {
    for (j in seq(i + 1, nIndiv)) {
      x <- Geno[i, ]
      y <- Geno[j, ]
      Diffs[i, j] <- mean((x - y)^2, na.rm = TRUE)
      Diffs[j, i] <- Diffs[i, j]
    }
  }
  Diffs
}

# V2 method to get .diffs matrix, only if V1 doesn't work
bed2diffs_v2 <- function(Geno) {
  nIndiv <- nrow(Geno)
  nSites <- ncol(Geno)
  Miss <- is.na(Geno)
  ## Impute NAs with the column means (= twice the allele frequencies)
  Mean <- matrix(colMeans(Geno, na.rm = TRUE), ## a row of means
                  nrow = nIndiv, ncol = nSites, byrow = TRUE) ## a matrix with nIndiv identical rows of

```

```

Mean[Miss == 0] <- 0 ## Set the means that correspond to observed genotypes to 0
Geno[Miss == 1] <- 0 ## Set the missing genotypes to 0 (used to be NA)
Geno <- Geno + Mean
## Compute similarities
Sim <- Geno %*% t(Geno) / nSites
SelfSim <- diag(Sim) ## self-similarities
vector1s <- rep(1, nIndiv) ## vector of 1s
## This chunk generates a `diffs` matrix
Diffs <- SelfSim %*% t(vector1s) + vector1s %*% t(SelfSim) - 2 * Sim
Diffs
}

geno_intron <- stratted.filt_introntab

# Get rid of non-biallelic loci
multi.loci_intron <- names(which(stratted.filt_introntab[, 2] != 2))
multi.cols_intron <- which(grep(paste0("^", multi.loci_intron, "\\\.\d+$"), collapse = "|"), colnames(geno_intron))
if (length(multi.cols_intron)) geno_intron <- geno_intron[, -multi.cols_intron]
nloci_intron <- dim(geno_intron)[2] / 2
dim(geno_intron)

## [1] 40 33300
stopifnot(identical(stratted.filt_introntab$type, 'codom'))

# bed2diffs functions
diffs.v1_intron <- bed2diffs_v1(geno_intron)
diffs.v2_intron <- bed2diffs_v2(geno_intron)
# Round to 6 digits
diffs.v1_intron <- round(diffs.v1_intron, digits = 6)
diffs.v2_intron <- round(diffs.v2_intron, digits = 6)

```

Check that the dissimilarity matrix has one positive eigenvalue and nIndiv-1 negative eigenvalues, as required by a full-rank Euclidean distance matrix. If the V1 method does not make a Euclidean matrix, you must use V2.

```

tail(sort(round(eigen(diffs.v1_intron)$values, digits = 2)))

## [1] -0.46 -0.45 -0.44 -0.44 -0.43 21.23
tail(sort(round(eigen(diffs.v2_intron)$values, digits = 2)))

## [1] -0.45 -0.44 -0.43 -0.43 -0.43 20.88
# Set suffix for EEMS input files
suf_intron <- "neutralintrondata-filt"

# This saves the file to directory
write.table(diffs.v1_intron, paste(suf_intron, ".v1.diffs", sep = ""),
            col.names = FALSE, row.names = FALSE, quote = FALSE)

```

## 2. datapath.coord

datapath.coord are the sample coordinates, two coordinates per sample, one sample per line. The sampling locations should be given in the same order as the rows and columns of the dissimilarity matrix.

Steps for generating datapath.coord are completed in RStudio.

```

## Get gps coordinates from previously created info matrix
x0R.info <- dplyr::filter(info)
gps_matrix <- select(x0R.info,c("Longitude","Latitude"))

#write .coord file
write.table(gps_matrix, paste(suf_intron,".v1.coord",sep=""),col.names = FALSE, row.names = FALSE,quote=FALSE)

```

### 3. datapath.outer

`datapath.outer` are the habitat coordinates, as a sequence of vertices that form a closed polygon. The habitat vertices should be listed counterclockwise and the first vertex should also be the last vertex, so that the outline is a closed ring. Otherwise, EEMS attempts to “correct” the polygon and prints a warning message.

`datapath.outer` is created manually in Excel,based on site coordinates gathered from Google Maps, copied into terminal using `nano`, and saved as the file with the appropriate extension.

`**runeems_snps` is then run in command-line following the steps documented in `NB_EEMS_OutlierHap.md`

Back in RStudio to plot `runeems_snps` outputs

```

# Install rEEMSpots
library(rEEMSpots)

# Plotting EEMS after running runeems_snps
path = "./NB_EEMS_Neutral/"
dirs = c(paste0(path,"neutralintrondata-D200-chain1"), paste0(path,"neutralintrondata-D300-chain1"), paste0(path,"neutralintrondata-D600-chain1"))

eems.plots(mcmcpath = c(paste0(path,"./neutralintrondata-D200-chain1"), paste0(path,"neutralintrondata-D300-chain1"), paste0(path,"neutralintrondata-D600-chain1")),
           longlat = T,add.grid=F,add.outline = T,add.demes = T,
           projection.in = "+proj=longlat +datum=WGS84",projection.out = "+proj=merc +datum=WGS84",
           add.map = T,add.abline = T, add.r.squared = T)

## Input projection: +proj=longlat +datum=WGS84
## Output projection: +proj=merc +datum=WGS84

## Loading rgdal (required by projection.in)
## Loading rworldmap (required by add.map)
## Loading rworldxtra (required by add.map)

## Using the default DarkOrange to Blue color scheme, with 'white' as the midpoint color.
## It combines two color schemes from the 'dichromat' package, which itself is based on
## a collection of color schemes for scientific data graphics:
## Light A and Bartlein PJ (2004). The End of the Rainbow? Color Schemes for Improved Data
## Graphics. EOS Transactions of the American Geophysical Union, 85(40), 385.
## See also http://geog.uoregon.edu/datagraphics/color\_scales.htm

## Using 'euclidean' distance to assign interpolation points to Voronoi tiles.

## Processing the following EEMS output directories :
## ./NB_EEMS_Neutral./neutralintrondata-D200-chain1./NB_EEMS_Neutral/neutralintrondata-D300-chain1./NB_EEMS_Neutral/neutralintrondata-D600-chain1

## Plotting effective migration surface (posterior mean of m rates)
## ./NB_EEMS_Neutral./neutralintrondata-D200-chain1
## ./NB_EEMS_Neutral/neutralintrondata-D300-chain1
## ./NB_EEMS_Neutral/neutralintrondata-D600-chain1

```

```

## Using the default DarkOrange to Blue color scheme, with 'white' as the midpoint color.
## It combines two color schemes from the 'dichromat' package, which itself is based on
## a collection of color schemes for scientific data graphics:
## Light A and Bartlein PJ (2004). The End of the Rainbow? Color Schemes for Improved Data
## Graphics. EOS Transactions of the American Geophysical Union, 85(40), 385.
## See also http://geog.uoregon.edu/datagraphics/color_scales.htm

## Plotting effective diversity surface (posterior mean of q rates)

## ./NB_EEMS_Neutral/./neutralintrondata-D200-chain1
## ./NB_EEMS_Neutral/neutralintrondata-D300-chain1
## ./NB_EEMS_Neutral/neutralintrondata-D600-chain1

## Using the default DarkOrange to Blue color scheme, with 'white' as the midpoint color.
## It combines two color schemes from the 'dichromat' package, which itself is based on
## a collection of color schemes for scientific data graphics:
## Light A and Bartlein PJ (2004). The End of the Rainbow? Color Schemes for Improved Data
## Graphics. EOS Transactions of the American Geophysical Union, 85(40), 385.
## See also http://geog.uoregon.edu/datagraphics/color_scales.htm

## Plotting posterior probability trace

## ./NB_EEMS_Neutral/./neutralintrondata-D200-chain1
## ./NB_EEMS_Neutral/neutralintrondata-D300-chain1
## ./NB_EEMS_Neutral/neutralintrondata-D600-chain1

## Plotting average dissimilarities within and between demes

## ./NB_EEMS_Neutral/./neutralintrondata-D200-chain1
## ./NB_EEMS_Neutral/neutralintrondata-D300-chain1

## Using 'euclidean' distance to assign interpolation points to Voronoi tiles.
##
##
##
## Using 'euclidean' distance to assign interpolation points to Voronoi tiles.

## EEMS results for at least two different population grids

```

## Pairwise Fst

```

fst.mat_introns <- pairwise.WCfst(hf.filt_introns)

gindF.fst.mat.triN_introns <- as.matrix(fst.mat_introns)
colnames(gindF.fst.mat.triN_introns) <- pop_order
rownames(gindF.fst.mat.triN_introns) <- pop_order

meltedN_introns <- melt(gindF.fst.mat.triN_introns, na.rm =TRUE)
round(gindF.fst.mat.triN_introns,4)

##          BIS      GB      NAR      PVD
## BIS      NA -0.0011  0.0089 -0.0015
## GB     -0.0011      NA  0.0109 -0.0003
## NAR    0.0089   0.0109      NA  0.0159
## PVD   -0.0015  -0.0003  0.0159      NA

```

```

summary(meltedN_intron$value)

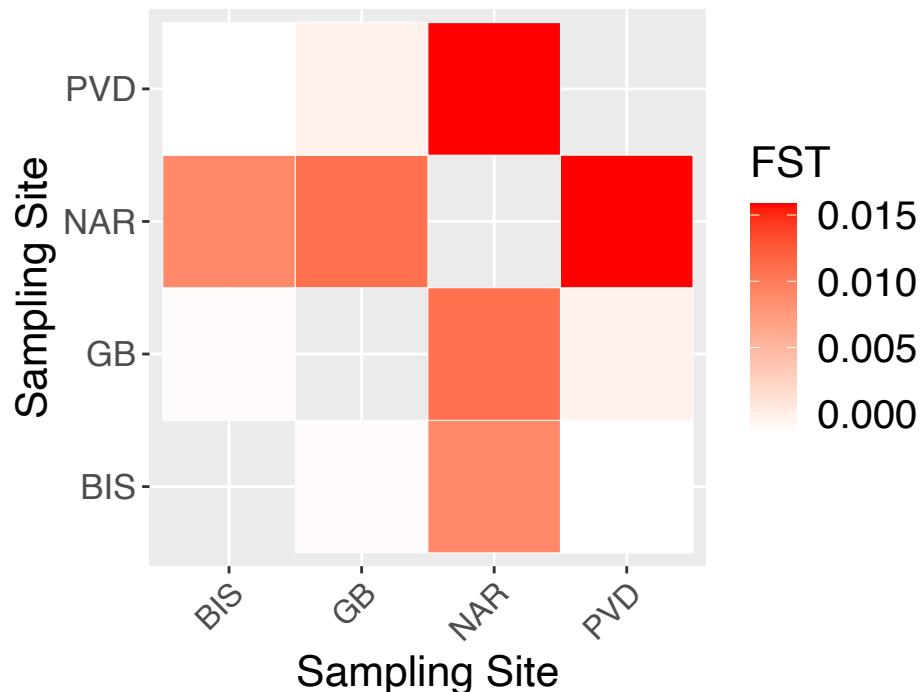
##      Min.    1st Qu.     Median      Mean    3rd Qu.     Max.
## -0.001470 -0.001132  0.004337  0.005479  0.010938  0.015864

#Plotting Pairwise fst
neutral_intron <- ggplot(data = meltedN_intron, aes(Var2, Var1, fill = value)) + geom_tile(color = "white")
  scale_fill_gradient(low = "white", high = "red", name="FST") +
  ggtitle(expression(atop("Pairwise FST, WC (1984) Neutral Intron SNPs", atop(italic("N = 40, L = 16,701")))))
  labs(x = "Sampling Site", y = "Sampling Site") +
  theme(axis.text.x = element_text(angle = 45, vjust = 1, size = 12, hjust = 1), axis.text.y = element_text(size = 12))
  theme(axis.title = element_text(size = 16), legend.text = element_text(size = 15), legend.title = element_text(size = 17)) +
  theme(plot.title = element_text(size = 17)) +
  coord_fixed()
neutral_intron

```

## Pairwise FST, WC (1984) Neutral Intron SNPs

$N = 40, L = 16,701$



## Genetic diversity (observed and expected heterozygosity)

```

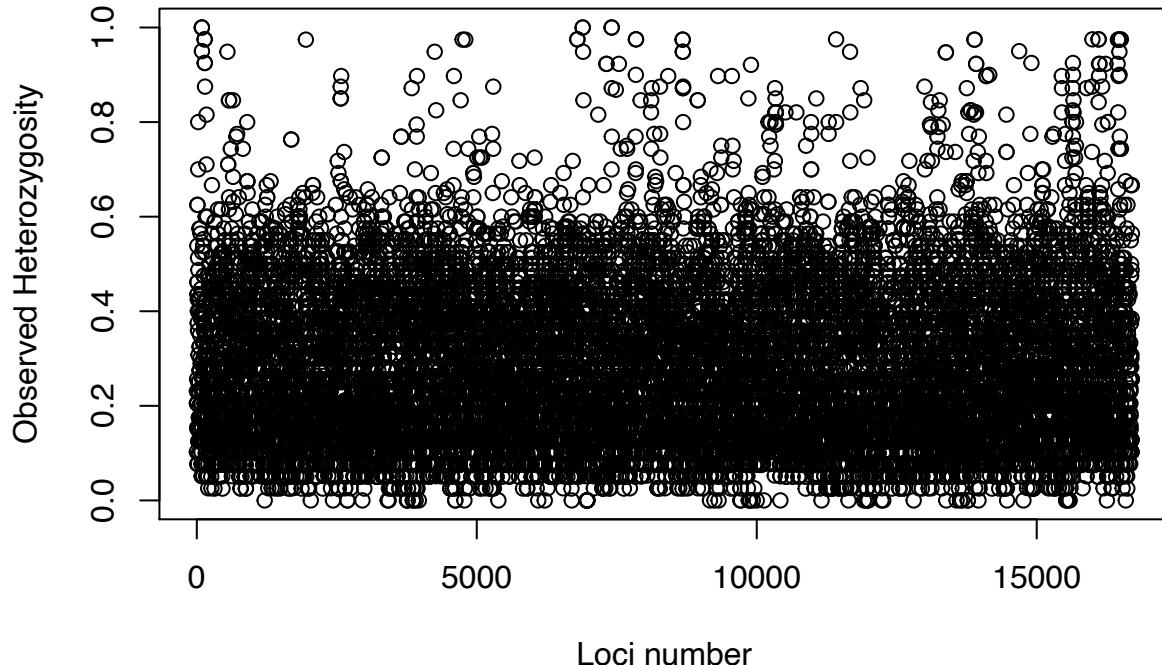
comb_intron <- summary(stratified.filt_intron)
names(comb_intron)

## [1] "n"          "n.by.pop"   "loc.n.all"  "pop.n.all"  "NA.perc"   "Hobs"
## [7] "Hexp"

plot(comb_intron$Hobs, xlab="Loci number", ylab="Observed Heterozygosity",
     main="Observed heterozygosity per locus")

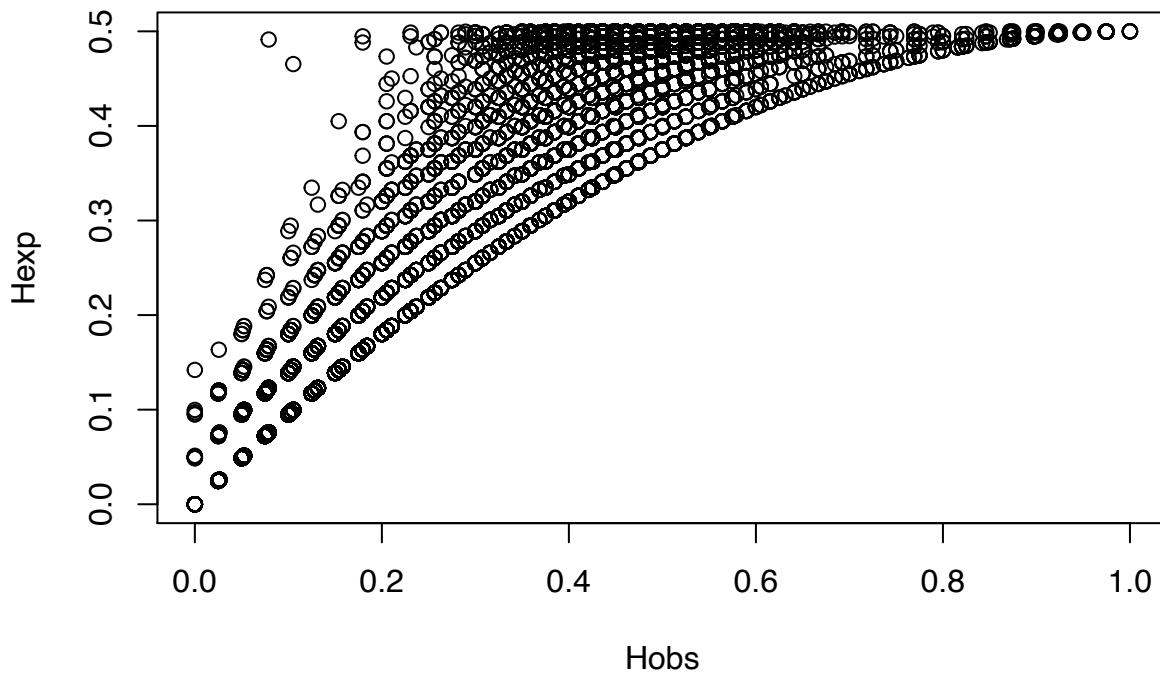
```

### Observed heterozygosity per locus



```
plot(comb_intron$Hobs, comb_intron$Hexp, xlab="Hobs", ylab="Hexp",
     main="Expected heterozygosity as a function of observed heterozygosity per locus")
```

## Expected heterozygosity as a function of observed heterozygosity per I



```
bartlett.test(list(comb_intron$Hexp, comb_intron$Hobs)) # a test : H0: Hexp = Hobs

##
##  Bartlett test of homogeneity of variances
##
## data: list(comb_intron$Hexp, comb_intron$Hobs)
## Bartlett's K-squared = 399.93, df = 1, p-value < 2.2e-16

Significant difference between Observed and expected heterozygosity.

basicstat_intron <- basic.stats(hf.filt_intron, diploid = TRUE, digits = 3)

as.data.frame(basicstat_intron$overall)

##      basicstat_intron$overall
## Ho          0.277
## Hs          0.273
## Ht          0.274
## Dst         0.001
## Htp         0.275
## Dstp        0.002
## Fst         0.004
## Fstp        0.005
## Fis        -0.015
## Dest        0.002

# get bootstrap confidence values for Fis
boot_intron <- boot.ppfis(hf.filt_intron, nboot = 1000)
```

```

boot5_introns <- boot.ppfis(hf.filt_introns, nboot = 1000, quant = 0.5)

# combine all pop statistics
colnames(basicstat_introns$Ho) <- pop_order
Ho_introns <- colMeans(basicstat_introns$Ho, na.rm = T)
He_introns <- colMeans(basicstat_introns$Hs, na.rm = T)
Fis_introns <- boot5_introns$fis.ci$ll
y_introns <- cbind(pop_order, Ho_introns, He_introns, Fis_introns, boot_introns$fis.ci, latitude, longitude, distance)
y_introns

##      pop_order Ho_introns He_introns Fis_introns      ll      hl latitude longitude
## BIS      BIS 0.2789041 0.2758841 -0.0109 -0.0169 -0.0054 41.545 -71.431
## GB       GB 0.2814673 0.2737831 -0.0281 -0.0343 -0.0216 41.654 -71.445
## NAR      NAR 0.2772171 0.2705006 -0.0247 -0.0315 -0.0186 41.505 -71.453
## PVD      PVD 0.2705747 0.2719077  0.0047 -0.0014  0.0112 41.816 -71.391
##      distance sewage temperature salinity pH Chlor_a DO
## BIS      4.76   8.82        23     30 7.9    4.9 8.2
## GB       0.47  14.60        24     28 7.4   18.8 5.7
## NAR      15.41   2.03        25     18 7.6    4.6 7.0
## PVD      1.49  59.86        23     25 7.4    8.1 4.9

summary(He_introns)

##      Min. 1st Qu. Median Mean 3rd Qu. Max.
##  0.2705 0.2716 0.2728 0.2730 0.2743 0.2759

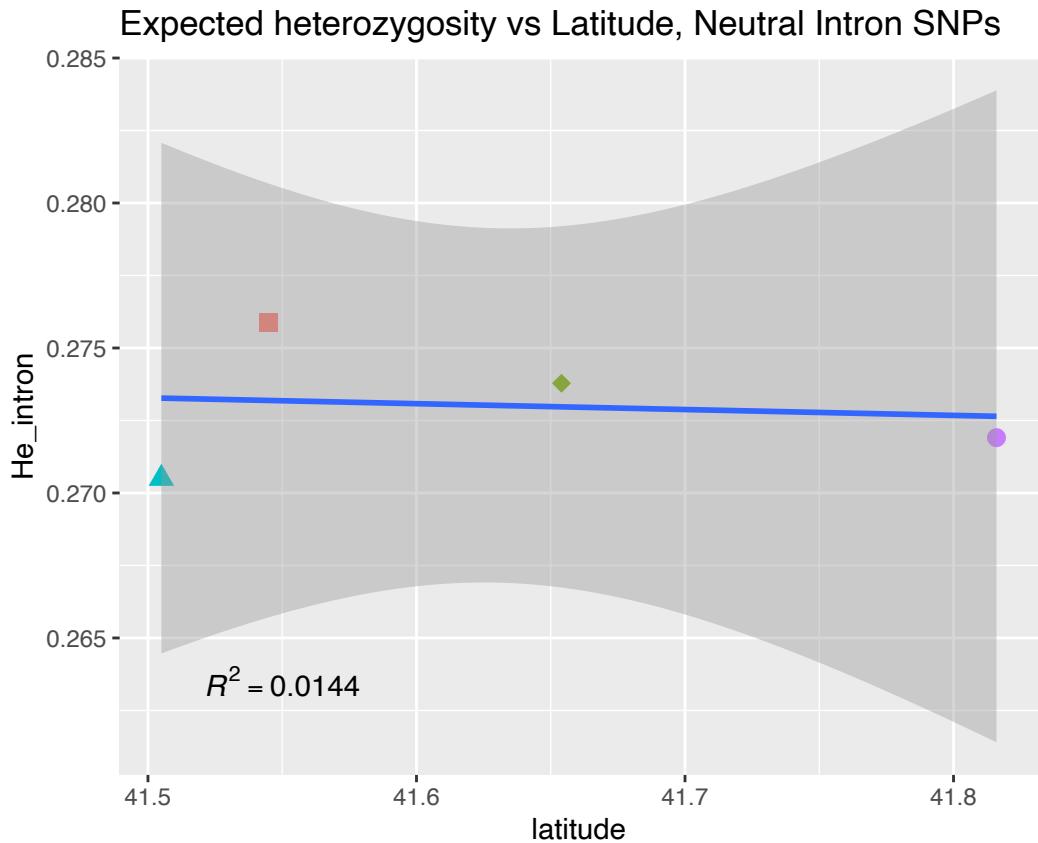
summary(Fis_introns)

##      Min. 1st Qu. Median Mean 3rd Qu. Max.
## -0.02810 -0.02555 -0.01780 -0.01475 -0.00700 0.00470

# Plot He vs Latitude
R2_introns = round(summary(lm(y_introns$He_introns ~ y_introns$latitude))$r.squared, 4)
ggplot(y_introns, aes(x = latitude, y = He_introns)) + geom_point(aes(shape=pop_order, color=pop_order), size=10) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggtitle("Expected heterozygosity vs Latitude, Neutral Intron SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R2_introns), x=41.55, y=0.2635, parse=T) +
  scale_x_continuous()

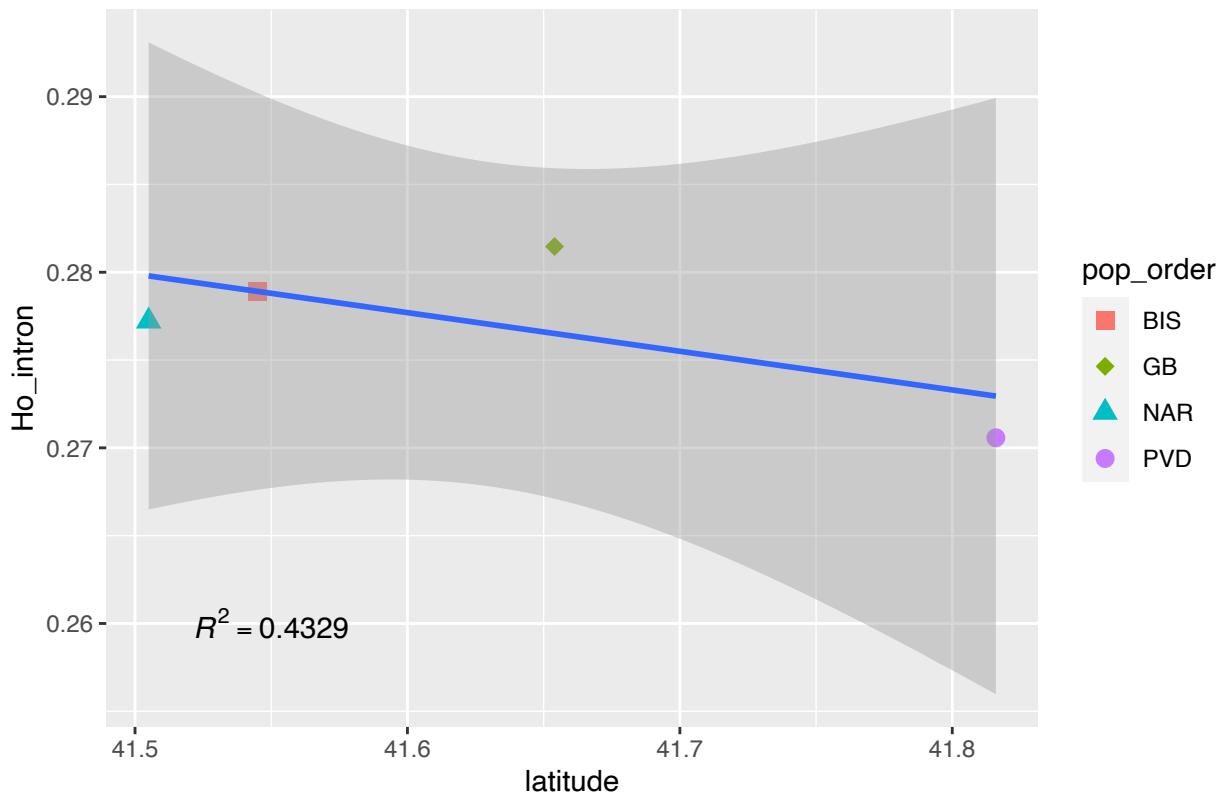
## `geom_smooth()` using formula 'y ~ x'

```



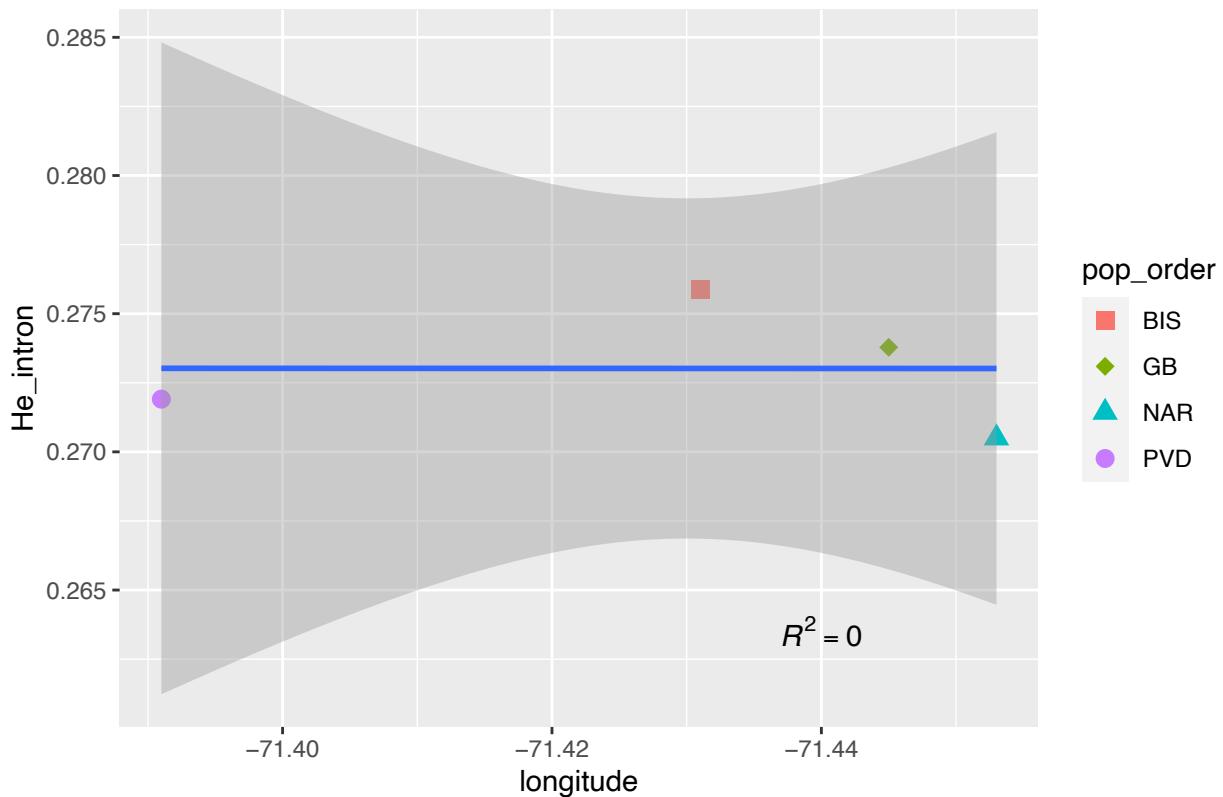
```
# Plot Ho vs Latitude
R2_intron = round(summary(lm(y_intron$Ho_intron ~ y_intron$latitude))$r.squared, 4)
ggplot(y_intron, aes(x = latitude, y = Ho_intron)) + geom_point(aes(shape=pop_order, color=pop_order),
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  gtitle("Observed heterozygosity vs Latitude, Neutral Intron SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R2_intron), x=41.55, y=0.26, parse=T) +
  scale_x_continuous()
## `geom_smooth()` using formula 'y ~ x'
```

### Observed heterozygosity vs Latitude, Neutral Intron SNPs



```
#Plot He vs Longitude
R3_intron = round(summary(lm(y_intron$He_intron ~ y_intron$longitude))$r.squared, 4)
ggplot(y_intron, aes(x = longitude, y = He_intron)) + geom_point(aes(shape=pop_order, color=pop_order),
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggttitle("Expected heterozygosity vs Longitude, Neutral Intron SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R3_intron), x=-71.44, y=0.2635, parse=T) +
  scale_x_reverse()
## `geom_smooth()` using formula 'y ~ x'
```

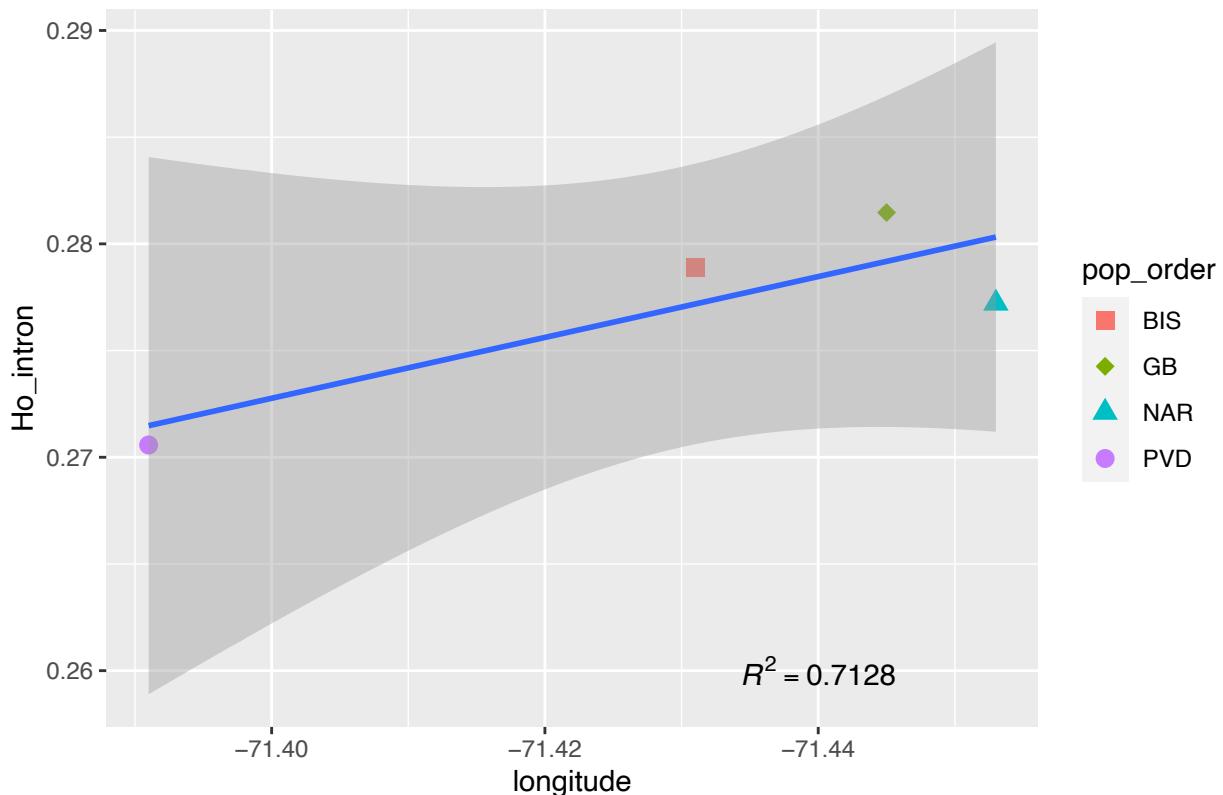
### Expected heterozygosity vs Longitude, Neutral Intron SNPs



```
#Plot Ho vs Longitude
R3_introns = round(summary(lm(y_introns$Ho_intron ~ y_introns$longitude))$r.squared, 4)
ggplot(y_introns, aes(x = longitude, y = Ho_intron)) + geom_point(aes(shape=pop_order, color=pop_order),
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  gtitle("Observed heterozygosity vs Longitude, Neutral Intron SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R3_introns), x=-71.44, y=0.26, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

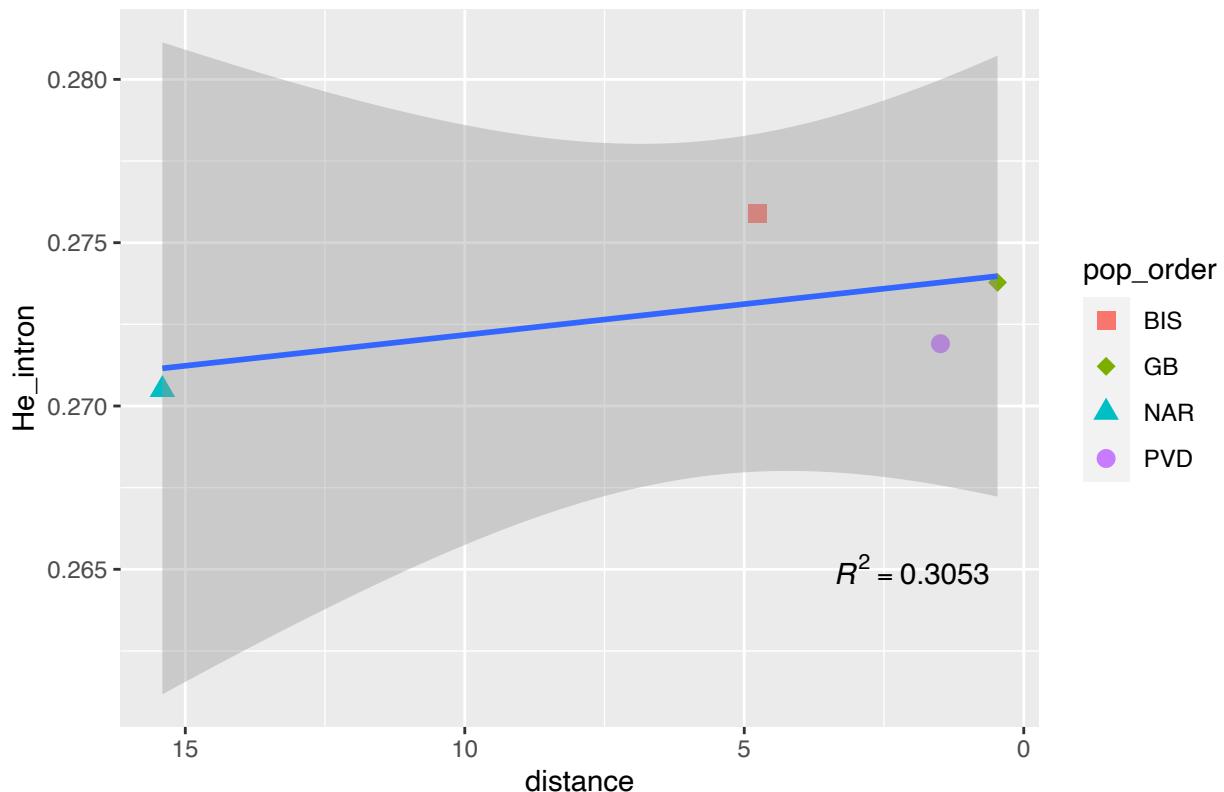
### Observed heterozygosity vs Longitude, Neutral Intron SNPs



```
#Plot He vs Distance from sewage outflow
R4_introns = round(summary(lm(y_introns$He_intron ~ y_introns$distance))$r.squared, 4)
ggplot(y_introns, aes(x = distance, y = He_intron)) + geom_point(aes(shape=pop_order, color=pop_order),
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggtitle("Expected heterozygosity vs Distance, Neutral Intron SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R4_introns), x=2, y=0.265, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

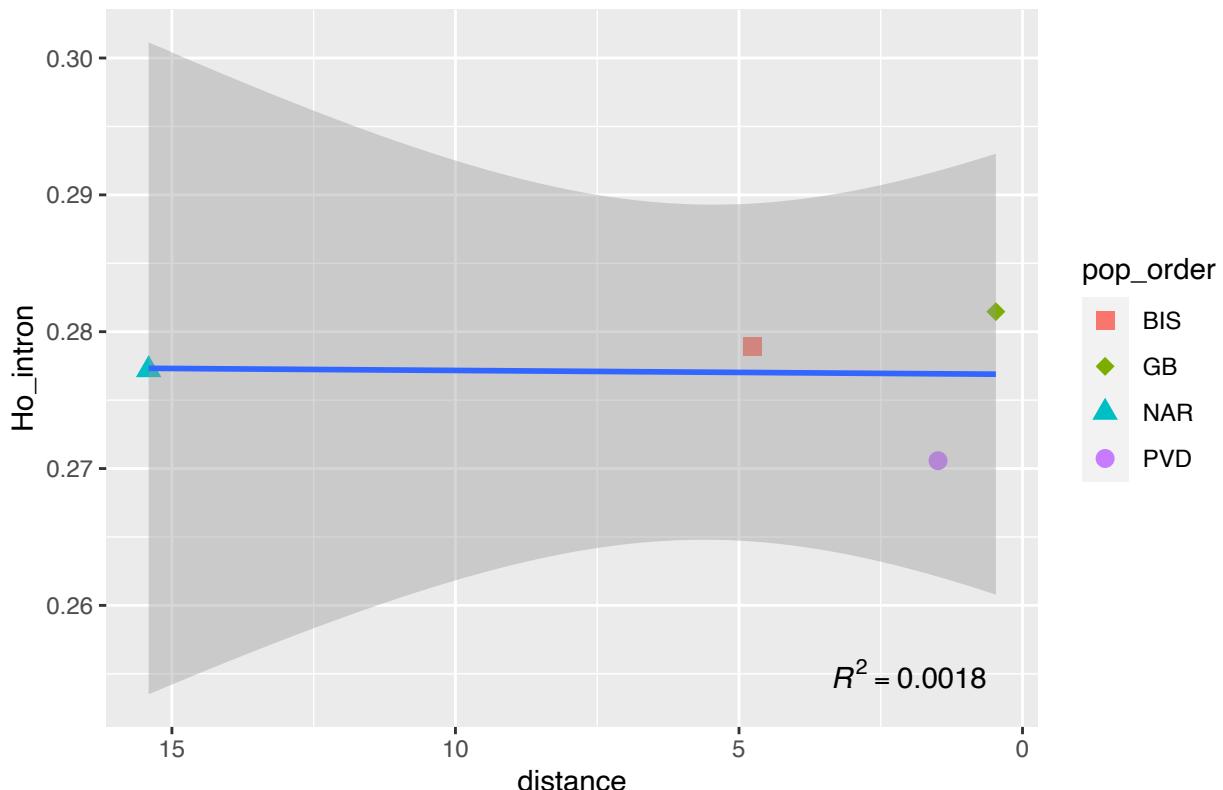
### Expected heterozygosity vs Distance, Neutral Intron SNPs



```
#Plot Ho vs Distance from sewage outflow
R4_introns = round(summary(lm(y_introns$Ho_intron ~ y_introns$distance))$r.squared, 4)
ggplot(y_introns, aes(x = distance, y = Ho_intron)) + geom_point(aes(shape=pop_order, color=pop_order),
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  gtitle("Observed heterozygosity vs Distance, Neutral Intron SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R4_introns), x=2, y=0.255, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

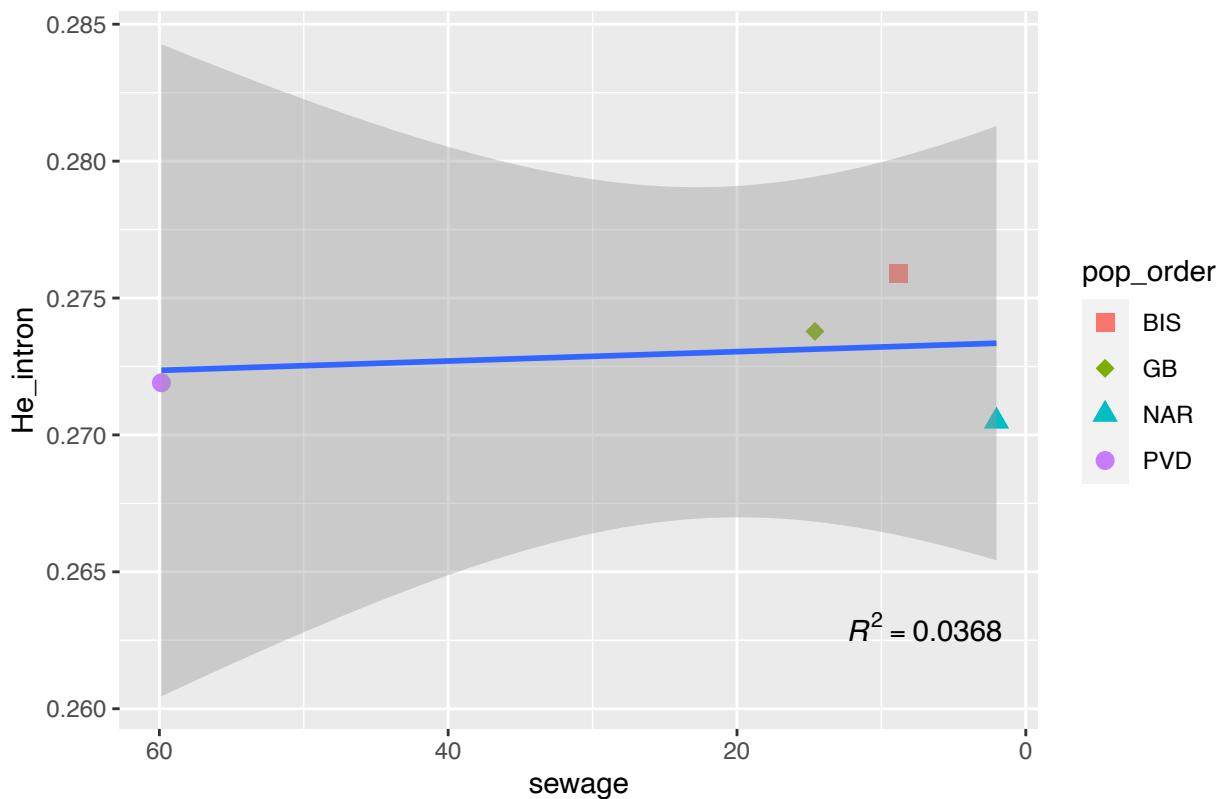
### Observed heterozygosity vs Distance, Neutral Intron SNPs



```
#Plot He vs Sewage effluent
R_intront = round(summary(lm(y_intront$He_intron ~ y_intront$sewage))$r.squared, 4)
ggplot(y_intront, aes(x = sewage, y = He_intron)) + geom_point(aes(shape=pop_order, color=pop_order), size=10) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggttitle("Expected heterozygosity vs Sewage Effluent, Neutral Intron SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R_intront), x=7, y=0.263, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

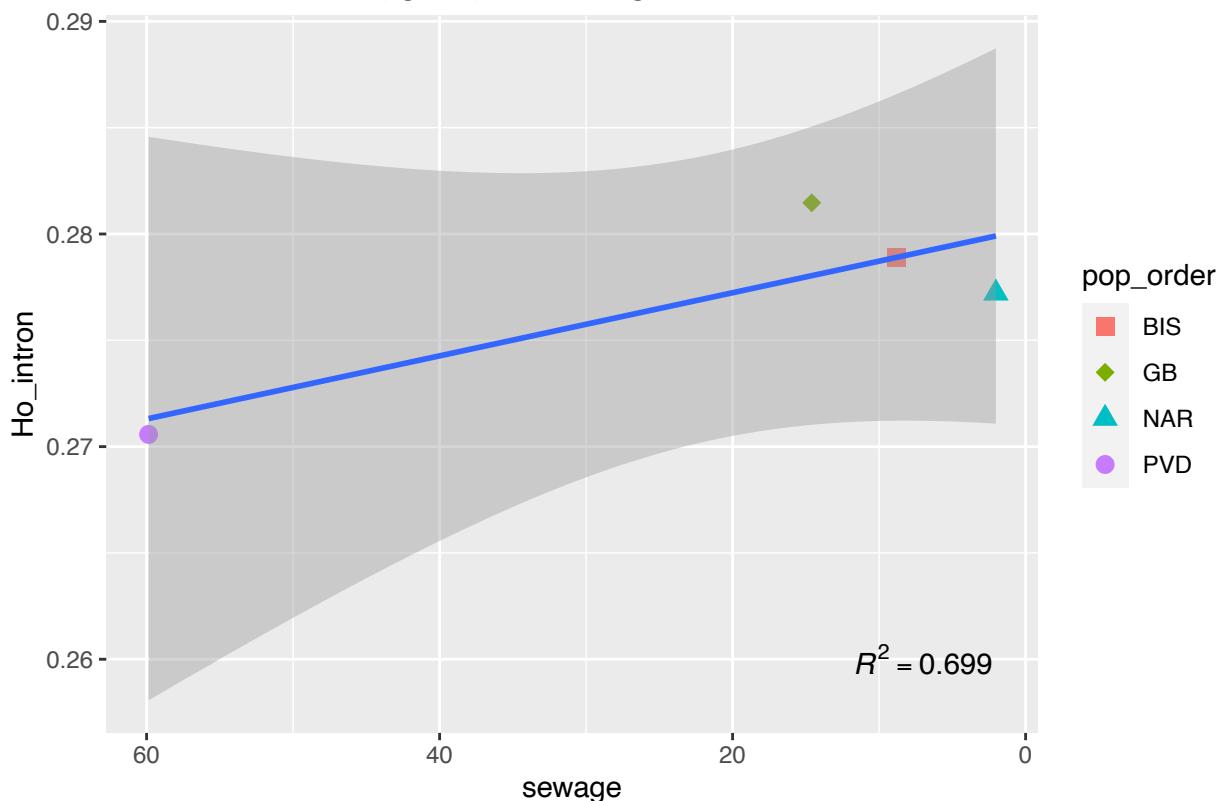
### Expected heterozygosity vs Sewage Effluent, Neutral Intron SNPs



```
#Plot Ho vs Sewage effluent
R_intront = round(summary(lm(y_intront$Ho_intron ~ y_intront$sewage))$r.squared, 4)
ggplot(y_intront, aes(x = sewage, y = Ho_intron)) + geom_point(aes(shape=pop_order, color=pop_order), size=10) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggttitle("Observed heterozygosity vs Sewage Effluent, Neutral Intron SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R_intront), x=7, y=0.26, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

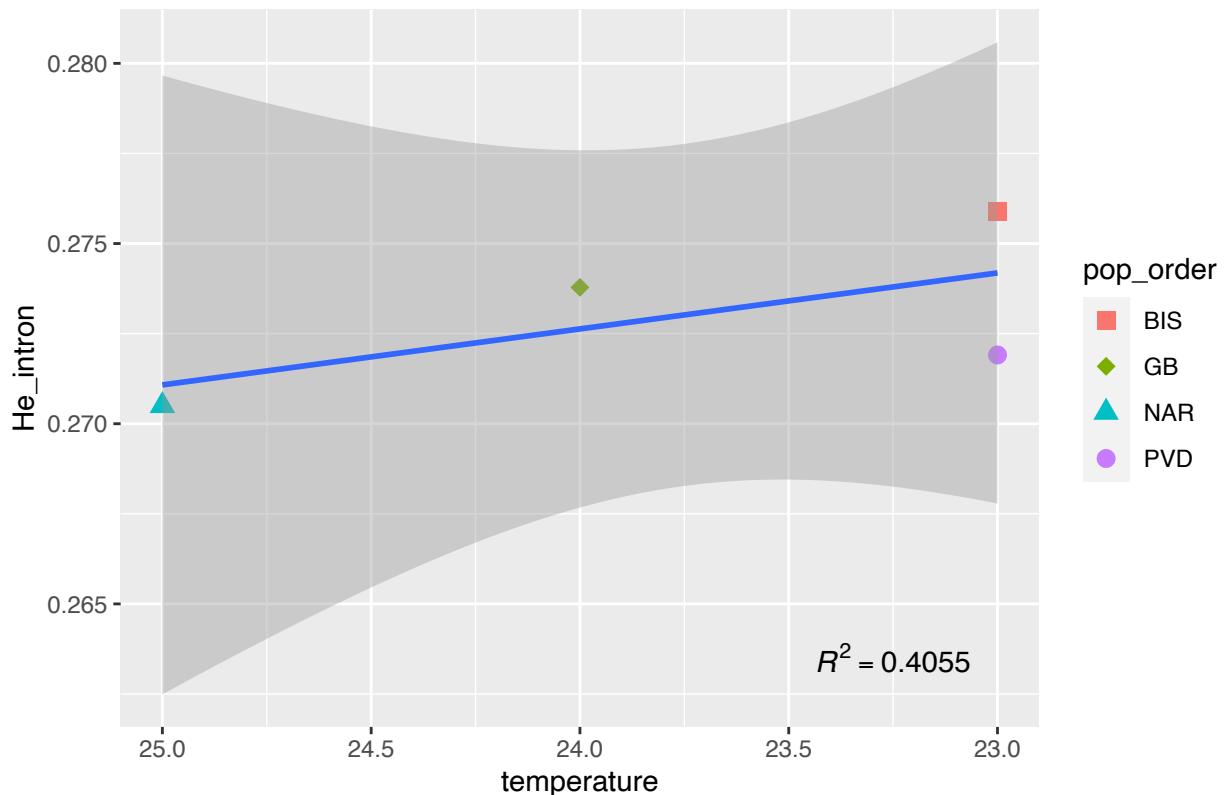
### Observed heterozygosity vs Sewage Effluent, Neutral Intron SNPs



```
#Plot He vs Temperature
R5_intron = round(summary(lm(y_intron$He_intron ~ y_intron$temperature))$r.squared, 4)
ggplot(y_intron, aes(x = temperature, y = He_intron)) + geom_point(aes(shape=pop_order, color=pop_order),
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggttitle("Expected heterozygosity vs Temperature, Neutral Intron SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R5_intron), x=23.25, y=0.2635, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

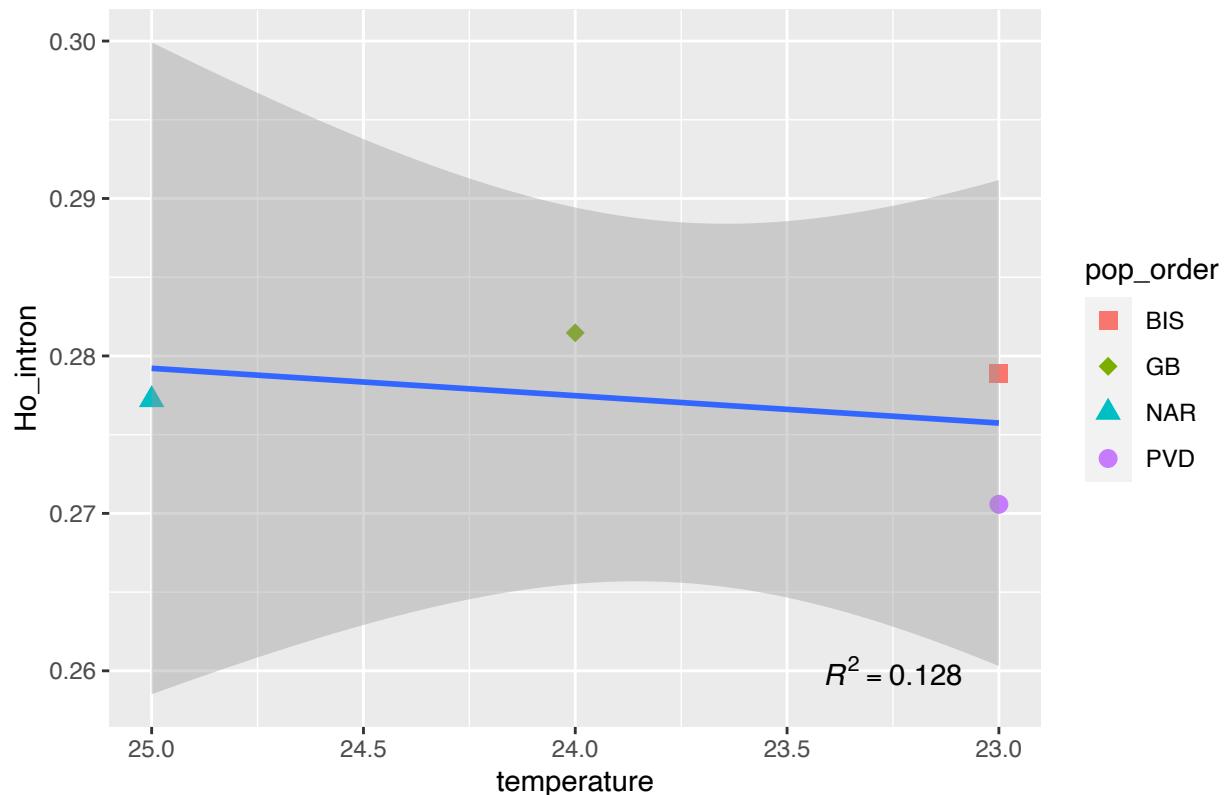
### Expected heterozygosity vs Temperature, Neutral Intron SNPs



```
#Plot Ho vs Temperature
R5_introns = round(summary(lm(y_introns$Ho_intron ~ y_introns$temperature))$r.squared, 4)
ggplot(y_introns, aes(x = temperature, y = Ho_intron)) + geom_point(aes(shape=pop_order, color=pop_order))
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  gtitle("Observed heterozygosity vs Temperature, Neutral Intron SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R5_introns), x=23.25, y=0.26, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

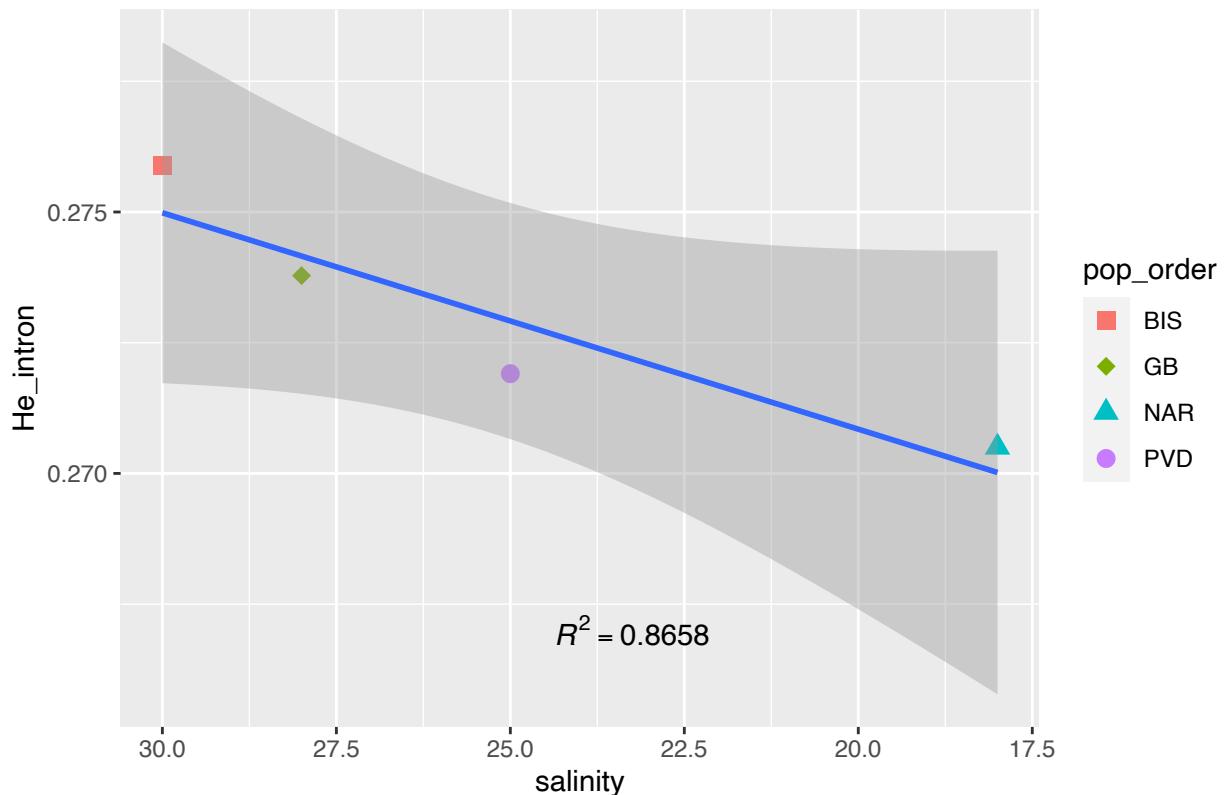
### Observed heterozygosity vs Temperature, Neutral Intron SNPs



```
#Plot He vs Salinity
R6_introns = round(summary(lm(He_intron ~ y_introns$salinity))$r.squared, 4)
ggplot(y_introns, aes(x = salinity, y = He_intron)) + geom_point(aes(shape=pop_order, color=pop_order),
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggtitle("Expected heterozygosity vs Salinity, Neutral Intron SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R6_introns), x=23.25, y=0.267, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

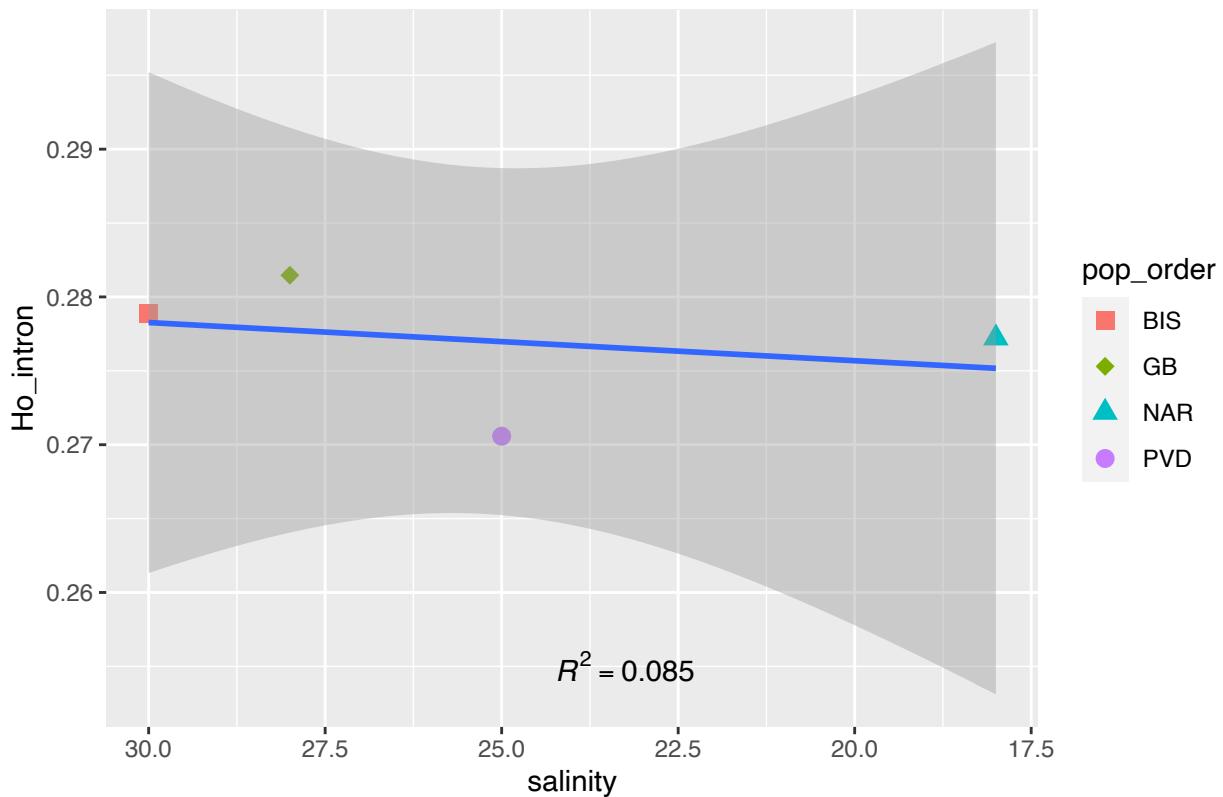
## Expected heterozygosity vs Salinity, Neutral Intron SNPs



```
#Plot Ho vs Salinity
R6_introns = round(summary(lm(y_introns$Ho_intron ~ y_introns$salinity))$r.squared, 4)
ggplot(y_introns, aes(x = salinity, y = Ho_intron)) + geom_point(aes(shape=pop_order, color=pop_order),
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  gtitle("Observed heterozygosity vs Salinity, Neutral Intron SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R6_introns), x=23.25, y=0.255, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

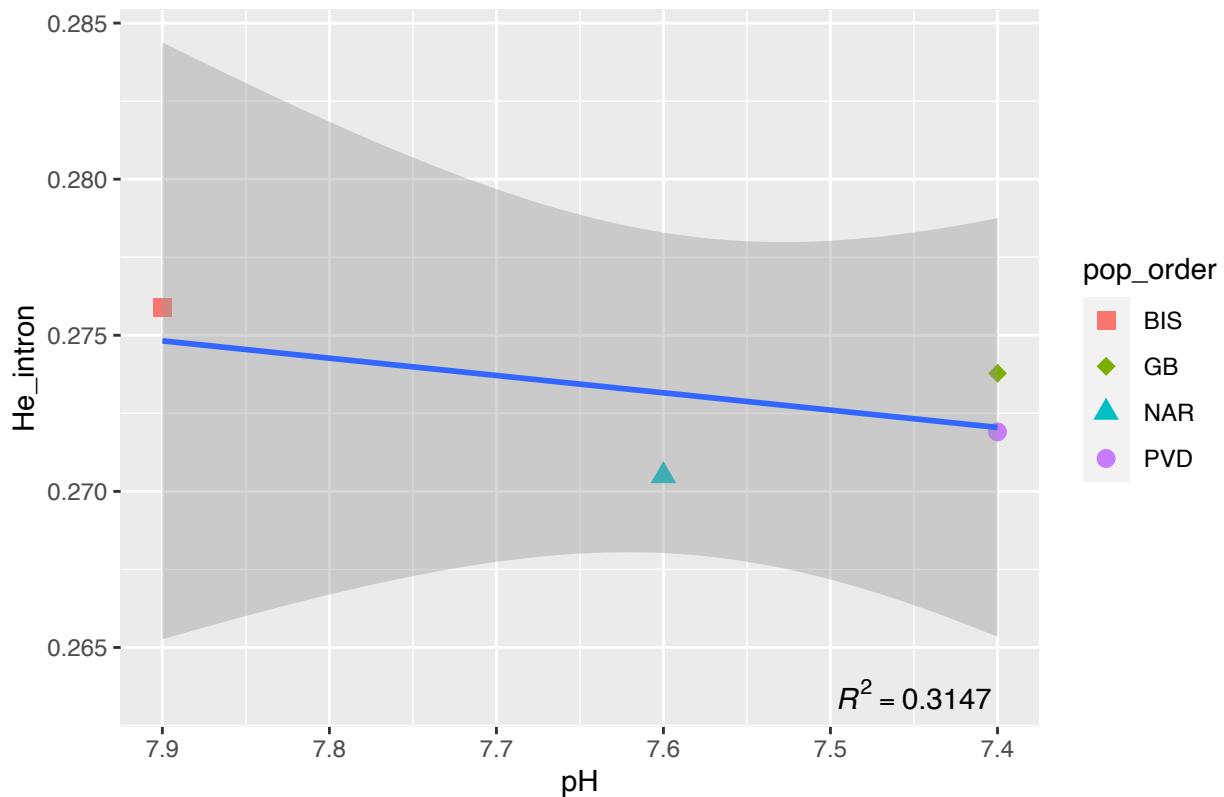
## Observed heterozygosity vs Salinity, Neutral Intron SNPs



```
#Plot He vs pH
R7_intron = round(summary(lm(y_intron$He_intron ~ y_intron$pH))$r.squared, 4)
ggplot(y_intron, aes(x = pH, y = He_intron)) + geom_point(aes(shape=pop_order, color=pop_order), size =
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggtitle("Expected heterozygosity vs pH, Neutral Intron SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R7_intron), x=7.45, y=0.2635, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

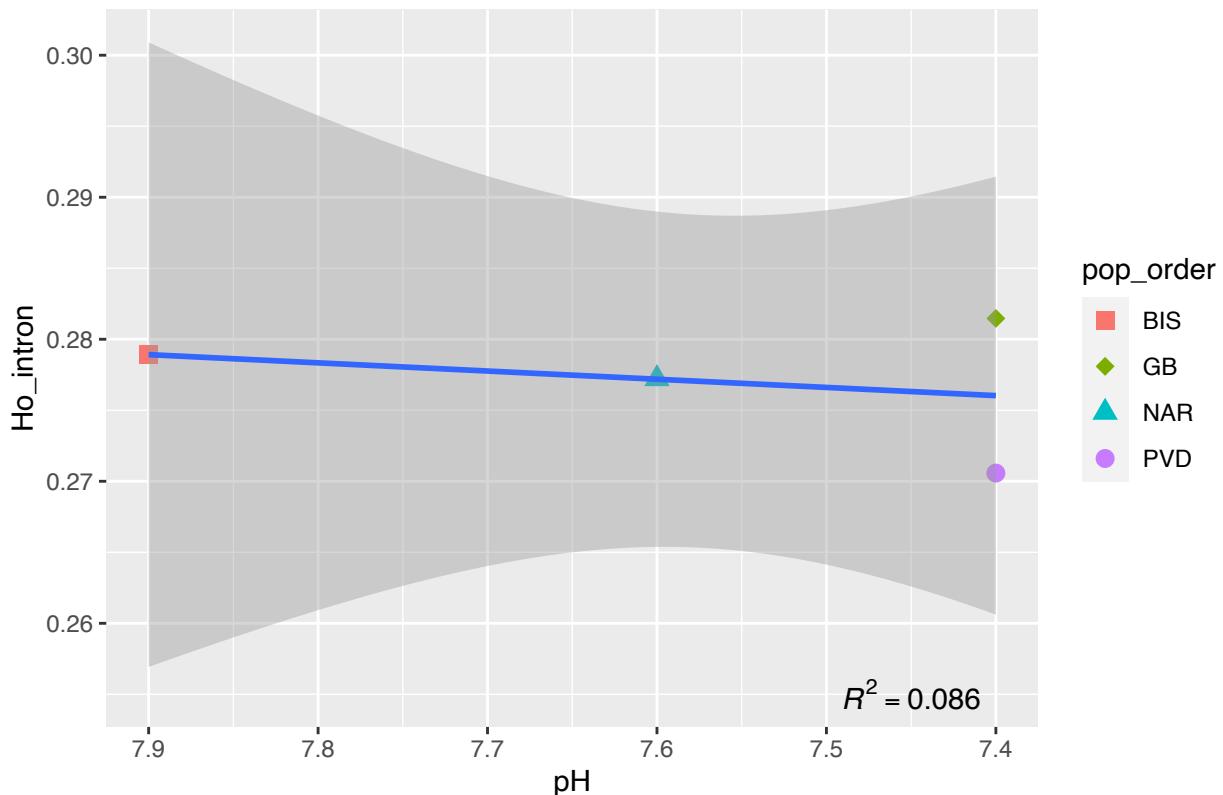
### Expected heterozygosity vs pH, Neutral Intron SNPs



```
#Plot Ho vs pH
R7_introns = round(summary(lm(y_introns$Ho_intron ~ y_introns$pH))$r.squared, 4)
ggplot(y_introns, aes(x = pH, y = Ho_intron)) + geom_point(aes(shape=pop_order, color=pop_order), size =
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggtitle("Observed heterozygosity vs pH, Neutral Intron SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R7_introns), x=7.45, y=0.255, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

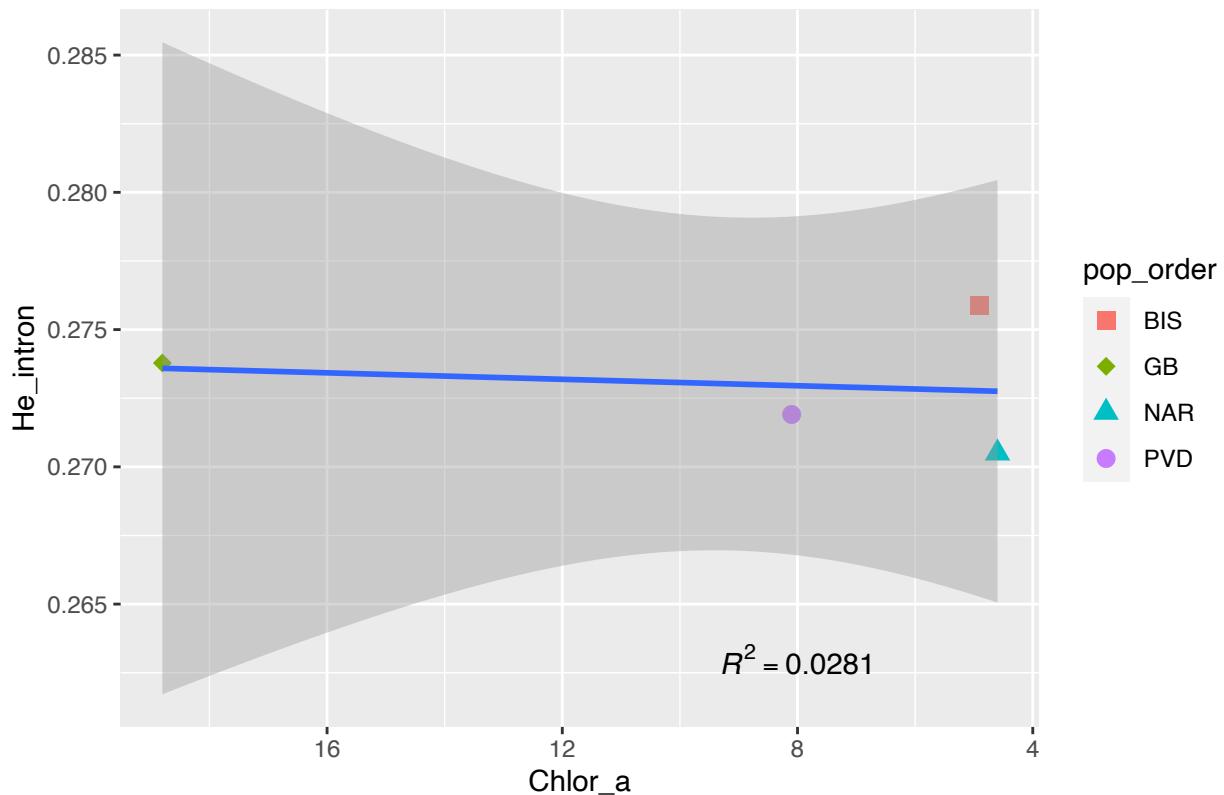
### Observed heterozygosity vs pH, Neutral Intron SNPs



```
#Plot He vs Chlorophyll a
R8_introns = round(summary(lm(y_introns$He_intron ~ y_introns$Chlor_a))$r.squared, 4)
ggplot(y_introns, aes(x = Chlor_a, y = He_intron)) + geom_point(aes(shape=pop_order, color=pop_order), size=10) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggtitle("Expected heterozygosity vs Chlorophyll a, Neutral Intron SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R8_introns), x=8, y=0.263, parse=T) +
  scale_x_reverse()

## `geom_smooth()` using formula 'y ~ x'
```

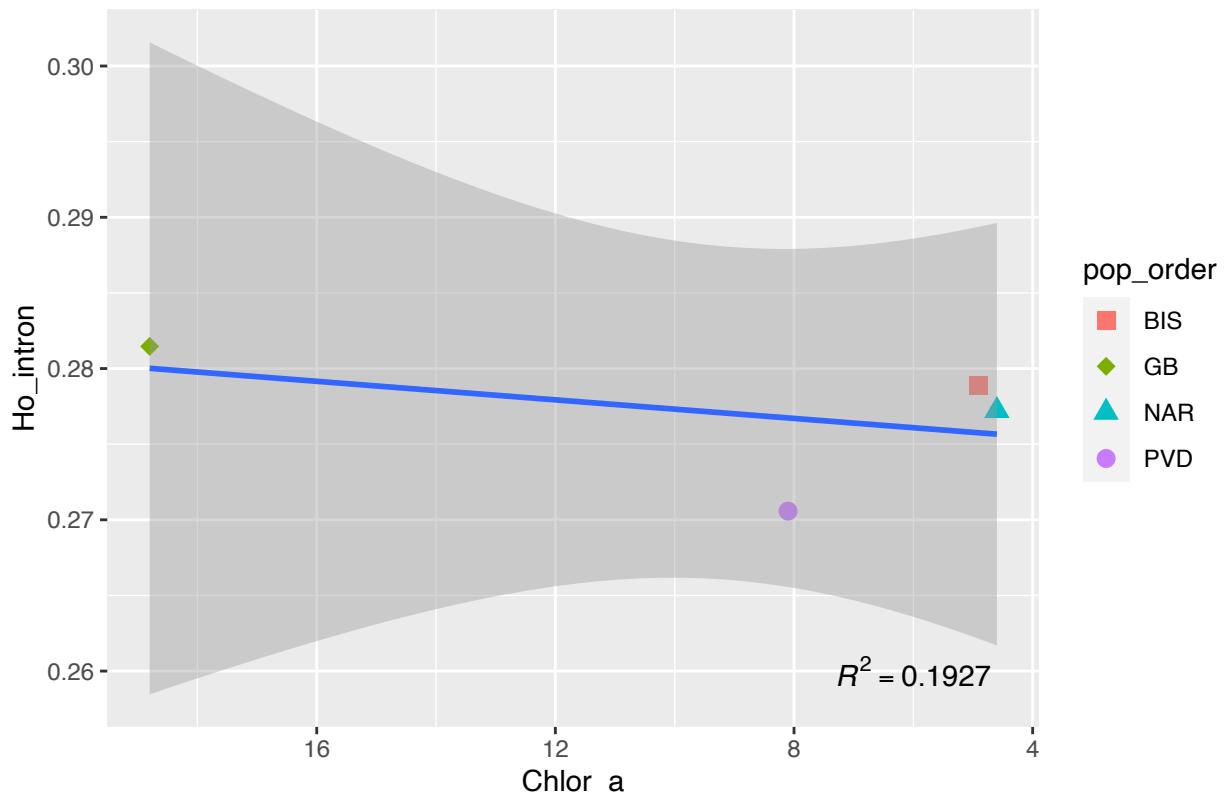
### Expected heterozygosity vs Chlorophyll a, Neutral Intron SNPs



```
#Plot Ho vs Chlorophyll a
R8_introns = round(summary(lm(y_introns$Ho_intron ~ y_introns$Chlor_a))$r.squared, 4)
ggplot(y_introns, aes(x = Chlor_a, y = Ho_intron)) + geom_point(aes(shape=pop_order, color=pop_order), size=10) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggtitle("Observed heterozygosity vs Chlorophyll a, Neutral Intron SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R8_introns), x=6, y=0.26, parse=T) +
  scale_x_reverse()

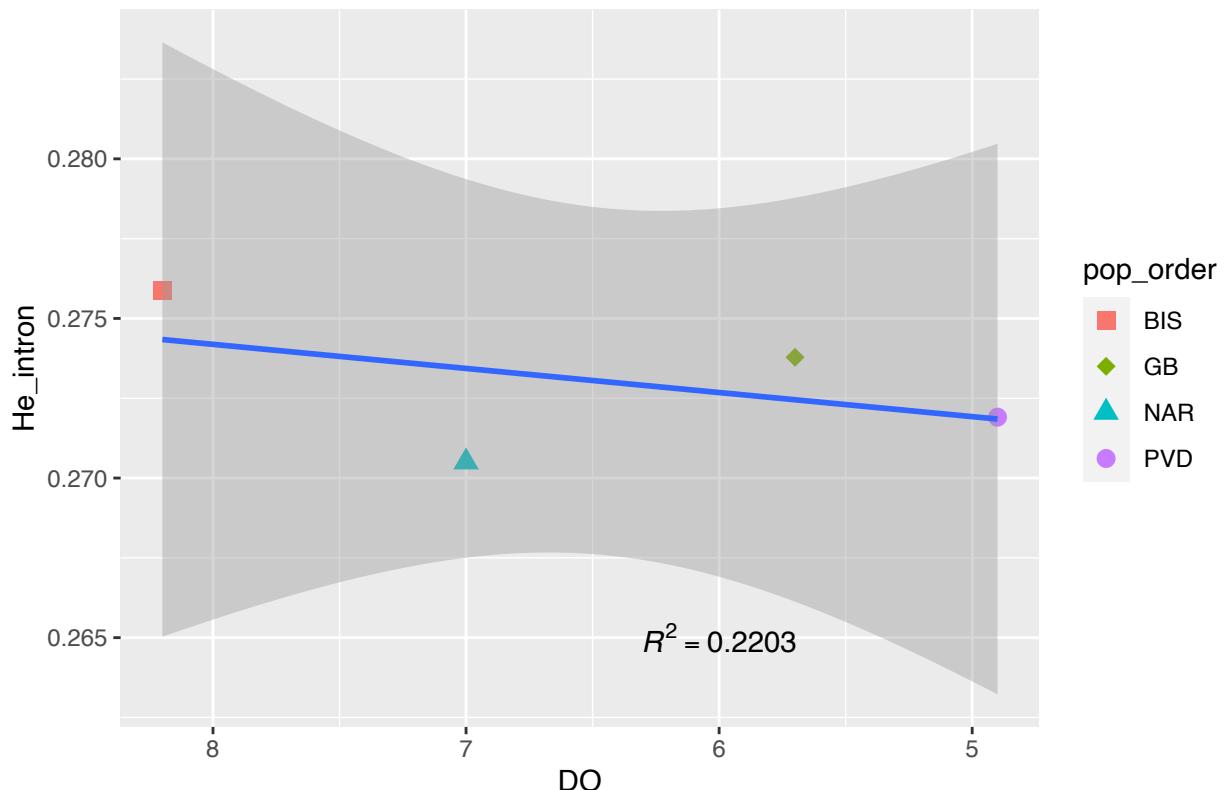
## `geom_smooth()` using formula 'y ~ x'
```

### Observed heterozygosity vs Chlorophyll a, Neutral Intron SNPs



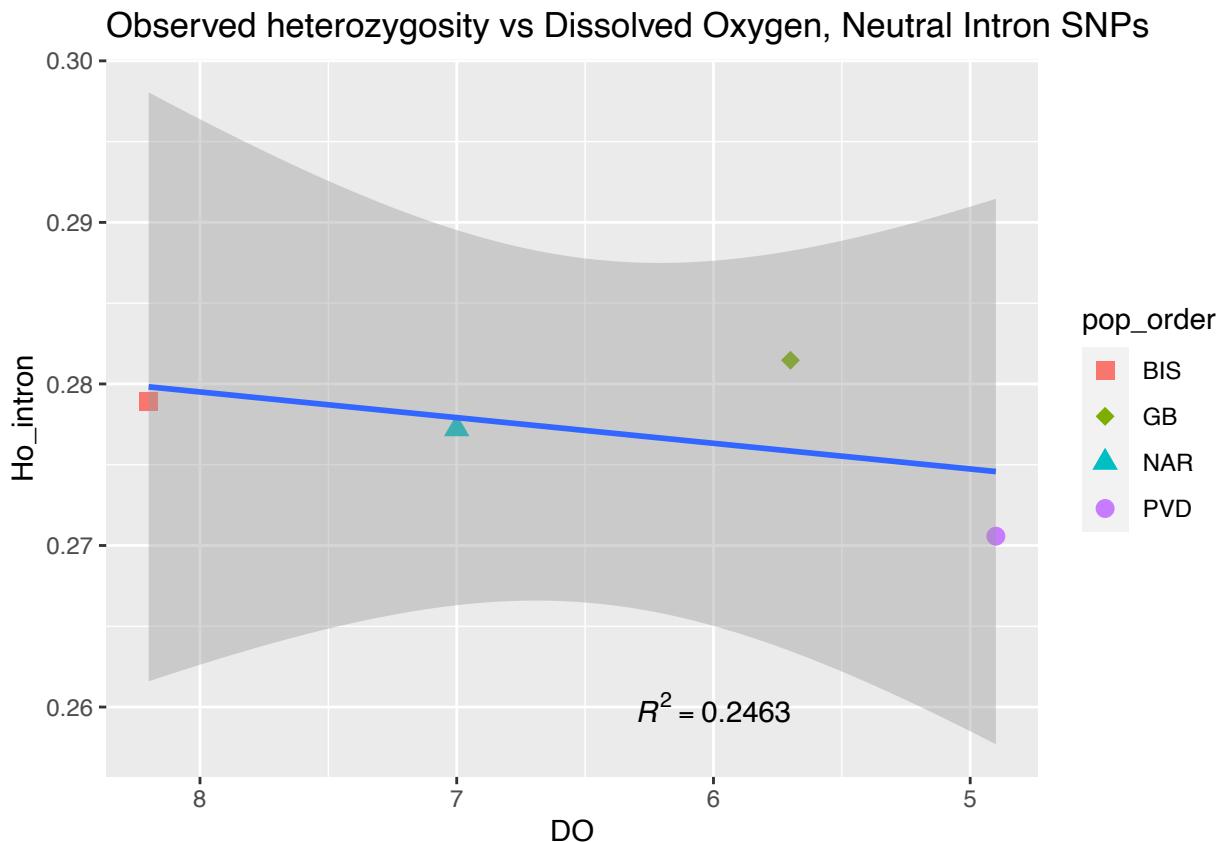
```
#Plot He vs Dissolved Oxygen
R9_introns = round(summary(lm(y_introns$He_intron ~ y_introns$D0))$r.squared, 4)
ggplot(y_introns, aes(x = D0, y = He_intron)) + geom_point(aes(shape=pop_order, color=pop_order), size =
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggtitle("Expected heterozygosity vs Dissolved Oxygen, Neutral Intron SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R9_introns), x=6, y=0.265, parse=T) +
  scale_x_reverse()
## `geom_smooth()` using formula 'y ~ x'
```

## Expected heterozygosity vs Dissolved Oxygen, Neutral Intron SNPs



```
#Plot Ho vs Dissolved Oxygen
R9_introns = round(summary(lm(y_introns$Ho_intron ~ y_introns$DO))$r.squared, 4)
ggplot(y_introns, aes(x = DO, y = Ho_intron)) + geom_point(aes(shape=pop_order, color=pop_order), size = 10) +
  scale_shape_manual(values = c(15,18,17,16)) +
  geom_smooth(method=lm) +
  ggtitle("Observed heterozygosity vs Dissolved Oxygen, Neutral Intron SNPs") +
  annotate(geom = "text", label=paste("italic(R^2)==",R9_introns), x=6, y=0.26, parse=T) +
  scale_x_reverse()

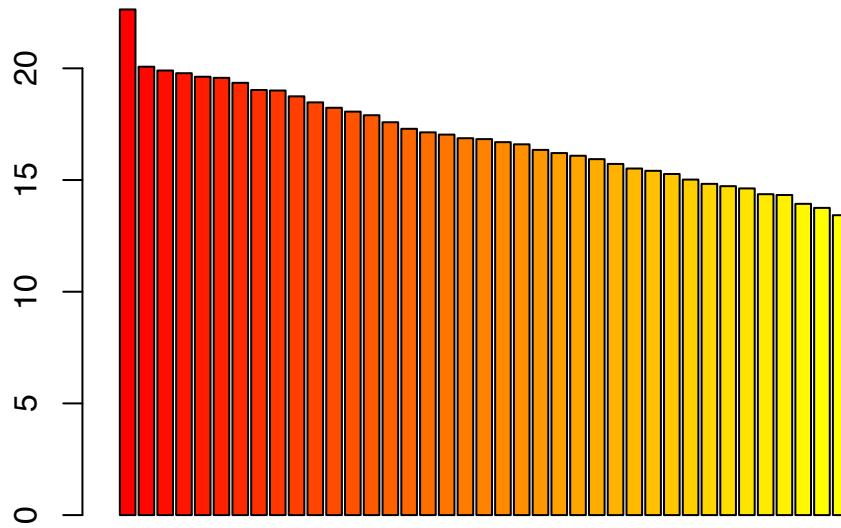
## `geom_smooth()` using formula 'y ~ x'
```



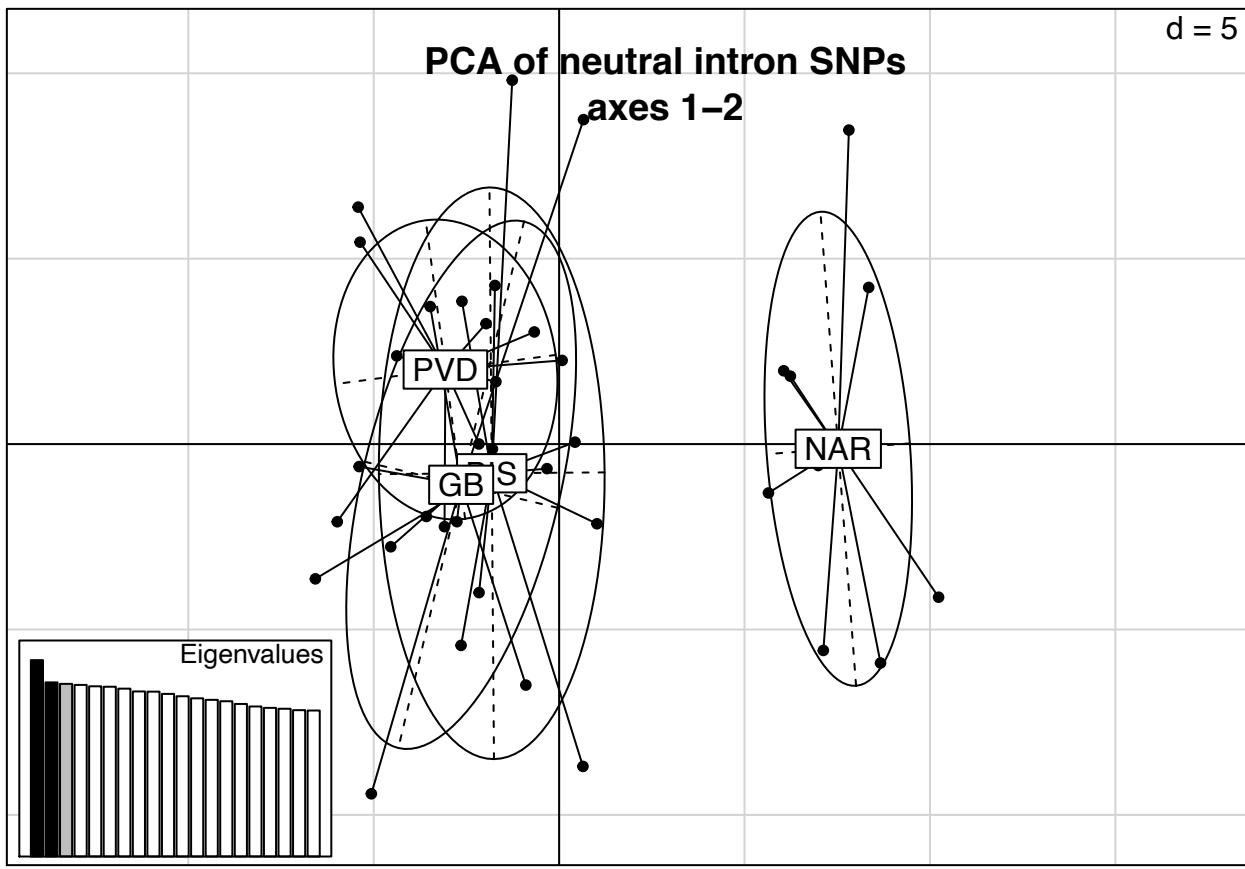
## PCA

```
X_introne <- tab(stratified.u_introne, freq = TRUE, NA.method = "mean")
pca1_introne <- dudi.pca(X_introne, scale = FALSE, scannf = FALSE, nf = 3)
barplot(pca1_introne$eig[1:50], main = "PCA eigenvalues", col = heat.colors(50))
```

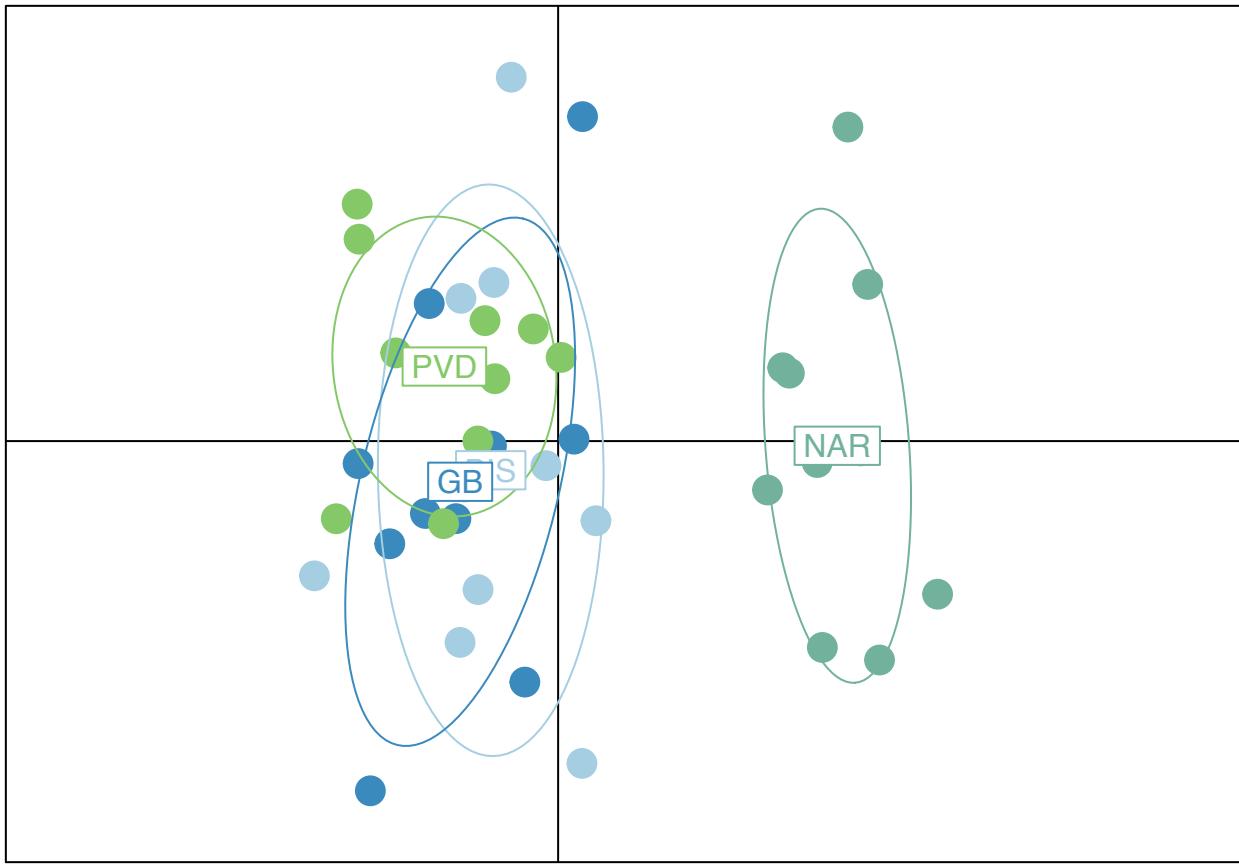
## PCA eigenvalues



```
s.class(pca1_intron$li, pop(stratated.filt_intron))
title("PCA of neutral intron SNPs\naxes 1-2")
add.scatter.eig(pca1_intron$eig[1:20], 3,1,2)
```



```
col <- funky(15)
s.class(pca1_introns$li, pop(stratified.filt_introns), xax=1, yax=2, col=col, axesell=FALSE, cstar=0, cpoint=
```



## Discriminant Analysis of Principal Components (DAPC)

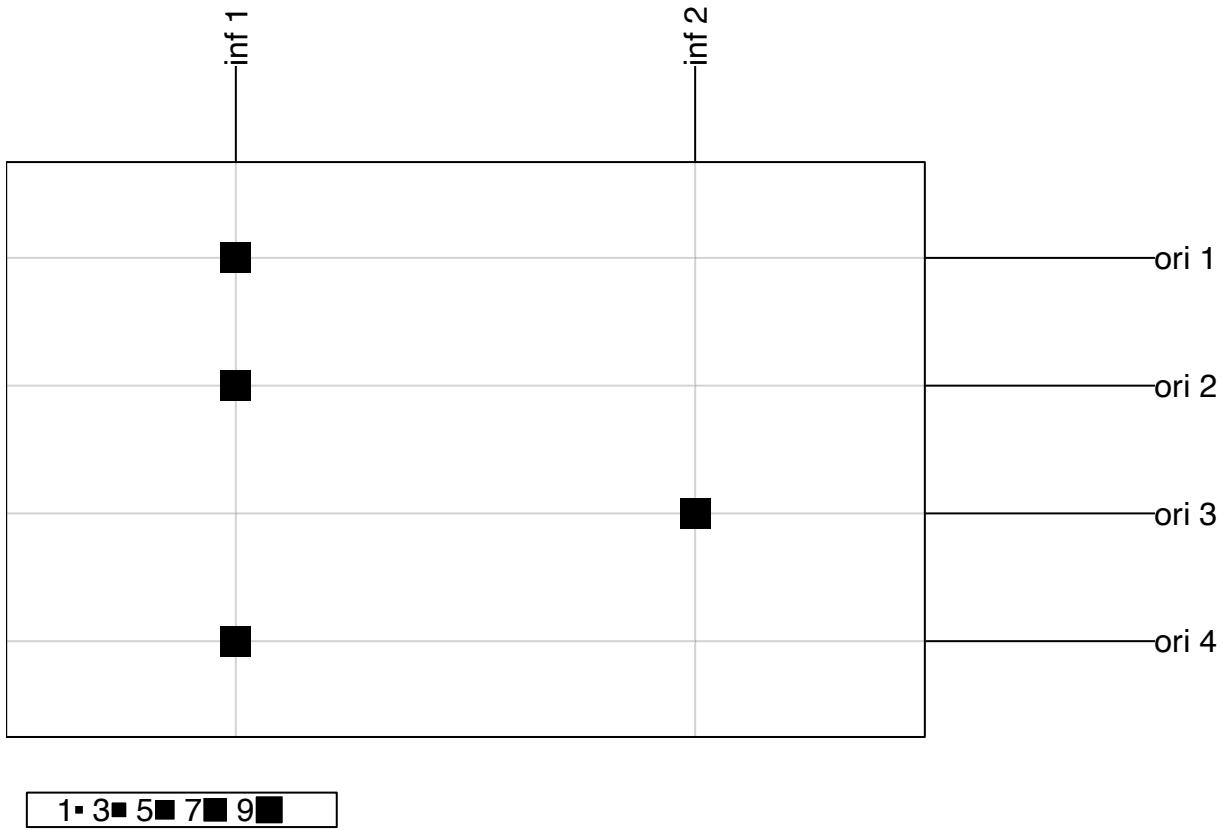
The first line of the code block below is interactive. In the console, it will tell you:

1. Choose the number PCs to retain ( $\geq 1$ ): I picked 1 2. Choose the number discriminant functions to retain ( $\geq 1$ ): I picked 2

```
grp_introns <- find.clusters(stratified.u_introns, n.pca = 1, choose.n.clust = FALSE)
table(pop(stratified.u_introns), grp_introns$grp)
```

```
##
##      1 2
##  BIS 10 0
##  GB  10 0
##  NAR  0 10
##  PVD 10 0

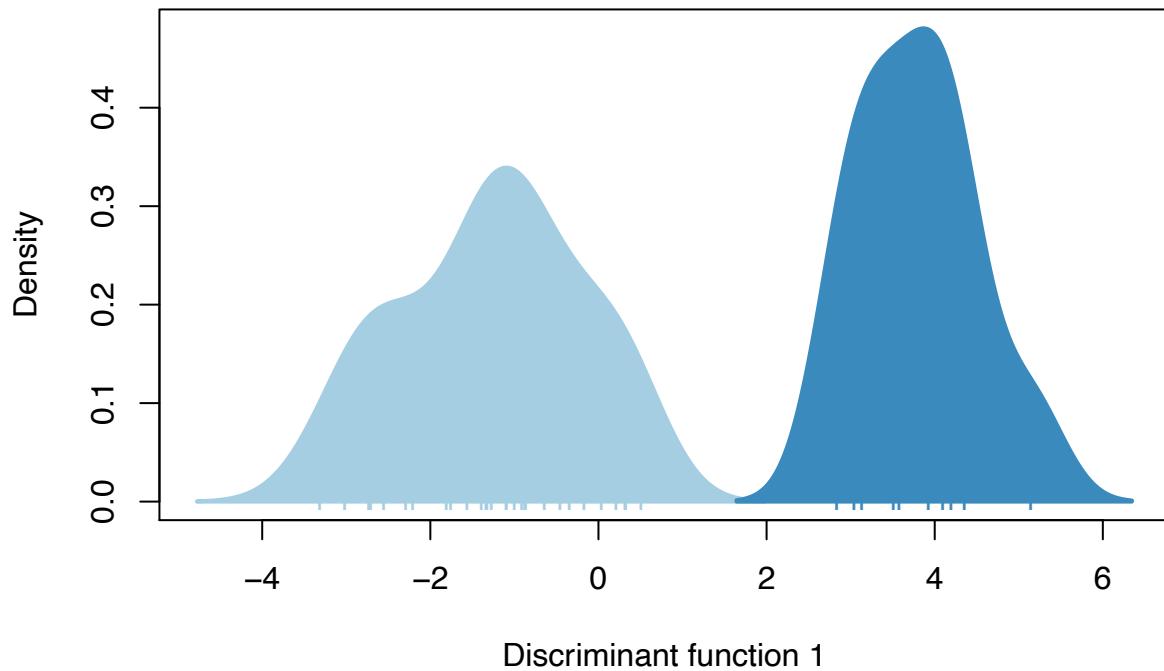
table.value(table(pop(stratified.u_introns), grp_introns$grp), col.lab= paste("inf", 1:2), row.lab= paste("or"))
```



Again, the first line of the code block below is interactive. In the console, it will tell you:

1. Choose the number PCs to retain ( $\geq 1$ ): I picked 1
2. Choose the number discriminant functions to retain ( $\geq 1$ ): I picked 1

```
dapc1_introns <- dapc(stratified.u_introns, grp_introns$grp, n.pca = 1, n.da= 1)
scatter(dapc1_introns, col=col, bg="white", solid=1)
```

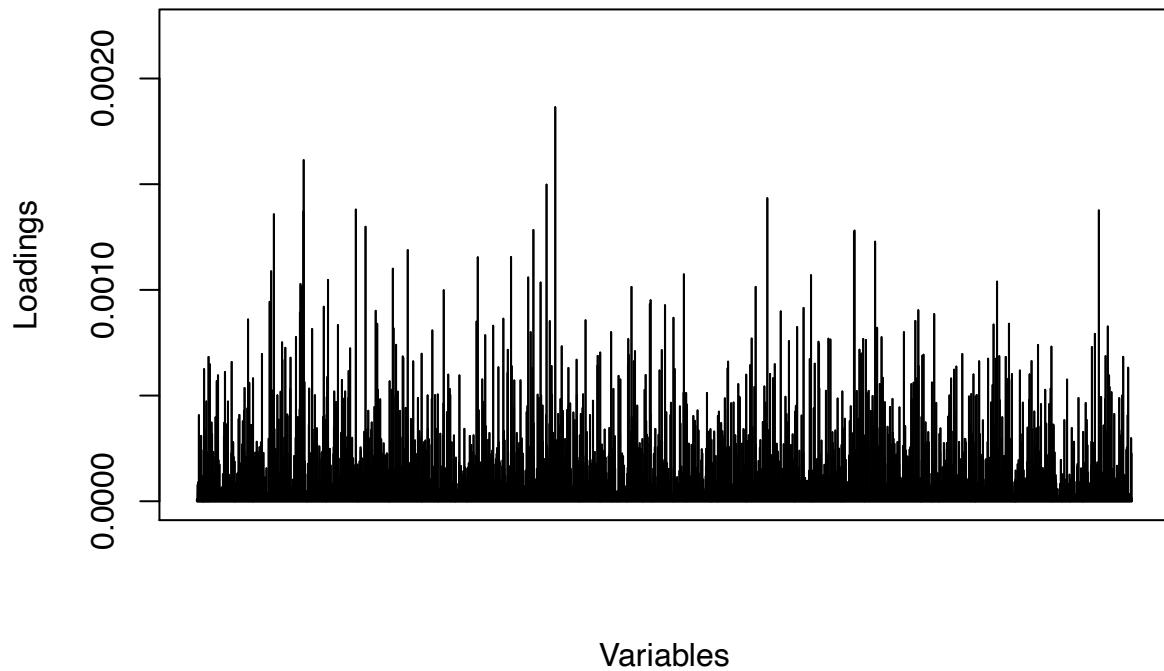


Again, the first line of the code block below is interactive. In the console, it will tell you:

1. Choose the number PCs to retain ( $\geq 1$ ): I picked 1 2. Choose the number discriminant functions to retain ( $\geq 1$ ): I picked 1

```
contrib_intron <- loadingplot(dapc1_intron$var.contr, axis=1, thres=.01, lab.jitter=1)
```

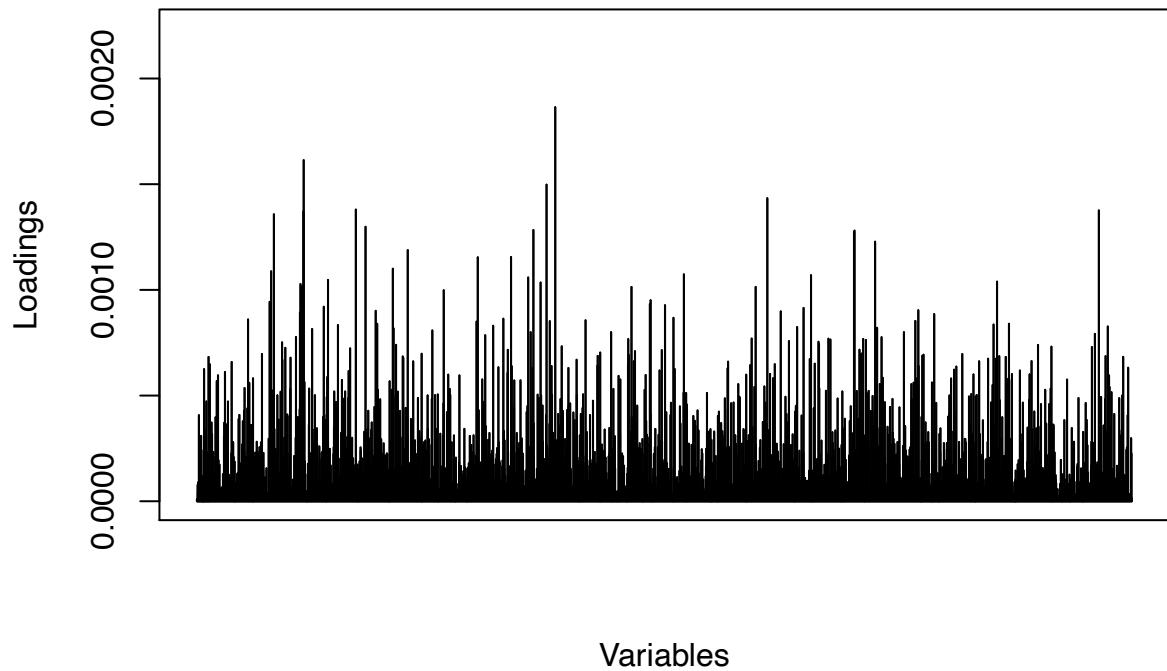
## Loading plot



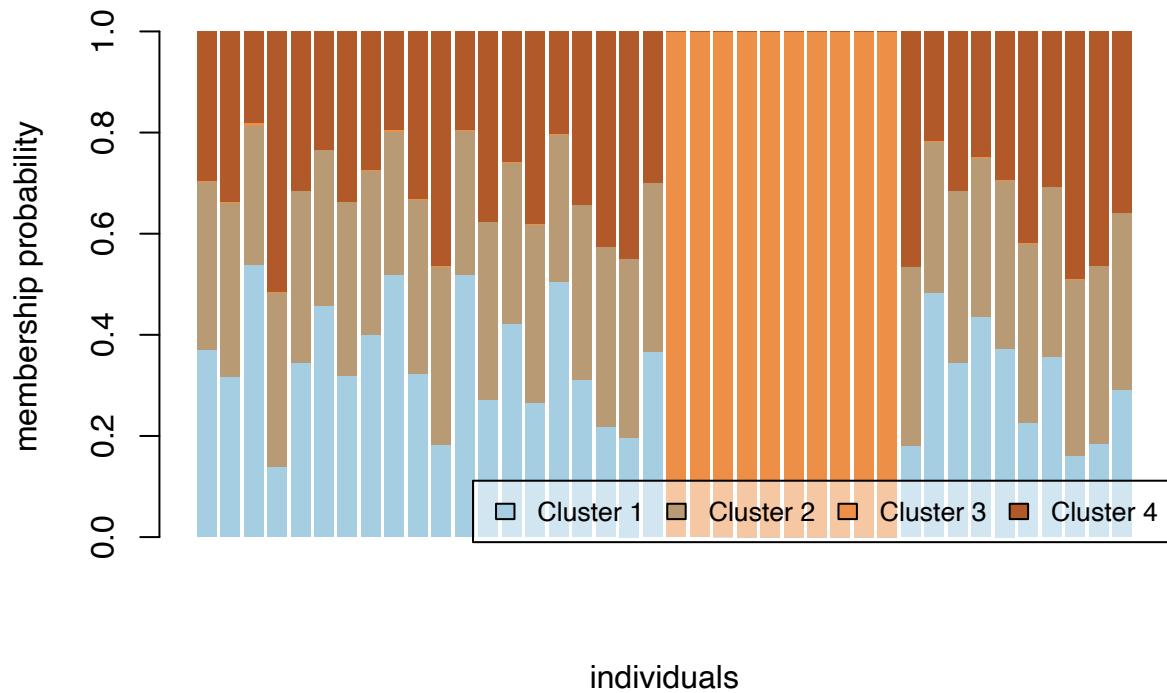
```
contrib_intron
```

```
## NULL  
setPop(rad.filt_intron) <- ~Library  
dapc1_intrон <- dapc(stratified.u_intrон, pop(stratified.u_intrон), n.pca = 1, n.da = 1)  
contrib_intron <- loadingplot(dapc1_intrон$var.contr, axis=1, thres=.05, lab.jitter=1)
```

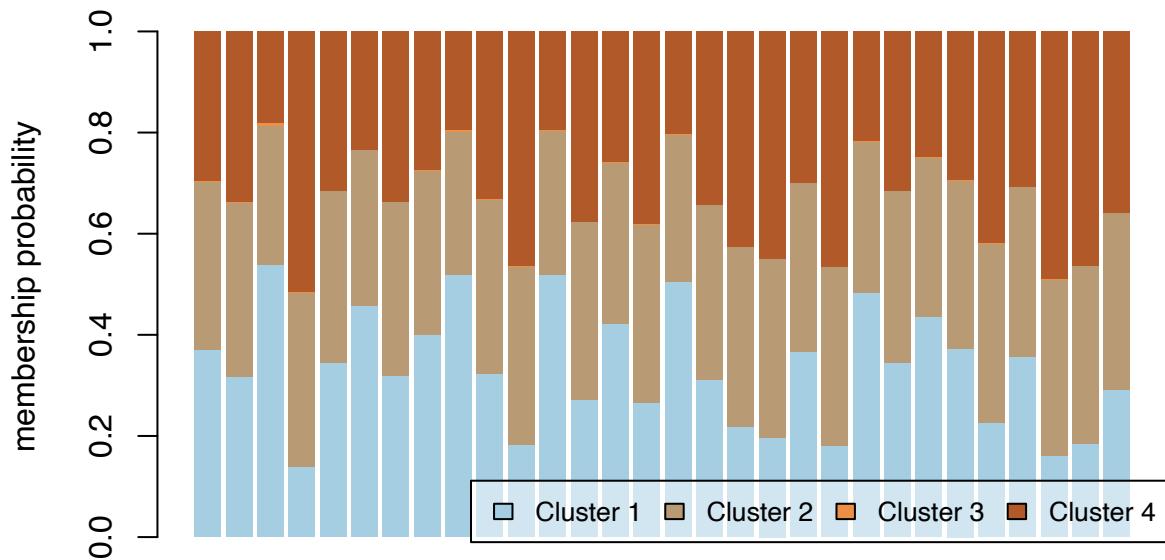
## Loading plot



```
#Structure Like
compoplot(dapc1_intrn, posi="bottomright", txt.leg=paste("Cluster", 1:4), lab="", ncol=1, xlab="individual")
```



```
temp_intron <- which(apply(dapc1_intron$posterior, 1, function(e) all(e<0.9)))
compoplot(dapc1_intron, subset=temp_intron, posi="bottomright", txt.leg=paste("Cluster", 1:4), ncol=2)
```



## PCAviz

```

NA.afDraw<- function(ind){
  ind.mat <- ind@tab
  new.mat <- ind.mat
  af = colSums(ind.mat[,seq(1,ncol(ind.mat)-1,2)],na.rm = TRUE) /
    (2*apply(ind.mat[,seq(1,ncol(ind.mat)-1,2)],2,function(x) sum(!is.na(x))))
  af.Draw <- function(geno, af){
    new <- function(geno,af){
      if(is.na(geno)){
        newA = rbinom(1,2,af)
      } else {newA <- geno}
      return(newA)
    }
    new.row <- mapply(geno,af,FUN = new)
    return(new.row)
  }

  new.mat[,seq(1,ncol(ind.mat)-1,2)] <- t(apply(ind.mat[,seq(1,ncol(ind.mat)-1,2)],1,af.Draw,af))
  new.mat[,seq(2,ncol(ind.mat),2)] <- 2-new.mat[,seq(1,ncol(ind.mat)-1,2)]
  new.ind <- ind
  new.ind@tab <- new.mat
  return(new.ind)
}

```

```

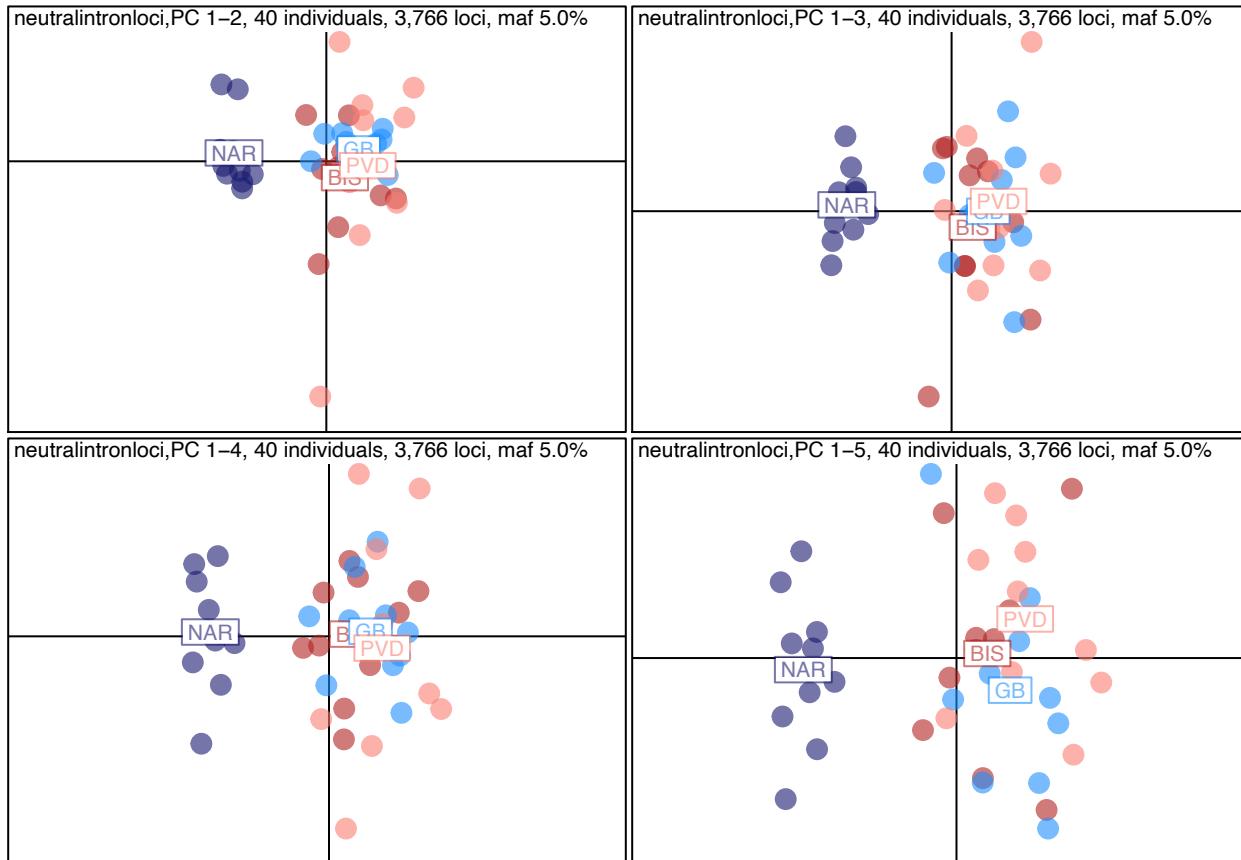
u.na_intron <- NA.afDraw(stratified.u_intron)

pca_intron <- dudi.pca(u.na_intron,center=TRUE,scale=TRUE,scannf = F, nf = 30)

col18 <- funky(length(unique(u.na_intron$strata$Population)))
#Colors that match the neutral Structure results
col6 <- c("firebrick","dodgerblue","midnightblue","salmon")
par(mfrow=c(2,2))

s.class(pca_intron$li, strata(u.na_intron)$Population,xax=1,yax=2,
         sub = "neutralintronloci,PC 1-2, 40 individuals, 3,766 loci, maf 5.0%",
         possub = "topleft",col=transp(col6,.6),axesell=FALSE,
         cstarc=0, cpoint=3, grid=FALSE, cellipse = 0)
s.class(pca_intron$li, strata(u.na_intron)$Population,xax=1,yax=3,
         sub = "neutralintronloci,PC 1-3, 40 individuals, 3,766 loci, maf 5.0%",
         possub = "topleft",col=transp(col6,.6),axesell=FALSE,
         cstarc=0, cpoint=3, grid=FALSE, cellipse = 0)
s.class(pca_intron$li, strata(u.na_intron)$Population,xax=1,yax=4,
         sub = "neutralintronloci,PC 1-4, 40 individuals, 3,766 loci, maf 5.0%",
         possub = "topleft",col=transp(col6,.6),axesell=FALSE,
         cstarc=0, cpoint=3, grid=FALSE, cellipse = 0)
s.class(pca_intron$li, strata(u.na_intron)$Population,xax=1,yax=5,
         sub = "neutralintronloci,PC 1-5, 40 individuals, 3,766 loci, maf 5.0%",
         possub = "topleft",col=transp(col6,.6),axesell=FALSE,
         cstarc=0, cpoint=3, grid=FALSE, cellipse = 0)

```



# Seascape Redundancy Analysis

This code follows that documented by Tom Jenkins.

## Prepare genetic data for redundancy analysis.

### Notes before execution:

1. Make sure all required R packages are installed.
2. Set working directory to the location of this R script.

```
# Load packages
library(adegenet)
library(poppr)

## Registered S3 method overwritten by 'pegas':
##   method      from
##   print.amova ade4

## This is poppr version 2.8.5. To get started, type package?poppr
## OMP parallel support: available

##
## Attaching package: 'poppr'

## The following object is masked from 'package:radiator':
##
##   private_alleles

library(dplyr)
library(reshape2)
library(ggplot2)
library(vcfR)

# Explore data
rad.filt

## /// GENIND OBJECT ///////////
##
## // 40 individuals; 74,520 loci; 148,794 alleles; size: 67.9 Mb
##
## // Basic content
##   @tab: 40 x 148794 matrix of allele counts
##   @loc.n.all: number of alleles per locus (range: 1-2)
##   @loc.fac: locus factor for the 148794 columns of @tab
##   @all.names: list of allele names for each locus
##   @ploidy: ploidy of each individual (range: 2-2)
##   @type: codom
##   @call: adegeneet::df2genind(X = t(x), sep = sep, pop = ..2, strata = ..1)
##
## // Optional content
##   @pop: population of each individual (group size range: 8-8)
##   @strata: a data frame with 12 columns ( Individual, Population, Latitude, Longitude, Distance, SE
nLoc(rad.filt) # number of loci

## [1] 74520
```

```

nPop(rad.filt) # number of sites
## [1] 5

nInd(rad.filt) # number of individuals
## [1] 40

summary(rad.filt$pop) # sample size
## 1 5 2 3 4
## 8 8 8 8 8

# Calculate allele frequencies for each site
allele_freqs = data.frame(rraf(rad.filt, by_pop = TRUE, correction = FALSE), check.names = FALSE)

# Keep only the first of the two alleles for each SNP (since p=1-q).
allele_freqs = allele_freqs[, seq(1, dim(allele_freqs)[2], 2)]

# Export allele frequencies
write.csv(allele_freqs, file = "all_allele_freqs.csv", row.names = TRUE)

```

---

## Calculate minor allele frequencies

---

```

# Separate genind object by site
site_list = seppop(rad.filt)
names(site_list)

## [1] "1" "5" "2" "3" "4"

# Calculate the minor allele frequency for each site
maf_list = lapply(site_list, FUN = minorAllele)

# Convert list to dataframe
maf = as.data.frame(maf_list) %>% t() %>% as.data.frame()

# Export minor allele frequencies
write.csv(maf, file = "minor_allele_freqs.csv", row.names = TRUE)

```

Prepare environmental data for redundancy analysis.

**Environmental variables:**

- Distance from sewage effluent source (km)

- Sewage Effluent (PW stats)
- Mean temperature (deg C)
- Mean Salinity (psu)
- Mean pH
- Mean Chlorophyll-a (ug/L)
- Mean Dissolved Oxygen (mg/L)

Environmental data for each population is saved in a `strata_pop` file that can be accessed here

```
# All environmental data was previously saved in strata file
strata_pop <- read.table("strata_pop", header=TRUE)
strata_pop

##   Population Latitude Longitude Distance      SE Temperature Salinity pH
## 1       BIS    41.545   -71.431     4.76  8.825        23     30 7.9
## 2       GB     41.654   -71.445     0.47 14.596        24     28 7.4
## 3       NAR    41.505   -71.453    15.41  2.027        25     18 7.6
## 4       PVD    41.816   -71.391     1.49 59.860        23     25 7.4
##   Chlorophylla DO
## 1           4.9 8.2
## 2           18.8 5.7
## 3           4.6 7.0
## 4           8.1 4.9

# Export data as a csv file
write.csv(strata_pop, file="environmental_data.csv", row.names = FALSE)
```

I also prepared spatial data for the redundancy analysis which is documented here.

Allele frequency, environmental, and spatial csv files are saved to your working directory and must be imported into the Rscript to run the redundancy analysis. Documentation of the RDA can be accessed here.

RDA and pRDA output plots can be viewed here.