

Q & A

Physalia last day

Stacks

Other calling variant tools ?

The matter of having variants of different sizes or to trim?

Pool-seq /low-coverage data

GEA:

- Baypass has a specific option for pool seq

http://www1.montpellier.inra.fr/CBGP/software/baypass/files/BayPass_manual_2.1.pdf

- Baypass /LFMM/RDA can all take a matrix of individuals (genotypes) or of populations (allelic frequencies)

	snps		snps
Ind1	0 1 1 2	Pop1	0.1 0.3 0.4 0.1
Ind2	2 1 2 0	Pop2	0.5 0.5 0.1 0.2
Ind3	2 1 1 1	Pop3	0.2 0.3 0.2 0.1

Baypass

about making independant runs

What we did

- Run baypass once
- Use 1 CPU!
- Take the value of xtx (or BF) from this run
- Keep as outliers SNPs with xtx (or BF) above the 99% of Xtx from simulated values
- Look at outliers SNPs that were shared with RDA (*but remember that RDA and Baypass works differently*)

Recommended Practices for your dataset

- Run baypass 3 to 5 times with a different seed
- Use 5 to 10 CPU (nthreads) if available
- Take median value of xtx (or BF) for each SNP
- Keep as outliers SNPs with xtx (or BF) above the 99,99...% of Xtx (or BF) from simulated values – Avoid considering BF below 3 (look at Jeffrey's rule)
- Look at outliers SNPs that were shared with any other method of genotype-environment association

RDA

how to interpret the triplot?
advanced options for geographic variables?

(prepared with the help of
Dr. Martin Laporte)

A super good vignette to understand and do Rda analysis:

https://popgen.nescent.org/2018-03-27_RDA_GEA.html

Population Genetics in R

Users ▾

Package Developers ▾

Contribute! ▾

Useful Links

Detecting multilocus adaptation using Redundancy Analysis (RDA)

- [Introduction](#)
- [Assumptions](#)
- [Data & packages](#)
- [Analysis](#)
- [Conclusions](#)
- [Contributors](#)
- [References](#)
- [Session Information](#)

Introduction

The purpose of this vignette is to illustrate the use of **Redundancy Analysis (RDA)** as a genotype-environment association (GEA) method to detect loci under selection (Forester et al., 2018). RDA is a multivariate ordination technique that can be used to analyze many loci and environmental predictors simultaneously. RDA determines how groups of loci covary in response to the multivariate environment, and can detect processes that result in weak, multilocus molecular signatures (Rellstab et al., 2015; Forester et al., 2018).

RDA is a two-step analysis in which genetic and environmental data are analyzed using multivariate linear regression, producing a matrix of fitted values. Then PCA of the fitted values is used to produce canonical axes, which are linear combinations of the predictors (Legendre & Legendre, 2012). RDA can be used to analyze genomic data derived from both individual and population-based sampling designs.

Assumptions

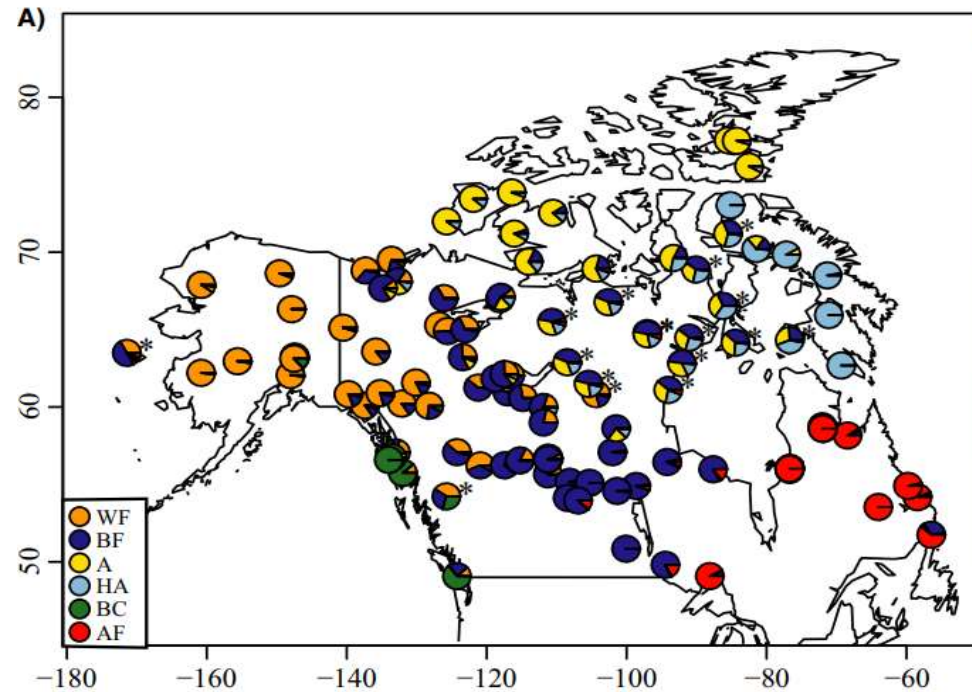
RDA is a linear model and so assumes a linear dependence between the response variables (genotypes) and the explanatory variables (environmental predictors). Additional detail can be found in Legendre & Legendre (2012). We also recommend Borcard et al. (2011) for details on the implementation and interpretation of RDA using the `vegan` package (Oksanen et al, 2017).

Contributors

- Brenna R. Forester (Author)
- Martin Laporte (reviewer)
- Stéphanie Manel (reviewer)

RDA

Multivariate
associations:

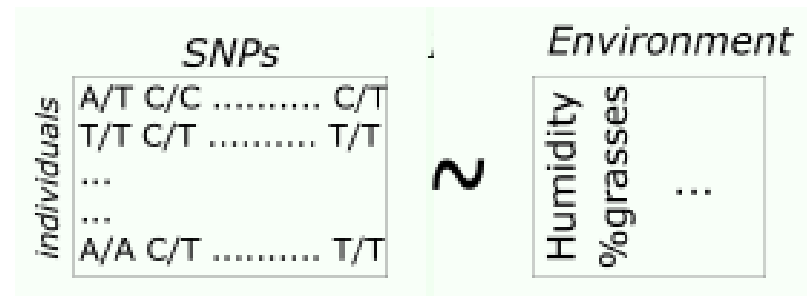


94 wolves
42 597 SNPs

species

sites

In community
ecology (package
vegan!)

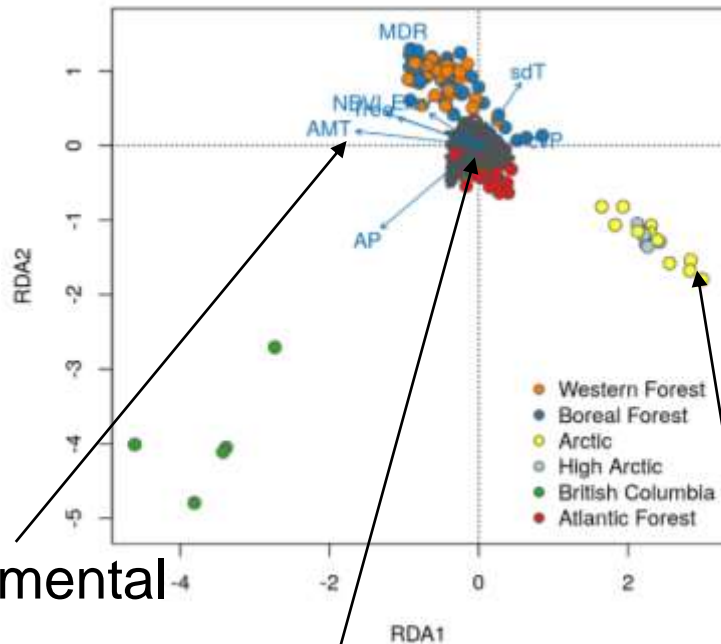


Forester et al 2018 Mol Ecol

RDA

```
points(X.rda, display="sites")
```

Triplot RDA (individual centered)



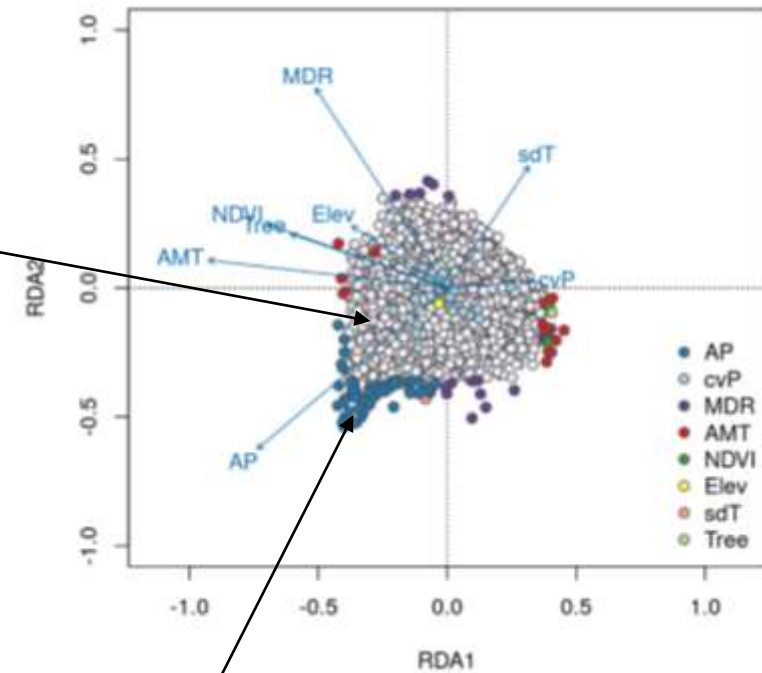
Environmental variable

Genetic markers

Individual

```
points(X.rda, display="species")
```

Triplot RDA (SNPs centered)



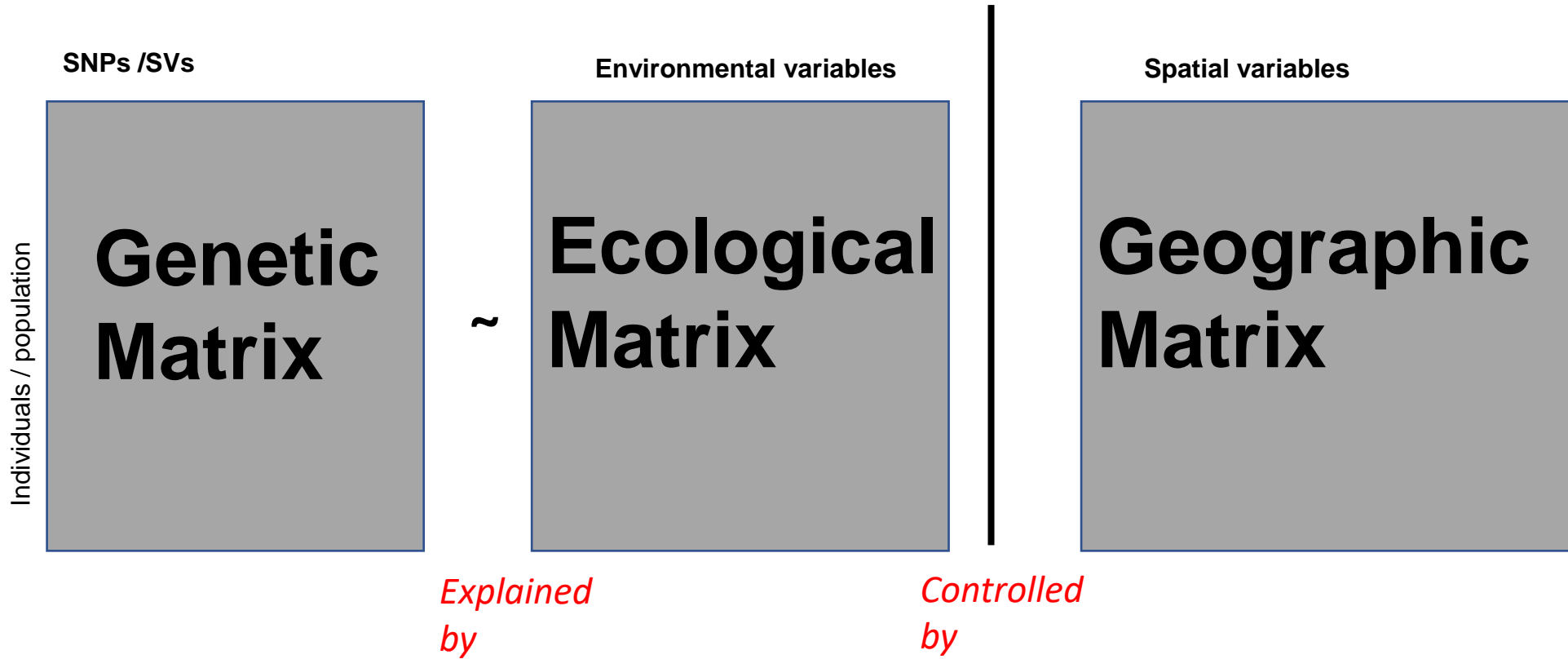
Neutral marker

Outlier marker putatively associated to *AP* variation

Forester et al 2018 Mol Ecol

Use the contribution of genetic markers along the different axis to detect putatively-selected loci

RDA

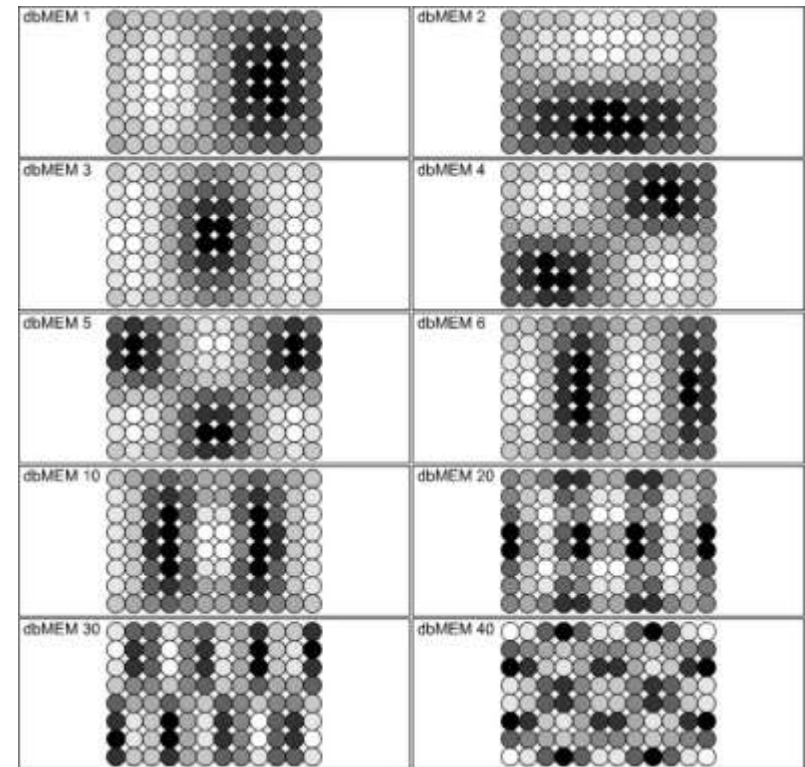


<https://doi.org/10.1016/B978-0-444-53868-0.50014-9>

RDA

$$\boxed{G} \sim \boxed{E} \mid \boxed{S} \quad \begin{array}{l} \text{Latitude + Longitude} \\ \text{or} \\ \text{Spatial eigenvectors} \\ = \text{db-MEM} \end{array}$$

Spatial-eigen vectors are a way to reduce a distance matrix between samples/populations
-> not necessarily neutral
-> describe different possible spatial combination



More information:
Legendre & Legendre

<https://doi.org/10.1016/B978-0-444-53868-0.50014-9>

Climatic Variables

how to extract them from databases?

Databases for climatic variables

<https://www.worldclim.org/>

<http://www.marspec.org/>
(with useful tutorials)

WORLDCLIM: R will gather the data itself

```
location_GPS<- read.delim("location_GPS.txt")
r <- getData("worldclim",var="bio",res=2.5)
div=10 #precision of the data

#1 is mean temp, 12 is annual precipitations, et...
Annual_mean_temp<-r[[1]]
variable<-paste0("bio1")

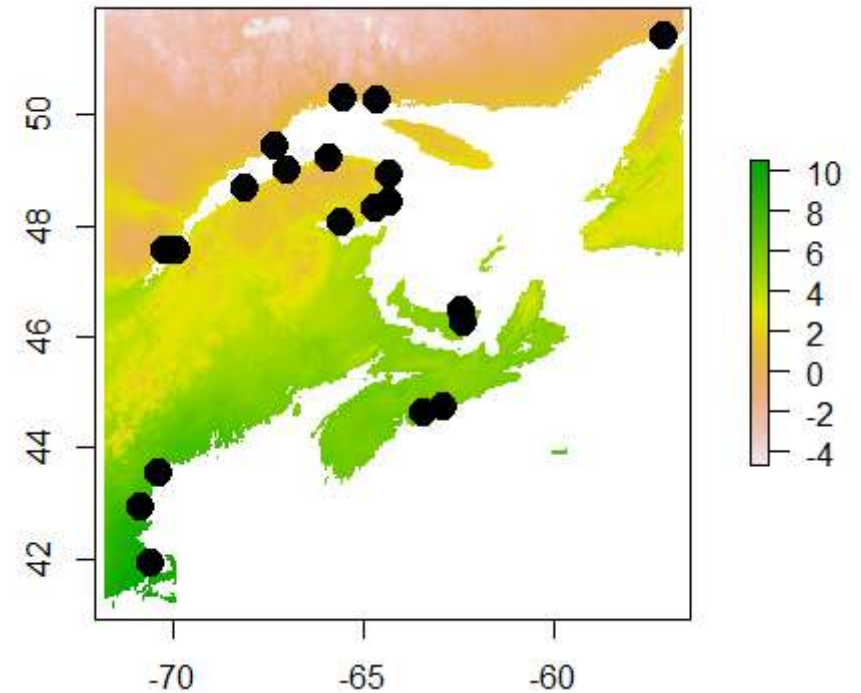
#make a plot of the area
aoi_area <- extent(min (location_GPS$GPS_EW)-1,max (location_GPS$GPS_EW)+0.5,min (location_GPS$GPS_NS)-1,max (location_GPS$GPS_NS)+1)
plot((crop(Annual_mean_temp, aoi_area)/div))
points(location_GPS$GPS_EW,location_GPS$GPS_NS, pch=19, col=1, cex=2)

# to get data round a point of your choice like pop 1
i=1
#determine the coordinates around your point
long_min<-floor(location_GPS$GPS_EW[i]*10)/10
long_max<-ceiling(location_GPS$GPS_EW[i]*10)/10
lat_min<-floor(location_GPS$GPS_NS[i]*10)/10
lat_max<-ceiling(location_GPS$GPS_NS[i]*10)/10

#prepare the area
aoi <- extent(long_min, long_max, lat_min, lat_max)

#get the value of the layer in the area
Annual_mean_temp.crop <- crop(Annual_mean_temp,aoi)
mean_value_i<-mean(Annual_mean_temp.crop@data@values, na.rm=T)/div
range_value_i<-(range(Annual_mean_temp.crop@data@values, na.rm=T)[2]-range(Annual_mean_temp.crop@data@values, na.rm=T)[1])/div

#print value
location_GPS[i,]
mean_value_i
range_value_i
```



MARSPEC Download data

I'll drop a tutorial on the github page of the course within day3