# Genotype-environment associations

E. Anne Chambers
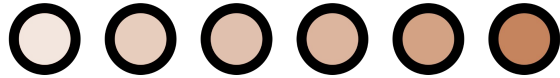March 21, 2024
eachambers@berkeley.edu

# Genotype-environment association (GEA) methods

Allele frequencies

**Local adaptation!**

*Enviro. 1*    *Enviro. 2*
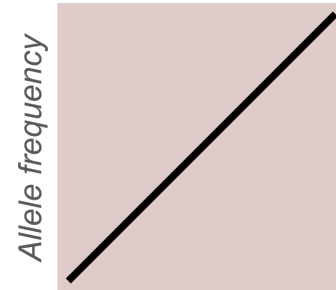
**Neutral processes**
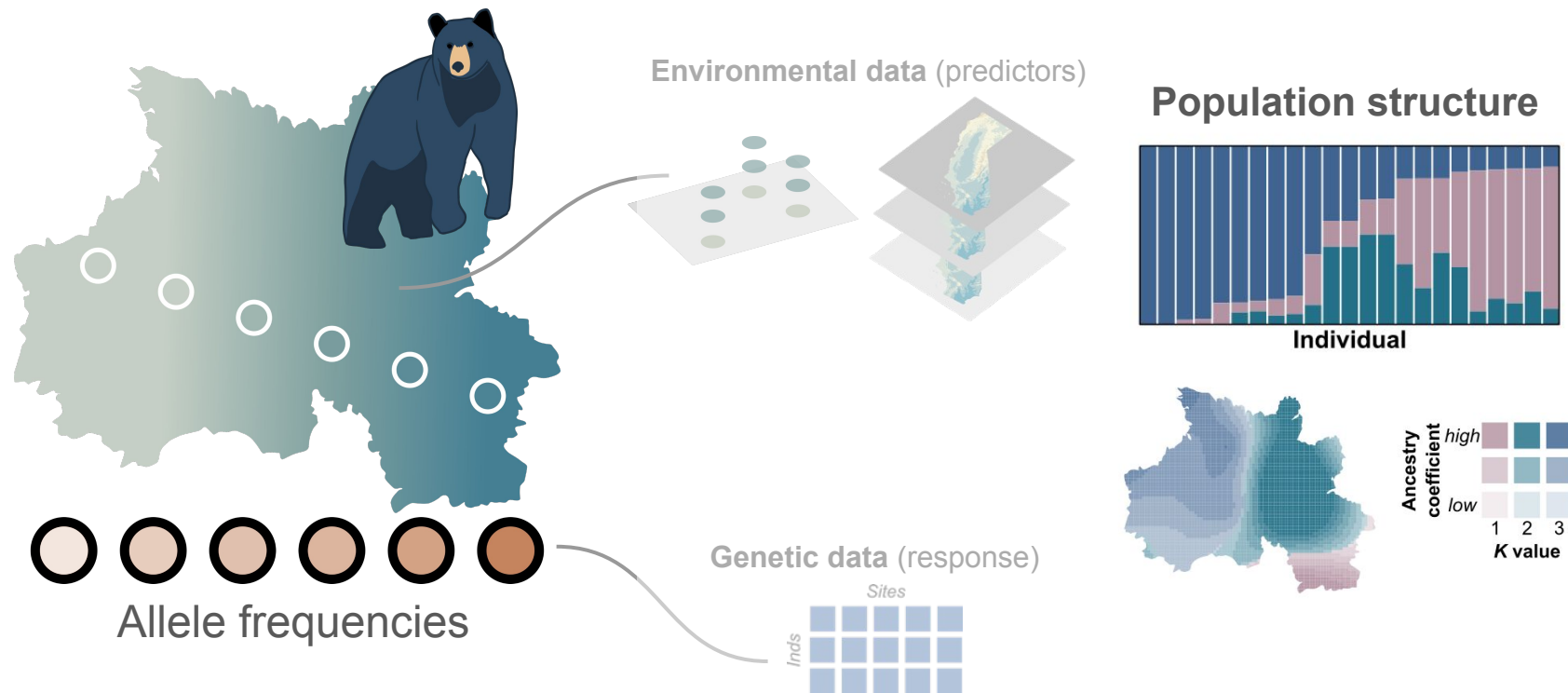
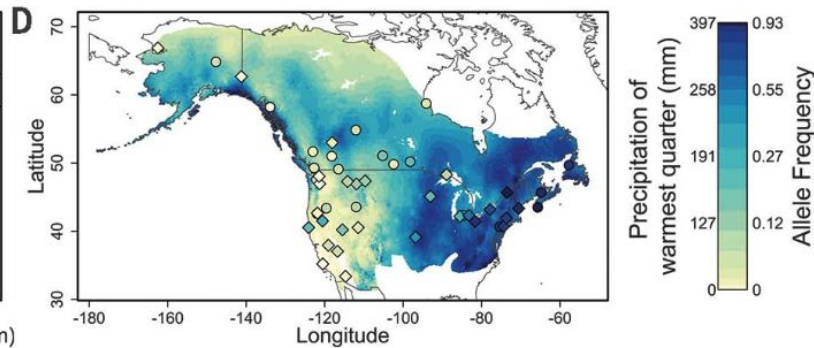*Allele frequency*

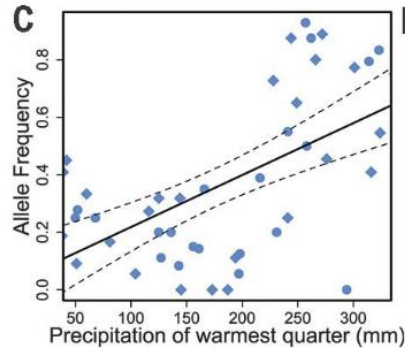*Geographic distance*

# Genotype-environment association (GEA) methods



**Environmental data** (predictors)

**Genetic data** (response)

Sites

Inds

Allele frequencies

# Genotype-environment association (GEA) methods



Environmental data (predictors)

Population structure

Individual

Allele frequencies

Genetic data (response)

Sites

Inds

Ancestry coefficient

high

low

1  2  3

K value

# What questions can we answer using GEA?



Bay et al. (2018) *Science*

# What questions can we answer using GEA?



(a) PC1/PC2 plot with environmental variable axes: Precipitation Driest Quarter (Bio17), Mean Temp Coldest Quarter (Bio11), Temp Seasonality (BIO4), Max Temp Warmest Month (BIO5)

(b) Map of latitude vs longitude

Ruegg et al. (2018) *Ecology Letters*

# What questions can we answer using GEA?



Capblancq & Forester (2021) *Methods Ecol. Evol.*

# What questions can we answer using GEA?



Bay et al. (2018) *Science*; Capblancq & Forester (2021) *Methods Ecol. Evol.*

# Different types of GEA

- BayEnv/BayPASS/BayeScEnv
- Redundancy analysis (RDA)
- Latent factor mixed models (LFMM)
- GLMM
- Gradient forest
- SAM/SamBada
- Weighted Z-analysis (WZA)

# Different types of GEA

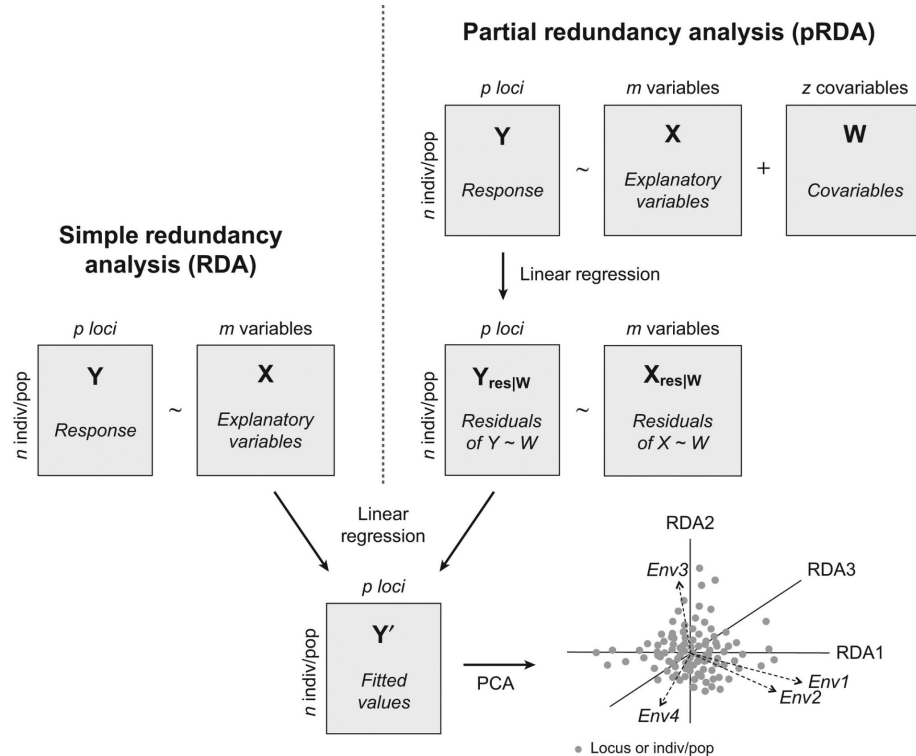| Method | Spatially explicit? | Accounts for neutral structure? | Individual- or population-based sampling? | Other tags |
|---|---|---|---|---|
| BayEnv/BayPASS | No | Yes | Population | Slow, Bayesian, linear |
| RDA | Optional | Optional | Both | Fast, ordination, linear |
| LFMM | No | Optional | Both | Fast, linear |
| GLMM | No | Optional | Both | Slow, linear |
| GF | Yes | No | Both | Nonlinear, map, machine learning |
| SAM/SamBada | No | No | Individual | Logistic |

Table: Anusha P. Bishop

# GEA: the logistics

*Some considerations:*

- May want to minimize **missing data** so as not to bias results; if lots of data are imputed double-check the relationship  between the strength of the association and % missingness (per site)

- Prune out sites that are in **linkage disequilibrium**

- Set a reasonable **MAF threshold**

- **Environmental data**: use realistic layers that you think are affecting your study species!
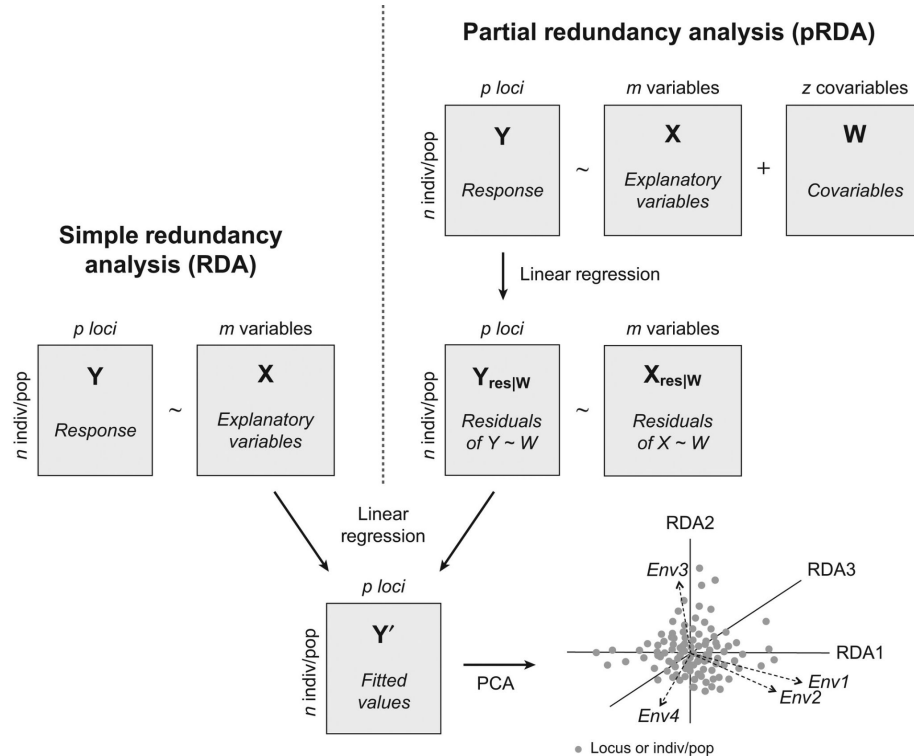
# Today's methods

## Redundancy analysis (RDA)



Outlier detection: RDadapt or Z-scores

Capblancq & Forester (2021) *Methods Ecol. Evol.*

# Today's methods

## Redundancy analysis (RDA)



## Latent factor mixed models (LFMM)

$$Y = XB^T + W + E$$

Latent matrix

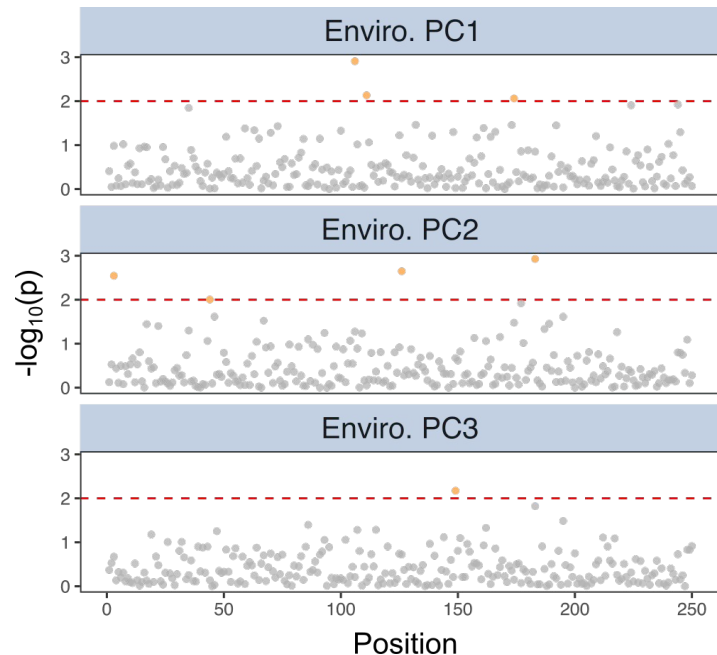Capblancq & Forester (2021) *Methods Ecol. Evol.*; Frichot et al. (2013) *Mol. Biol. Evol.*

# Redundancy analysis (RDA)
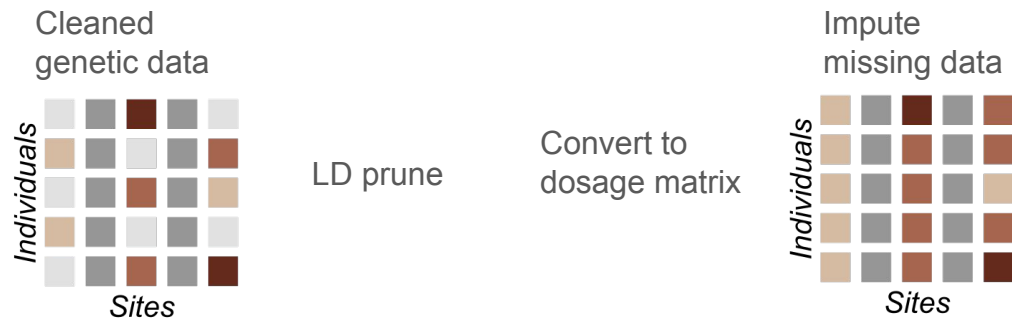
# Latent factor mixed models (LFMM)



**These methods can't accept missing data!** Two choices with different tradeoffs…
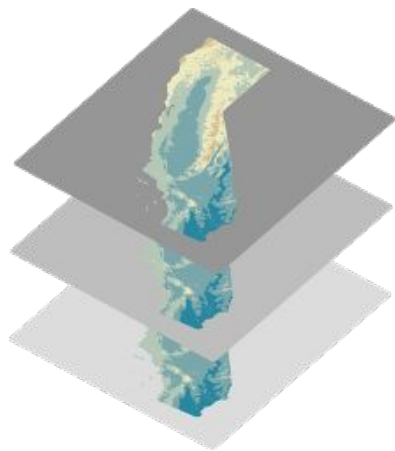
**GEA: the logistics**

**Steps to perform GEA:**

1. Gather genetic data
2. Prune out sites that are in linkage disequilibrium
3. Convert to dosage matrix
4. Impute missing values
5. Gather environmental data (or harvest from online given your sampling)
6. Extract environmental data for each sampling locality
7. Decide on model (and covariables)
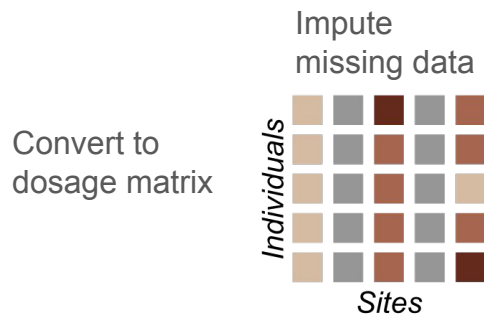8. Run the GEA method!
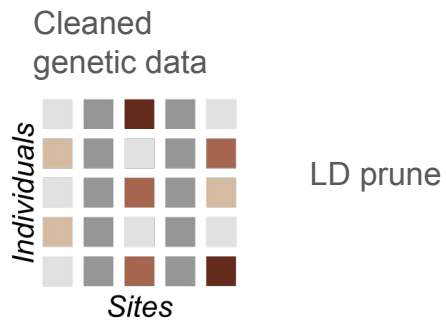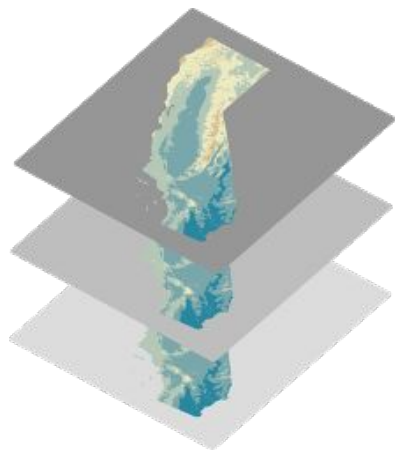
**Genetic data processing**

Cleaned genetic data

*Individuals*

*Sites*

LD prune

Convert to dosage matrix

Impute missing data

*Individuals*

*Sites*

**Environmental and geo data processing**

Enviro data layers

Genetic data processing

Cleaned genetic data

*Individuals*

*Sites*

LD prune

Convert to dosage matrix

Impute missing data

*Individuals*

*Sites*

Environmental and geo data processing

Enviro data layers

Enviro. PC2

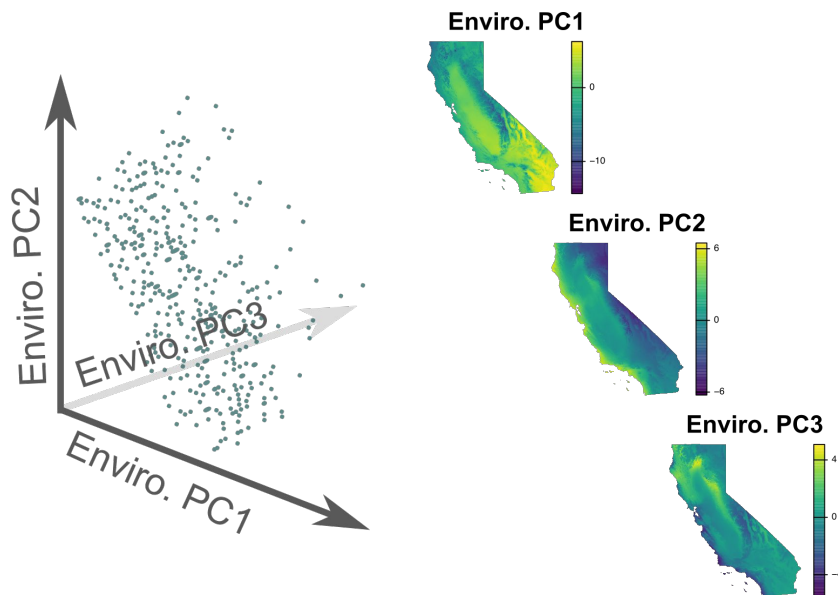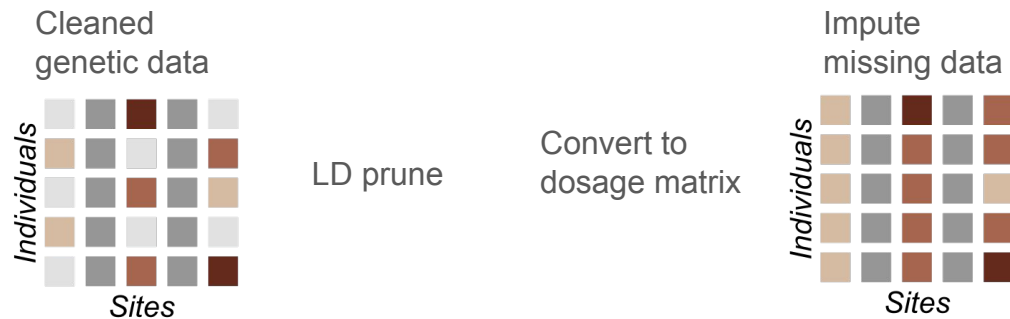Enviro. PC3

Enviro. PC1

Enviro. PC1

Enviro. PC2

Enviro. PC3

Genetic data processing

Cleaned
genetic data

*Individuals*

*Sites*

LD prune

Convert to
dosage matrix

Impute
missing data

*Individuals*

*Sites*

Environmental and geo data processing

Enviro data layers

Test for collinearity

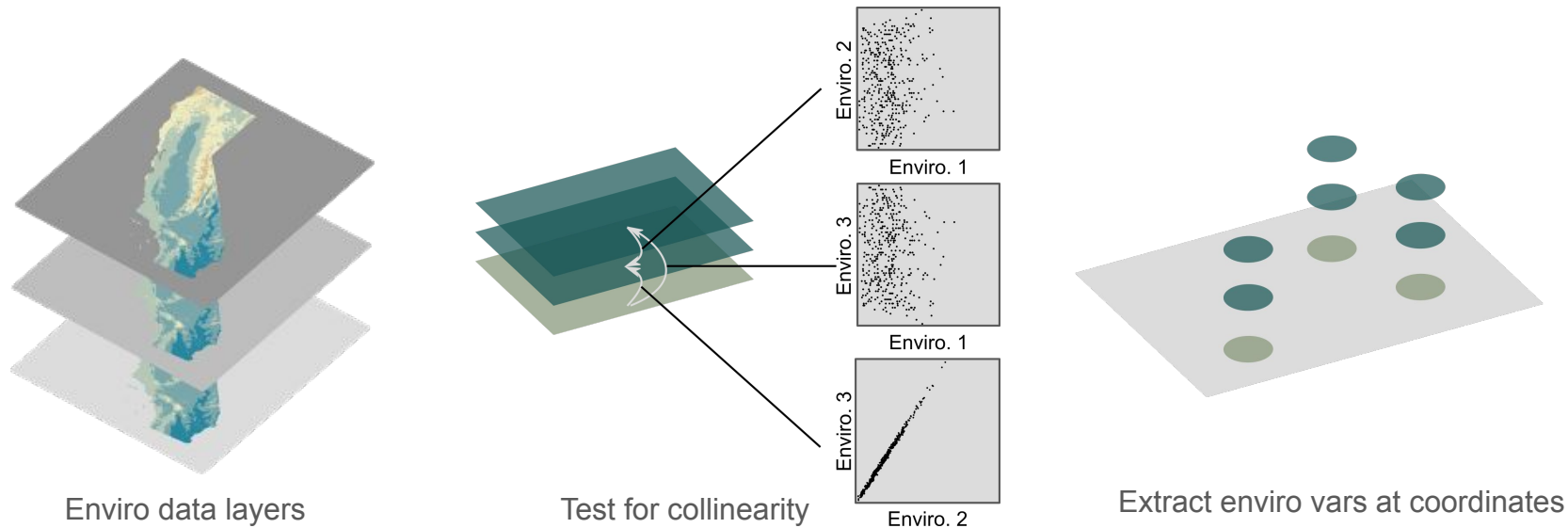Enviro. 2

Enviro. 1

Enviro. 3

Enviro. 1

Enviro. 3

Enviro. 2

Genetic data processing

Cleaned
genetic data

*Individuals*

*Sites*

LD prune

Convert to
dosage matrix

Impute
missing data

*Individuals*

*Sites*

Environmental and geo data processing

Enviro data layers

Test for collinearity

Enviro. 2

Enviro. 1

Enviro. 3

Enviro. 1

Enviro. 3

Enviro. 2

Extract enviro vars at coordinates

| Genetic diversity | wingen | |
| Population structure | TESS | |
| IBD and IBE | MMRR | GDM |
| Adaptive genetic variation | LFMM | RDA |

**algatr**

https://github.com/TheWangLab/algatr

Chambers et al. (2024) *Mol. Ecol. Res.*

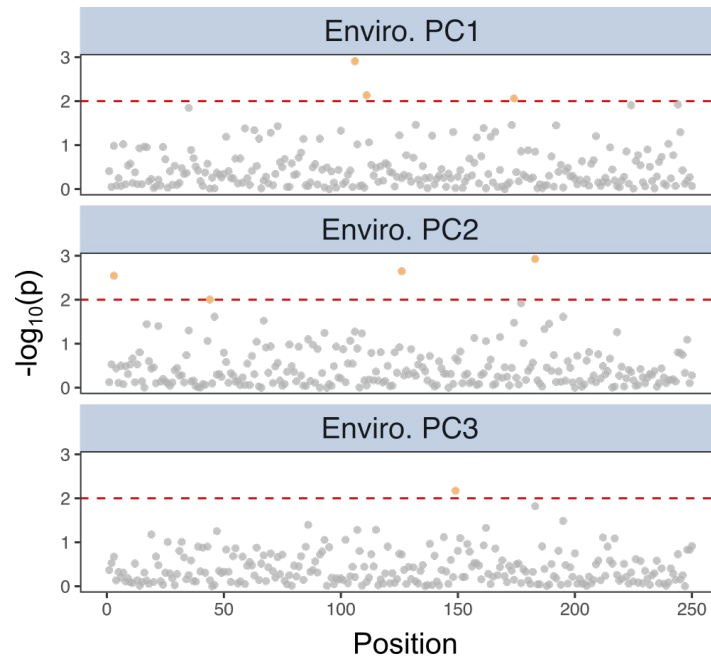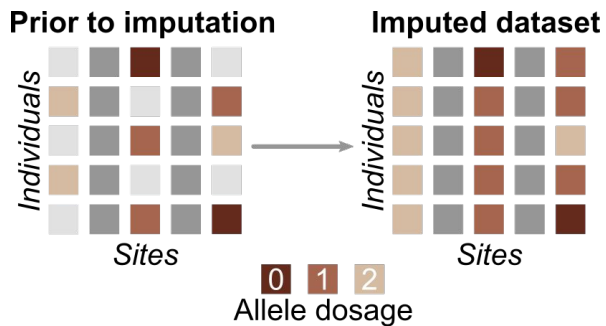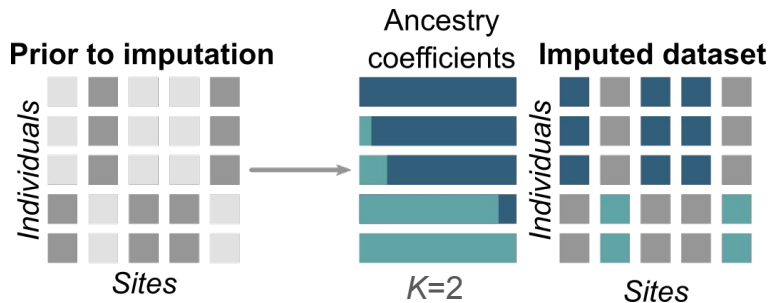**Redundancy analysis (RDA)**

**Latent factor mixed models (LFMM)**

These methods can't accept missing data! Two choices with different tradeoffs…

# Imputation for GEA methods

**Median-based**

**Prior to imputation** → **Imputed dataset**

Individuals / Sites

0 1 2
Allele dosage

**Structure-based**

**Prior to imputation** → Ancestry coefficients / **Imputed dataset**

Individuals / Sites

*K*=2

Present
Absent

# EXERCISE 1

# Example dataset: Bouzid et al. (2022)

# Example dataset: Bouzid et al. (2022)



Pacific
Northwest

Western
Sierra Nevada

Central California

Southern
California

NV

UT

CA

MX

AZ

**Genetic data:**
53 individuals (53 localities)
Individual-based sampling
1,000 SNPs (ddRAD data)

**Environmental data:**
We're going to gather some for this dataset!

# Exercise 1

`?vcf_to_dosage`

# Exercise 1

`?vcf_to_dosage`

# Exercise 1

`?vcf_to_dosage`



vcf_to_dosage {algatr}                                    R Documentation

## Convert a vcf to a dosage matrix

**Description**

Convert a vcf to a dosage matrix

**Usage**

```
vcf_to_dosage(x)
```

**Arguments**

x  can either be an object of class 'vcfR' or a path to a .vcf file



```
# (2) Process genetic data ----------------------------------------------

# Convert the loaded vcf to a dosage matrix using `vcf_to_dosage()`.
###### * YOUR CODE HERE * ######

###### *Q2a*: Do your genetic data have missing values? How do you know?
###### * YOUR ANSWER/CODE HERE * ######
```

# Exercise 1

**?vcf_to_dosage**



vcf_to_dosage {algatr}                                      R Documentation

## Convert a vcf to a dosage matrix

**Description**

Convert a vcf to a dosage matrix

**Usage**

```
vcf_to_dosage(x)
```

**Arguments**

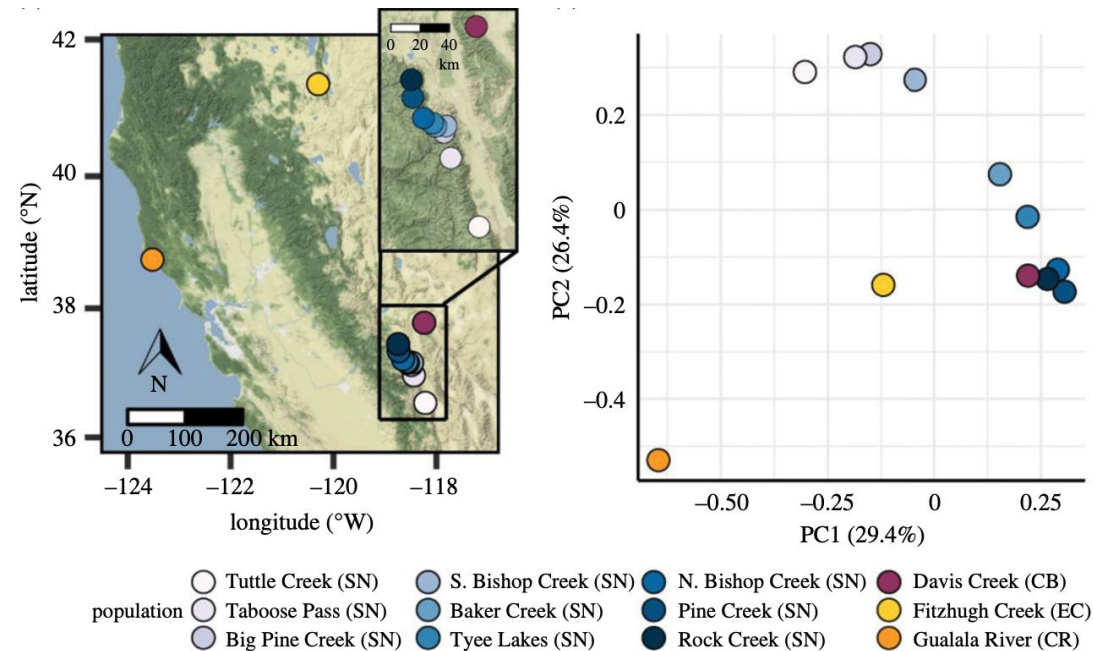x  can either be an object of class 'vcfR' or a path to a .vcf file

```
# (2) Process genetic data ---------------------------------------------

# Convert the loaded vcf to a dosage matrix using `vcf_to_dosage()`.
###### * YOUR CODE HERE * ######

###### *Q2a*: Do your genetic data have missing values? How do you know?
###### * YOUR ANSWER/CODE HERE * ######
```

# Exercise 1

1. Load the example dataset
2. Process genetic data:
   a. Convert vcf to dosage using `vcf_to_dosage()`
   b. Impute missing values using structure-based imputation using `str_impute()`
3. Process environmental data:
   a. Extract environmental values using coordinates using `raster::extract()`
4. Run simple RDA using `rda_run()`
5. Run partial RDA, correcting for geodist using four PCs using `rda_run()`
6. Get outliers using `rda_getoutliers()`
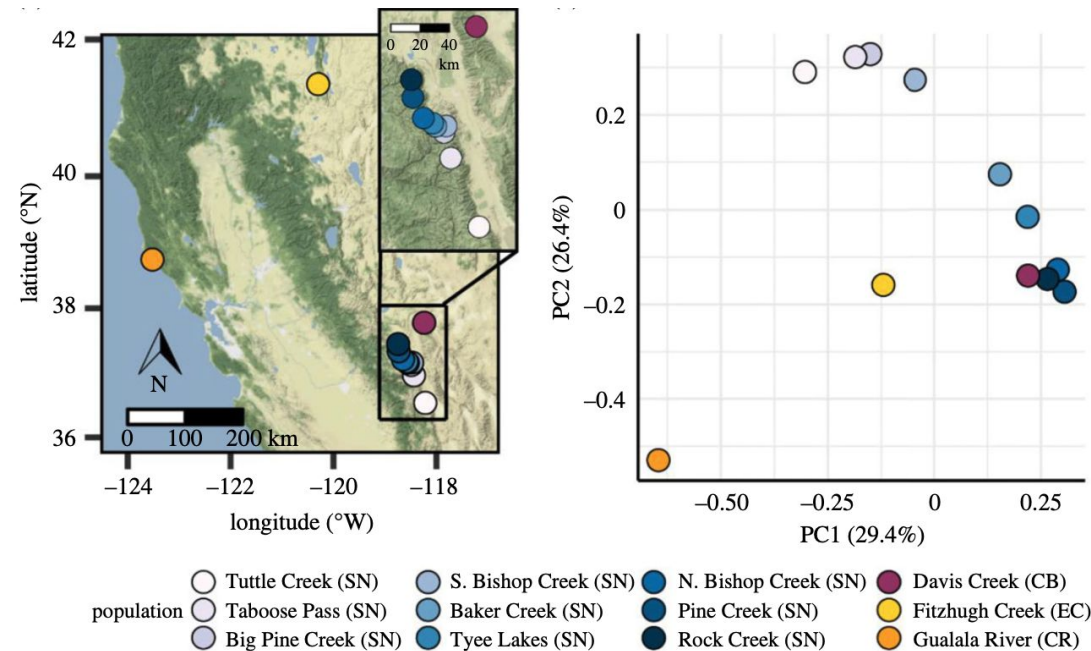7. Build a Manhattan plot and an RDA biplot of results using `rda_plot()`

# EXERCISE 2

# Example dataset: Keller et al. (2023)



Multi-locus genomic signatures of local adaptation to snow across the landscape in California populations of a willow leaf beetle

Abigail G. Keller[1], Elizabeth P. Dahlhoff[2], Ryan Bracewell[3], Kamalakar Chatla[1], Doris Bachtrog[1], Nathan E. Rank[4] and Caroline M. Williams[1]

# Example dataset: Keller et al. (2023)



**Genetic data:**
175 individuals (12 populations)
Site-based sampling
22,323 SNPs (WGS data)

**Environmental data:**
Point data for each population

population
- Tuttle Creek (SN)
- Taboose Pass (SN)
- Big Pine Creek (SN)
- S. Bishop Creek (SN)
- Baker Creek (SN)
- Tyee Lakes (SN)
- N. Bishop Creek (SN)
- Pine Creek (SN)
- Rock Creek (SN)
- Davis Creek (CB)
- Fitzhugh Creek (EC)
- Gualala River (CR)

# Example dataset: Keller et al. (2023)



https://github.com/eachambers/GEA_tutorial/tree/main/Data

**Genetic data:**
175 individuals (12 populations)
Site-based sampling
22,323 SNPs (WGS data)

**Environmental data:**
Point data for each population

# Exercise 2

1. Import and process data using the tidyverse
2. Impute missing genetic data using the median using `simple_impute()`
3. Perform two types of *K* selection to determine how many latent factors you want to use with `select_K()`
4. Run LFMM using `lfmm_run()`
5. Get summary statistics with `lfmm_table()`
6. Make a Manhattan plot of the results using `lfmm_manhattanplot()`