

# Linear Models in R

Helene Wagner, University of Toronto

## Goal: Test for adaptation in Western white pine

- Estimate trait from common garden experiment
- Assess trait heritability
- Compare trait variation to neutral markers (SNPs)
- Explain trait variation by 'climate' and 'geography'

Trait: d13C

(related to water use efficiency)



## Methodological Challenges

1. Design of Western white pine study
2. Introducing SNP data
3. Specifying linear models
  - What type of linear model to fit?
  - Model formulas in R
4. Video 2: Under the hood
  - Degrees of freedom
  - Beyond p-values: effect size

Image source: wikipedia.com

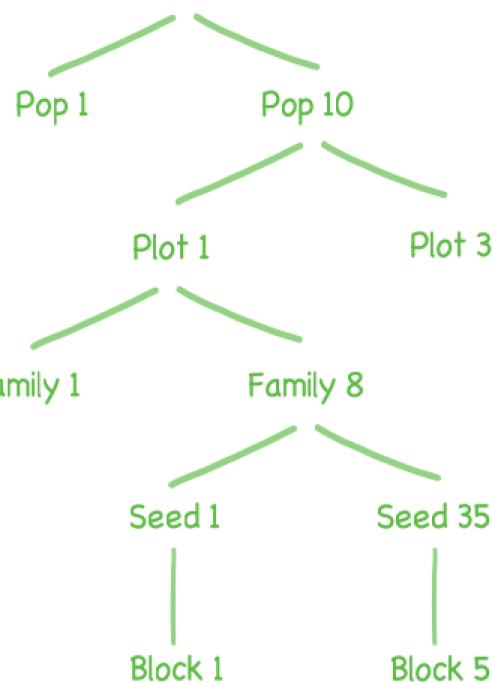
Residual analysis (Week 4)

Spatial linear models (Week 7)

Model selection (Week 12)

# Study Design

**10 Populations**



**3 Plots per Population**

- Spatial coordinates
- 7 bioclimatic variables

**2 - 8 Trees per Plot**

- 164 SNP markers

**5 - 35 Seeds per Family**

**5 Blocks in Common Garden**

- $\delta^{13}\text{C}$  of each seedling

**Seedling trait:**

- 3 - 6 seedlings per family
- common garden
- trait:  $\delta^{13}\text{C}$
- no molecular markers

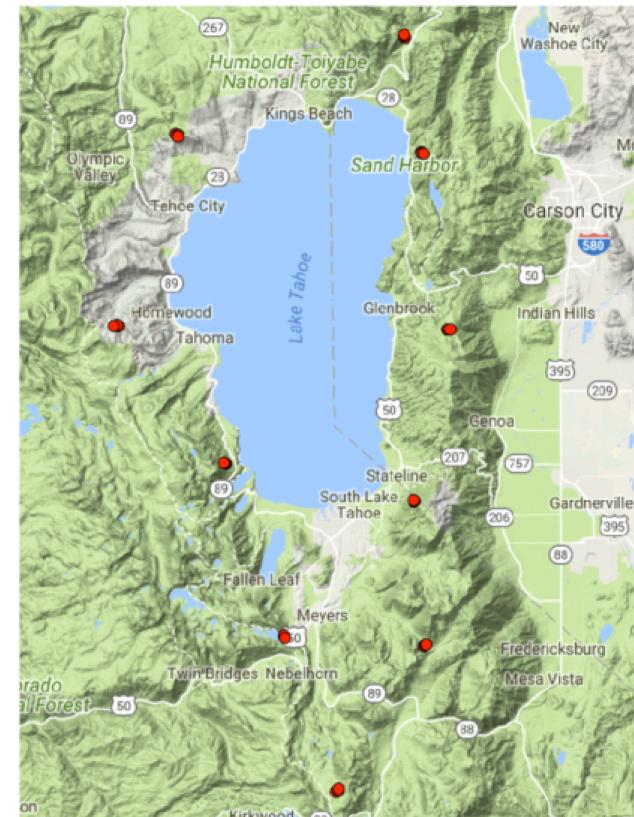


**Estimate heritable trait:**

- Trait =  $G + E + G:E$
- no trait measurement for mother
- estimate from offspring trait
- account for hierarchical sampling
- account for blocks



**Trait - Environment Association:** Variation in heritable trait explained by Climate or Geography?



**Genetic differentiation:**

- Trait: calculate  $Q_{ST}$
- SNPs: calculate  $F_{ST}$
- Compare  $Q_{ST}$  to  $F_{ST}$

# Introducing SNP Data

| family | population | snp102.Plmn | snp106.Plmn |
|--------|------------|-------------|-------------|
| <int>  | <fctr>     | <fctr>      | <fctr>      |
| 59     | blk cyn    | CC          | AC          |
| 60     | blk cyn    | AC          | AA          |
| 61     | blk cyn    | CC          | AA          |
| 63     | blk cyn    | CC          | AA          |
| 64     | blk cyn    | AC          | AA          |
| 65     | blk cyn    | AC          | AA          |
| 67     | blk cyn    | AC          | AA          |
| 69     | blk cyn    | CC          | AA          |
| 72     | blk cyn    | AC          | AA          |
| 73     | blk cyn    | AC          | NA          |

1-10 of 157 rows... Previous 1 2 3 4 5 6 ... 16

## SNP markers:

- co-dominant
- genotype or haplotype?
- wildtype vs. mutant?
- monomorphic?
- many markers!
- neutral or not?

## Here:

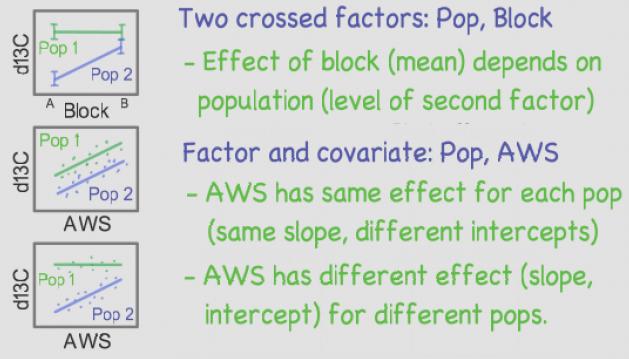
- 164 SNP markers
- monomorphic: 4
- 0 - 20% missing values
- unrelated to trait of interest

# What Type of Linear Model to Fit?

## 1. Variable types?

| Predictors X: | Response Y:             |                     |
|---------------|-------------------------|---------------------|
|               | Quantitative            | Binary              |
|               | Regression (dummy vars) | Logistic Regression |
|               | t - Test                | Odds Ratio          |
| Categorical   | ANOVA (covariates)      | Log-linear Model    |

## 5. Interactions?



## 2. LM, LMM, or GLMM?

| Method: | Conditions: |          |
|---------|-------------|----------|
|         | Normality   | Balanced |
|         | LM          | ✓        |
|         | LMM         | ✓        |
| GLMM    | ✗           | ✗        |

## 4. Fixed or random factors?

Warning: Different philosophies!

|        | Replicate study  | # Levels    |
|--------|------------------|-------------|
| Random | Different levels | Many (> 5?) |
| Fixed  | Same levels      | Often few   |

- Fixed effect: fits one parameter per level
- Random effect: fits 2 parameters

Example:

- Y: seedling d13C
- Hierarchical sampling: Families = trees within plots within populations
- Blocking: common garden with 5 blocks
- Covariate (plot level): AWS (soil available water supply)

## 3. Nested or crossed?

Nested factors: hierarchical sampling



Crossed factors: randomized experiment

|          | Block 1 | Block 2 | Block 3 |
|----------|---------|---------|---------|
| Family 1 | 1       | 1       | 1       |
| Family 2 | 1       | 1       | 2       |
| Family 3 | 1       | 1       | 0       |

balanced      unbalanced

# Specify Model Formula in R

## Factors only:

- Global mean:
- Population mean:
- Additive effects:
- Interaction:
- Hierarchical model:

## With covariate (AWS, standardized):

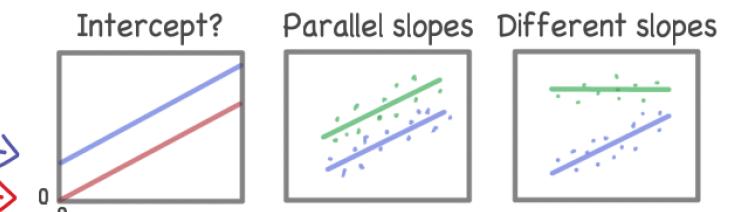
- Simple regression:
- No intercept:
- Parallel slopes:
- Different slopes:

## Fixed effects only:

$d_{13}C \sim 1$   
 $d_{13}C \sim Pop$   
 $d_{13}C \sim Pop + Block$   
 $d_{13}C \sim Pop * Block$   
 $d_{13}C \sim Pop/Family$

## With random effects:

$d_{13}C \sim (1 | Pop)$   
 $d_{13}C \sim (1 | Pop) + (1 | Block)$   
 $d_{13}C \sim (1 | Pop/Family)$



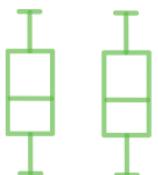
$d_{13}C \sim AWS + (1 | Block)$   
 $d_{13}C \sim AWS + (AWS | Block)$

| Interaction             | Nested              | Intercept   | Sum        | Quadratic |
|-------------------------|---------------------|-------------|------------|-----------|
| $A * B = A + B + A : B$ | $A / B = A + A : B$ | $A = 1 + A$ | $I(A + B)$ | $I(A^2)$  |

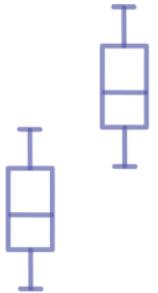
# Degrees of Freedom

## ANOVA F-test (LM)

F follows F-distribution with df(model), df(residual)



$$F = \frac{\text{Variability between groups}}{\text{Variability within groups}} = 1$$



$$F = \frac{\text{Variability between groups}}{\text{Variability within groups}} \gg 1$$

## Assumptions

- Same variability in all groups
- Same group size (balanced)
- Residuals normally distributed
- Observations are independent

## Degrees of freedom

$$\text{df(total)} = (n - 1)$$

- n = # independent observations
- one df used for global mean

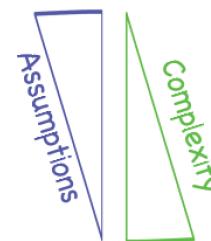
$$\text{df(model)} = (k - 1)$$

- k = number of parameters fitted

$$\text{df(residual)} = n - k$$

## Alternative tests?

- Wald chi-square test
- Likelihood ratio test
- Conditional F-tests (balanced)
- Cond. F-tests with df correction
- MCMC or bootstrap tests



## Non-independence?

### Study design:

- Hierarchical sampling
- Blocking
- Paired samples
- Repeated measures

### Autocorrelation:

- Spatial
- Temporal
- Phylogenetic

### Co-ancestry:

- Kinship
- Population history
- Phylogeography

# Beyond p-values: Effect Size

$$\text{Eta}^2 = \frac{\text{Variability between groups}}{\text{Total variability}} = \frac{\text{Explained variance}}{\text{Total variance}} = R^2$$

## Variance Components

How much does each factor explain?

$$\text{var}(y) = \frac{\sum (y_i - \bar{y})^2}{(n-1)} = \frac{SS}{df_{\text{total}}}$$

|          | SS   | Eta-squared     |
|----------|------|-----------------|
| Factor A | 0.35 | 0.35/0.9 = 0.39 |
| Residual | 0.55 | 0.55/0.9 = 0.61 |
| Total    | 0.90 | 1.00            |

## Mixed Models

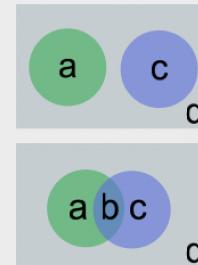
Avoid bias in variance estimates!

|     | Model Comparison | Model Fitting |
|-----|------------------|---------------|
| LM  | LS               | LS            |
| LMM | ML               | REML          |

Best model: lowest AIC      Best predictor: highest SS

## Covariates

What if predictors are correlated?



Crossed factors, balanced:

- Predictors uncorrelated
- No overlap of SS

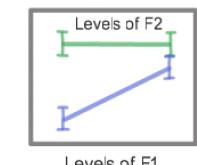
Correlated predictors:

- Overlap of SS
- Results depend on model

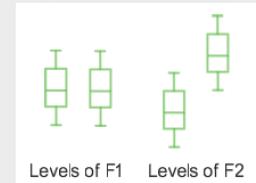
## Interactions

Which effects to interpret?

F1:F2 significant



No interaction



If 'F1 : F2' significant, don't interpret F1, F2!

## Type I, II and III SS

The order of effects matters!

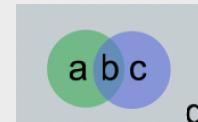


- Type I: for nested factors?
- Type II: balanced, additive?
- Type III: crossed, interaction?

|       | Type I  | Type II | Type III |
|-------|---------|---------|----------|
| F1    | a+d+f+g | a+f     | a        |
| F2    | b+e     | b+e     | b        |
| F1:F2 | c       | c       | c        |

## Variation Partitioning

Shared vs. unique explained variance



- a: Climate alone (7 vars)
- b: Shared variance
- c: Geography alone (2 vars)

Partial regression (accounting for Z):

- Take residuals of Y ~ Z
- Take residuals of X ~ Z
- Fit model with residuals on both sides