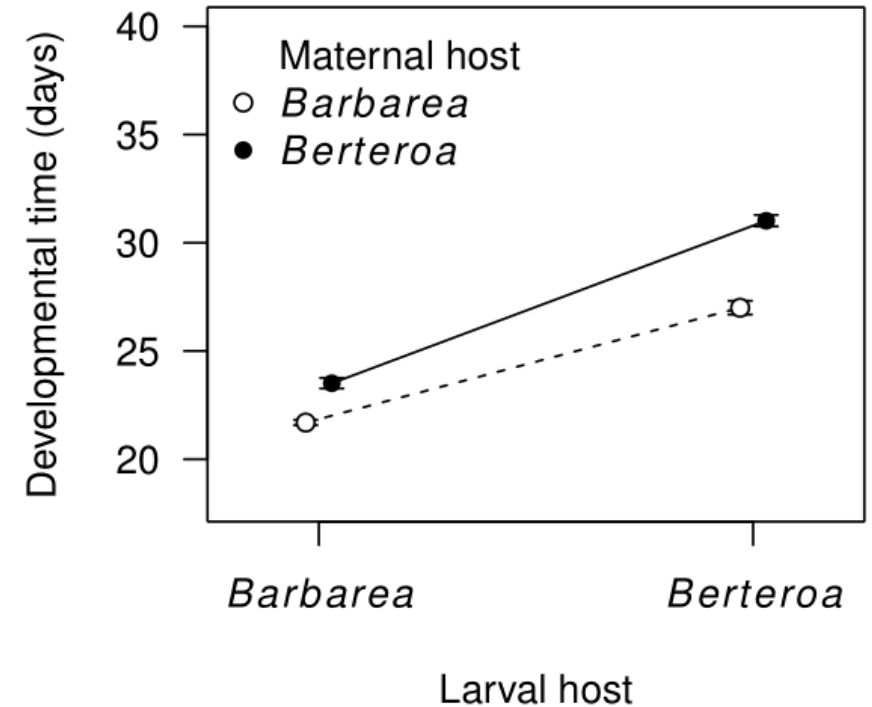


Discussion of exercise 3

- ANOVA analysis of factorial experiment

- “Interaction plots” illustrate the interactive effect of two factors
- Here, the effect of the maternal host is slightly stronger when the larval host is *Berteroa*, than when the larval host is *Barbarea*

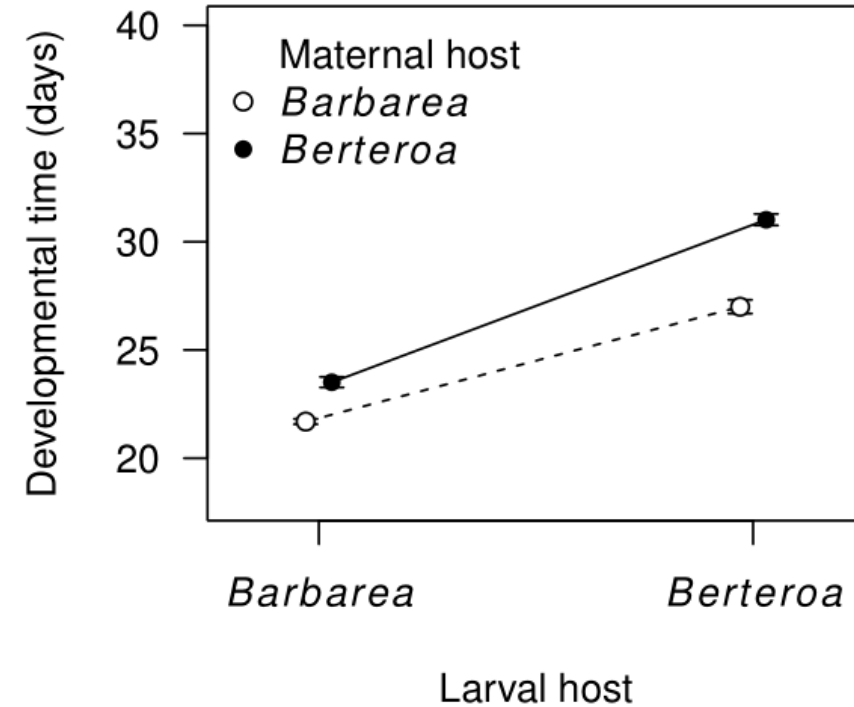


Suggested analysis methods

Methods To assess differences in development time between larvae grown on *Barbarea* and *Berteroa*, and between larvae whose mothers were grown on the same two hosts, we fitted a linear model with development time as response variable, and larval and maternal hosts as predictors, and performed an analysis of variance based on the fitted model. To assess transgenerational effects, we also includes the interaction term between maternal and larval host. Thus, in R syntax, our model took the form

$$DevelopmentTime \sim LarvalHost * MaternalHost$$

Suggested results



Results The larvae developed 22.1% faster when grown on Barbarea than when grown on Berteroa (mean development time = 22.6 and 29.0 days, respectively, $F_{1,283} = 765.21$, Fig. 1). Larvae whose mothers were grown on Barbarea developed 10.7% faster (mean development time = 24.3 and 27.3 days, respectively, $F_{1,283} = 177.90$). The difference in development time between larval host plants was slightly larger when the mother was grown on Berteroa than when the mother was grown on Barbarea (24.2% vs. 19.6% reduction in developmental time on Barbarea, respectively).

Processing and Analysis of Biological Data

BIOS15 2025

Lecture 4. ANCOVA and Multiple Regression

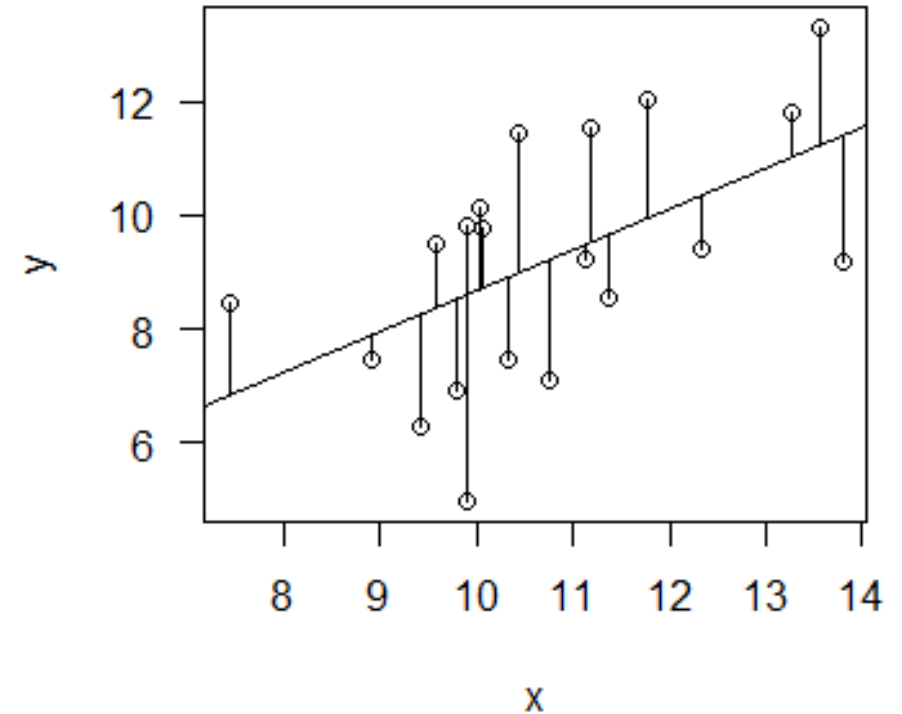
Øystein H. Opedal

$$y_i = \beta_0 + \sum_j \beta_j x_{ij} + \varepsilon_i$$



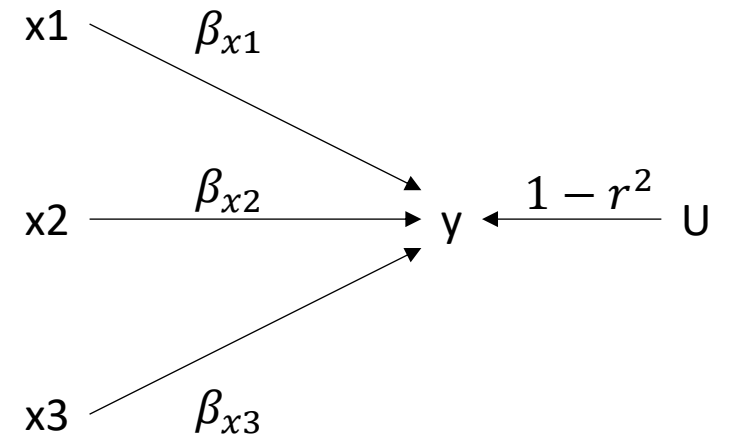
The linear model

- Most of the models we will work with in this course are linear models, that describe how a linear set of predictors relate to a response variable
- A key element of the model is the so-called linear predictor:
- $y_i = \beta_0 + \sum_j \beta_j x_{ij} + \varepsilon_i, \varepsilon \sim N(0, \sigma^2)$
- The term $\varepsilon \sim N(0, \sigma^2)$ means that the residuals (epsilon) are assumed to follow a normal distribution



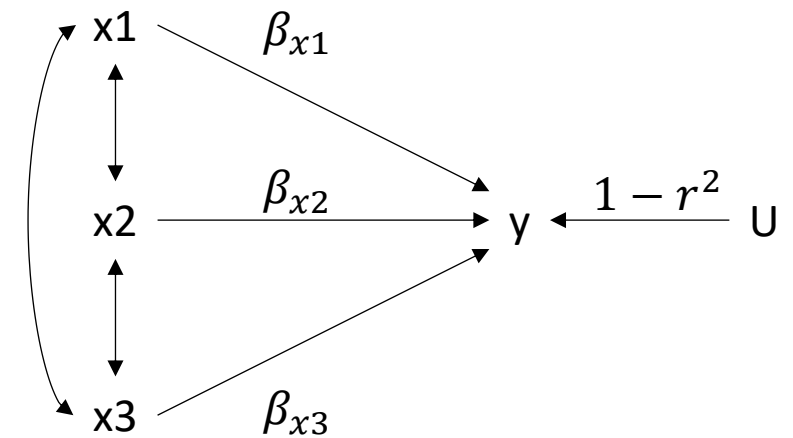
Multiple regression

- A linear model with multiple continuous predictors is called a multiple regression.
- Each slope is estimated while holding the other predictors constant, and are thus **marginal effects**.



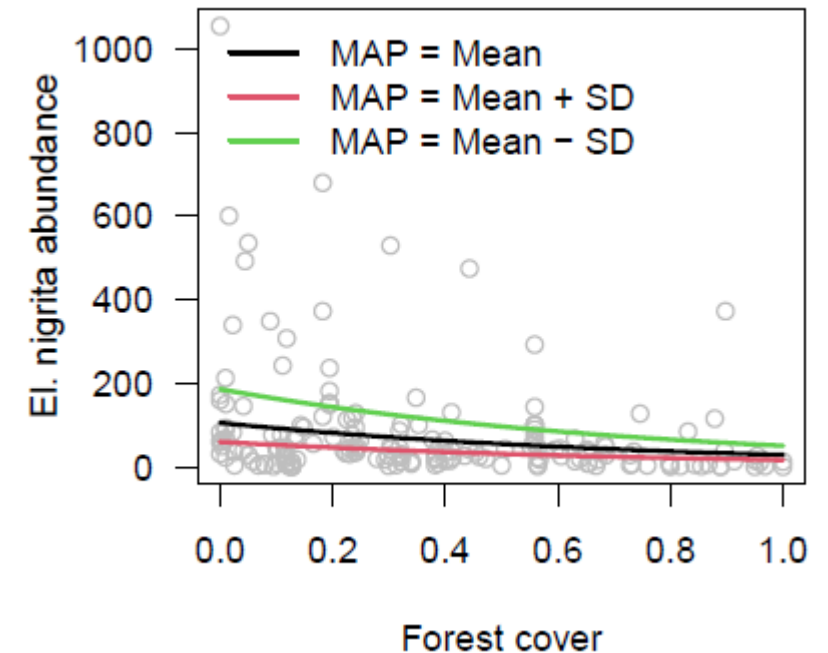
Multiple regression

- A linear model with multiple continuous predictors is called a multiple regression.
- Each slope is estimated while holding the other predictors constant, and are thus **marginal effects**.
- The **net effect** of a predictor (as detected in a univariate model) will also include indirect effects of other variables



Multiple regression

- A linear model with multiple continuous predictors is called a multiple regression.
- Each slope is estimated while holding the other predictors constant, and are thus **marginal effects**.
- The **net effect** of a predictor (as detected in a univariate model) will also include indirect effects of other variables



Example: multivariate selection gradients

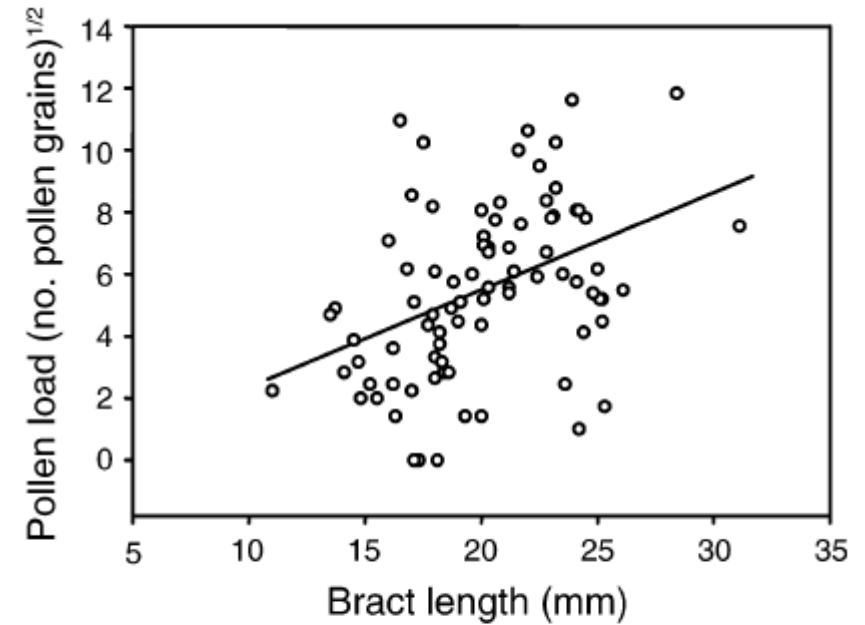
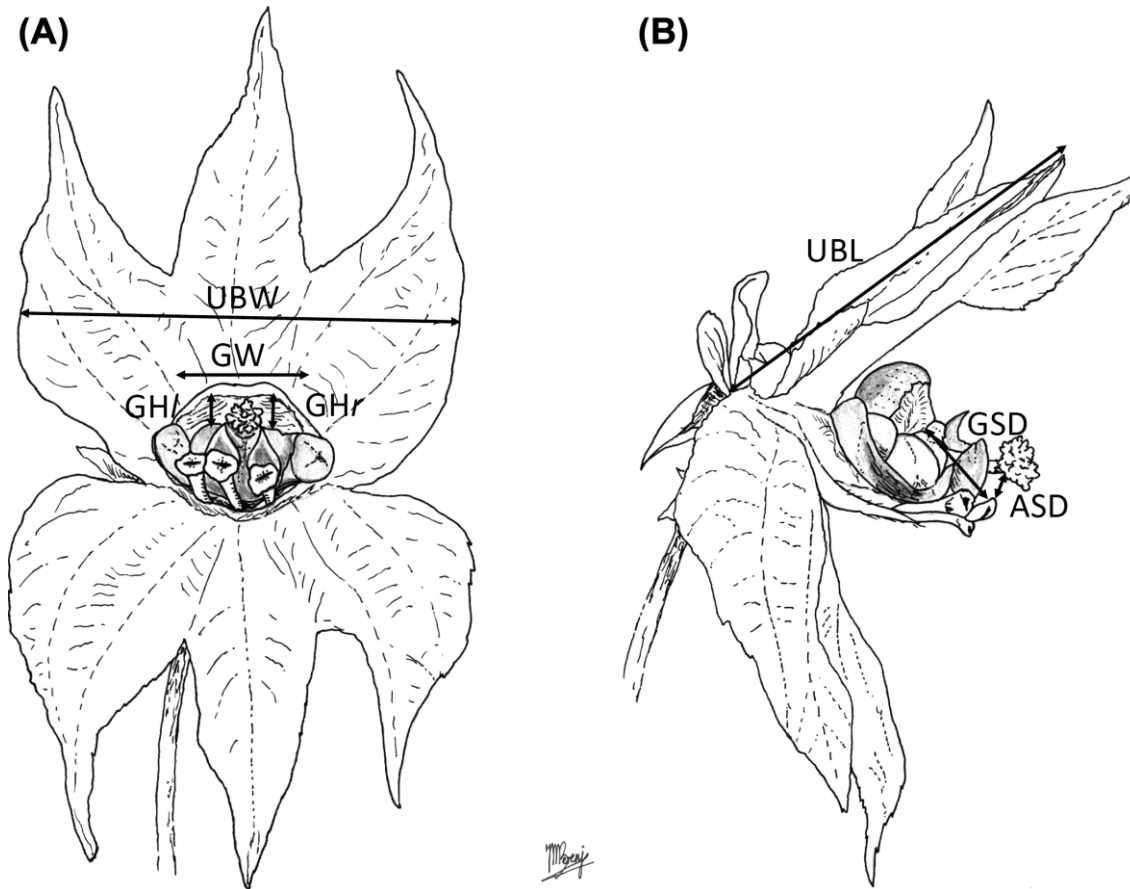
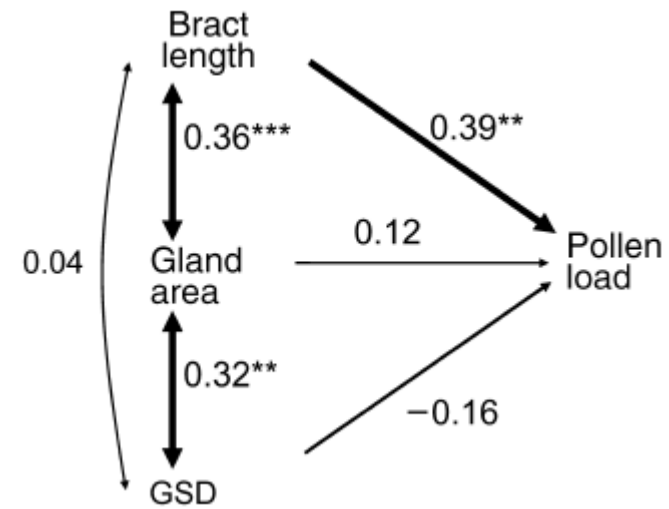
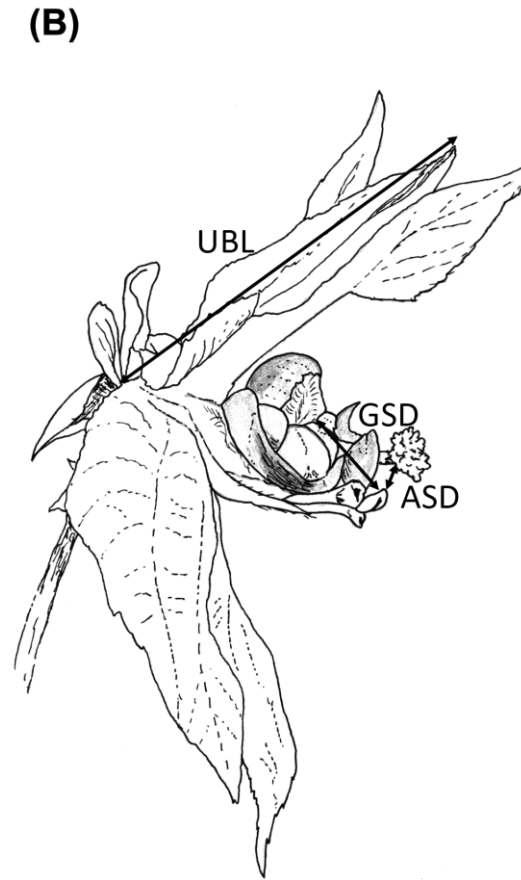
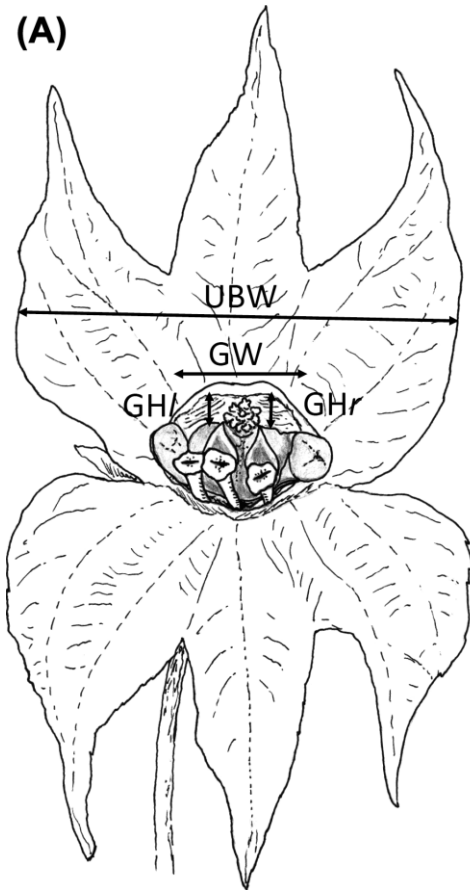


FIG. 4. Effect of natural variation in bract length on the amount of pollen found on the stigmas at the end of the female phase (i.e., pollen transported to stigmas by pollinators). The coefficients (\pm SE) for the regression line are: intercept = -0.88 ± 1.55 , slope = 0.31 ± 0.07 ; $R^2 = 0.18$. See Fig. 5 for significance testing.

Example: multivariate selection gradients



Multiple-regression model in R

- The parameter estimates from a multiple-regression model are marginal effects, i.e. the effect with all other variables held constant

```
##
## Call:
## lm(formula = y ~ x1 + x2)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -8.4276 -2.7240 -0.0065  2.7041  9.7580
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.48722    1.34745   0.362   0.718
## x1           0.64178    0.13246   4.845 2.56e-06 ***
## x2           2.18446    0.06422  34.017 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.618 on 197 degrees of freedom
## Multiple R-squared:  0.8683, Adjusted R-squared:  0.8669
## F-statistic: 649.3 on 2 and 197 DF, p-value: < 2.2e-16
```

Multiple-regression model in R

- The parameter estimates from a multiple-regression model are marginal effects, i.e. the effect with all other variables held constant
- If we standardize the predictors, we can compare the strength of effects across variables (for example in units of standard deviations)

```
##  
## Call:  
## lm(formula = y ~ x1_z + x2_z)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max   
## -8.4276 -2.7240 -0.0065  2.7041  9.7580   
##  
## Coefficients:  
##              Estimate Std. Error t value Pr(>|t|)      
## (Intercept)  19.4090     0.2558   75.866 < 2e-16 ***  
## x1_z          1.2683     0.2618    4.845 2.56e-06 ***  
## x2_z          8.9047     0.2618   34.017 < 2e-16 ***  
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
##  
## Residual standard error: 3.618 on 197 degrees of freedom  
## Multiple R-squared:  0.8683, Adjusted R-squared:  0.8669   
## F-statistic: 649.3 on 2 and 197 DF,  p-value: < 2.2e-16
```

$$z = \frac{x - \bar{x}}{\sigma(x)}$$

Overfitting and multicollinearity

- When we increase the number of variables in the model, we risk problems with “overfitting”, that is fitting a model that explains much variance, but makes poor predictions for independent data
- If the independent (predictor) variables are strongly correlated, this can lead to imprecise estimates (multicollinearity).
- Thus, we often want to select the simplest, most parsimonious model we can (cf. “Occam’s Razor”).
- We can quantify these effects through variance inflation factors, or through cross-validation.

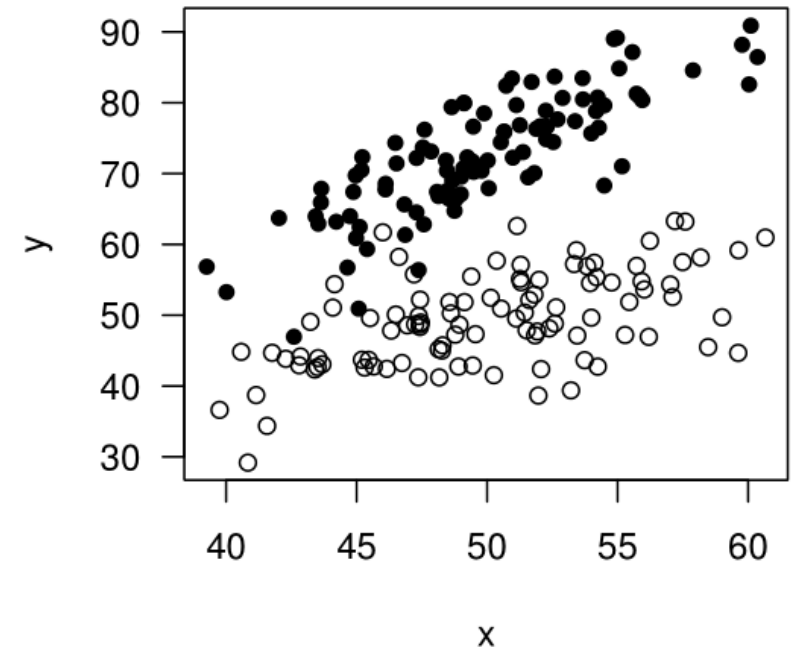
Variance inflation factors

- Variance inflation factors quantify potentially problematic multicollinearity with values greater than 3(5) considered problematic
- Can be quantified by the ratio $1/(1 - r^2)$, where the r^2 is for a model predicting the focal variable, with all other predictors as explanatory variables

$$VIF_i = \frac{1}{1 - r_i^2}$$

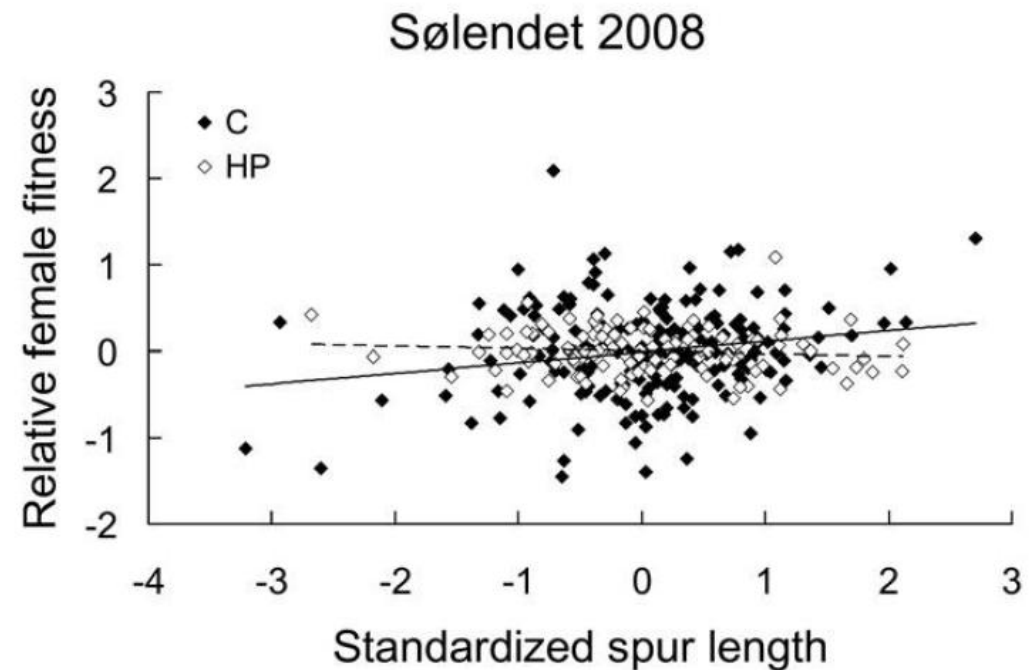
Analysis of Covariance (ANCOVA)

- Analyses of covariance are linear models with both continuous and categorical predictors
- We can use these models to assess whether the slope of a regression differs between groups (e.g. treatments)
- Residuals are assumed to be normally distributed within each group, and variances are assumed to be equal



Example: pollinator-mediated selection

- Analysis of covariance can be used to assess differences in slopes between experimental treatments
- In this case, to show that the relationship between floral spur length and fitness is less steep, as expected, when plants are hand-pollinated



ANCOVA model in R

- The ANOVA table gives the **sums of squares** associated with each predictor, i.e. the sum of square deviations from the predicted value (mean).
- The interaction term tests for heterogeneity of slopes, and the main effect of the grouping variable tests for different intercepts

```
## Analysis of Variance Table
##
## Response: y
##           Df  Sum Sq Mean Sq F value    Pr(>F)
## x           1   4910.5   4910.5  174.019 < 2.2e-16 ***
## gr          1 27641.3 27641.3  979.564 < 2.2e-16 ***
## x:gr         1   849.9   849.9   30.121 1.246e-07 ***
## Residuals 196   5530.7    28.2
```

ANCOVA model in R

- As for all linear models, the summary table gives the parameter estimates, their standard errors, and other model statistics
- In an ANCOVA model, the intercept gives the intercept for the categorical reference level (first level of the factor “gr”, here “Female”).
- The parameter “grMale” gives the contrast between the male and female intercepts.

```
##
## Call:
## lm(formula = y ~ x * gr)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -13.9024  -2.9997   0.0212   3.4958  15.3626
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  12.4340     5.3751   2.313  0.0217 *
## x             0.7371     0.1069   6.897 7.12e-11 ***
## grMale      -21.2230     8.1867  -2.592  0.0102 *
## x:grMale      0.8960     0.1633   5.488 1.25e-07 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 5.312 on 196 degrees of freedom
## Multiple R-squared:  0.8579, Adjusted R-squared:  0.8558
## F-statistic: 394.6 on 3 and 196 DF,  p-value: < 2.2e-16
```

Overview of linear models

- Continuous covariates: (multiple) regression
- Categorical covariates: N-way ANOVA
- Continuous and categorical covariates: ANCOVA

