

# ESS 575: Bayesian Regression Lab

Team England

12 October, 2022

Team England:

- Caroline Blommel
- Carolyn Coyle
- Bryn Crosby
- George Woolsey

[cblommel@mail.colostate.edu](mailto:cblommel@mail.colostate.edu), [carolynm@mail.colostate.edu](mailto:carolynm@mail.colostate.edu), [brcrosby@rams.colostate.edu](mailto:brcrosby@rams.colostate.edu), [george.woolsey@colostate.edu](mailto:george.woolsey@colostate.edu)

## Setup

In this lab you will practice building regression models that are appropriate for a given data set.

- [Coexistence data](#)
- [Abundance data](#)

## Load the data

```
coexist_pth <- "https://nthobbs50.github.io/ESS575/content/labs/coexist.csv"
coverage_pth <- "https://nthobbs50.github.io/ESS575/content/labs/hesp_coverage.csv"

coexist <- read.csv(coexist_pth)
coverage <- read.csv(coverage_pth)
```

## Problem A.

Hein et. al (2012) investigated how environmental conditions influence coexistence between Arctic char (*Salmo trutta*) and pike (*Esox lucius*) in Swedish lakes. Pike were introduced to 151 lakes containing brown trout. Coexistence of the two species was recorded, as  $y_i = 1$  if both species were found in lake  $i$  and 0 otherwise. For each lake, five environmental conditions deemed relevant to coexistence patterns were observed:

- elevation
- upstream catchment area

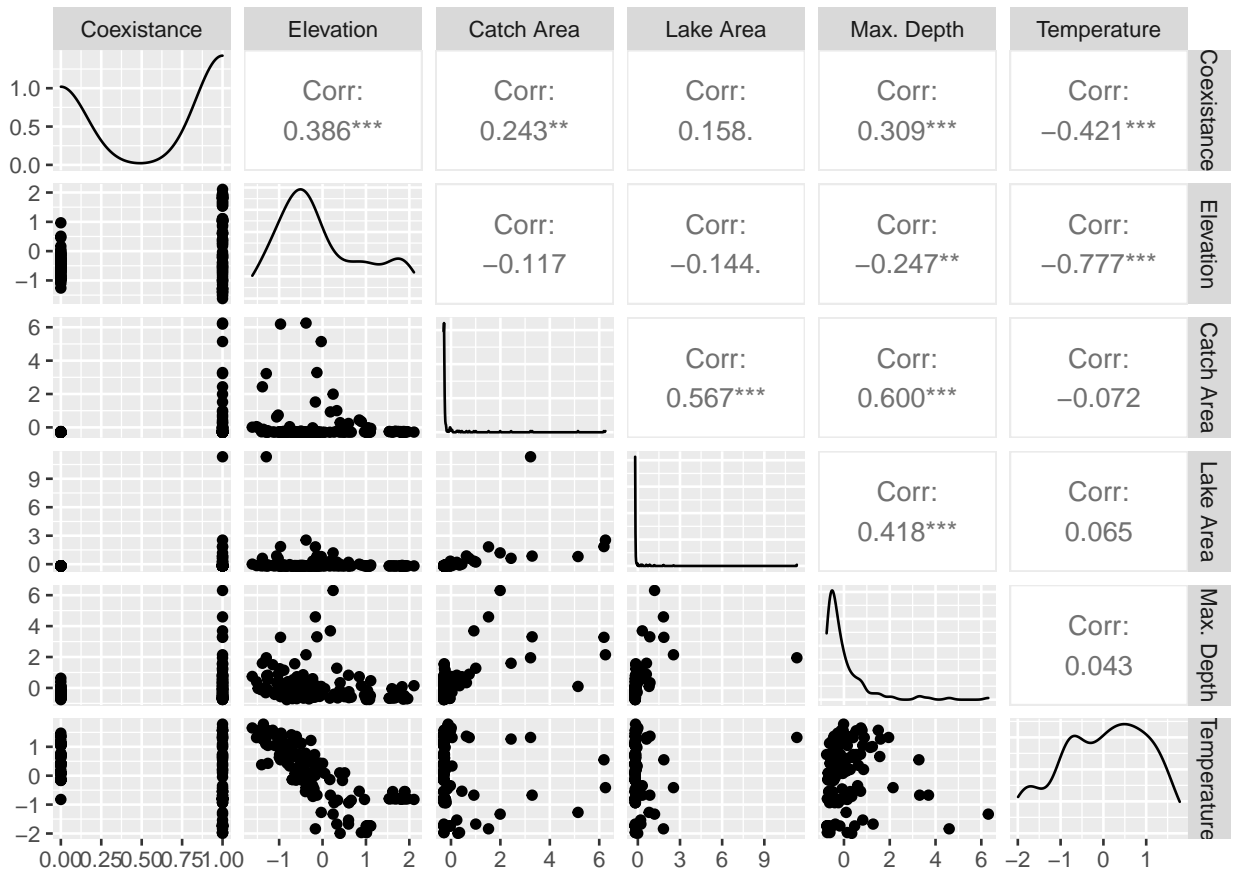
- maximum area
- maximum depth
- mean annual air temperature at outlet

The predictors have been standardized.

## Question 1

Plot the data using `pairs` or `GGally::ggpairs`. What do you notice? How might this impact your modeling choices?

```
GGally::ggpairs(
  data = coexist
  , columns = c("coexist", "elev", "catcharea", "lakearea", "maxdepth", "temp1")
  , columnLabels = c("Coexistence", "Elevation", "Catch Area", "Lake Area", "Max. Depth", "Temperature")
)
```



This grid of plots given by `GGally` shows the empirical density (a.k.a. marginal distribution) of each variable on the diagonal, the scatterplot of points for pairs of variables on the lower triangle, and the Pearson correlation between variables in the upper right triangle.

We notice that the `coexist` outcome variable takes on two values (0 and 1). We also notice that there is one relatively large lake in the data set. A binary data model is appropriate given the support of the data.

## Question 2

We seek to understand the relationship between the probability of coexistence and the five environmental covariates. Write out a reasonable data model.

$$y_i \sim \text{Bernoulli}(p_i)$$

where:

$$p_i = g(\beta_0, \beta_1, \beta_2, \beta_3, \beta_4, \beta_5, x_i) = \text{inverse logit}(\beta_0 + \beta_1 x_{1i} + \cdots + \beta_5 x_{5i}) = \frac{\exp(\beta_0 + \beta_1 x_{1i} + \cdots + \beta_5 x_{5i})}{1 + \exp(\beta_0 + \beta_1 x_{1i} + \cdots + \beta_5 x_{5i})}$$

## Question 3

Assume we know very little about the impact of the environmental covariates on coexistence. Specify appropriate prior distributions for all unknown parameters. Explain your choice of priors.

Assuming vague priors on the intercept and slope can be accomplished by setting the variance  $\sigma^2 = 2.7$  normally distributed with a mean of 0; e.g.,  $\beta_0 \sim \text{normal}(0, 2.7)$ ,  $\beta_1 \sim \text{normal}(0, 2.7)$ ,  $\cdots$ ,  $\beta_5 \sim \text{normal}(0, 2.7)$ . Such that:

$$[\beta_0, \beta_1, \cdots, \beta_5 \mid \mathbf{y}] \propto \prod_{i=1}^n \text{Bernoulli}(y_i \mid g(\beta_0, \beta_1, \cdots, \beta_5, x_i)) \times \text{normal}(\beta_0 \mid 0, 2.7) \times \text{normal}(\beta_1 \mid 0, 2.7) \times \cdots \times \text{normal}(\beta_5 \mid 0, 2.7)$$

## Question 4

Write the expression for the posterior in terms of the joint distribution of the parameters and data.

$$[\beta_0, \beta_1, \cdots, \beta_5 \mid \mathbf{y}] \propto \prod_{i=1}^n \text{Bernoulli}(y_i \mid g(\beta_0, \beta_1, \cdots, \beta_5, x_i)) \times \text{normal}(\beta_0 \mid 0, 2.7) \times \text{normal}(\beta_1 \mid 0, 2.7) \times \cdots \times \text{normal}(\beta_5 \mid 0, 2.7)$$

where:

$$p_i = g(\beta_0, \beta_1, \beta_2, \beta_3, \beta_4, \beta_5, x_i) = \text{inverse logit}(\beta_0 + \beta_1 x_{1i} + \cdots + \beta_5 x_{5i}) = \frac{\exp(\beta_0 + \beta_1 x_{1i} + \cdots + \beta_5 x_{5i})}{1 + \exp(\beta_0 + \beta_1 x_{1i} + \cdots + \beta_5 x_{5i})}$$

## Question 5

Write out the JAGS code for the model.

```
# define explanatory variables
exp_vars <- c("elev", "catcharea", "lakearea", "maxdepth", "temp1")
# list of data
data = list(
  n = nrow(coexist) # n is required in the JAGS program to index the for structure
  , y = as.double(coexist$coexist)
  , x_mtrx = as.matrix(
```

```

    coexist %>%
      dplyr::mutate(intercept = 1) %>%
      dplyr::select(c("intercept", exp_vars)) %>%
      #the execution of JAGS is about 5 times faster on double precision than on integers.
      dplyr::mutate_if(is.numeric, as.double)
  )
)

## JAGS Model
model{
  # priors
  b0 ~ dnorm(0, (1/2.7))
  b1 ~ dnorm(0, (1/2.7))
  b2 ~ dnorm(0, (1/2.7))
  b3 ~ dnorm(0, (1/2.7))
  b4 ~ dnorm(0, (1/2.7))
  b5 ~ dnorm(0, (1/2.7))
  # likelihood
  for (i in 1:n) {
    p[i] <- ilogit(
      b0
      + b1*x_mtrx[i,2] # or b1*x_mtrx[i,"elev"]
      + b2*x_mtrx[i,3] # or b2*x_mtrx[i,"catcharea"]
      + b3*x_mtrx[i,4]
      + b4*x_mtrx[i,5]
      + b5*x_mtrx[i,6]
    )
    y[i] ~ dbern(p[i])
  }
}

```

## Problem B

Kembel and Cahill Jr. (2011) collected data from temperate grassland plant communities in Alberta, Canada. Twenty-seven plots are established, and the slope, slope position, aspect, and relative moisture of each plot is recorded. For each plot, the abundance of several species is recorded, interpreted as the proportion of land area covered by the species. The land coverage by needle-and-thread grass (*Hesperostipa comata* ssp. *comata*) is considered here.

### Question 1

Plot the data using `pairs` or `GGally::ggpairs`. What do you notice? How might this impact your modeling choices?

```

my_vars <- c("coverage", "slope", "aspect", "slope_position", "rel_moisture")
# plot
coverage %>% names

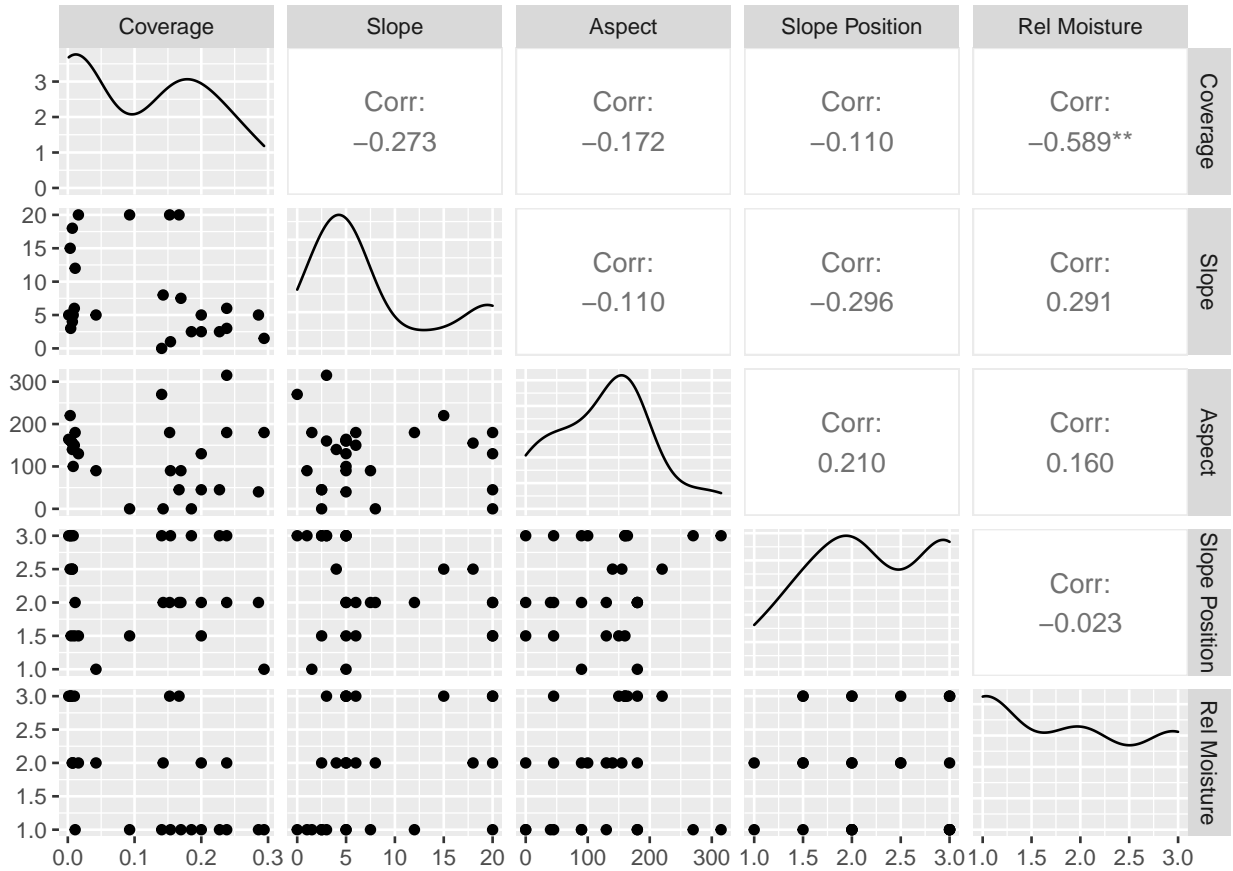
```

```

## [1] "coverage"      "slope"         "aspect"        "slope_position"
## [5] "rel_moisture"

```

```
GGally::ggpairs(
  data = coverage
  , columns = my_vars
  , columnLabels = my_vars %>% stringr::str_replace_all("[:punct:]", " ") %>% stringr::str_to_title()
)
```



Slope position and relative moisture appear to be discrete data (not continuous). Slope should be constrained to values between and including 0-180 while aspect should be constrained to values between and including 0-360. Coverage is a proportion and can take on any value between and including 0 and 1.

## Question 2

Let  $y_i$  represent the proportion of land covered by *Hesperostipa comata* ssp. *comata*. What deterministic model would you use to predict the mean of  $y_i$  as a function of four covariates  $x_1 \cdots x_4$ ?

$$y_i \sim \text{beta}(\alpha, \beta)$$

with a mean ( $\mu$ ) function of:

$$\mu_i = g(\beta_0, \beta_1, \dots, \beta_4, x_i) = \text{inverse logit}(\beta_0 + \beta_1 x_{1i} + \dots + \beta_4 x_{4i}) = \frac{\exp(\beta_0 + \beta_1 x_{1i} + \dots + \beta_4 x_{4i})}{1 + \exp(\beta_0 + \beta_1 x_{1i} + \dots + \beta_4 x_{4i})}$$

### Question 3

Write the likelihood for an observation,  $y_i$ .

The likelihood function for  $y_i$  is:

$$y_i \sim \text{beta}(a_i, b_i)$$

where:

$$a_i = \frac{\mu_i^2 - \mu_i^3 - \mu_i \sigma^2}{\sigma^2}$$

$$b_i = \frac{\mu_i - 2\mu_i^2 + \mu_i^3 - \sigma^2 + \mu_i \sigma^2}{\sigma^2}$$

### Question 4

What would you use for vague priors on the coefficients to assure that the prior on  $\mu_i$  is vague?

Assuming vague priors on the intercept and slope can be accomplished by setting the variance  $\sigma^2 = 2.7$  normally distributed with a mean of 0; e.g.,  $\beta_0 \sim \text{normal}(0, 2.7)$ ,  $\beta_1 \sim \text{normal}(0, 2.7)$ ,  $\dots$ ,  $\beta_4 \sim \text{normal}(0, 2.7)$ . Such that:

$$[\beta_0, \beta_1, \dots, \beta_4, \sigma^2 \mid \mathbf{y}] \propto \prod_{i=1}^n \text{beta}(y_i \mid a_i, b_i) \times \text{normal}(\beta_0 \mid 0, 2.7) \times \dots \times \text{normal}(\beta_4 \mid 0, 2.7) \times \text{uniform}(\sigma \mid 0, 100)$$

## Question 5

Express the posterior distribution as proportional to joint distribution for your model.

$$y_i \sim \text{beta}(a_i, b_i)$$

$$a_i = \frac{\mu_i^2 - \mu_i^3 - \mu_i \sigma^2}{\sigma^2}$$

$$b_i = \frac{\mu_i - 2\mu_i^2 + \mu_i^3 - \sigma^2 + \mu_i \sigma^2}{\sigma^2}$$

$$\mu_i = g(\beta_0, \beta_1, \dots, \beta_4, x_i) = \text{inverse logit}(\beta_0 + \beta_1 x_{1i} + \dots + \beta_4 x_{4i}) = \frac{\exp(\beta_0 + \beta_1 x_{1i} + \dots + \beta_4 x_{4i})}{1 + \exp(\beta_0 + \beta_1 x_{1i} + \dots + \beta_4 x_{4i})}$$

$$[\beta_0, \beta_1, \dots, \beta_4, \sigma^2 \mid \mathbf{y}] \propto \prod_{i=1}^n \text{beta}(y_i \mid a_i, b_i) \times \text{normal}(\beta_0 \mid 0, 2.7) \times \dots \times \text{normal}(\beta_4 \mid 0, 2.7) \times \text{uniform}(\sigma \mid 0, 100)$$

asdf