# Bayesian Regression Bayesian Models for Ecologists

Becky Tang

#### June 07, 2024



#### Outline

- Be able to write proper Bayesian regression models for different types of data.
- Appreciate one-to-one relationship between math and JAGS code.
- Be able to interpret coefficients of general linear models.

# A great follow-up

This book should be in your library:



Data Analysis Using Regression and Multilevel/Hierarchical Models

ANDREW GELMAN JENNIFER HILL

#### The general Bayesian set-up

Recall that the posterior distribution of the unobserved quantities conditional on the observed ones is proportional to their joint distribution:

$$[\boldsymbol{\theta}, \sigma^2 | \mathbf{y}] \propto [\boldsymbol{\theta}, \sigma^2, \mathbf{y}].$$

The joint distribution can be factored into a likelihood and priors for simple Bayesian models:

$$\begin{bmatrix} \boldsymbol{\theta}, \boldsymbol{\sigma}^2, \boldsymbol{y} \end{bmatrix} = \begin{bmatrix} \boldsymbol{y} \mid \boldsymbol{\theta}, \boldsymbol{\sigma}^2 \end{bmatrix} \begin{bmatrix} \boldsymbol{\theta}, \boldsymbol{\sigma}^2 \end{bmatrix} \stackrel{\textit{ind.}}{=} \begin{bmatrix} \boldsymbol{y} \mid \boldsymbol{\theta}, \boldsymbol{\sigma}^2 \end{bmatrix} \begin{bmatrix} \boldsymbol{\theta} \end{bmatrix} \begin{bmatrix} \boldsymbol{\sigma}^2 \end{bmatrix}$$

A deterministic model of an ecological or socioenvironmental process is embedded in the likelihood like this:

$$[\boldsymbol{\theta}, \boldsymbol{\sigma}^2, \mathbf{y}] \propto [\mathbf{y} \mid g(\boldsymbol{\theta}, x), \boldsymbol{\sigma}^2] [\boldsymbol{\theta}] [\boldsymbol{\sigma}^2]$$

#### Simple Bayesian regression models

We use likelihood to connect the underlying process to data:

$$\underbrace{[y_i \mid \mu_i, \sigma^2]}_{\text{stochastic model}}, \quad i = 1, \dots, n$$

We formulate the deterministic model:

$$\mu_i = \underline{g(\beta, \mathbf{x}_i)}, \quad i = 1, \dots, n$$

deterministic model

where  $\beta$  is a vector of regression coefficients and  $\mathbf{x}_i$  is a vector of predictor variables.

Assuming conditional independence of the data,

$$\left[\beta,\sigma^2 \mid \mathbf{y}\right] \propto \prod_{i=1}^n \left[y_i \mid g(\beta,x_i),\sigma^2\right] \times \left[\beta\right] \left[\sigma^2\right]$$

We choose appropriate deterministic functions (linear or non-linear) and appropriate probability distributions to compose specific models.

#### Identical notation

$$y_i = g(eta, x_i) + arepsilon_i$$
  
 $arepsilon_i \sim \mathsf{Normal}(0, \sigma^2)$ 

is the same as:

$$y_i|\beta, \sigma^2 \sim \text{Normal}(g(\beta, x_i), \sigma^2),$$

but the second notation is much more flexible because it generalizes to distributions that do not have additive errors.

## You don't have to be Normal!

Assuming you have one predictor x:

Data (y-values)	Distribution	"Mean" function	Link
continuous, real valued	Normal	$\mu = eta_0 + eta_1  imes$	NA (i.e. identity)
discrete, strictly positive	Poisson	$\mu=e^{eta_0+eta_1 imes}$	$\log(\mu) = \beta_0 + \beta_1 x$
0 or 1	Bernoulli	$\mu = \frac{\exp(\beta_0 + \beta_1 x)}{\exp(\beta_0 + \beta_1 x) + 1}$	$\operatorname{logit}(\mu) = \log\left(\frac{\mu}{1-\mu}\right) = \beta_0 + \beta_1 x$
[0,1]	Beta	$\mu = \frac{\exp(\beta_0 + \beta_1 x)}{\exp(\beta_0 + \beta_1 x) + 1}$	$\operatorname{logit}(\mu) = \log\left(\frac{\mu}{1-\mu}\right) = \beta_0 + \beta_1 x$
continuous, strictly positive, variance ↑ as a f(mean)	lognormal	$\mu=e^{eta_0+eta_1 imes}$	$\log(\mu) = eta_0 + eta_1  imes$
continuous, strictly positive, constant variance	Gamma	$\mu=e^{eta_0+eta_1 imes}$	$\log(\mu) = \beta_0 + \beta_1 x$

#### Continuous and real valued data

Suppose you have collected some continuous data  $\mathbf{y} = (-10.7, -4.3, \cdots, 49)$  at *n* sites, along with a predictor  $x_i$  measured at each site *i* which you believe is likely to affect these measurements. Write a model regressing *y* on *x* as follows:

- Choose a specific stochastic and deterministic model.
- Specify (vague) priors for your parameters.
- Write out the DAG and express posterior distribution as proportional to joint distribution for your model.
- Write the JAGS code for the model.
- Interpret the coefficients of your model.

Normal data, continuous and real valued Stochastic model:

$$y_i | \beta_0, \beta_1, \sigma^2 \stackrel{ind.}{\sim} \mathsf{Normal}(g(\beta_0, \beta_1, x_i), \sigma^2)$$
  $i = 1, ..., n$ 

Deterministic model:

$$\mu_i = g(\beta_0, \beta_1, x_i) = \beta_0 + \beta_1 x_i$$

Priors:

$$\begin{array}{ll} \beta_0 \sim ? & [\beta_0] = ? \\ \beta_1 \sim ? & [\beta_1] = ? \\ \sigma^2 \sim ? & [\sigma^2] = ? \end{array}$$

# Normal data, continuous and real valued DAG:

Posterior distribution:

$$\begin{split} \begin{bmatrix} \beta_0, \beta_1, \sigma \mid \mathbf{y} \end{bmatrix} & \propto & [\beta_0, \beta_1, \sigma, \mathbf{y}] \\ & \propto & [\mathbf{y} \mid \beta_0, \beta_1, \sigma] [\beta_0] [\beta_1] [\sigma] \\ & \propto & \prod_{i=1}^n \operatorname{Normal}(y_i \mid g(\beta_0, \beta_1, x_i), \sigma^2) \\ & \times \operatorname{Normal}(\beta_0 \mid 0, 1000) \operatorname{Normal}(\beta_1 \mid 0, 1000) \\ & \times \operatorname{uniform}(\sigma \mid 0, 100) \\ & g(\beta_0, \beta_1, x_i) & = & \beta_0 + \beta_1 x_i \end{split}$$

# Normal data, continuous and real valued JAGS code for the model:

```
b0 ~ dnorm(0, .001)
b1 ~ dnorm(0, .001)
sigma ~ dunif(0, 100)
tau <- 1/sigma^2
for (i in 1:length(y)){
    mu[i] <- b0 + b1 * x[i]
    y[i] ~ dnorm(mu[i], tau)
}</pre>
```

Interpretation:

- $\beta_0$  : expected outcome when x = 0
- $\beta_1$  : average change in the outcome for a one-unit change in x
- $\sigma$  : std deviation of the outcomes about their respective means

#### Counts, discrete and non-negative

You have collected some count data (y = 12, 17, 1, 0, 31, ..., 25) at *n* sites, along with a covariate  $x_i$  at each location which you believe is likely to affect these counts. Write a model regressing y on x.

- Choose a specific stochastic and deterministic model.
- Specify (vague) priors for your parameters.
- Write out the DAG and express posterior distribution as proportional to joint distribution for your model.
- Write the JAGS code for the model.
- Interpret the coefficients of your model.

Stochastic model:

$$y_i|\beta_0,\beta_1 \stackrel{ind.}{\sim} \mathsf{Poisson}(g(\beta_0,\beta_1,x_i))$$
  $i=1,...,n$ 

Deterministic model:

$$\mu_i = g(\beta_0, \beta_1, x_i) = e^{\beta_0 + \beta_1 x_i}$$

Priors:

DAG:

Posterior distribution:

$$\begin{bmatrix} \beta_0, \beta_1 \mid \mathbf{y} \end{bmatrix} \propto \prod_{i=1}^n \operatorname{Poisson}(y_i \mid g(\beta_0, \beta_1, x_i)) \\ \times \operatorname{Normal}(\beta_0 \mid 0, 1000) \operatorname{Normal}(\beta_1 \mid 0, 1000) \\ g(\beta_0, \beta_1, x_i) = e^{\beta_0 + \beta_1 x_i}$$

JAGS code:

```
b0 ~ dnorm(0, .001)
b1 ~ dnorm(0, .001)
for(i in 1:length(y)){
    log(mu[i]) <- b0 + b1 * x[i]
    y[i] ~ dpois(mu[i])
}</pre>
```

or

```
mu[i] <- exp(b0 + b1 * x[i])
y[i] ~ dpois(mu[i])</pre>
```

$$\begin{aligned} \begin{bmatrix} \beta_0, \beta_1 \mid \mathbf{y} \end{bmatrix} & \propto & \prod_{i=1}^n \operatorname{Poisson}(y_i \mid g(\beta_0, \beta_1, x_i)) \\ & \times \operatorname{Normal}(\beta_0 \mid 0, 1000) \operatorname{Normal}(\beta_1 \mid 0, 1000) \\ \mu_i &= g(\beta_0, \beta_1, x_i) &= e^{\beta_0 + \beta_1 x_i} = e^{\beta_0} e^{\beta_1 x_i} \end{aligned}$$

Interpretation: (Exponentiate coeff. and report multiplicative change in mean counts.)

- $e^{\beta_0}$ : average count when x = 0
- $e^{\beta_1}$ : multiplicative change in the mean count per one unit change in x

For example: "Mean western toad juvenile abundance is reduced by a factor of 5.1 (95% CI: 3.4, 10.8) per unit change in UV-B radiation."

What was the estimate of  $\beta_1$ ?





What happens when we want to relate p to a predictor x, where x can be any value on the real line? How do we connect x to  $p \in [0,1]$ ?

• odds:  $\frac{p}{1-p} \in [0,\infty)$ • log odds: log(odds) = log  $\left(\frac{p}{1-p}\right) \in (-\infty,\infty)$ 

Moving between probability and log odds and relating to x:

- $\operatorname{logit}(p) = \log\left(\frac{p}{1-p}\right) = x$ 
  - input to logit() is probability p, output is log odds x
- Inverting the above, we obtain  $p = \frac{e^x}{e^x + 1} = \text{inverse logit}(x) = \exp(x)$ 
  - input to inverse logit() is x = log odds, output is probability p

Inverse logit mapping: input is log odds =  $\log\left(\frac{p}{1-p}\right)$ , output is probability



log(odds)

You have collected some binary data (y = 1, 0, 0, 1, 1, 0, 1, ..., 1) at *n* sites, along with a covariate  $x_i$  measured at each location which you believe is likely to affect these counts. Write a model regressing y on x.

- Choose a specific stochastic and deterministic model.
- Specify (vague) priors for your parameters.
- Write out the DAG and express posterior distribution as proportional to joint distribution for your model.
- Write the JAGS code for the model.
- Interpret the coefficients of your model.

Stochastic model:

$$y_i | \beta_0, \beta_1 i n d$$
. Bernoulli $(g(\beta_0, \beta_1, x_i))$   $i = 1, ..., n$ 

Deterministic model:

$$\mu_i = p_i = g(eta_0, eta_1, x_i) = rac{e^{eta_0 + eta_1 x_i}}{e^{eta_0 + eta_1 x_i} + 1}$$

Priors:

Choosing reasonable flat priors on logit intercept Imagine for now that we have no predictor, so  $p_i = \mu_i = \frac{e^{\beta_0}}{e^{\beta_0}+1}$ .

If we use the same Normal(0, large variance) prior as before, what is the induced prior for the success probability p?



## Choosing reasonable flat priors on logit intercept

Now instead consider the prior  $\beta_0 \sim \text{Normal}(0, 2.7)$ .



Choosing reasonable flat priors on logit effects



Returning to case with a single predictor, the posterior distribution is:

$$\begin{bmatrix} \beta_0, \beta_1 \mid \mathbf{y} \end{bmatrix} \propto \prod_{i=1}^n \operatorname{Bernoulli}(y_i \mid g(\beta_0, \beta_1, x_i)) \\ \times \operatorname{Normal}(\beta_0 \mid 0, 2.7) \operatorname{Normal}(\beta_1 \mid 0, 2.7) \\ g(\beta_0, \beta_1, x_i) = \frac{e^{\beta_0 + \beta_1 x_i}}{e^{\beta_0 + \beta_1 x_i} + 1}$$

JAGS code for the model:

```
b0 ~ dnorm(0, 1/2.7)
b1 ~ dnorm(0, 1/2.7)
for(i in 1:length(y)){
    logit(p[i]) <- b0 + b1 * x[i]
    y[i] ~ dbern(p[i])
}</pre>
```

or

```
p[i] <- inv.logit(b0 + b1 * x[i])
y[i] ~ dbern(p[i])</pre>
```

$$\begin{bmatrix} \beta_0, \beta_1 \mid \mathbf{y} \end{bmatrix} \propto \prod_{i=1}^n \operatorname{Bernoulli}(y_i \mid g(\beta_0, \beta_1, x_i)) \\ \times \operatorname{Normal}(\beta_0 \mid 0, 2.7) \operatorname{Normal}(\beta_1 \mid 0, 2.7) \\ p_i = g(\beta_0, \beta_1, x_i) = \frac{e^{\beta_0 + \beta_1 x_i}}{e^{\beta_0 + \beta_1 x_i} + 1} \iff \frac{p_i}{1 - p_i} = e^{\beta_0 + \beta_1 x_i} = e^{\beta_0} e^{\beta_1 x_i}$$

Interpretation: (Exponentiate coef. and report odds and odds ratios.)

- $e^{\beta_0}$  : odds when x = 0
- $e^{eta_1}$  : multiplicative change in the odds for a one unit change in x

For example: "the odds of detecting weevils in upland willow stems were 3.2 (95% CI: 2.3, 4.8) times greater than detecting them in riparian willow stems." What was the estimate of  $\beta_1$ ? (Might be helpful to identify/define  $x_i$ ).

# Nonlinear regression





# Centering and standardizing

The remainder of the slides apply to all of the general linear models, but we will use a simple linear model for Normally distributed data as an example.

#### Centering predictor data

Rather than the usual linear model  $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$  where  $\varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$ , consider the following linear model:

$$y_i = \beta_0 + \beta_1 (x_i - \bar{x}) + \varepsilon_i$$

where  $\bar{x} = \sum_{i=1}^{n} x_i$  is the sample mean of the predictor.

Why complicate things?

- To reduce autocorrelation in MCMC chain and speed convergence.
- To make the intercept more easily interpretable.

Centering predictor data

$$\begin{bmatrix} \beta_0, \beta_1, \sigma \mid \mathbf{y} \end{bmatrix} \propto \prod_{i=1}^n \operatorname{Normal}(y_i \mid g(\beta_0, \beta_1, x_i), \sigma^2) \times \\ \operatorname{Normal}(\beta_0 \mid 0, 1000) \operatorname{Normal}(\beta_1 \mid 0, 1000) \times \\ \operatorname{uniform}(\sigma \mid 0, 100) \\ g(\beta_0, \beta_1, x_i) = \beta_0 + \beta_1(x_i - \bar{x})$$

```
b0 ~ dnorm(0, .001)
b1 ~ dnorm(0, .001)
sigma ~ dunif(0, 100)
tau <- 1/sigma^2
xBar <- mean(x)
for (i in 1:length(y)){
    mu[i] <- b0 + b1 * (x[i] - xBar)
    y[i] ~ dnorm(mu[i], tau)
}
b0_UC <- b0 - b1 * xBar</pre>
```



What is the interpretation of  $\beta_0$  in this model? What is the interpretation of  $\beta_1$  in this model?

#### Recovering uncentered parameters



$$egin{array}{rcl} eta_0^* &=& eta_0 - eta_1 ar x \ eta_1^* &=& eta_1 \end{array}$$

- For this to work properly, all the coefficients in the model must be *added*.
- Slopes will not be the same if there is an interaction term or quadratic. In these cases, back transforming is not simple.

#### Standardizing predictor data

$$y_i = \beta_0 + \beta_1 \left(\frac{x_i - \bar{x}}{s_x}\right)$$

where  $s_x$  is sample standard deviation of predictor  $(s_x^2 = \frac{1}{n-1}\sum_{i=1}^n (x_i - \bar{x})^2)$ 

Why complicate things?

- To reduce autocorrelation in MCMC chain and speed convergence.
- To make the intercept more easily interpretable.
- To make coefficients more easily comparable

Standardizing predictor data

$$\begin{bmatrix} \beta_0, \beta_1, \sigma \mid \mathbf{y} \end{bmatrix} \propto \prod_{i=1}^n \operatorname{Normal}(y_i \mid g(\beta_0, \beta_1, x_i), \sigma^2) \times \\ \operatorname{Normal}(\beta_0 \mid 0, 1000) \operatorname{Normal}(\beta_1 \mid 0, 1000) \times \\ \operatorname{uniform}(\sigma \mid 0, 100) \\ g(\beta_0, \beta_1, x_i) = \beta_0 + \beta_1 \left(\frac{x_i - \bar{x}}{s_x}\right) \\ \end{bmatrix}$$
  
b0 ~ dnorm(0, .001)  
b1 ~ dnorm(0, .001)  
sigma ~ dunif(0, .100)

```
b1 ~ dnorm(0, .001)
sigma ~ dunif(0, 100)
tau <- 1/sigma^2
xBar <- mean(x)
xSD <- sd(x)
for (i in 1:length(y)){
    mu[i] <- b0 + b1 * ((x[i] - xBar)/xSD
    y[i] ~ dnorm(mu[i], tau)
}</pre>
```

#### Recovering unstandardized parameters

$$y_i = \beta_0 + \beta_1 \left(\frac{x_i - \bar{x}}{s_x}\right)$$
$$= \beta_0 + \frac{\beta_1}{s_x} x_i - \frac{\beta_1 \bar{x}}{s_x}$$
$$= \beta_0^* + \beta_1^* x_i$$
$$\beta_0^* = \beta_0 - \frac{\beta_1 \bar{x}}{s_x}$$
$$\beta_1^* = \frac{\beta_1}{s_x}$$

- This only works if there are not squared values or interactions.
- It is fine to make predictions using  $\hat{y}_i = \beta_0 + \beta_1 \frac{x_i \bar{x}}{s_x}$  and plot  $\hat{y}_i$  against  $x_i$  and the observed  $y_i$

lognormal, data continuous and > 0 (log link)

$$\begin{bmatrix} \beta_0, \beta_1, \sigma \mid \mathbf{y} \end{bmatrix} \propto \prod_{i=1}^n \operatorname{lognormal}(y_i \mid \log(g(\beta_0, \beta_1, x_i)), \sigma^2) \\ \times \operatorname{Normal}(\beta_0 \mid 0, 1000) \operatorname{Normal}(\beta_1 \mid 0, 1000) \\ \times \operatorname{uniform}(\sigma \mid 0, 5) \\ g(\beta_0, \beta_1, x_i) = e^{\beta_0 + \beta_1 x_i}$$

Talk about the interpretation of  $\sigma$ .

```
b0 ~ dnorm(0, .001)
b1 ~ dnorm(0, .001)
sigma ~ dunif(0, 5)
tau <- 1/sigma^2
for(i in 1:length(y)){
    mu[i] <- exp(b0 + b1 * x[i])
    y[i] ~ dlnorm(log(mu[i]), tau)
}
```

lognormal, data continuous and > 0 (not log link)

$$\begin{bmatrix} \beta_0, \beta_1, \sigma \mid \mathbf{y} \end{bmatrix} \propto \prod_{i=2}^n \operatorname{lognormal}(y_i \mid \log(g(\beta_0, \beta_1, y_{i-1})), \sigma^2) \\ \times \operatorname{Normal}(\beta_0 \mid 0, 1000) \operatorname{Normal}(\beta_1 \mid 0, 1000) \\ \times \operatorname{uniform}(\sigma \mid 0, 5) \operatorname{uniform}(y_1 \mid 1, 1E6) \\ g(\beta_0, \beta_1, y_{i-1}) = y_{i-1}e^{\beta_0 + \beta_1 y_{i-1}}$$

```
b0 ~ dnorm(0, .001)
b1 ~ dnorm(0, .001)
sigma ~ dunif(0, 5); tau <- 1/sigma^2
y[1] ~ dunif(1, 1E6)
for(i in 2:length(y)){
    mu[i] <- y[i-1] * exp(b0 + b1 * y[i-1])
    y[i] ~ dlnorm(log(mu[i]), tau)
}
```