

Today

- Tidy data
 - reshaping with tidyr
 - data manipulation with dplyr
- Ants GLM, Bayesian
 - ulam: equations, incl priors
 - working with posterior samples
- Model matrix

Independent project

- Homework: ideas submitted today
- I'll look at your ideas tomorrow
- Questions?
- Also chat after class

Tidy data

country	year	cases	population
Afghanistan	1999	745	19987071
Afghanistan	2000	2666	20095360
Brazil	1999	30737	172006362
Brazil	2000	80488	174004898
China	1999	210258	1272015272
China	2000	210706	1280425583

**Variables
in
columns**

country	year	cases	population
Afghanistan	1999	745	19987071
Afghanistan	2000	2666	20095360
Brazil	1999	30737	172006362
Brazil	2000	80488	174004898
China	1999	210258	1272015272
China	2000	210706	1280425583

**Observations
in
rows**

country	year	cases	population
Afghanistan	99	745	19987071
Afghanistan	00	2666	20095360
Brazil	99	30737	172006362
Brazil	00	80488	174004898
China	99	210258	1272015272
China	00	210706	1280425583

**Values
in
cells**

`pivot_longer()` - tidy a variable that is in multiple columns
`pivot_wider()` - tidy an observation that is in multiple rows

Base R: `reshape()`, `stack()`, `unstack()`, `strsplit()`, `paste()`

dplyr - working with data

`filter()` - pick observations by their values

`select()` - pick columns by name

`arrange()` - reorder rows

`mutate()` - create new variables from existing variables

`summarize()` - collapse values to a summary statistic

`group_by()` - all the above split by group

`|>` pipe to combine (or `%>%`)

Base R: `subset()`, `order()`, `sort()`, `table()`, `aggregate()`, `|>`

tibbles (tables with NZ accent)

Data frame; optimized printing for **large** datasets.
You probably spent a lot of time collecting data.
Wouldn't you want to spend **a few minutes** to inspect each row and column? That's how you **catch errors**.

Convert to dataframe

```
my_tibble |> data.frame()
```

Keep tibble, but print all the data by default

```
options(tibble.width=Inf)
```

```
options(tibble.print_max=Inf)
```

```
options(max.print=1500)
```

```
?print.tbl - see options
```

Other ways to inspect tibbles

```
View(my_tibble) – only works in Rstudio
```

```
glimpse(my_tibble) - summary
```

dplyr vs base

How many trees with known mortality status are missing a diameter in 2013?

```
tree_dat |>
  filter(year == 2013) |>
  filter(!is.na(mortality)) |>
  mutate(diam_missing=is.na(diam)) |>
  summarize(sum(diam_missing))

sum(is.na(subset(tree_dat, year == 2013 &
  !is.na(mortality))$diam))
```

dplyr, complex queries

How many species observed per habitat fragment in each year per treatment?

```
beetle_species |>
  group_by(species, year, fragment,
            treatment) |>
  summarize(mean_abun = mean(abundance)) |>
  mutate(present = mean_abun > 0) |>
  group_by(year, fragment, treatment) |>
  summarize(richness = sum(present)) |>
  group_by(year, treatment) |>
  summarize(mean_richness = mean(richness))
```

Most common models

Normal
+
Identity link

$$y_i \sim \text{Normal}(\mu_i, \sigma)$$

$$\mu_i = \beta_0 + \beta_1 x_i$$

Poisson
+
Log link

$$y_i \sim \text{Poisson}(\mu_i)$$

$$\log(\mu_i) = \beta_0 + \beta_1 x_i$$

Binomial
+
Logit link

$$y_i \sim \text{Binomial}(\mu_i, n)$$

$$\log\left(\frac{\mu_i}{1 - \mu_i}\right) = \beta_0 + \beta_1 x_i$$

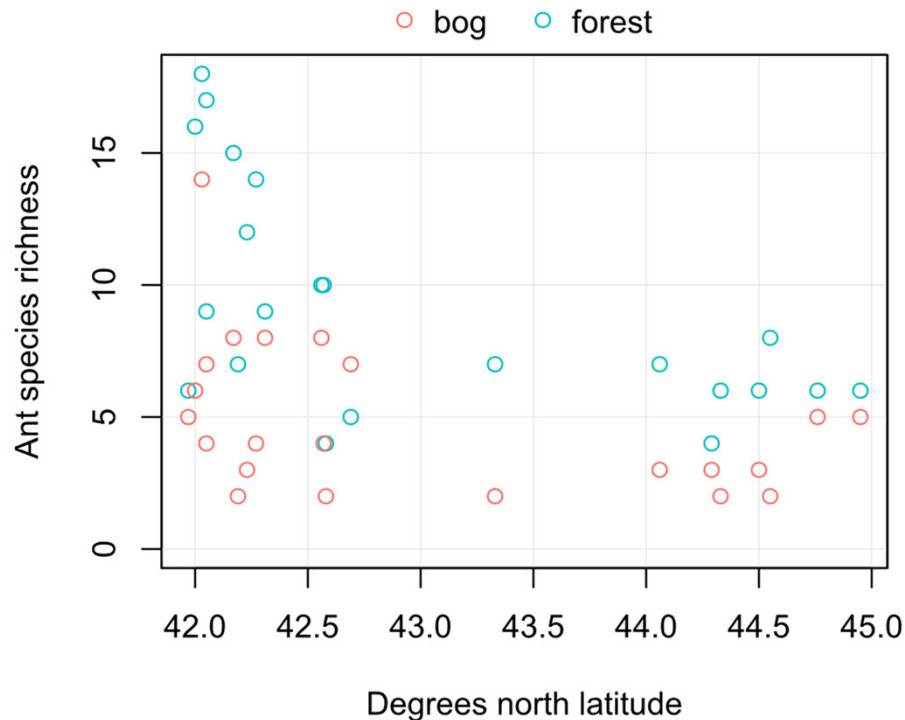
Key properties:

y : $-\infty$ to ∞ , continuous
 μ : $-\infty$ to ∞ , continuous

y : 0 to ∞ , discrete, integer
 μ : 0 to ∞ , continuous

y : 0, 1, discrete, binary
 μ : 0 to 1, probability

Which makes the most sense?



Poisson
+
Log link

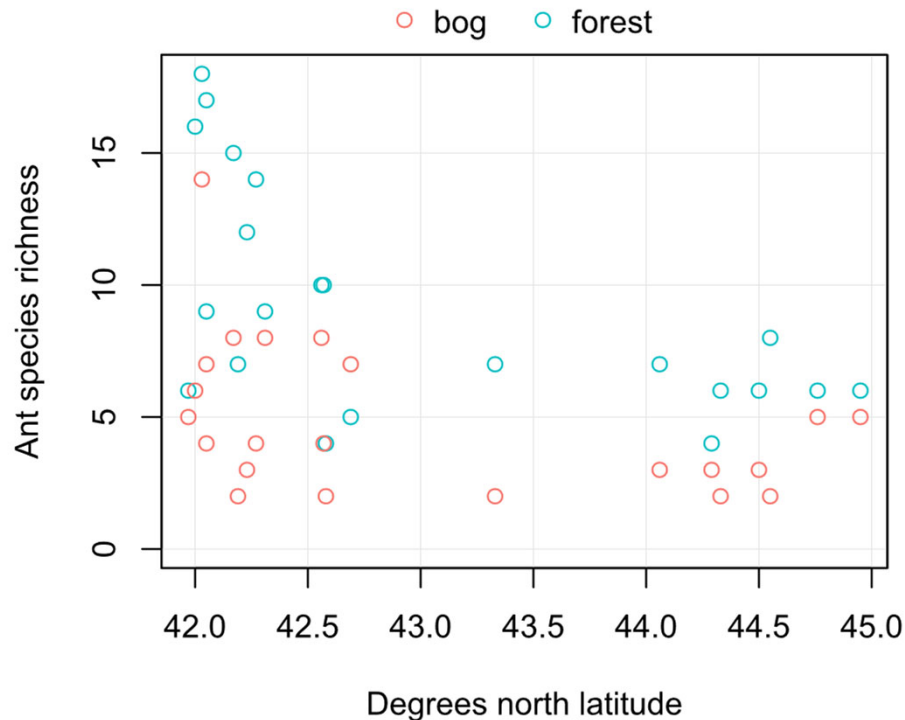
$$y_i \sim \text{Poisson}(\mu_i)$$

$$\log(\mu_i) = \beta_0 + \beta_1 x_i + \dots$$

$$\eta_i = \beta_0 + \beta_1 x_i + \dots$$

$$\mu_i = e^{\eta_i}$$

Write the full linear predictor



Poisson
+
Log link

$$y_i \sim \text{Poisson}(\mu_i)$$

$$\log(\mu_i) = \beta_0 + \beta_1 x_i + \dots$$

$$\eta_i = \beta_0 + \beta_1 x_i + \dots$$

$$\mu_i = e^{\eta_i}$$

$$\eta_i = \beta_0 + \beta_1 \text{forest}_i + \beta_2 \text{latitude}_i + \beta_3 \text{forest}_i \times \text{latitude}_i$$

Model matrix (design matrix)

Data

habitat	latitude	richness
forest	42	16
forest	42.56	10
forest	43.33	7
forest	44.76	6
bog	42.17	8
bog	42.57	4
bog	44.06	3
bog	44.95	5

$$\eta_i = \beta_0 + \beta_1 \text{forest}_i + \beta_2 \text{latitude}_i + \beta_3 \text{forest}_i \times \text{latitude}_i$$

Model matrix (design matrix)

Data

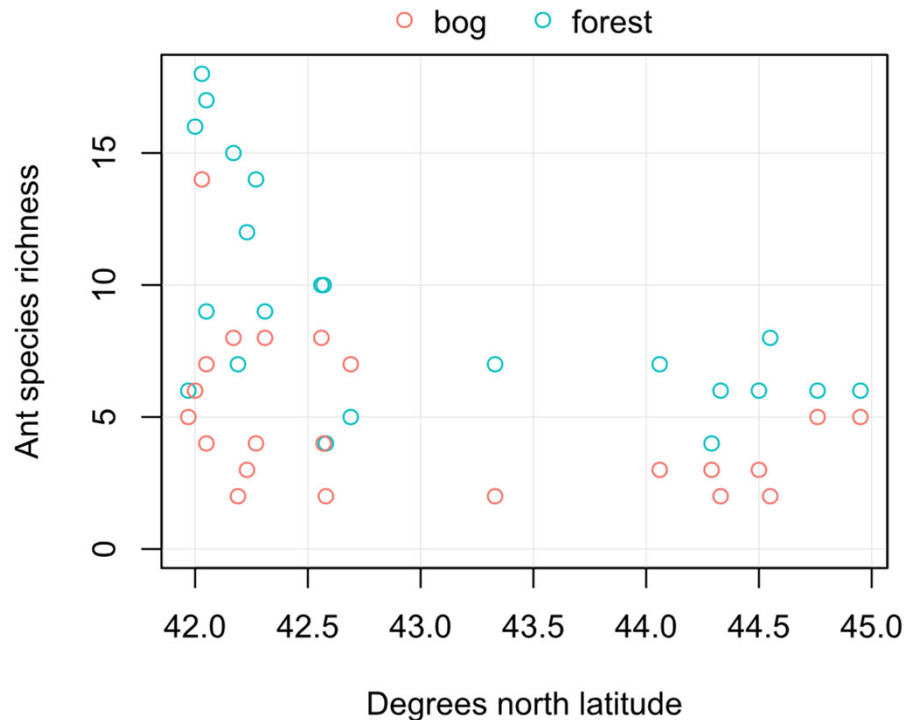
habitat	latitude	richness
forest	42	16
forest	42.56	10
forest	43.33	7
forest	44.76	6
bog	42.17	8
bog	42.57	4
bog	44.06	3
bog	44.95	5

Model matrix

intercept	forest	latitude	forest:latitude
1	1	42	42
1	1	42.56	42.56
1	1	43.33	43.33
1	1	44.76	44.76
1	0	42.17	0
1	0	42.57	0
1	0	44.06	0
1	0	44.95	0

$$\eta_i = \beta_0 \text{intercept}_i + \beta_1 \text{forest}_i + \beta_2 \text{latitude}_i + \beta_3 \text{forest}_i \times \text{latitude}_i$$

Bayesian GLM: Priors?



Poisson

+

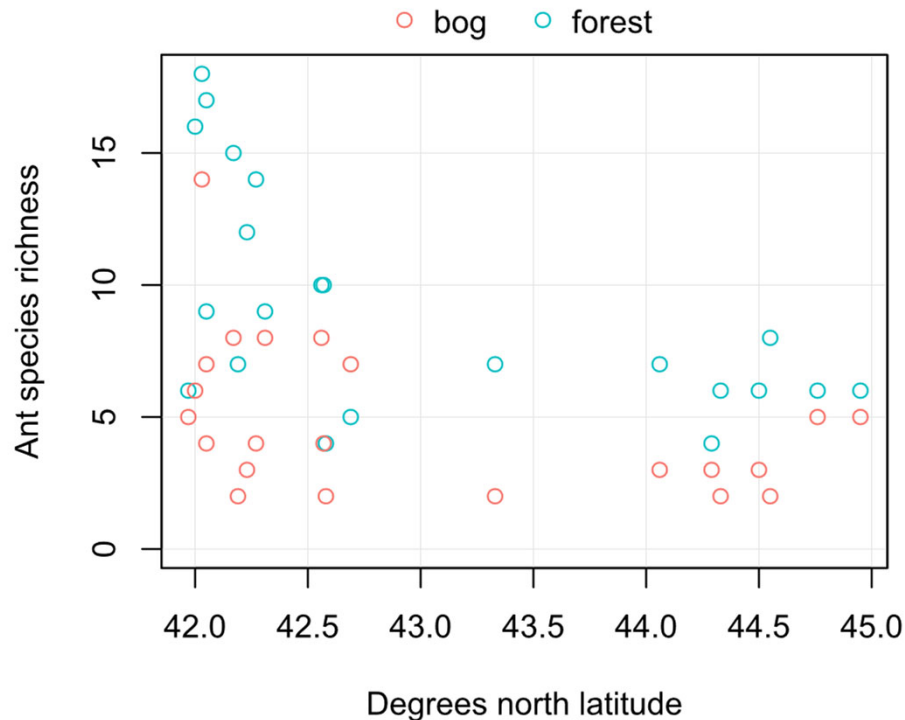
Log link

$$y_i \sim \text{Poisson}(\mu_i)$$

$$\log(\mu_i) = \eta_i$$

$$\eta_i = \beta_0 + \beta_1 \text{forest}_i + \beta_2 \text{latitude}_i + \beta_3 \text{forest}_i \times \text{latitude}_i$$

Bayesian GLM: Code



Poisson

+

Log link

$$y_i \sim \text{Poisson}(\mu_i)$$

$$\log(\mu_i) = \eta_i$$

$$\eta_i = \beta_0 + \beta_1 \text{forest}_i + \beta_2 \text{latitude}_i + \beta_3 \text{forest}_i \times \text{latitude}_i$$

Write the code for `ulam()`