

STA 360/602L: MODULE 7.2

METROPOLIS IN ACTION

DR. OLANREWAJU MICHAEL AKANDE

COUNT DATA

- We will use the Metropolis sampler on count data with predictors, so let's first do some general review.
- Suppose you have count data as your response variable.
- For example, we may want to explain the number of c-sections carried out in hospitals using potential predictors such as hospital type, (that is, private vs public), location, size of the hospital, etc.
- The models we have covered so far are not (completely) adequate for count data with predictors.
- Of course there are instances where linear regression, with some transformations (especially taking logs) on the response variable, might still work reasonably well for count data.
- That's not the focus here, so we won't cover that.

POISSON REGRESSION

- As we have seen so far, a good distribution for modeling count data with no limit on the total number of counts is the **Poisson distribution**.
- As a reminder, the Poisson pmf is given by

$$\Pr[Y = y|\lambda] = \frac{\lambda^y e^{-\lambda}}{y!}; \quad y = 0, 1, 2, \dots; \quad \lambda > 0.$$

- Remember that

$$\mathbb{E}[Y = y] = \mathbb{V}[Y = y] = \lambda.$$

- When our data fails this assumption, we may have what is known as **over-dispersion** and may want to consider the **Negative Binomial distribution** instead (actually easy to fit within the Bayesian framework!).
- With predictors, index λ with i , so that each λ_i is a function of \mathbf{X} . Therefore, the **random component** of the glm is

$$p(y_i|\lambda_i) = \text{Poisson}(\lambda_i); \quad i = 1, \dots, n.$$

POISSON REGRESSION

- We must ensure that $\lambda_i > 0$ at any value of \mathbf{X} , therefore, we need a **link function** that enforces this. A natural choice is

$$\log(\lambda_i) = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}.$$

- Combining these pieces give us our full mathematical representation for the **Poisson regression**.
- Clearly, λ_i has a natural interpretation as the "expected count", and

$$\lambda_i = e^{\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}}$$

so the e^{β_j} 's are **multiplicative effects** on the expected counts.

- For the frequentist version, in **R**, use the `glm` command but set the option `family = "poisson"`.

ANALYSIS OF HORSESHOE CRABS

- We have data from a study of nesting horseshoe crabs (J. Brockmann, *Ethology*, 102: 1–21, 1996). The data has been discussed in Agresti (2002).
- Each female horseshoe crab in the study had a male crab attached to her in her nest.
- The study investigated factors that affect whether the female crab had any other males, called satellites, residing nearby her.
- The response outcome for each female crab is her number of satellites.
- We have several factors (including the female crab's color, spine condition, weight, and carapace width) which may influence the presence/absence of satellite males.
- The data is called `hcrabs` in the R package `rsq`.

ANALYSIS OF HORSESHOE CRABS

- Let's fit the Poisson regression model to the data. In vector form, we have

$$y_i \sim \text{Poisson}(\lambda_i); \quad i = 1, \dots, n;$$

$$\log[\lambda_i] = \boldsymbol{\beta}^T \mathbf{x}_i$$

where y_i is the number of satellites for female crab i , and \mathbf{x}_i contains the intercept and female crab i 's

- color;
 - spine condition;
 - weight; and
 - carapace width.
- Suppose we specify a normal prior for $\boldsymbol{\beta} = (\beta_0, \beta_1, \beta_2, \dots, \beta_{p-1})$,
 $\pi(\boldsymbol{\beta}) = \mathcal{N}_p(\boldsymbol{\beta}_0, \boldsymbol{\Sigma}_0)$.
 - Can you write down the posterior for $\boldsymbol{\beta}$? Can you sample directly from it?

ANALYSIS OF HORSESHOE CRABS

- We can use Metropolis to generate samples from the posterior.
- First, we need a "symmetric" proposal density $\beta^* \sim g[\beta^* | \beta^{(s)}]$; a reasonable choice is usually a multivariate normal centered on $\beta^{(s)}$.
- What about the variance of the proposal density? We can use the variance of the ols estimate, that is, $\hat{\sigma}^2 (\mathbf{X}^T \mathbf{X})^{-1}$, which we can scale using δ , to tune the acceptance ratio.
- Here, $\hat{\sigma}^2$ is calculated as the sample variance of $\log[y_i + c]$, for some small constant c , to avoid problems when $y_i = 0$.
- So we have $g[\beta^* | \beta^{(s)}] = \mathcal{N}_p \left(\beta^{(s)}, \delta \hat{\sigma}^2 (\mathbf{X}^T \mathbf{X})^{-1} \right)$.
- Finally, since we do not have any information apriori about β , let's set the prior for it to be $\pi(\beta) = \mathcal{N}_p(\beta_0 = \mathbf{0}, \Sigma_0 = \mathbf{I})$.

ANALYSIS OF HORSESHOE CRABS

- The Metropolis algorithm for this model is:

1. Given a current $\beta^{(s)}$, generate a candidate value

$$\beta^* \sim g[\beta^* | \beta^{(s)}] = \mathcal{N}_p \left(\beta^{(s)}, \delta \hat{\sigma}^2 (\mathbf{X}^T \mathbf{X})^{-1} \right).$$

2. Compute the acceptance ratio

$$\begin{aligned} r &= \frac{\pi(\beta^* | Y)}{\pi(\beta^{(s)} | Y)} = \frac{\pi(\beta^*) \cdot p(Y | \beta^*)}{\pi(\beta^{(s)}) \cdot p(Y | \beta^{(s)})} \\ &= \frac{\mathcal{N}_p(\beta^* | \beta_0 = \mathbf{0}, \Sigma_0 = \mathbf{I}) \cdot \prod_{i=1}^n \text{Poisson} \left(Y_i | \lambda_i = \exp \left\{ (\beta^*)^T \mathbf{x}_i \right\} \right)}{\mathcal{N}_p(\beta^{(s)} | \beta_0 = \mathbf{0}, \Sigma_0 = \mathbf{I}) \cdot \prod_{i=1}^n \text{Poisson} \left(Y_i | \lambda_i = \exp \left\{ (\beta^{(s)})^T \mathbf{x}_i \right\} \right)}. \end{aligned}$$

3. Sample $u \sim U(0, 1)$ and set

$$\beta^{(s+1)} = \begin{cases} \beta^* & \text{if } u < r \\ \beta^{(s)} & \text{if otherwise} \end{cases}.$$

MOVE TO THE R SCRIPT **HERE.**

WHAT'S NEXT?

MOVE ON TO THE READINGS FOR THE NEXT MODULE!