

ADVANCED BAYESIAN MODELING

Loglinear Regression Modeling

Data Model

$$X_i = (x_{i1}, \dots, x_{ik}), \quad y_i \quad i = 1, \dots, n$$

X has k explanatory variables (usually including an intercept)

$$y_i \mid \beta, X_i \sim \text{indep. Poisson}(\lambda_i)$$

$$\log \lambda_i = X_i \beta \quad \lambda_i = e^{X_i \beta}$$

This is a **(Poisson) loglinear regression model**.

Poisson models are appropriate for responses that are counts with no obvious maximum value, especially counts of rare, unrelated events, such as:

- ▶ Cases of a rare genetic disease
- ▶ Shark attacks on swimmers
- ▶ Numbers of lottery winners

If all counts are large enough (hundreds or more), a transformed normal-theory model might be a good approximation.

If some counts are small (especially, if some are zero), normal-theory models are not appropriate.

Priors

Conjugate (gamma-type) priors exist for special cases, but are not very useful in general.

Flat priors on coefficients (in β) may be safe, provided data are not too “sparse” (too many zeros). Otherwise, may need to make priors informative.

We consider normal priors, which may be diffuse (large variance).

As in ordinary regression, standardizing quantitative explanatory variables is recommended.

Overdispersion

Mean and variance of Poisson are equal:

$$E(y_i \mid \beta, X_i) = \lambda_i \qquad \text{var}(y_i \mid \beta, X_i) = \lambda_i$$

Common problem: Data exhibit variances that do not equal the means – typically larger than they should be: **overdispersion**.

Possible causes:

- ▶ Explanatory variables are missing
- ▶ Events being counted are dependent
- ▶ Events of different rates are combined into a single count

The **chi-square discrepancy** measures overdispersion:

$$T(y, \theta, X) = \sum_{i=1}^n \frac{(y_i - \lambda_i)^2}{\lambda_i}$$

(When the model is correct and the λ_i s are large enough, it has an approximate χ_n^2 distribution.)

When there is overdispersion, it tends to be larger.

Can be used to compute a posterior predictive p -value for overdispersion:

$$\Pr(T(y^{\text{rep}}, \theta, X) \geq T(y, \theta, X) \mid y)$$

One possible way to account for overdispersion:

Add a separate random effect to the linear combination of each observation.

$$\log \lambda_i = X_i \beta + \epsilon_i \quad \epsilon_i \sim \text{iid } N(0, \sigma_\epsilon^2)$$

Then add a reasonably noninformative prior on σ_ϵ^2 (e.g. flat on σ_ϵ).

(See BDA3, Sec. 16.4, for example.)

Alternative: Use a *negative binomial* response distribution, which naturally allows for overdispersion.

(BDA3, Sec. 17.2)