

Generalized Linear Models

Matt Farr

- Science is too complex to understand → must reduce complexity
- Describe and predict systems/processes
- Model is set of assumptions about scientific systems/processes
- Simplification of reality

- Description of a system or process composed of variables, typically written in algebra
- Algebra: explicit and great transparency
- Simple example is the equation for a line: $y=mx+b$
- Format for linear relationships, where the expected response (y) can be treated as the result of explanatory variables (x) whose effects (mx + b) are additive.

- Stochasticity or randomness
- Statistical models: a description of a system/process composed of variables but where one or more **random variables** are related to other variables
- Explicitly acknowledge stochasticity in system/processes
$$\text{Response} = \text{deterministic part} + \text{stochastic part}$$
- Parametric statistical models (generalized linear models): a description of a system/process using probability distributions thought to have produced the data

GENERALIZED LINEAR MODEL

Response Data = Deterministic Part + Stochastic Part

$$y = \alpha + \beta \cdot x + \epsilon$$

$$\epsilon \sim \text{Normal}(0, \sigma^2)$$

Response Data \sim Distribution(Parameters)

$$y \sim \text{Normal}(\mu, \sigma^2)$$

Link(Mean Response) = Linear Model

$$I(\mu) = \alpha + \beta \cdot x$$

$$\mu = \alpha + \beta \cdot x$$

- Types of responses (y)
 - Continuous (body mass)
 - Binary (heads/tails; dead/survived)
 - Categorical (rolling a die, nationality; geographic location)
 - Counts (number of birds in a quadrant)

- Continuous (body mass)
Random part: $y \sim \text{Normal}(\mu, \sigma^2)$
Deterministic part: $\mu = \alpha + \beta \cdot x$
- Binary: Binomial distribution (coin flip; sex; alive/dead)
Random part: $y \sim \text{Binomial}(p, N)$
Deterministic part: $\text{logit}(p) = \alpha + \beta \cdot x$
- Categorical (rolling a die; nationality; age)
Random part: $y \sim \text{Multinomial}(\boldsymbol{\pi}, N)$
Deterministic part: $\text{logit}(\boldsymbol{\pi}) = \alpha + \beta \cdot x$
- Counts (number of birds in a quadrant)
Random part: $y \sim \text{Poisson}(\lambda)$
Deterministic part: $\log(\lambda) = \alpha + \beta \cdot x$

- Deterministic part of the linear model
- A matrix of explanatory variables
- Manipulation of the design matrix that leads to classical statistical models
- Column for each predictor variable
- Means parameterization or an effects parameterization (intercept)
- Row for each data point (individuals, sites, time period, etc.)

$$y_{li} \sim \text{Normal}(\mu_{li}, \sigma^2)$$

$$I(\mu_{li}) = \beta_F \cdot \text{females}_{li} + \beta_M \cdot \text{males}_{li}$$

$$(\begin{bmatrix} \mu_1 & \mu_2 & \mu_3 & \mu_4 & \mu_5 & \mu_6 \end{bmatrix}) = (\begin{bmatrix} 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 \end{bmatrix}) * (\begin{bmatrix} \beta_F & \beta_M \end{bmatrix})$$

DESIGN MATRIX

$$y_{\downarrow i} \sim \text{Normal}(\mu_{\downarrow i}, \sigma^2)$$

$$I(\mu_{\downarrow i}) = \beta_{\downarrow F} \cdot \text{females}_{\downarrow i} + \beta_{\downarrow M} \cdot \text{males}_{\downarrow i}$$

$$\mu_{\downarrow 1} = \beta_{\downarrow F} \cdot 1 + \beta_{\downarrow M} \cdot 0$$

$$\mu_{\downarrow 2} = \beta_{\downarrow F} \cdot 1 + \beta_{\downarrow M} \cdot 0$$

$$\mu_{\downarrow 3} = \beta_{\downarrow F} \cdot 1 + \beta_{\downarrow M} \cdot 0$$

$$\mu_{\downarrow 4} = \beta_{\downarrow F} \cdot 0 + \beta_{\downarrow M} \cdot 1$$

$$\mu_{\downarrow 5} = \beta_{\downarrow F} \cdot 0 + \beta_{\downarrow M} \cdot 1$$

$$\mu_{\downarrow 6} = \beta_{\downarrow F} \cdot 0 + \beta_{\downarrow M} \cdot 1$$

DESIGN MATRIX

$$y_{\downarrow i} \sim \text{Normal}(\mu_{\downarrow i}, \sigma^2)$$

$$I(\mu_{\downarrow i}) = \beta_{\downarrow F} \cdot \text{females}_{\downarrow i} + \beta_{\downarrow M} \cdot \text{males}_{\downarrow i}$$

$$\mu_{\downarrow 1} = \beta_{\downarrow F}$$

$$\mu_{\downarrow 2} = \beta_{\downarrow F}$$

$$\mu_{\downarrow 3} = \beta_{\downarrow F}$$

$$\mu_{\downarrow 4} = \beta_{\downarrow M}$$

$$\mu_{\downarrow 5} = \beta_{\downarrow M}$$

$$\mu_{\downarrow 6} = \beta_{\downarrow M}$$

$\beta_{\downarrow F}$ = mean
of females

$\beta_{\downarrow M}$ = mean
of males

DESIGN MATRIX

$$y_i = \beta_F \cdot \text{females}_i + \beta_M \cdot \text{males}_i + \epsilon_i$$

$$\epsilon_i \sim \text{Normal}(0, \sigma^2)$$

$$\begin{pmatrix} y_1 & y_2 & y_3 & y_4 & y_5 & y_6 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 \end{pmatrix} \cdot \begin{pmatrix} \beta_F \\ \beta_M \end{pmatrix} + \begin{pmatrix} \epsilon_1 & \epsilon_2 & \epsilon_3 & \epsilon_4 & \epsilon_5 & \epsilon_6 \end{pmatrix}$$

MODEL OF THE MEAN

$$y_i \sim \text{Normal}(\mu, \sigma^2)$$

$$y_i = \mu + \epsilon_i$$

$$\epsilon_i \sim \text{Normal}(0, \sigma^2)$$

$$(y_1 \ y_2 \ y_3 \ y_4 \ y_5 \ y_6) = (1 \ 1 \ 1 \ 1 \ 1 \ 1) * (\mu) + (\epsilon_1 \ \epsilon_2 \ \epsilon_3 \ \epsilon_4 \ \epsilon_5 \ \epsilon_6)$$

- Parameters are unknown!!
- Frequentist analysis
 - Least-squares (LM function in R)
 - Maximum likelihood
 - Method of moments
- Bayesian analysis

- Data Exploration
- Response explanation
- Model structure & development
- Data formatting
- Model analysis
- Model results & interpretation

Exercise:

<https://raw.githubusercontent.com/farrmt/Rworkshop/master/GLMfun.csv>

References

- Kéry, Marc. 2010. Introduction to WinBUGS for ecologists: A Bayesian approach to regression, ANOVA, mixed models, and related analyses. Academic Press, Boston.