

Poisson Model with an Offset

General Principles

When we want to model count data, where the counts are observed over different periods or areas of exposure, we use a *Poisson model with an offset*. This is a type of generalized linear model used for modeling count data and contingency tables.

An *offset* is a predictor variable with a coefficient that is fixed at 1. It is used to account for the “exposure” variable, which represents the opportunity for an event to occur. For instance, if we are counting the number of sick individuals in different cities, the population of each city would be the exposure variable. A city with a larger population is expected to have more sick individuals. The offset accounts for this by essentially modeling the rate of events per unit of exposure.

Considerations

Note

- The dependent variable in a Poisson regression must be a non-negative count.
- The exposure variable used as an offset cannot contain zeros.
- A key assumption of the Poisson distribution is that the mean and variance of the count variable are equal. If the variance is greater than the mean, a condition known as overdispersion, a Negative Binomial regression might be more appropriate.
- The logarithm of the exposure variable is typically used as the offset. This is because Poisson regression models the logarithm of the expected count. By including the log of the exposure as an offset, we are effectively modeling the rate.

Example

Below is an example of code that demonstrates a Bayesian Poisson regression with an offset. The data consists of the number of elephant aggressions (**agressions**), the age of the elephants (**age**), and the number of years they have been observed (**years_obs**). The goal is to model the rate of aggressions per year, accounting for the age of the elephants.

Python

```
from main import*

# Setup device-----
m = bi(platform='cpu')

# Import Data & Data Manipulation -----
# Import
from importlib.resources import files
data_path = files('BI.resources.data') / 'elephants.csv'
m.data(data_path, sep=',')
m.log(['years_obs']) # Log transform the exposure variable
m.data_to_model(['agressions', 'age', "log_years_obs"]) # Send to model

# Define model -----
def model(agressions, age, log_years_obs):
    a = m.bi.dist.normal(0, 1, name = 'a')
    b = m.bi.dist.normal(0, 0.5, name = 'b')
    log_lambda = a + b * age + log_years_obs # Add offset to the linear model
    m.poisson(m.bi.numpy.exp(log_lambda), obs=agressions)

# Run mcmc -----
m.fit(model) # Optimize model parameters through MCMC sampling

# Summary -----
m.summary()
```

R

```
library(BI)
m=importbi(platform='cpu')
```

```

# Load csv file
m$data(paste(system.file(package = "BI"),"/data/elephants.csv", sep = ''), sep=',')
m$log(list('years_obs')) # Log transform the exposure variable
m$data_to_model(list('agressions', 'age', 'log_years_obs')) # Send to model

# Define model -----
model <- function(agressions, age, log_years_obs){
  # Parameter prior distributions
  a = bi.dist.normal(0, 1, name = 'a')
  b = bi.dist.normal(0, 0.5, name = 'b')
  # Likelihood
  log_lambda = a + b * age + log_years_obs # Add offset to the linear model
  m$poisson(exp(log_lambda), obs=agressions)
}

# Run mcmc -----
m$fit(model) # Optimize model parameters through MCMC sampling

# Summary -----
m$summary() # Get posterior distributions

```

Mathematical Details

Frequentist formulation

We model the relationship between the independent variables (X) and the expected count (λ) using the following equation:

$$\log(\lambda_i) = \alpha + \beta X_i + \log(\text{exposure}_i)$$

Where:

- λ_i is the expected count for observation i .
- α is the intercept term.
- β is the regression coefficient for the independent variable.
- X_i is the value of the independent variable for observation i .
- $\log(\text{exposure}_i)$ is the offset, which is the natural logarithm of the exposure for observation i .

The number of observed counts Y_i is assumed to follow a Poisson distribution with mean λ_i :

$$Y_i \sim \text{Poisson}(\lambda_i)$$

Bayesian formulation

In the Bayesian framework, we assign prior distributions to the model parameters. The model can be expressed as:

$$Y_i \sim \text{Poisson}(\lambda_i)$$

$$\log(\lambda_i) = \alpha + \beta X_i + \log(\text{exposure}_i)$$

$$\alpha \sim \text{Normal}(0, 1)$$

$$\beta \sim \text{Normal}(0, 1)$$

Where:

- Y_i is the observed count for observation i .
- λ_i is the expected count.
- α is the intercept with a unit-normal prior.
- β is the slope coefficient with a unit-normal prior.
- X_i is the independent variable.
- $\log(\text{exposure}_i)$ is the offset.

Notes

i Note

- The use of an offset is crucial when the goal is to compare rates of events rather than absolute counts.
- It is a common practice to use the natural logarithm of the exposure variable as the offset.

Reference(s)