

# Poisson Regression

## Motivation and Set-up:

In order to model the expected *count* or rate of some process, we use a Poisson model.

- Disease occurrence
- Claims made (Insurance)

If *exposure* is constant across all subjects, the mean count can be modeled directly. However usually the **rate** is modeled. E.g. cases / 100 patient-years or events / month, etc.

**Set-up:** for  $i = 1, \dots, n$  i.i.d samples

- $Y_i$  = event count
- $\mathbf{x}_i$  = covariate vector
- $T_i$  = exposure;  $\log T_i$  is the offset
- rate =  $\frac{E[\text{count}]}{\text{exposure}} = \frac{\mu_i}{T_i} = \lambda_i$
- $Y_i \sim \text{Poisson}(\lambda_i = e^{\mathbf{x}_i^T \boldsymbol{\beta}})$
- $P(Y_i | \mathbf{x}_i, T_i) = \frac{e^{\lambda_i} \lambda_i^{Y_i}}{Y_i!}$

## Rate Model

- $\log(\lambda_i) = \mathbf{x}_i^T \boldsymbol{\beta}$
- $\log E[\frac{Y_i}{T_i}] = \mathbf{x}_i^T \boldsymbol{\beta}$
- $\log \underbrace{E[Y_i]}_{\mu_i} = \log T_i + \mathbf{x}_i^T \boldsymbol{\beta}$

## GLM details:

- log link;  $\log \lambda_i = \mathbf{x}_i^T \boldsymbol{\beta}$
- Mean function:  
$$\mu_i = T_i \lambda_i \iff \frac{\mu_i}{T_i} = \lambda_i = \mathbf{x}_i^T \boldsymbol{\beta}$$
- Variance function  $v(\mu_i) = \mu_i = T_i \lambda_i$
- $a(\phi) = 1$  – usually, check dispersion section

## Interpretation

Consider the following toy example

$$\log(\lambda_i) = \beta_0 + \beta_1(x_i - \bar{x})$$

Coefficient interpretations

- $\beta_0$ 
  - if  $x_i = \bar{x}$  then  $\beta_0 = \log(\frac{\mu_i}{T_i}) = \log(\lambda_i)$
  - so  $\beta_0$  is the log rate when  $x_i = \bar{x}$
- $\beta_1$ 
  - $\log \lambda_a - \log \lambda_b = \beta_0 + \beta_1(x_a + 1) - \beta_0 - \beta_1(x_b)$   
 $= \log \frac{\lambda_a}{\lambda_b} = \beta_1$
  - So  $\beta_1$  is the log rate ratio for a one unit increase in  $x$

## Likelihood Function and Derivatives

$$\begin{aligned}
 L(\beta) &\propto \prod_{i=1}^n \frac{e^{\lambda_i} \lambda_i^{Y_i}}{Y_i!} \\
 &\propto \prod_{i=1}^n \frac{\exp\{e^{\mathbf{x}_i^T \beta}\} \exp\{\mathbf{x}_i^T \beta\}^{Y_i}}{Y_i!} \\
 \Rightarrow l(\beta) &= \sum_{i=1}^n -T_i e^{\mathbf{x}_i^T \beta} + Y_i \log T_i + Y_i \mathbf{x}_i^T \beta \\
 U(\beta) &= \mathbf{X}^T (\mathbf{Y} - \boldsymbol{\mu}) = \sum_{i=1}^n \mathbf{x}_i (y_i - T_i e^{\mathbf{x}_i^T \beta}) \\
 J(\beta) &= \mathbf{X}^T \mathbf{V} \mathbf{X} = \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T T_i \exp\{\mathbf{x}_i^T \beta\} \\
 \mathbf{V} &= \text{diag}(v(\mu_i)) = \text{diag}(T_i e^{\mathbf{x}_i^T \beta})
 \end{aligned}$$

## Poisson Regression Example

**Context:** Analyzing the coronary heart disease of 3,154 males aged 40-50 in a prospective cohort design. We're interested in modeling the number of CHD cases. Risk factors recorded include smoking, blood pressure and behavior type (A and B).

We first fit a main effects model without an offset.

- lack of offset ruins interpretation
  - not all men had same exposure
  - if all men *had* had the same exposure, this model would be fine
- including offset restores interpretative value
- Changing the exposure from Person Years to Person Years / 1000 only changes the **intercept**

$$\beta_0 = \beta_0^* - \log(1000)$$

$$\beta_0^* = \beta_0 + \log(1000)$$