



# ETC3580: Advanced Statistical Modelling

Week 5: Count responses

# Outline

- 1 Poisson regression
- 2 Quasi-Poisson regression
- 3 Negative binomial regression
- 4 Zero inflated count models

# Poisson distribution

Let  $Y$  = number of events in given time interval. If events independent, and prob of event proportional to length of interval, then  $Y$  is Poisson distributed.

## Poisson( $\mu$ ) distribution

$$P(Y = y) = \frac{e^{-\mu} \mu^y}{y!}$$

- $E(Y) = \text{Var}(Y) = \mu$
- If  $Y \sim B(n, p)$ , then  $Y \approx \text{Poisson}(np)$  for small  $p/n$ .
- If  $Y \sim \text{Poisson}(\mu)$ , then  $Y \approx N(\mu, \mu)$  for large  $\mu$ .
- $\text{Poisson}(\mu_1) + \text{Poisson}(\mu_2) \sim \text{Poisson}(\mu_1 + \mu_2)$ .

# Poisson distribution

# Regression with count data

Suppose response  $Y$  is a count  $(0,1,2,\dots)$ .

- If count is bounded and bound is small, use binomial regression.
- If min count is large, use normal approximation.
- Otherwise, use Poisson or negative binomial.

# Poisson regression

$$y_i \sim \text{Poisson}(\mu_i)$$

$$\log(\mu_i) = \eta_i = \beta_0 + \beta_1 x_{i,1} + \cdots + \beta_q x_{i,q}$$

- Log link function forces positive mean.

# Poisson regression

$$y_i \sim \text{Poisson}(\mu_i)$$

$$\log(\mu_i) = \eta_i = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_q x_{i,q}$$

- Log link function forces positive mean.
- Likelihood:  $L = \prod_{i=1}^n \frac{\exp^{-\mu_i} \mu_i^{y_i}}{y_i!}$

# Poisson regression

$$y_i \sim \text{Poisson}(\mu_i)$$

$$\log(\mu_i) = \eta_i = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_q x_{i,q}$$

- Log link function forces positive mean.

- Likelihood: 
$$L = \prod_{i=1}^n \frac{\exp^{-\mu_i} \mu_i^{y_i}}{y_i!}$$

$$\begin{aligned} \log L &= \sum_{i=1}^n [-\mu_i + y_i \log(\mu_i) - \log(y_i!)] \\ &= \sum_{i=1}^n [-\exp(\mathbf{x}_i' \boldsymbol{\beta}) + y_i \mathbf{x}_i' \boldsymbol{\beta} - \log(y_i!)] \end{aligned}$$



# Poisson regression

$$y_i \sim \text{Poisson}(\mu_i)$$

$$\log(\mu_i) = \eta_i = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_q x_{i,q}$$

- Log link function forces positive mean.

- Likelihood: 
$$L = \prod_{i=1}^n \frac{\exp^{-\mu_i} \mu_i^{y_i}}{y_i!}$$

$$\begin{aligned} \log L &= \sum_{i=1}^n [-\mu_i + y_i \log(\mu_i) - \log(y_i!)] \\ &= \sum_{i=1}^n [-\exp(\mathbf{x}_i' \boldsymbol{\beta}) + y_i \mathbf{x}_i' \boldsymbol{\beta} - \log(y_i!)] \end{aligned}$$

```
fit <- glm(y ~ x1 + x2,  
           family='poisson', data)
```

# Deviance

$$\begin{aligned} D &= -2 \log L = -2 \sum_{i=1}^n [-\hat{\mu}_i + y_i \log(\hat{\mu}_i) - \log(y_i!)] \\ &= 2 \sum_{i=1}^n (y_i \log(y_i/\hat{\mu}_i) - (y_i - \hat{\mu}_i)) \end{aligned}$$

# Deviance

$$\begin{aligned} D &= -2 \log L = -2 \sum_{i=1}^n [-\hat{\mu}_i + y_i \log(\hat{\mu}_i) - \log(y_i!)] \\ &= 2 \sum_{i=1}^n (y_i \log(y_i/\hat{\mu}_i) - (y_i - \hat{\mu}_i)) \end{aligned}$$

- Check distributional assumptions by comparing  $D$  against  $\chi^2$
- Compare changes in deviance using a  $\chi^2$  test as for binomial regression.
- Use profile likelihood to find confidence intervals for parameters.
- Common for a model to be “over-dispersed”.

# Dispersed Poisson model

When model is over-dispersed (variance too large):

- estimates of  $\beta$  consistent, but standard errors incorrect.
- could correct model using quasi-Poisson model or negative binomial model.

# Outline

- 1 Poisson regression
- 2 Quasi-Poisson regression
- 3 Negative binomial regression
- 4 Zero inflated count models

# Quasi-Poisson models

Assumes:

$$Y \in \{0, 1, 2, \dots\} \quad E(Y) = \mu, \quad \text{Var}(Y) = \phi\mu$$

```
fit <- glm(y ~ x1 + x2,  
  family='quasipoisson', data)
```

- Use  $F$ -tests not  $\chi^2$  tests when using quasi-Poisson models
- Overdispersion parameter  $\phi$  represents the variance inflation

# Outline

- 1 Poisson regression
- 2 Quasi-Poisson regression
- 3 Negative binomial regression
- 4 Zero inflated count models

# Negative binomial distribution

In series of independent trials, each with prob  $p$  of success, let  $Z$  be number of trials until  $k$ th success.

$$P(Z = z) = \binom{z-1}{k-1} p^k (1-p)^{z-k}, \quad z = k, k+1, \dots$$

- $k = 1$  gives the *geometric* distribution.
- The NegBin distribution also arises when  $Y \sim \text{Poisson}(\lambda)$  and  $\lambda\theta \sim \gamma$  for some constant  $\theta$ .
- For negative binomial regression, we model  $Y = Z - k$ .
- $E(Y) = \mu = k(1-p)/p$  and  $\text{Var}(Y) = \mu + \mu^2/k$ .
- $p = k/(k + \mu)$



# Negative binomial regression

$Y_i \sim \text{NegBin} - k$  with mean  $\mu_i$  and variance  $\mu_i + \mu_i^2/k$ .

$$\eta_i = \log \left( \frac{\mu_i}{\mu_i + k} \right) = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_q x_{i,q}$$

- $k$  (the “dispersion” parameter) is usually estimated along with the coefficients by MLE:

$$L = \prod_{i=1}^n \binom{y_i + k - 1}{k - 1} p_i^k (1 - p_i)^{y_i}$$

$$\begin{aligned} \log L = & \sum_{i=1}^n \left( y_i \log \left( \frac{\mu_i}{\mu_i + k} \right) - k \log(1 + \mu_i/k) \right. \\ & \left. + \sum_{j=0}^{y_i-1} \log(j + k) - \log(y_i!) \right) \end{aligned}$$

# Negative binomial regression

```
# Specified k
```

```
fit <- glm(y ~ x1 + x2,  
  family=negative.binomial(k), data)
```

```
# Estimated k
```

```
fit <- MASS::glm.nb(y ~ x1 + x2, data)
```

# Outline

- 1 Poisson regression
- 2 Quasi-Poisson regression
- 3 Negative binomial regression
- 4 Zero inflated count models

# Zero inflated count models

## Examples of zero-inflated data:

- Number of insurance claims for each account
- Number of arrests for criminal offences for each individual
- Number of articles written by PhD students

Over-dispersed models do not deal adequately with this type of data.

# Zero-inflated count models

## Solution 1: Hurdle model

- Model for probability of zero (logistic).
- Model for non-zero counts (truncated Poisson).

$$P(Y = 0) = p$$

$$P(Y = j) = \frac{1 - p}{1 - f(0)} f(j), \quad j > 0$$

- $p$  is probability of zero;  $f$  is Poisson probability.
- Two sets of coefficients for the two parts of the model.

```
fit <- pscl::hurdle(y ~ x1 + x2, data)
```

# Zero-inflated count models

## Solution 2: Mixture model

- Model for probability of always zero (logistic).
- Model for counts (regular Poisson).

$$P(Y = 0) = p + (1 - p)f(0)$$

$$P(Y = j) = (1 - p)f(j), \quad j > 0$$

- $p$  is probability of zero;  $f$  is Poisson probability.
- Two sets of coefficients for the two parts of the model.

```
fit <- pscl::zeroinfl(y ~ x1 + x2, data)
```

# Zero-inflated count models

- Often difficult to select between these – what makes most sense for the application?
- Can have different predictors for the two sub-models:

```
fit <- zeroinfl(y ~ x1 + x2 | x3, data)
```

count model before the | and zero model after.