

POISSON REGRESSION

If Y_i is a count variable, with values in $\mathbb{N}_0 = \{0, 1, 2, \dots\}$, assuming a Gaussian distribution is not adequate

The most common distribution for a count variable is the Poisson.

Recall that:

$$Y \sim \text{Poisson}(\lambda) \quad \lambda > 0$$

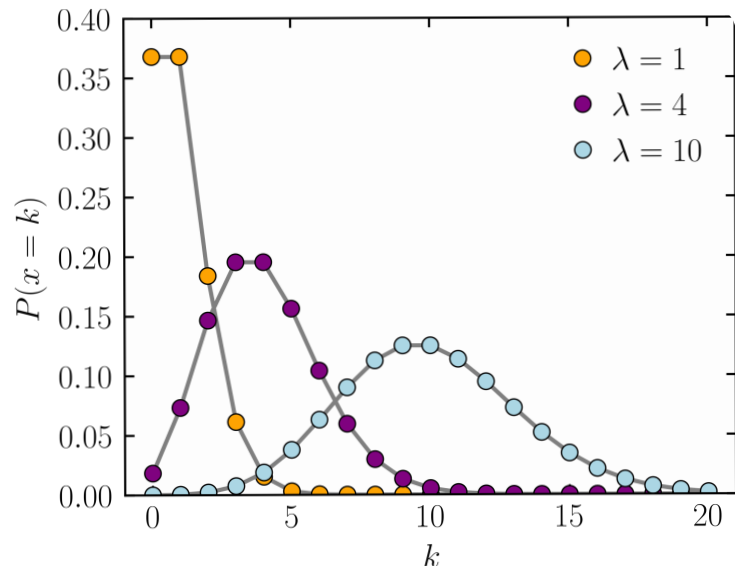
• parameter space: $\Theta = (0, +\infty)$

• support: $\mathcal{Y} = \mathbb{N}_0 = \{0, 1, 2, \dots\}$

$$\text{• probability mass function } p(y; \lambda) = \mathbb{P}(Y=y) = \frac{e^{-\lambda} \lambda^y}{y!}$$

$$\text{• moments: } \mathbb{E}[Y] = \lambda$$

$$\text{var}(Y) = \lambda$$



POISSON REGRESSION: ASSUMPTIONS

1. $Y_i \sim \text{Poisson}(\lambda_i)$ independent for $i=1, \dots, n$

$$2. \eta_i = \mathbf{x}_i^T \underline{\beta}$$

3. $\log(\lambda_i) = \eta_i$ LOGARITHMIC LINK FUNCTION "log-linear model"

Remarks:

• the log link allows mapping the linear predictor $\eta_i = \mathbf{x}_i^T \underline{\beta} \in \mathbb{R}$ to $(0, +\infty)$, the parameter space of λ_i

$$\text{indeed } \log(\lambda_i) = \eta_i \Rightarrow \lambda_i = e^{\eta_i} = e^{\mathbf{x}_i^T \underline{\beta}} > 0$$

We could also use other link functions, however, the log link leads to better theoretical properties (it is the "canonical" link).

• non-constant variance: the Poisson distribution assumes that $\text{var}(Y_i) = \mathbb{E}[Y_i]$

Hence $\text{var}(Y_i) = \lambda_i$ (different between units, by construction).

The distribution of Y_i hence is

$$\mathbb{P}(Y_i = y_i) = \frac{e^{-\lambda_i} \lambda_i^{y_i}}{y_i!}$$

$$\log(\lambda_i) = \mathbf{x}_i^T \underline{\beta} \Rightarrow \lambda_i = e^{\mathbf{x}_i^T \underline{\beta}}$$

$$= \frac{e^{-e^{\mathbf{x}_i^T \underline{\beta}}} e^{\mathbf{x}_i^T \underline{\beta} y_i}}{y_i!}$$

INTERPRETATION OF THE MODEL PARAMETERS

Let's study the mean $\mathbb{E}[Y]$ for two individuals i and k with all the covariates equal except the j -th one, for which we assume $x_{kj} = x_{ij} + 1$

i.e., $x_{ih} = x_{kh}$ for $h=1, \dots, p$ $h \neq j$, $x_{ij} = x_{kj} + 1$

For individual i we get

$$\mathbb{E}[Y_i] = \lambda_i = e^{\mathbf{x}_i^T \underline{\beta}} = \exp\{\beta_1 + \beta_2 x_{i2} + \dots + \beta_{j-1} x_{i,j-1} + \beta_j x_{ij} + \beta_{j+1} x_{i,j+1} + \dots + \beta_p x_{ip}\}$$

For individual k we get

$$\mathbb{E}[Y_k] = \lambda_k = e^{\mathbf{x}_k^T \underline{\beta}} = \exp\{\beta_1 + \beta_2 x_{k2} + \dots + \beta_{j-1} x_{k,j-1} + \beta_j x_{kj} + \beta_{j+1} x_{k,j+1} + \dots + \beta_p x_{kp}\}$$

$$= \exp\{\beta_1 + \beta_2 x_{k2} + \dots + \beta_{j-1} x_{k,j-1} + \beta_j (x_{ij} + 1) + \beta_{j+1} x_{k,j+1} + \dots + \beta_p x_{kp}\}$$

If we study the RATIO

$$\frac{\mathbb{E}[Y_k]}{\mathbb{E}[Y_i]} = \frac{\lambda_k}{\lambda_i} = \frac{\exp\{\beta_1 + \beta_2 x_{k2} + \dots + \beta_{j-1} x_{k,j-1} + \beta_j (x_{ij} + 1) + \beta_{j+1} x_{k,j+1} + \dots + \beta_p x_{kp}\}}{\exp\{\beta_1 + \beta_2 x_{i2} + \dots + \beta_{j-1} x_{i,j-1} + \beta_j x_{ij} + \beta_{j+1} x_{i,j+1} + \dots + \beta_p x_{ip}\}}$$

$$= \exp\{\cancel{\beta_1} + \cancel{\beta_2 x_{i2}} + \dots + \cancel{\beta_{j-1} x_{i,j-1}} + \beta_j (x_{ij} + 1) + \cancel{\beta_{j+1} x_{i,j+1}} + \dots + \cancel{\beta_p x_{ip}} - \cancel{\beta_1} - \cancel{\beta_2 x_{i2}} - \dots - \cancel{\beta_{j-1} x_{i,j-1}} - \beta_j x_{ij} - \cancel{\beta_{j+1} x_{i,j+1}} - \dots - \cancel{\beta_p x_{ip}}\}$$

$$= \exp\{\beta_j (x_{ij} + 1) - \beta_j x_{ij}\}$$

$$= \exp\{\beta_j x_{ij} + \beta_j - \beta_j x_{ij}\} = \exp\{\beta_j\}$$

all terms except the j -th simplify since we assumed $x_{ih} = x_{kh}$ for $h \neq j$.

$$\Rightarrow \frac{\lambda_k}{\lambda_i} = e^{\beta_j}$$

$$\Rightarrow \beta_j = \log \frac{\lambda_k}{\lambda_i} = \log \lambda_k - \log \lambda_i = \log \mathbb{E}[Y | x_j = x_{ij} + 1] - \log \mathbb{E}[Y | x_j = x_{ij}]$$

The parameter β_j represents the DIFFERENCE IN THE LOG OF THE EXPECTED COUNTS IF WE INCREASE x_j OF 1 UNIT, WHILE KEEPING THE OTHER COVARIATES FIXED.

$$\text{or, if we write } e^{\beta_j} = \frac{\lambda_k}{\lambda_i} \Rightarrow \lambda_k = \lambda_i \cdot e^{\beta_j} \Rightarrow \mathbb{E}[Y | x_j = x_{ij} + 1] = \mathbb{E}[Y | x_j = x_{ij}] \cdot e^{\beta_j}$$

The expected counts change of a MULTIPLICATIVE FACTOR e^{β_j} if we increase the j -th covariate of 1 unit, while keeping the other covariates fixed.

ESTIMATION

data (y_1, \dots, y_n) from $Y_i \sim \text{Pois}(\lambda_i) = \text{Pois}(e^{\mathbf{x}_i^T \underline{\beta}})$ indep.

joint density

$$p(y_1, \dots, y_n) = \prod_{i=1}^n p(y_i) = \prod_{i=1}^n \frac{e^{-\lambda_i} \lambda_i^{y_i}}{y_i!} = \prod_{i=1}^n \frac{e^{-e^{\mathbf{x}_i^T \underline{\beta}}} e^{\mathbf{x}_i^T \underline{\beta} y_i}}{y_i!} = \frac{e^{-\sum_{i=1}^n e^{\mathbf{x}_i^T \underline{\beta}}} e^{\sum_{i=1}^n \mathbf{x}_i^T \underline{\beta} y_i}}{\prod_{i=1}^n y_i!}$$

likelihood

$$L(\underline{\beta}) \propto \prod_{i=1}^n p(y_i; \underline{\beta}) \propto e^{-\sum_{i=1}^n e^{\mathbf{x}_i^T \underline{\beta}}} e^{\sum_{i=1}^n \mathbf{x}_i^T \underline{\beta} y_i}$$

log-likelihood

$$\ell(\underline{\beta}) = -\sum_{i=1}^n e^{\mathbf{x}_i^T \underline{\beta}} + \sum_{i=1}^n y_i \mathbf{x}_i^T \underline{\beta}$$

score function

$$\ell_{\mathbf{x}}(\underline{\beta}) = \left\{ \frac{\partial \ell(\underline{\beta})}{\partial \beta_r} \right\}_{r=1, \dots, p}$$

$$\frac{\partial \ell(\underline{\beta})}{\partial \beta_r} = -\sum_{i=1}^n x_{ir} e^{\mathbf{x}_i^T \underline{\beta}} + \sum_{i=1}^n y_i x_{ir} = \sum_{i=1}^n x_{ir} (y_i - e^{\mathbf{x}_i^T \underline{\beta}})$$

Hence the score function can be written as a function of the entire vector $\underline{\beta}$ as:

$$\frac{\partial \ell(\underline{\beta})}{\partial \underline{\beta}} = -\sum_{i=1}^n \mathbf{x}_i e^{\mathbf{x}_i^T \underline{\beta}} + \sum_{i=1}^n \mathbf{x}_i y_i = \sum_{i=1}^n \mathbf{x}_i (y_i - e^{\mathbf{x}_i^T \underline{\beta}}) = \mathbf{X}^T (\underline{y} - \underline{\lambda})$$

The MLE $\hat{\underline{\beta}}$ is the solution of the equation $\ell_{\mathbf{x}}(\underline{\beta}) = 0$

$$\Rightarrow \text{solution of } \mathbf{X}^T (\underline{y} - \underline{\lambda}) = 0$$

$$\mathbf{X}^T (\underline{y} - e^{\mathbf{X} \underline{\beta}}) = 0$$

it resembles the normal equations in the Gaussian LM however, here $\underline{\lambda}$ is a non-linear function of $\underline{\beta}$

This equation does not have an analytical solution: the maximum is found numerically using iterative optimisation methods.

Hence we do not have a closed-form expression for the MLE $\hat{\underline{\beta}}$.

Remark:

notice that, similarly to the LM, since $\hat{\underline{\beta}}$ is the solution of the equation, we obtain

$$\mathbf{X}^T (\underline{y} - e^{\mathbf{X} \hat{\underline{\beta}}}) = 0$$

$$\Rightarrow \mathbf{X}^T (\underline{y} - \hat{\underline{\lambda}}) = 0$$

$$\begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_p^T \end{bmatrix}_{p \times n} \cdot (\underline{y} - \hat{\underline{\lambda}}) = \begin{bmatrix} \mathbf{x}_1^T (\underline{y} - \hat{\underline{\lambda}}) \\ \vdots \\ \mathbf{x}_p^T (\underline{y} - \hat{\underline{\lambda}}) \end{bmatrix} = 0$$

If the model includes the intercept $\Rightarrow \mathbf{x}_1 = \mathbf{1}_n$

$$\Rightarrow \mathbf{x}_1^T (\underline{y} - \hat{\underline{\lambda}}) = \mathbf{1}^T (\underline{y} - \hat{\underline{\lambda}}) = \sum_{i=1}^n (y_i - \hat{\lambda}_i) = 0$$

$$\text{second derivative } \ell_{\mathbf{x}\mathbf{x}}(\underline{\beta}) = \left\{ \frac{\partial^2 \ell(\underline{\beta})}{\partial \beta_r \partial \beta_s} \right\}_{r,s=1, \dots, p} = -\sum_{i=1}^n x_{ir} x_{is} e^{\mathbf{x}_i^T \underline{\beta}}$$

$$= -\sum_{i=1}^n x_{ir} x_{is} \lambda_i$$

In matrix form we get $\ell_{\mathbf{x}\mathbf{x}}(\underline{\beta}) = -\mathbf{X}^T \mathbf{U} \mathbf{X}$ with \mathbf{U} an $n \times n$ diagonal matrix

$$\mathbf{U} = \begin{bmatrix} \lambda_1 & 0 & 0 & \dots & 0 \\ 0 & \lambda_2 & 0 & \dots & 0 \\ 0 & 0 & \ddots & & \vdots \\ \vdots & \vdots & & \lambda_{n-1} & 0 \\ 0 & 0 & \dots & 0 & \lambda_n \end{bmatrix} = \text{diag}\{\lambda_1, \dots, \lambda_n\} = \text{diag}\{e^{\mathbf{x}_1^T \underline{\beta}}, \dots, e^{\mathbf{x}_n^T \underline{\beta}}\}$$

\Rightarrow it is a function of $\underline{\beta} \Rightarrow \mathbf{U} = \mathbf{U}(\underline{\beta})$

The OBSERVED INFORMATION evaluated at the MLE $\hat{\underline{\beta}}$ is

$$\mathbf{j}(\hat{\underline{\beta}}) = -\ell_{\mathbf{x}\mathbf{x}}(\underline{\beta}) \big|_{\underline{\beta}=\hat{\underline{\beta}}} = \mathbf{X}^T \mathbf{U}(\hat{\underline{\beta}}) \mathbf{X}$$

where $\mathbf{U}(\hat{\underline{\beta}}) = \text{diag}\{e^{\mathbf{x}_1^T \hat{\underline{\beta}}}, \dots, e^{\mathbf{x}_n^T \hat{\underline{\beta}}}\}$

INFERENCE

inference here is based on APPROXIMATE distributions

Remarks:

- notation: we write "Y approximately distributed as (some distribution $p(y)$)" as " $Y \sim p(y)$ "
- approximations get better with n (large samples \rightarrow better approximation)

DISTRIBUTION of the MAXIMUM LIKELIHOOD ESTIMATOR of the REGRESSION PARAMETERS

$$\hat{\underline{\beta}} \sim N_p(\underline{\beta}, \mathbf{j}(\hat{\underline{\beta}})^{-1})$$

the marginal distribution for the j -th element is $\hat{\beta}_j \sim N(\beta_j, [\mathbf{j}(\hat{\underline{\beta}})^{-1}]_{jj})$ $j=1, \dots, p$

CONFIDENCE INTERVAL FOR β_j

A pivotal quantity is

$$\frac{\hat{\beta}_j - \beta_j}{\sqrt{[\mathbf{j}(\hat{\underline{\beta}})^{-1}]_{jj}}} \sim N(0, 1)$$

a confidence interval with level $(1-\alpha)$ for β_j ($j=1, \dots, p$) can be obtained as

$$\mathbb{P}\left(-z_{1-\frac{\alpha}{2}} < \frac{\hat{\beta}_j - \beta_j}{\sqrt{[\mathbf{j}(\hat{\underline{\beta}})^{-1}]_{jj}}} < z_{1-\frac{\alpha}{2}}\right) = 1-\alpha$$

Gaussian is symmetric $z_{\frac{\alpha}{2}} = -z_{1-\frac{\alpha}{2}}$

quantile of level $1-\frac{\alpha}{2}$ of a $N(0, 1)$

with the data:

$$\hat{\beta}_j - \sqrt{[\mathbf{j}(\hat{\underline{\beta}})^{-1}]_{jj}} \cdot z_{1-\frac{\alpha}{2}} < \beta_j < \hat{\beta}_j + \sqrt{[\mathbf{j}(\hat{\underline{\beta}})^{-1}]_{jj}} \cdot z_{1-\frac{\alpha}{2}}$$

$$\Rightarrow \beta_j \in \hat{\beta}_j \pm z_{1-\frac{\alpha}{2}} \sqrt{[\mathbf{j}(\hat{\underline{\beta}})^{-1}]_{jj}}$$

TEST ABOUT β_j

$$\text{consider the test } \begin{cases} H_0: \beta_j = b_j \\ H_1: \beta_j \neq b_j \end{cases}$$

We can use the test statistic

$$Z_j = \frac{\hat{\beta}_j - b_j}{\sqrt{[\mathbf{j}(\hat{\underline{\beta}})^{-1}]_{jj}}} \sim N(0, 1) \text{ under } H_0$$

the observed value of the test is z_j^{obs}

if we use a fixed significance level α

$$\alpha = \mathbb{P}_{H_0}(|Z_j| > z_{1-\frac{\alpha}{2}})$$

\rightarrow reject region is $R = (-\infty, -z_{1-\frac{\alpha}{2}}) \cup (z_{1-\frac{\alpha}{2}}, +\infty)$



if we use the observed significance level

$$\text{the p-value is } \omega^{\text{obs}} = \mathbb{P}_{H_0}(|Z_j| \geq |z_j^{\text{obs}}|) = 2(1 - \Phi(|z_j^{\text{obs}}|))$$

