

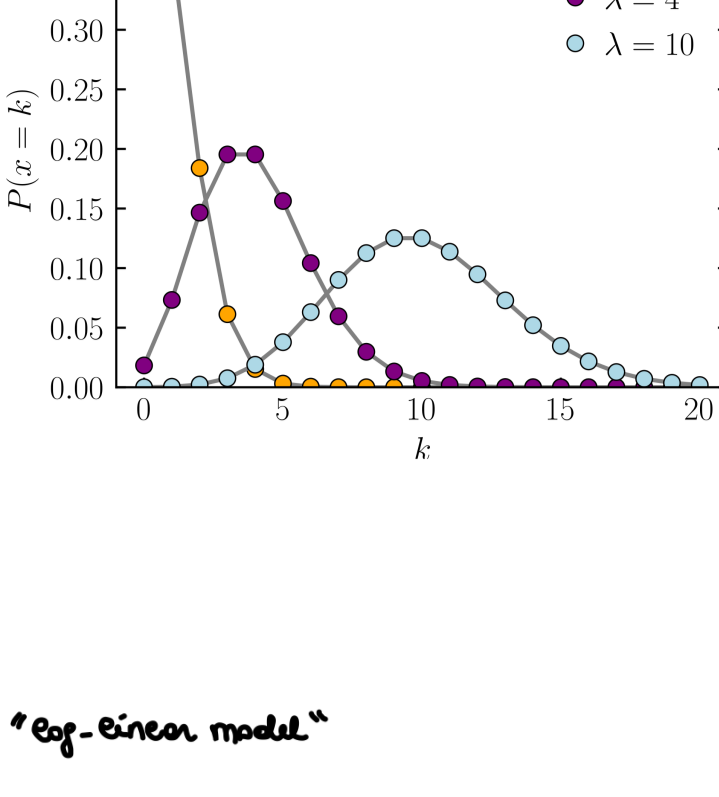
POISSON REGRESSION

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:

- $Y \sim \text{Poisson}(\mu)$
- parameter space: $\mu > 0 \Rightarrow (0, +\infty)$
- support: $\mathcal{Y} = \mathbb{N}_0 = \{0, 1, 2, \dots\}$
- probability mass function $p(y; \mu) = P(Y=y) = \frac{e^{-\mu} \mu^y}{y!}$
- moments: $E[Y] = \mu$
 $\text{var}(Y) = \mu$



POISSON REGRESSION: ASSUMPTIONS

- $Y_i \sim \text{Poisson}(\mu_i)$ independent for $i=1, \dots, n$
- $\eta_i = \beta^T \mathbf{x}_i$
- $g(\mu_i) = \eta_i$ with $g = \text{Exp}$ LOGARITHMIC LINK FUNCTION "log-linear model"

Remarks:

- the log link allows mapping the linear predictor $\eta_i = \mathbf{x}_i^T \beta \in \mathbb{R}$ to \mathbb{R}^+ , the parameter space of μ_i
indeed $\text{Exp}(\mu_i) = \eta_i \Rightarrow \mu_i = e^{\eta_i} = e^{\mathbf{x}_i^T \beta}$
- 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 $\text{var}(Y_i) = E[Y_i]$
Hence $\text{var}(Y_i) = \mu_i$ (different between units, by construction).

The model is

$$P(Y_i = y_i) = \frac{e^{-\mu_i} \mu_i^{y_i}}{y_i!} \quad \text{Exp}(\mu_i) = \mathbf{x}_i^T \beta \Rightarrow \mu_i = e^{\mathbf{x}_i^T \beta}$$

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

INTERPRETATION of the regression parameters

Let's study the mean μ_i at two values x_j and $x_j + 1$ of the j -th covariate

at x_j we obtain $\mu_1 = \exp\{\beta_0 + \beta_1 x_1 + \dots + \beta_j x_j + \dots + \beta_p x_p\}$

at $(x_j + 1)$ $\mu_2 = \exp\{\beta_0 + \beta_1 x_1 + \dots + \beta_j (x_j + 1) + \dots + \beta_p x_p\}$

$$\frac{\mu_2}{\mu_1} = \frac{\exp\{\beta_0 + \beta_1 x_1 + \dots + \beta_j (x_j + 1) + \dots + \beta_p x_p\}}{\exp\{\beta_0 + \beta_1 x_1 + \dots + \beta_j x_j + \dots + \beta_p x_p\}}$$

$$= \exp\{\beta_0 + \beta_1 x_1 + \dots + \beta_j (x_j + 1) + \dots + \beta_p x_p - \beta_0 - \beta_1 x_1 - \dots - \beta_j x_j - \dots - \beta_p x_p\}$$

$$= \exp\{\beta_j\} = e^{\beta_j}$$

$$\Rightarrow \beta_j = \text{Log} \frac{\mu_2}{\mu_1} = \text{Log} \mu_2 - \text{Log} \mu_1 = \text{Log} E(Y | x_j = x_j + 1) - \text{Log} E(Y | x_j = x_j)$$

The parameter β_j represents the difference of the Logs of the expected counts if we increase x_j of 1 unit, while keeping the other predictors fixed.

Or, if we write: $e^{\beta_j} = \frac{\mu_2}{\mu_1} \Rightarrow \mu_2 = \mu_1 \cdot e^{\beta_j}$

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.

ESTIMATE

likelihood

$$L(\beta) \propto \prod_{i=1}^n p(y_i | \beta) = \prod_{i=1}^n \frac{e^{-\mu_i} \mu_i^{y_i}}{y_i!} \quad \mu_i = e^{\mathbf{x}_i^T \beta}$$

$$\propto e^{-\sum_{i=1}^n \mu_i} \prod_{i=1}^n \mu_i^{y_i}$$

$$\ell(\beta) = -\sum_{i=1}^n \mu_i + \sum_{i=1}^n y_i \text{Log} \mu_i = -\sum_{i=1}^n e^{\mathbf{x}_i^T \beta} + \sum_{i=1}^n y_i \mathbf{x}_i^T \beta \quad (\text{log-likelihood})$$

$$\ell_{\beta}(\beta) = \left\{ \frac{\partial \ell(\beta)}{\partial \beta_r} \right\}_{r=0, \dots, p} \quad (\text{score function})$$

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

$$= \sum_{i=1}^n x_{ir} (y_i - \mu_i)$$

for the vector β

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

$p \times n \quad n \times 1$

The MLE $\hat{\beta}$ is the solution of $\mathbf{X}^T (\mathbf{y} - \boldsymbol{\mu}) = \mathbf{0} \rightarrow$ it resembles the normal equations, $\mathbf{X}^T (\mathbf{y} - e^{\boldsymbol{\mu}}) = \mathbf{0}$ but here $\boldsymbol{\mu}$ is a nonlinear function of β .

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 $\hat{\beta}$.

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

$$\mathbf{X}^T (\mathbf{y} - \hat{\boldsymbol{\mu}}) = \mathbf{0} \Rightarrow \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_n^T \end{bmatrix} \cdot (\mathbf{y} - \hat{\boldsymbol{\mu}}) = \begin{bmatrix} \mathbf{x}_1^T (\mathbf{y} - \hat{\boldsymbol{\mu}}) \\ \vdots \\ \mathbf{x}_n^T (\mathbf{y} - \hat{\boldsymbol{\mu}}) \end{bmatrix} = \mathbf{0}$$

If the model includes the intercept $\Rightarrow \mathbf{x}_2 = \mathbf{1}_n \Rightarrow \mathbf{x}_2^T (\mathbf{y} - \hat{\boldsymbol{\mu}}) = \mathbf{1}_n^T (\mathbf{y} - \hat{\boldsymbol{\mu}}) = \sum_{i=1}^n (y_i - \hat{\mu}_i) = 0$

second derivative $\ell_{\beta\beta}(\beta) = \left\{ \frac{\partial^2 \ell(\beta)}{\partial \beta_r \partial \beta_s} \right\}_{r,s=0, \dots, p} = -\sum_{i=1}^n x_{ir} x_{is} e^{\mathbf{x}_i^T \beta}$
 $= -\sum_{i=1}^n x_{ir} x_{is} \mu_i < 0 \Rightarrow \hat{\beta}$ is a max

Hence the matrix is $\ell_{\beta\beta}(\beta) = -\mathbf{X}^T \mathbf{U} \mathbf{X}$ with $\mathbf{U} = \text{diag}\{\mu_1, \dots, \mu_n\} = \text{diag}\{e^{\mathbf{x}_1^T \beta}, \dots, e^{\mathbf{x}_n^T \beta}\} = \mathbf{U}(\beta)$

observed information evaluated at the MLE $\hat{\beta}$ is

$$j(\hat{\beta}) = -\text{Cov}(\beta) \Big|_{\beta=\hat{\beta}} = \mathbf{X}^T \mathbf{U}(\hat{\beta}) \mathbf{X} \quad \text{where } \mathbf{U}(\hat{\beta}) = \text{diag}\{e^{\mathbf{x}_1^T \hat{\beta}}, \dots, e^{\mathbf{x}_n^T \hat{\beta}}\}$$

INFERENCE

inference here is based on approximate distributions
(we write "approximately distributed as" with $\hat{\sim}$)
(we write "approximately get better for large n ")

DISTRIBUTION of the MAXIMUM LIKELIHOOD ESTIMATOR of the REGRESSION PARAMETERS

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

the marginal is $\hat{\beta}_j \hat{\sim} N(\beta_j, [j(\hat{\beta})^{-1}]_{jj})$

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

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

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

A test $H_0: \beta_j = b_j$ vs $H_1: \beta_j \neq b_j$ is performed as

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

the p-value is $\alpha^{\text{obs}} = P_{H_0}(|Z_j| \geq |z_j^{\text{obs}}|) = 2(1 - \Phi(|z_j^{\text{obs}}|))$

TEST for comparing NESTED MODELS

We have a model $Y_i \sim \text{Pois}(\mu_i)$ ($i=1, \dots, n$)

with $\text{Log}(\mu_i) = \mathbf{x}_i^T \beta = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \beta_{p+1} x_{i,p+1} + \dots + \beta_p x_{ip}$
we call it the "full" model (it is the proposed model). under H_0 this is 0

we want to test

$$\begin{cases} H_0: \beta_{p+1} = \dots = \beta_p = 0 \\ H_1: \bar{H}_0 \end{cases}$$

We can partition the vector $\beta = \begin{bmatrix} \beta^{(0)} \\ \beta^{(1)} \end{bmatrix}$ $\beta^{(0)} \in \mathbb{R}^p$
 $\beta^{(1)} \in \mathbb{R}^{p-p}$

$$\begin{cases} H_0: \beta^{(1)} = \mathbf{0} \\ H_1: \beta^{(1)} \neq \mathbf{0} \end{cases}$$

under H_0 we have the "restricted model"

$$Y_i \sim \text{Pois}(\mu_i) \text{ with } \text{Log}(\mu_i) = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$$

To compare two nested models we use the LIKELIHOOD RATIO TEST: it compares the maximum of the likelihood of the full and of the restricted model:

$$W = 2 \log \frac{\hat{\ell}(\text{model})}{\hat{\ell}(\text{restricted})} = 2 \{ \hat{\ell}(\text{model}) - \hat{\ell}(\text{restricted}) \} \hat{\sim} \chi_{p-p_0}^2 \text{ under } H_0$$

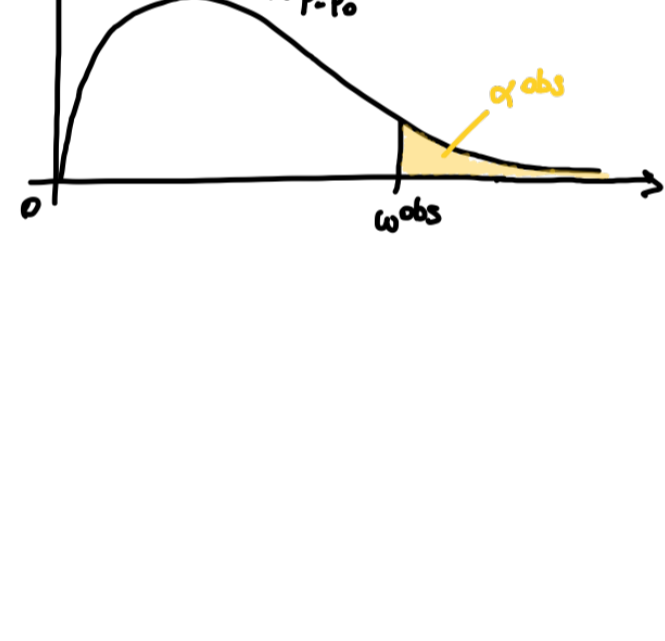
$(\# \text{ covariates under } H_1) - (\# \text{ covariates under } H_0)$

We can denote with $\hat{\beta} = (\hat{\beta}^{(0)}, \hat{\beta}^{(1)})$ the MLE under H_1 (full model) and with $\hat{\beta}^{(0)} = (\hat{\beta}^{(0)}, \mathbf{0})$ the MLE under H_0 (restricted). Then,

$$W = 2 \{ \ell(\hat{\beta}^{(0)}, \hat{\beta}^{(1)}) - \ell(\hat{\beta}^{(0)}, \mathbf{0}) \} \hat{\sim} \chi_{p-p_0}^2 \text{ under } H_0$$

with the data we compute w^{obs}

and the p-value is $\alpha^{\text{obs}} = P_{H_0}(W \geq w^{\text{obs}}) = P(\chi_{p-p_0}^2 \geq w^{\text{obs}})$



TEST about the OVERALL SIGNIFICANCE

Similarly to the LM, we can test

$$\begin{cases} H_0: \beta_1 = \beta_2 = \dots = \beta_p = 0 \\ H_1: \bar{H}_0 \end{cases}$$

We can use the test for nested models by setting $\beta_0 = 1$.

In this case we compare the full model with a model with only the intercept ("null model").

Under H_0 : $Y_i \sim \text{Pois}(\mu_i)$ indep $i=1, \dots, n$

$$\mu_i = e^{\beta_0} = \mu$$

$$L(\mu) = \prod_{i=1}^n \frac{e^{-\mu} \mu^{y_i}}{y_i!} \propto e^{-n\mu} \mu^{\sum y_i}$$

$$\ell(\mu) = -n\mu + \sum_{i=1}^n y_i \text{Log} \mu = -n\mu + n\bar{y} \text{Log} \mu$$

$$\ell_{\mu}(\mu) = -n + \frac{n\bar{y}}{\mu}$$

$$\ell_{\mu}(\mu) = 0 \Rightarrow -n\mu = -n\bar{y} \Rightarrow \hat{\mu} = \bar{y} \quad \text{MLE under the restricted ("null") model}$$

Moreover, since $\text{Log}(\mu_i) = \text{Log} \mu = \beta_0$ (one-to-one correspondence: bijective function)
we automatically obtain $\hat{\beta}_0 = \text{Log} \hat{\mu} = \text{Log} \bar{y}$

$$\ell_{\beta\beta}(\mu) = -\frac{n\bar{y}}{\mu^2} \frac{1}{\mu^2} = -\frac{n\bar{y}}{\mu^3} < 0 \quad \text{it's a max}$$

Hence $\hat{\ell}(\text{restricted}) = \ell(\hat{\mu}) = -n\bar{y} + n\bar{y} \text{Log} \bar{y} = n(\bar{y} \text{Log} \bar{y} - \bar{y})$

Under H_1 we have the model with p covariates

$$\mu_i = \exp\{\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}\}$$

we estimate $\hat{\beta}$ numerically, and we compute $\hat{\ell}(\text{model}) = \ell(\hat{\beta}) = -\sum_{i=1}^n \mu_i + \sum_{i=1}^n y_i \text{Log} \mu_i = -\sum_{i=1}^n e^{\mathbf{x}_i^T \hat{\beta}} + \sum_{i=1}^n y_i \mathbf{x}_i^T \hat{\beta}$

The likelihood ratio test in this case is

$$W = 2 \{ \hat{\ell}(\text{model}) - \hat{\ell}(\text{restricted}) \} = 2 \{ \ell(\hat{\beta}) - \ell(\bar{y}) \} \hat{\sim} \chi_{p-1}^2 \text{ under } H_0$$

with the data:

$$w^{\text{obs}} = 2 \left\{ -\sum_{i=1}^n \mu_i + \sum_{i=1}^n y_i \text{Log} \mu_i - n(\bar{y} \text{Log} \bar{y} - \bar{y}) \right\}$$

$$= 2 \left\{ \sum_{i=1}^n y_i \text{Log} \frac{\mu_i}{\bar{y}} - \sum_{i=1}^n \mu_i \frac{\mu_i}{\bar{y}} + n\bar{y} \right\}$$

$$= 2 \left\{ \sum_{i=1}^n y_i \text{Log} \frac{\mu_i}{\bar{y}} - \sum_{i=1}^n \frac{\mu_i^2}{\bar{y}} + n\bar{y} \right\}$$

reject H_0 if $w^{\text{obs}} \geq \chi_{p-1, 1-\alpha}^2$ (quantile of level $1-\alpha$ of a χ_{p-1}^2)

TEST about the GOODNESS OF FIT of the model

we first need to introduce the concept of "SATURATED MODEL". This is the most elaborated model one can consider, i.e., a model with n parameters (one for each observation).

We want to estimate a model with n parameters using n observations \Rightarrow we obtain $\hat{\mu}_i = y_i$ $i=1, \dots, n$
a model with a perfect fit (perfect but useless: interpolation, there is no simplification \rightarrow the model is keeping all the erratic variability of the data and is not highlighting the underlying systematic behavior)

$$\begin{cases} \ell(\mu_i) = -\mu_i + y_i \text{Log} \mu_i \\ \Rightarrow \ell_{\mu_i}(\mu_i) = -1 + \frac{y_i}{\mu_i} \Rightarrow \frac{y_i}{\mu_i} - 1 = 0 \Rightarrow \hat{\mu}_i = y_i \\ \Rightarrow \ell(\hat{\mu}_1, \dots, \hat{\mu}_n) = \sum_{i=1}^n y_i \text{Log} y_i - \sum_{i=1}^n y_i \end{cases}$$

The ("full") model with p parameters can be compared with the saturated model using the likelihood ratio test. For this particular case, this quantity is called DEVNANCE (or "residual deviance")

$$D = \text{deviance}(\text{model}) = 2 \{ \hat{\ell}(\text{saturated}) - \hat{\ell}(\text{model}) \}$$

Since $\hat{\mu}_i = y_i$ for all i in the saturated model, the expected likelihood evaluated at $\hat{\mu}_i$ is always $\hat{\ell}(\text{saturated}) = \sum_{i=1}^n y_i \text{Log} y_i - \sum_{i=1}^n y_i$

Hence, $D = 2 \{ \hat{\ell}(\text{saturated}) - \ell(\hat{\mu}) \} = 2 \left\{ \sum_{i=1}^n y_i \text{Log} y_i - \sum_{i=1}^n y_i - \sum_{i=1}^n (y_i \text{Log} \mu_i - \mu_i) \right\}$
estimate of μ in the model $\hat{\mu}$ is a transformation of β
 $= 2 \left\{ \sum_{i=1}^n y_i \text{Log} \frac{y_i}{\mu_i} - \sum_{i=1}^n (y_i - \mu_i) \right\} = 2 \sum_{i=1}^n y_i \text{Log} \frac{y_i}{\mu_i}$
does not depend on β
if the model includes the intercept

Since the saturated model has a perfect fit, for sure $\hat{\ell}(\text{saturated}) > \hat{\ell}(\text{model}) \Rightarrow D > 0$.
Moreover, if the model with p covariates fits the data well, $\hat{\ell}(\text{model})$ will not be "too far" from $\hat{\ell}(\text{saturated})$.

\rightarrow A good model will have a small deviance

\rightarrow WE DO NOT HAVE A DISTRIBUTION FOR THE DEVNANCE

When the LR test is used to compare the saturated model we lose the approximation to a χ^2 distribution

We can't do formal tests, a deviance $< n-p$ is generally ok. It indicates that the model fits the data well: it does not lack too much accuracy compared to the "perfect" saturated model.

Notice that deviance = 0 means that you fit the data perfectly, but with p parameters instead of n .

The residual deviance is more useful when used to compare different models (on the same data), with the same number of covariates p (but different covariates).

RELATIONSHIP between SATURATED, PROPOSED and NULL MODEL

Notice that the saturated model and the null model (with only the intercept) are the two extreme cases:

- SATURATED model: n parameters
 - PROPOSED model: p parameters
 - NULL model: 1 parameter
- $\left. \begin{array}{l} \text{nested} \\ \text{nested} \end{array} \right\}$

When you estimate a gem in \mathbb{R}^p , the output returns 2 quantities:

- "Residual deviance" and "Null deviance"

likelihood ratio test:

between saturated and proposed (the proposed model can be seen as a "restricted" model w.r.t. the saturated)

$$2 \{ \hat{\ell}(\text{saturated}) - \hat{\ell}(\text{model}) \} = D(\text{model}) \quad \text{"residual deviance"}$$

between saturated and null (also the null model can be seen as a restricted model w.r.t. the saturated)

$$2 \{ \hat{\ell}(\text{saturated}) - \hat{\ell}(\text{null}) \} = D(\text{null}) \quad \text{"null deviance"}$$

between model and null (test of overall significance)

$$2 \{ \hat{\ell}(\text{model}) - \hat{\ell}(\text{null}) \} = 2 \{ \hat{\ell}(\text{model}) + \hat{\ell}(\text{saturated}) - \hat{\ell}(\text{saturated}) - \hat{\ell}(\text{null}) \}$$

$$= 2 \{ [\hat{\ell}(\text{saturated}) - \hat{\ell}(\text{null})] - [\hat{\ell}(\text{saturated}) - \hat{\ell}(\text{model})] \}$$

$$= D(\text{null}) - D(\text{model})$$

\rightarrow we can write the LR test in terms of difference of the deviances

MODEL CHECKING: RESIDUALS

In the linear model we had $\hat{y} = \hat{\mu} + (y - \hat{\mu}) = \hat{\mu} + e$ and we studied the expected behavior of the residuals when the model assumptions are valid, to compare it with the observed behavior after fitting the model.

In the linear model the residuals were the "sample counterpart" of the errors: here we do not have them since now it is not so clear how to define residuals, several versions have been proposed.

• Pearson's residuals: they are the analogous version of the standardized residuals in the LM.

$$e_i = \frac{y_i - \hat{\mu}_i}{\sqrt{V(\hat{\mu}_i)}} \quad i=1, \dots, n$$

For the Poisson, we have $V(\mu) = \mu \Rightarrow e_i = \frac{y_i - \hat{\mu}_i}{\sqrt{\hat{\mu}_i}} \quad i=1, \dots, n$

They have approximately zero mean and constant variance.