

**MBA  
USP  
ESALQ**

**SUPERVISED MACHINE LEARNING:  
REGRESSION MODELS FOR  
COUNT DATA**

Prof. Dr. Luiz Paulo Fávero

\*The responsibility for trustworthiness, originality, and legality of the didactic content presented is responsibility of the professor.

The total or partial reproduction of this material without authorization is **prohibited**.

Law No. 9610/98



## **MODELS FOR COUNT DATA**

**Theoretical foundation, concepts and applications**

**Model specification and canonical connection functions**

**Models of Poisson and negative binomial types**

**Estimation of parameters by maximum likelihood**

**Identification of the phenomenon of overdispersion in the data**

**Zero-Inflated Models**

**Estimate in R**

# Generalized Linear Model (GLM)

$$\eta_i = \alpha + \beta_1 \cdot X_{1i} + \beta_2 \cdot X_{2i} + \dots \beta_k \cdot X_{ki}$$

Modelos lineares generalizados, características da variável dependente e funções de ligação canônica.

Modelo de Regressão	Característica da Variável Dependente	Distribuição	Função de Ligação Canônica ( $\eta$ )
Linear	Quantitativa	Normal	$\hat{Y}$
Com Transformação de Box-Cox	Quantitativa	Normal Após a Transformação	$\frac{\hat{Y}^\lambda - 1}{\lambda}$
Logística Binária	Qualitativa com 2 Categorias ( <i>Dummy</i> )	Bernoulli	$\ln\left(\frac{p}{1-p}\right)$
Logística Multinomial	Qualitativa $M$ ( $M > 2$ ) Categorias	Binomial	$\ln\left(\frac{p_m}{1-p_m}\right)$
Poisson	Quantitativa com Valores Inteiros e Não Negativos (Dados de Contagem)	Poisson	$\ln(\lambda_{poisson})$
Binomial Negativo	Quantitativa com Valores Inteiros e Não Negativos (Dados de Contagem)	Poisson-Gama	$\ln(\lambda_{bneg})$

Siméon Denis Poisson



(1781-1840)

## Models for Count Data

---

Poisson and negative binomial regression models are part of what is known as regression models for count data, and aim to analyze the behavior, according to predictor variables, of a determined dependent variable that is presented in the quantitative form with discrete and non-negative values. Exposure must also be defined (temporal, spatial, social unit, among others).

## Models for Count Data: Examples and Applications

---

- Evaluation of the number of times that a group of elderly patients goes to the doctor per year, according to the age of each of them, the gender and the characteristics of their health plans.
- A study on the number of public offers of actions that are carried out in a sample of developed and emerging countries in a given year, based on their economic performance, such as inflation, interest rate, gross domestic product, and foreign investment rate.

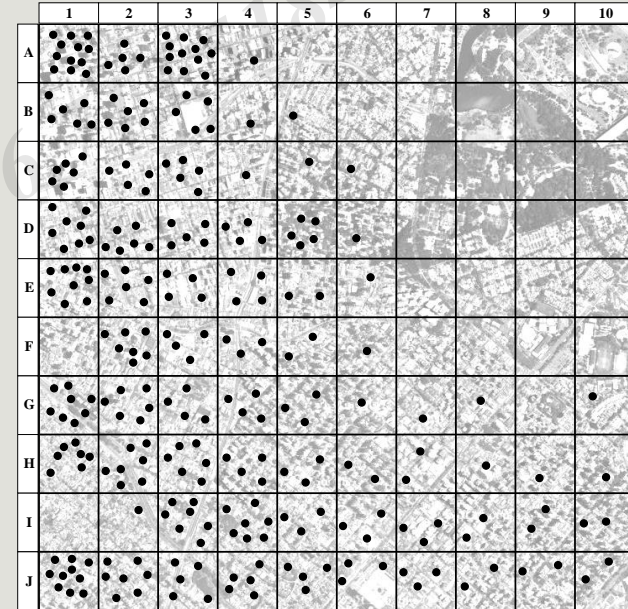
*Notice that the number of going to the doctor or the number of public offers of shares are the dependent variables in both cases, being represented by quantitative data that assume discrete, non-negative values, and with annual exposure. That is, they offer count data.*



# Models for Count Data: Examples and Applications



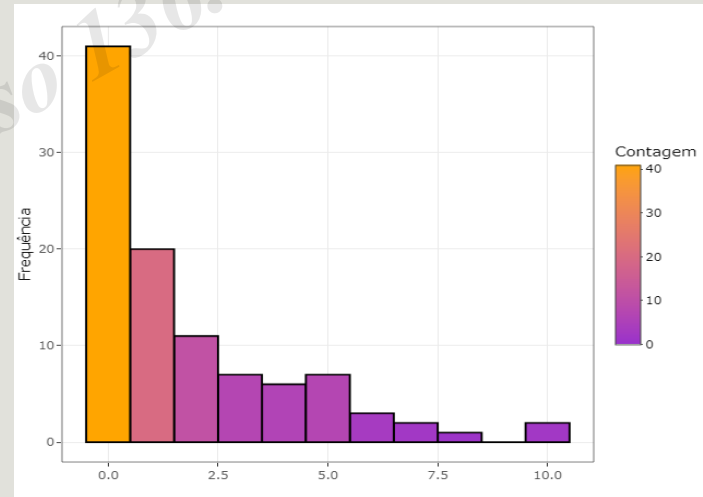
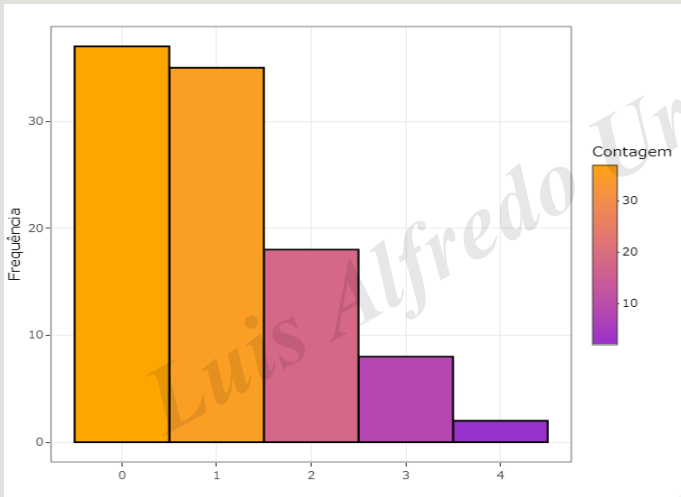
Environment



Property Market

## Poisson and Binomial Negative Distributions

$$\ln(\hat{Y}_i) = \alpha + \beta_1.X_{1i} + \beta_2.X_{2i} + \dots + \beta_k.X_{ki}$$



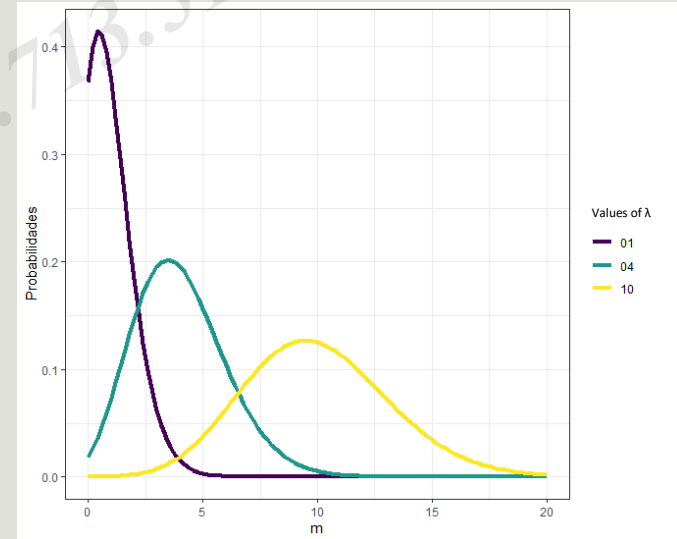


# Poisson Distribution

Certain observation  $i$  ( $i = 1, 2, \dots, n$ , in which  $n$  is the sample size) has the following probability of occurrence of a count  $m$  in a certain exposure (period, area, region, among other examples):

$$p(Y_i = m) = \frac{e^{-\lambda_i} \cdot \lambda_i^m}{m!}$$

in which  $\lambda$  is the expected number of occurrences or the estimated average incidence rate of the phenomenon under study for a given exposure.



# The Poisson Distribution and the Poisson Model

**Mean:** 
$$E(Y) = \sum_{m=0}^{\infty} m \cdot \frac{e^{-\lambda} \cdot \lambda^m}{m!} = \lambda \cdot \sum_{m=1}^{\infty} \frac{e^{-\lambda} \cdot \lambda^{m-1}}{(m-1)!} = \lambda \cdot 1 = \lambda$$

**Variance:** 
$$\begin{aligned} Var(Y) &= \sum_{m=0}^{\infty} m \cdot \frac{e^{-\lambda} \cdot \lambda^m}{m!} \cdot (m - \lambda)^2 = \sum_{m=0}^{\infty} m \cdot \frac{e^{-\lambda} \cdot \lambda^m}{m!} \cdot (m^2 - 2 \cdot m \cdot \lambda + \\ &\lambda^2) = \lambda^2 \cdot \sum_{m=2}^{\infty} \frac{e^{-\lambda} \cdot \lambda^{m-2}}{(m-2)!} + \lambda \cdot \sum_{m=1}^{\infty} \frac{e^{-\lambda} \cdot \lambda^{m-1}}{(m-1)!} - \lambda^2 = \lambda \end{aligned}$$

**General Model:**

$$\ln(\hat{Y}_i) = \ln(\lambda_{poisson_i}) = \alpha + \beta_1 \cdot X_{1i} + \beta_2 \cdot X_{2i} + \dots + \beta_k \cdot X_{ki}$$





## Overdispersion Test

---

$$Y_i^* = \frac{\left[ \left( Y_i - \lambda_{poisson_i} \right)^2 - Y_i \right]}{\lambda_{poisson_i}}$$

$$Y_i^* = \beta \cdot \lambda_{poisson_i}$$

Cameron and Trivedi (1990) state that if the phenomenon of overdispersion in the data occurs, the estimated parameter  $\beta$  of this **auxiliary model without intercept** will be statistically different from zero, at a given level of significance (5%, usually).

## Gamma-Poisson Distribution or Negative Binomial

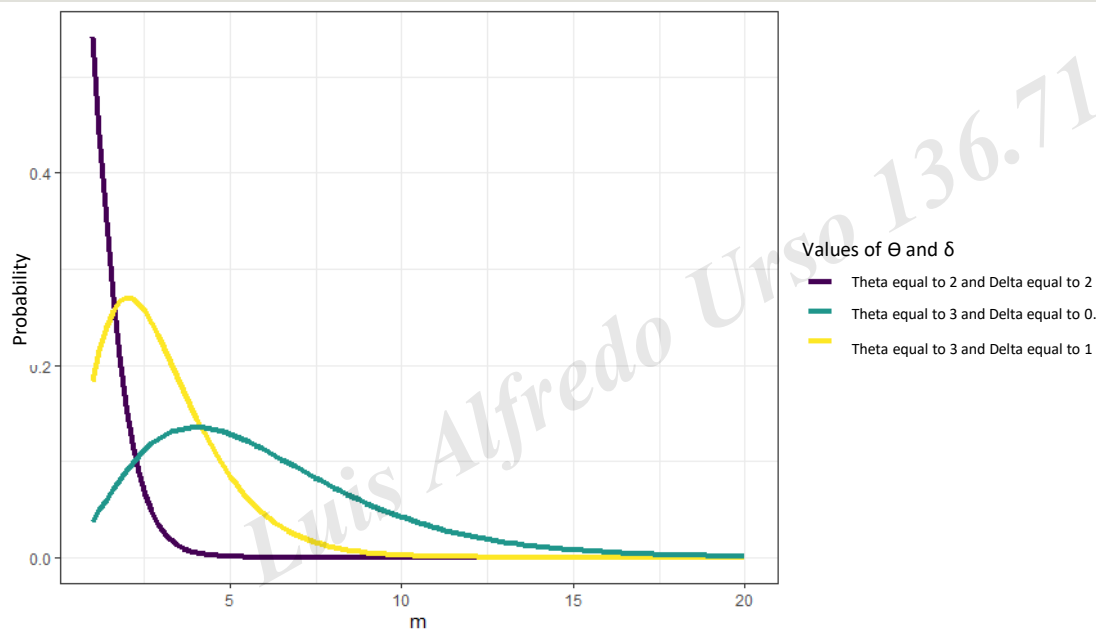
---

For a certain observation  $i$  ( $i = 1, 2, \dots, n$ , in which  $n$  is the sample size), the function of the probability distribution of the dependent variable  $Y$  will be given by:

$$p(Y_i = m) = \frac{\delta^\theta \cdot m_i^{\theta-1} \cdot e^{-m_i \cdot \delta}}{(\theta-1)!}$$

in which  $\theta$  is called a form parameter ( $\theta > 0$ ) and  $\delta$  is called a decay rate parameter ( $\delta > 0$ ).

# Gamma-Poisson Distribution or Negative Binomial



- Mean:

$$E(Y) = \lambda_{bneg}$$

- Variance:

$$Var(Y) = \lambda_{bneg} + \phi \cdot (\lambda_{bneg})^2$$

$$\phi = \frac{1}{\theta}$$

**NB2 Models**

## The Gamma-Poisson Model or Negative Binomial



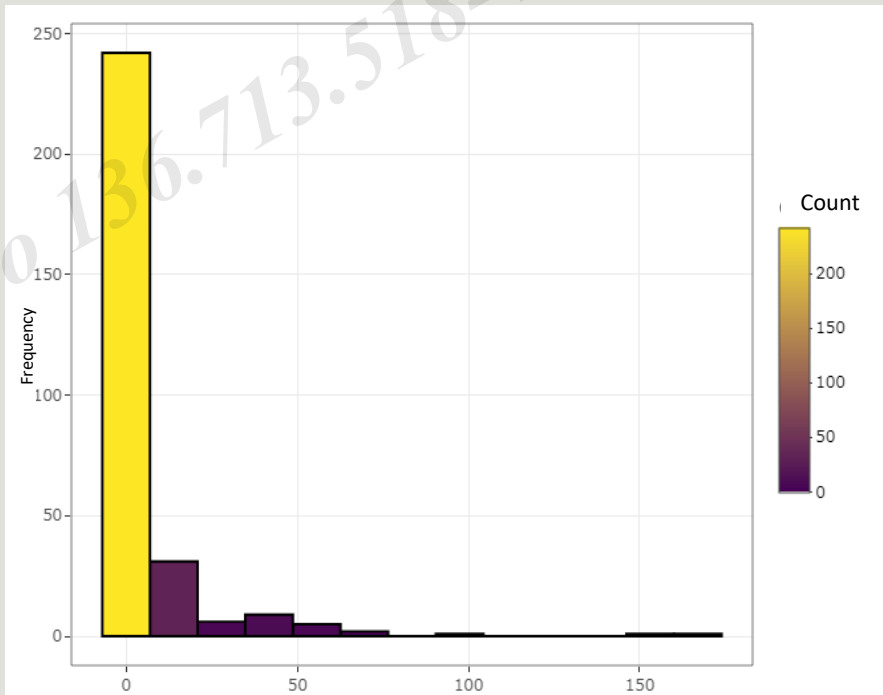
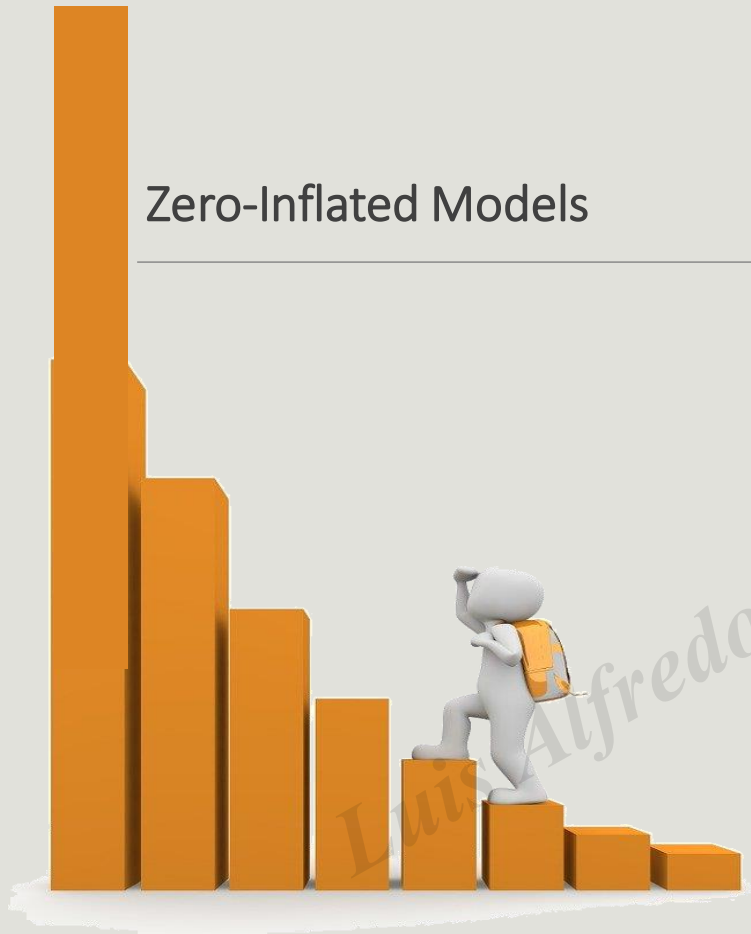
$$\ln(\hat{Y}_i) = \ln(\lambda_{bneg_i}) = \alpha + \beta_1.X_{1i} + \beta_2.X_{2i} + \dots + \beta_k.X_{ki}$$



A rustic wooden signpost with a weathered, brown surface is mounted on two wooden posts. The sign is shaped like a horizontal arrow pointing to the right. In the center of the sign, the words "EXCEL" and "R" are written in a large, white, sans-serif font, separated by a small gap. The background features a lush green field in the foreground, a dense forest of evergreen trees in the middle ground, and a mountain range with snow-capped peaks in the distance under a clear blue sky.

EXCEL R

## Zero-Inflated Models



## Model Choice

Verification	Modelo de Regressão para Dados de Contagem			
	Poisson	Negative Binomial	Zero-Inflated Poisson (ZIP)	Zero-Inflated Negative Binomial (ZINB)
Overdispersion in Data of the Dependent Variable	No	Yes	No	Yes
Excessive Number Of Zeros In The Dependent Variable	No	No	Yes	Yes



## Zero-Inflated Models

---

They are considered a combination of a model for count data and a model for binary data, since they are used to investigate the reasons that conduct to a certain number of occurrences (counts) of a phenomenon, as well as the reasons (or not) that conduct to the occurrence of this phenomenon, regardless of the amount of counts observed.

While a Poisson model inflated with zeros is estimated from the **combination of a Bernoulli distribution with a Poisson distribution**, a negative binomial model inflated with zeros is estimated by combining **a Bernoulli distribution with a Gamma-Poisson** distribution.

LAMBERT, D. Zero-inflated Poisson regression, with an application to defects in manufacturing. **Technometrics**, v. 34, n. 1, p. 1-14, 1992.

## Zero-Inflated Models

---

The definition of the existence of an excessive number of zeros in the dependent variable  $Y$  is elaborated through a specific test, known as **the Vuong** (1989) test, which will represent an important *output* to be analyzed in the estimation of regression models for count data, when there is a suspicion of existence of zero-inflation.

VUONG, Q. H. Likelihood ratio tests for model selection and non-nested hypotheses. **Econometrica**, v. 57, n. 2, p. 307-333, 1989.



## Zero-Inflated Models of the Poisson Type (ZIP)

---

Regarding specifically the **Poisson regression models inflated with zeros**, we can define that, while the **probability  $p$  of occurrence of no count** for a given observation  $i$  ( $i = 1, 2, \dots, n$ , in which  $n$  is the sample size), that is,  **$p(Y_i = 0)$** , is calculated taking into account the sum of a dichotomous component with a counting component, and therefore, the *plogit* probability of not occurring any count due to the dichotomous component, the **probability  $p$  of occurrence of a certain count  $m$**  ( $m = 1, 2, \dots$ ), i.e.,  **$p(Y_i = m)$** , follows the expression of the probability of the Poisson distribution, multiplied by  $(1 - p_{logit})$ .

## Zero-Inflated Models of the Poisson Type (ZIP)

$$\begin{cases} p(Y_i = 0) = p_{\text{logit}_i} + (1 - p_{\text{logit}_i}) \cdot e^{-\lambda_i} \\ p(Y_i = m) = (1 - p_{\text{logit}_i}) \cdot \frac{e^{-\lambda_i} \cdot \lambda_i^m}{m!}, \quad m = 1, 2, \dots \end{cases}$$

$$p_{\text{logit}_i} = \frac{1}{1 + e^{-(\gamma + \delta_1 \cdot W_{1i} + \delta_2 \cdot W_{2i} + \dots + \delta_q \cdot W_{qi})}}$$

$$\lambda_{\text{poisson}_i} = e^{(\alpha + \beta_1 \cdot X_{1i} + \beta_2 \cdot X_{2i} + \dots + \beta_k \cdot X_{ki})}$$

The zero-inflated Poisson regression models have two processes that generate zeros, one of them due to the binary distribution (in this case, the so-called structural zeros are generated) and the other one due to the Poisson distribution (in this situation, count data are generated, including the so-called sample zeros).





## Zero-Inflated Models of the Negative Binomial Type (ZINB)

---

Regarding the **zero-inflated negative binomial regression models**, we can define that, while the **probability  $p$  of occurrence of no count** for a given observation  $i$ , that is,  **$p(Y_i = 0)$** , is also calculated taking into account the sum of a dichotomous component with a count component, the **probability  $p$  of occurrence of a certain count  $m$**  ( $m = 1, 2, \dots$ ), that is,  **$p(Y_i = m)$** , follows now the expression of the probability of the Gamma-Poisson distribution.

## Zero-Inflated Models of the Negative Binomial Type (ZINB)

$$\begin{cases} p(Y_i = 0) = p_{logit_i} + (1 - p_{logit_i}) \cdot \left( \frac{1}{1 + \theta^{-1} \cdot \lambda_{bneg_i}} \right)^\theta \\ p(Y_i = m) = (1 - p_{logit_i}) \cdot \left[ \frac{\delta^\theta \cdot m_i^{\theta-1} \cdot e^{-m_i \cdot \delta}}{(\theta-1)!} \right], \quad m = 1, 2, \dots \end{cases}$$

$$p_{logit_i} = \frac{1}{1 + e^{-(\gamma + \delta_1 \cdot W_{1i} + \delta_2 \cdot W_{2i} + \dots + \delta_q \cdot W_{qi})}}$$

$$\lambda_{bneg_i} = e^{(\alpha + \beta_1 \cdot X_{1i} + \beta_2 \cdot X_{2i} + \dots + \beta_k \cdot X_{ki})}$$







THANK YOU VERY  
MUCH!

Prof. Dr. Luiz Pa 