



FEUP

Phd Program in Transportation

Transport Demand Modeling

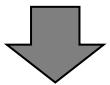
Filipe Moura

Generalized Linear Models

Why Generalized Linear Models?

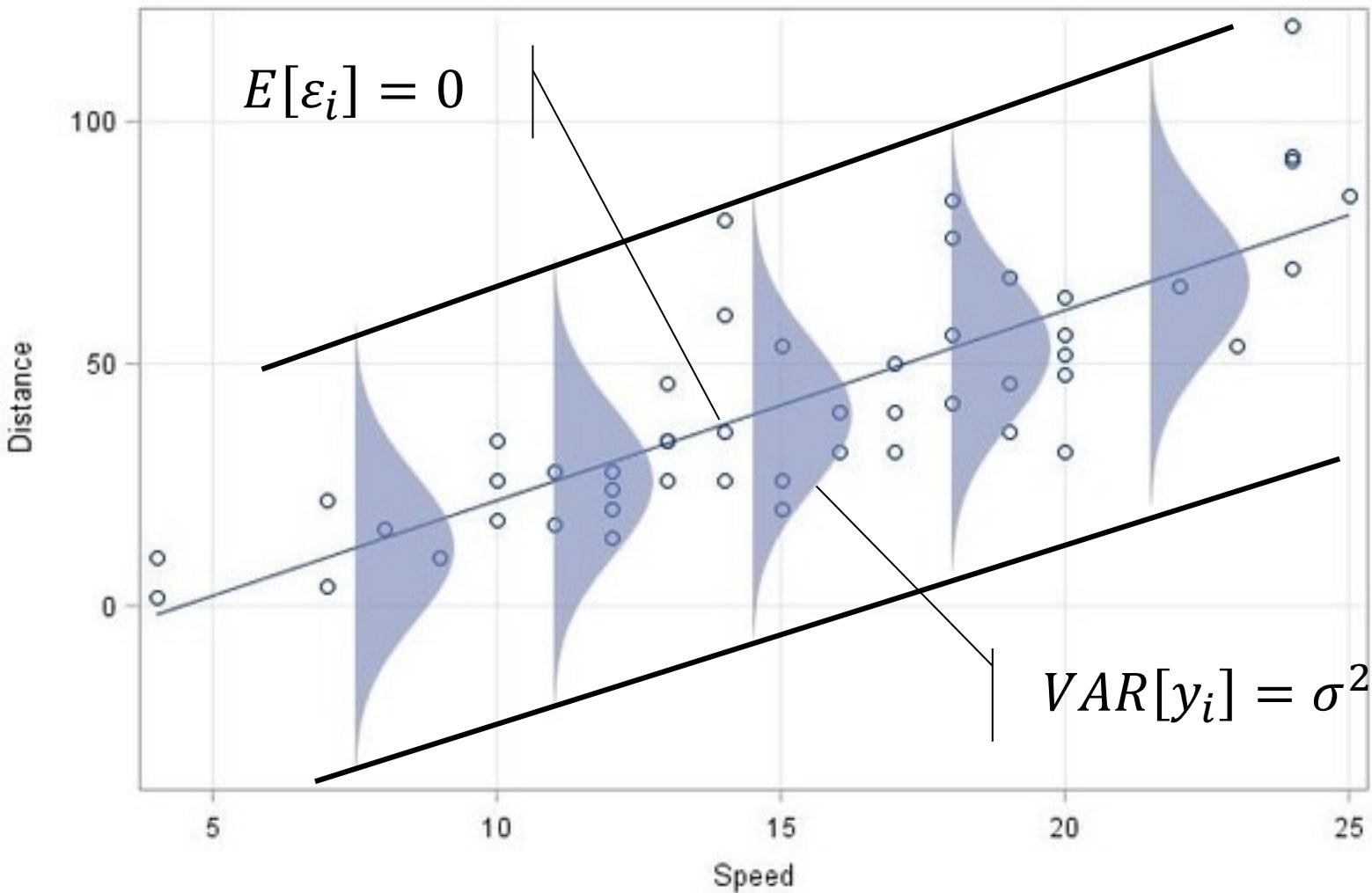
□ Why using GZLM?

- We shall see that these models extend the linear modelling framework to models where:
 - The dependent variable may not be continuous
 - The effect of independent variables may not be linear
 - The expected value of the errors terms might not be 0.



- GZLM unify all non linear models, used to explain the situation were the linear normal regression was not able to explain the relation under analysis
- GZLMs are most commonly used to model binary or count data, so we will focus on models for these types of data.

Errors are assumed to have constant variance and normally distributed.



When do GZLM come into play?

With MLR:

- $Y_i = BX + \varepsilon$, where X is a vector of predictors and B is a vector of coefficients β
- $E[Y_i] = \hat{B}X$ because $E[\varepsilon_i] = 0$ and $VAR[y_i] = \sigma^2$
- Etc.

When such conditions are not met,

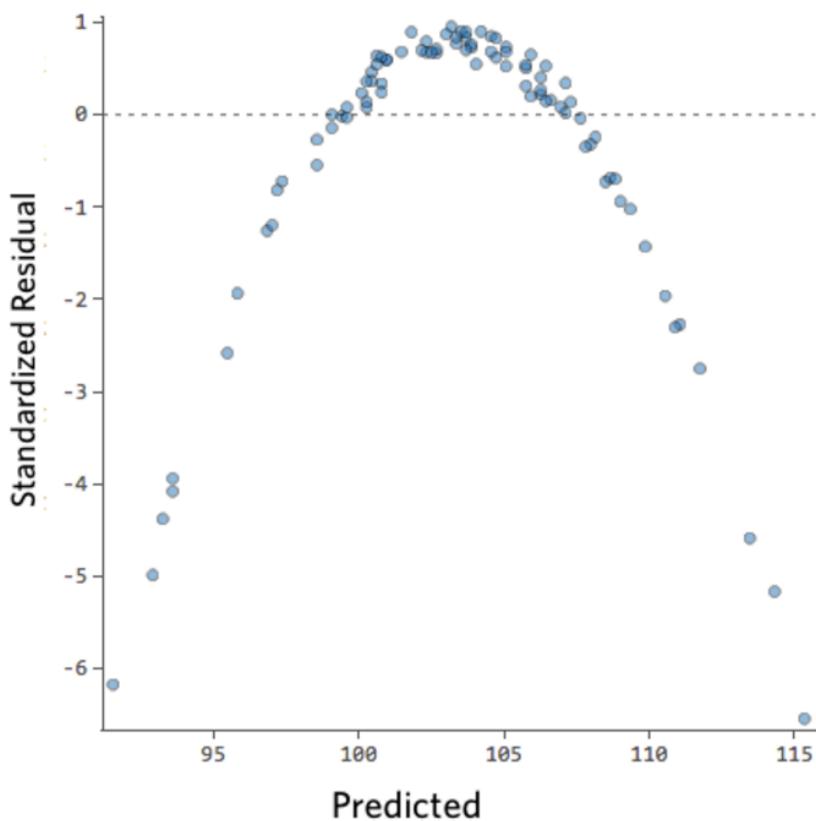
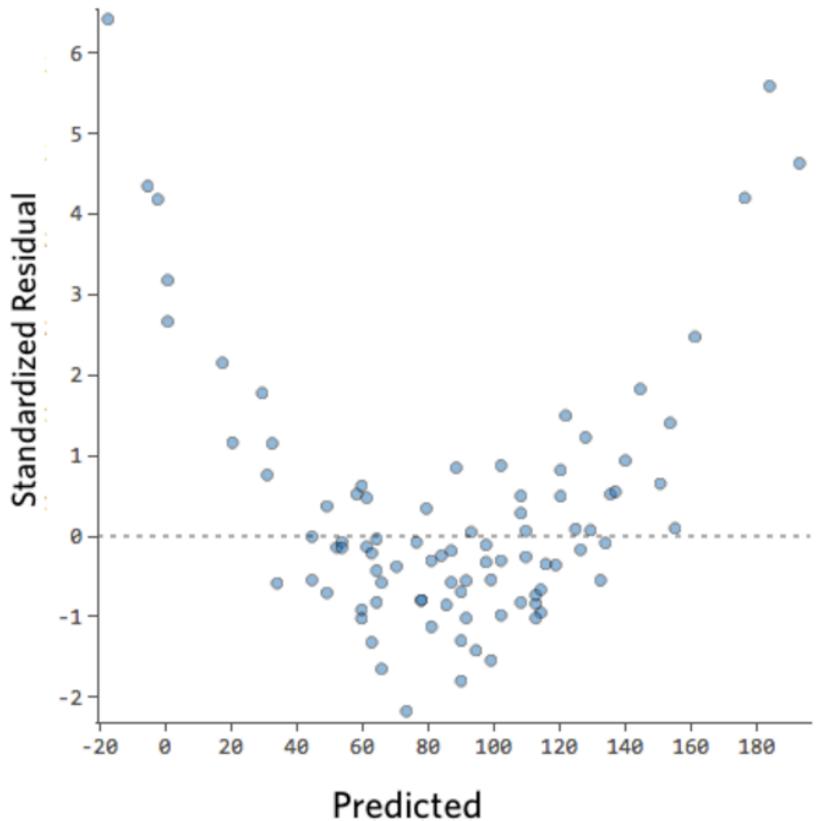
- When $E[\varepsilon_i] \neq 0$ or $VAR[Y_i] \neq \sigma^2$

you use GZLM where...

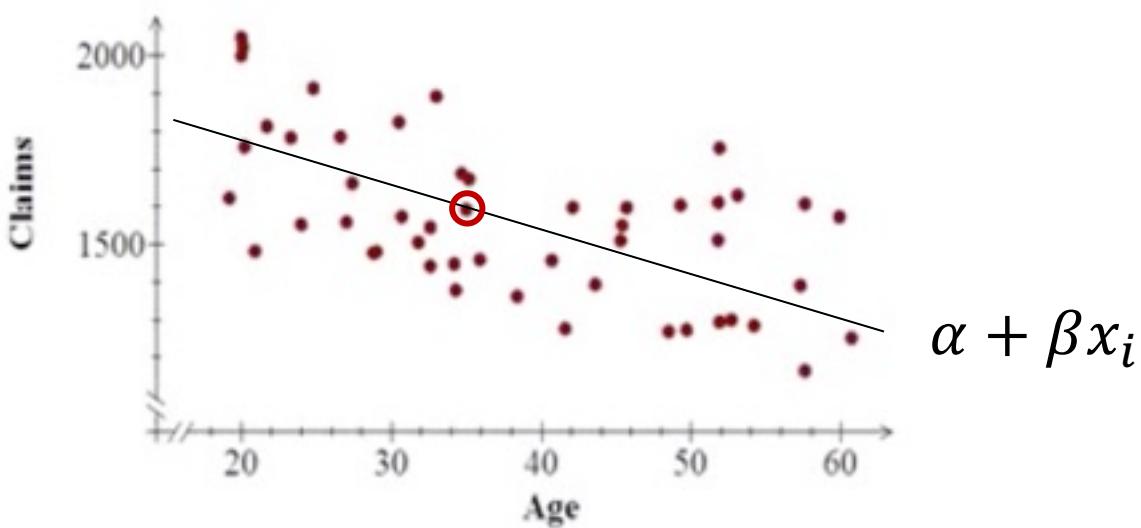
- The variation (probabilistic distribution) in the response variable Y_i can be explained in terms of the values of X
- We want to find some link function $g(.)$, that mediates the response variable (Y_i) and the regressors X_i , such that

$$E[g(Y_i)] = \hat{B}X, \text{ where } g(.) \text{ is the link function.}$$

Nonlinear



Recap on linear models and generalizing



- Distributions of the Y_i 's: $Y_i \sim N(\mu_i, \sigma^2)$
- Function of the explanatory variable, x_i 's: $\alpha + \beta x_i$
- Connection between explanatory variable and the distribution of Y_i :

$$\mu_i = E[Y_i] = \alpha + \beta x_i$$
- We will now generalize the distribution of the Y_i variables according to different distributions, besides the normal distribution

Components of the GZLM

□ Distribution of the Y_i 's

- Linear models: $Y_i \sim N(\mu_i, \sigma^2)$
- GZLM: $Y_i \sim \text{exponential family}$

□ Linear predictor = function of the covariates (explanatory variables)

- Linear models: $\eta_i = \alpha + \beta x_i$
- GZLM:
 - e.g. $\eta_i = \alpha + \beta x_i + \gamma z_i$
 - $\eta_i = \alpha + \beta x_i + \gamma x_i^2$

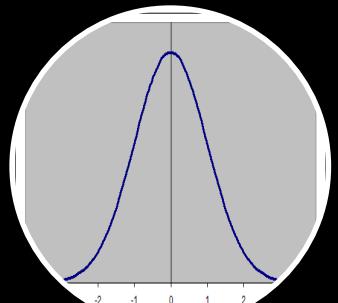
□ Link function = connection between the linear predictor

and $\mu_i = E[Y_i]$

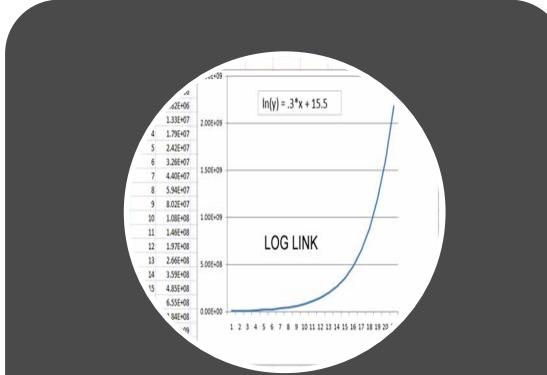
- Linear models: $\eta_i = \mu_i \Rightarrow \mu_i = \alpha + \beta x_i$
- GZLM:
 - e.g. $\eta_i = \ln(\mu_i) \Rightarrow \mu_i = e^{(\alpha + \beta x_i + \gamma x_i^2)}$

Structure of Generalized Linear Models

□ The overall structure of GLZM



Random
Component



Link
Function



Systematic
Component



Fundamental condition for using GZLM

- Response variable distribution must be a member of the **Exponential Family**
 - It corresponds to a y function that belongs to the exponential family with a single parameter θ and a probability distribution function (pdf) such as

$f(u, \theta) = s(u) \cdot t(\theta) \cdot \exp\{a(u) \cdot b(\theta)\}$, where s, t, a, b are known functions
or

$$f(u, \theta) = \exp\{a(u) \cdot b(\theta) + d(u) + c(\theta)\}$$

where $d(u) = \ln(s(u))$ and $c(\theta) = \ln(t(\theta))$

- When $a(u) = u$, the distribution is said to be in **canonical form**.
- $b(\theta)$ **is called the natural parameter** of the distribution function.
- For each function of the Exponential Family, one **parameter is of interest**.
The remaining are said **nuisance parameters**.

Members of the Exponential Family: Normal Distribution



- Normal distribution: $N(\mu, \sigma^2)$

$$f(u, \mu) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp^{-\frac{1}{2}\left(\frac{u-\mu}{\sigma}\right)^2}, \text{ with } -\infty \leq \mu \leq \infty$$

Random variable

Interest parameter

$$f(u, \mu) = \exp \left\{ u \cdot \frac{\mu}{\sigma^2} + \left[\frac{-\mu^2}{2\sigma^2} - \frac{1}{2} \ln(2\pi\sigma^2) \right] - \left(\frac{-u^2}{2\sigma^2} \right) \right\}$$

$a(u)$

$b(\mu)$

$c(\mu)$

$d(u)$

Natural parameter

Members of the Exponential Family: Binomial Distribution



- Binomial distribution: $\text{Bin}(n, p)$

$$f(u, p) = \binom{n}{u} \cdot p^u \cdot (1 - p)^{n-u}, n=0,1,2,\dots,n \text{ trials}$$

Random variableInterest parameter

$$= \binom{n}{u} \cdot \left(\frac{p}{1-p}\right)^u \cdot (1 - p)^n$$

$$f(u, \mu) = \exp \left\{ u \cdot \ln \left(\frac{p}{1-p} \right) + n \ln(1-p) - \ln \binom{n}{u} \right\}$$

$a(u)$ $b(p)$ $c(p)$ $d(u)$

↑
Natural parameter

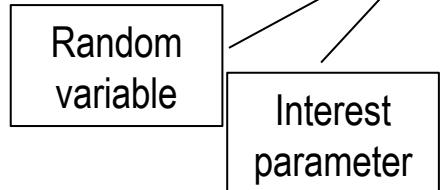
Members of the Exponential Family: Poisson Distribution



- Poisson distribution: $P(\lambda)$

$$f(u, \lambda) = \frac{e^{-\lambda} \cdot \lambda^u}{u!}, u=0,1,2,\dots, n \text{ observations}$$

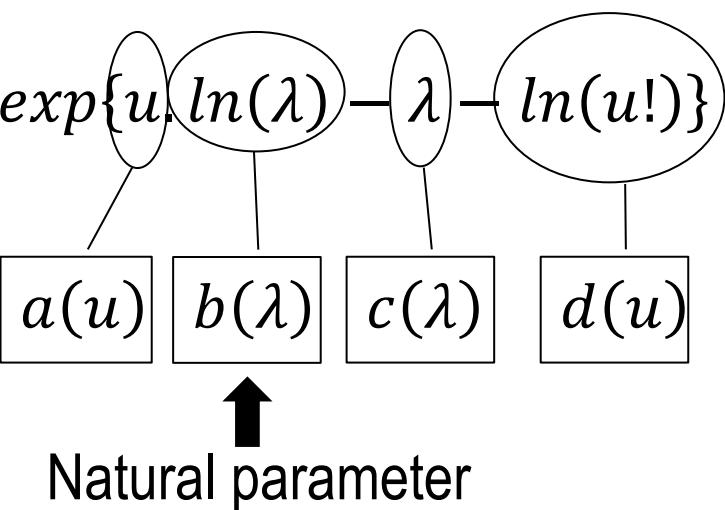
Random variable
Interest parameter



$$f(u, \mu) = \exp\{u \cdot \ln(\lambda) - \lambda - \ln(u!)\}$$

$a(u)$ $b(\lambda)$ $c(\lambda)$ $d(u)$

Natural parameter



Calibration of GZLM

- Suppose we have a set of independent observations, where
 - Y_i, X_i , for $i=1,2,\dots,n$ observations
and X_i is a vector of regressors = X_1, X_2, \dots, X_p
and Y_i is a response variable (dependent variable) we want to estimate and that belongs to some Exponential Family distribution where $a(Y)=Y$

- The joint pdf can be written as:

$$f(Y_1, Y_2, \dots, Y_n, \theta, \phi) = \prod_{i=1}^n \exp\{Y_i \cdot b(\theta_i) + c(\theta_i) + d(Y_i)\}$$

$$= \exp\left\{\sum_{i=1}^n Y_i \cdot \sum_{i=1}^n b(\theta_i) + \sum_{i=1}^n c(\theta_i) + \sum_{i=1}^n d(Y_i)\right\}$$

- The variation of Y_i can be explained in terms of the regressors X_i , based on the calibrated coefficients β
- NOT THE VALUES OF Y_i THEMSELVES!!!!

Calibration of GZLM

- For X_i is a vector of regressors = X_1, X_2, \dots, X_p
 - We hope to find a set of parameters $\beta=(\beta_1, \beta_2, \dots, \beta_p)$ that fits the regressor values to a link function $g(\cdot)$ that transforms the response variable (Y_i) such that, in the case of the normal distribution – $N(\mu, \sigma^2)$,

$$E[\mu_i] = \mu_i = g(\mu_i) = BX$$

Link function

- When the response variable Y_i follows a normal distribution, the GZLM is equal to the MLR

$$E[Y_i] = \mu_{y_i} = BX + E[\epsilon] = BX, \text{ because } E[\epsilon] = 0$$

Calibrating regressions

Maximum Likelihood Estimation (MLE) - Recap



- Obtain the likelihood:

$$L(\mu) = f(y_1) \cdot f(y_2) \cdots f(y_n)$$

- Log it – to make it easier to differentiate

$$\ln[L(\mu)]$$

- Differentiate and set the derivative equal to zero:

$$\frac{\delta}{\delta \mu} (\ln[L(\mu)]) = 0 \Rightarrow \hat{\mu} = \dots$$

- Check its maximum:

$$\frac{\delta^2}{\delta \mu^2} (\ln[L(\mu)]) < 0 \Rightarrow \text{max}$$

Calibrating GLZM

Maximum Likelihood Estimation (MLE)



- Obtain the likelihood:

$$L(\mu_1, \mu_2, \dots, \mu_n) = f(y_1) \cdot f(y_2) \cdots f(y_n)$$

- Log it – to make it easier to differentiate

$$\ln[L(\mu_1, \mu_2, \dots, \mu_n)]$$

- Use the link function to replace the $\mu_i's$:

$$\ln[L(\alpha, \beta, \gamma, \dots)]$$

- Differentiate and set the derivative equal to zero:

$$\frac{\delta}{\delta \alpha} (\ln[L(\alpha, \beta, \gamma, \dots)]) = 0 \Rightarrow \hat{\alpha} = \dots$$

$$\frac{\delta}{\delta \beta} (\ln[L(\alpha, \beta, \gamma, \dots)]) = 0 \Rightarrow \hat{\beta} = \dots$$

Link functions and Inverse functions

□ If the random component of Y_i follows a **normal distribution**

- The corresponding natural parameter is: $b(\mu_i) = \mu_i$
- The link function is: $g(\mu_i) = \mu_i$ (Identity link)
- Then: $g(\mu_i) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots = BX$
- And inversely: $\mu_i = BX$

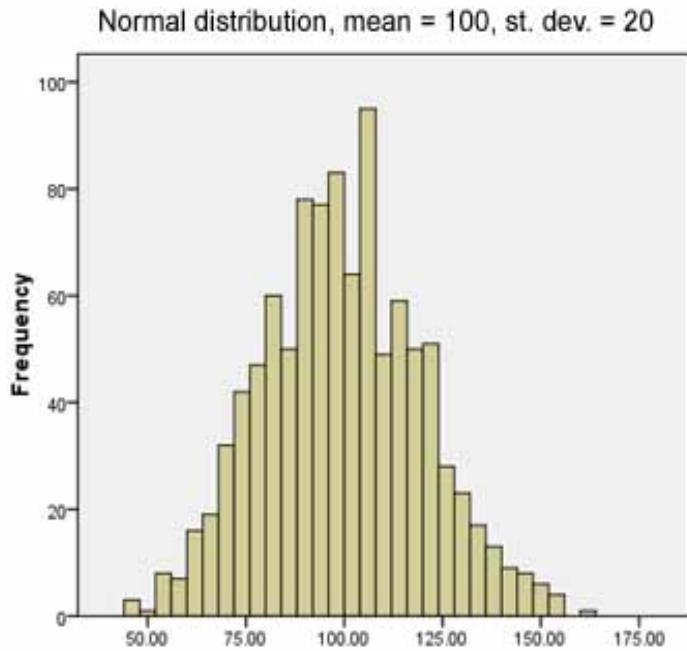
□ If the random component of Y_i follows a **binomial distribution**

- The corresponding natural parameter is: $b(\mu_i) = \ln\left(\frac{p_i}{1-p_i}\right)$
- The link function is: $g(\mu_i) = \ln\left(\frac{p_i}{1-p_i}\right)$ (logit link – or logistic)
- Then: $g(p_i) = \ln\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots = BX$
- And inversely: $E[Y_i] = p_i = \frac{\exp(BX)}{1+\exp(BX)}$

Generalized Linear Models

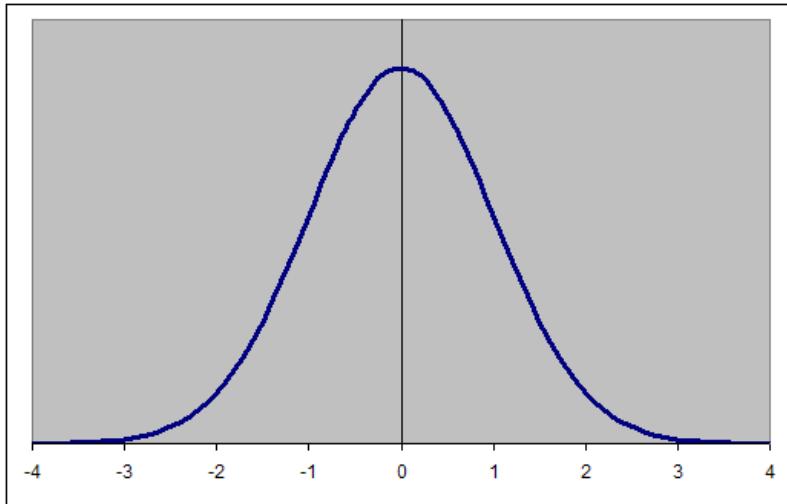
Normal distribution

- The distribution of the dependent variable has the form of the bell-shaped symmetrical curve centered in the mean.
- This implies the dependent variable is continuous.



$$f(x) = \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x-\mu)^2}{2\sigma^2}},$$

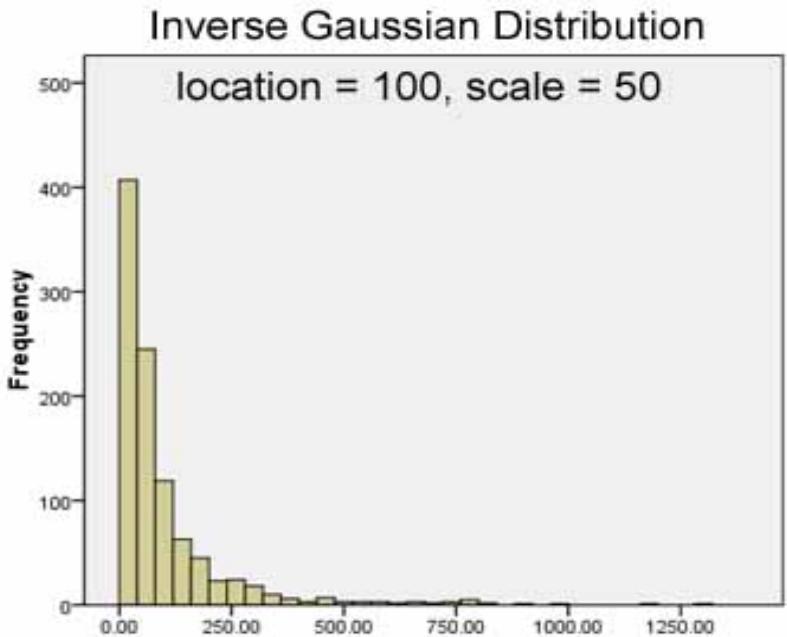
Normal distribution (standardized)



Generalized Linear Models

Inverse Gaussian (Wald Distribution)

- It is used for dependent variables that are **positively skewed** and have values always greater than 0.
- Values must be greater than 0 or are dropped. It has been used to model **diffusion processes, insurance claims**



If λ tends to infinity, the distribution becomes similar to a Normal distribution

$$f(x; \mu, \lambda) = \left[\frac{\lambda}{2\pi x^3} \right]^{1/2} \exp \frac{-\lambda(x - \mu)^2}{2\mu^2 x}$$

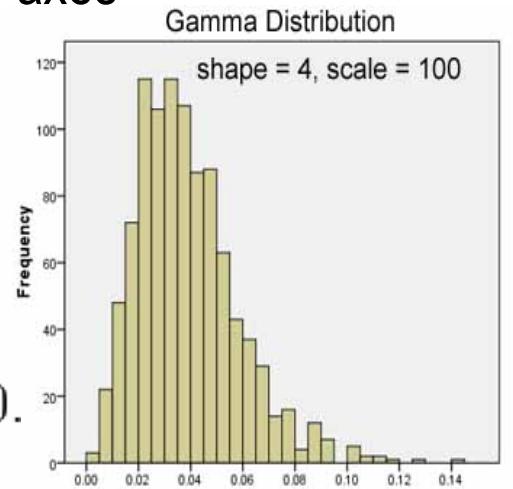
Generalized Linear Models

Gamma



- This is an **alternative for positively skewed** dependent variables. It is highly sensitive to the shape parameter.
 - When the shape parameter is greater than 1, the gamma distribution is bounded but positively skewed as shown in the figure below.
 - When the shape parameter is 1, the gamma distribution is exponentially declining.
 - When the shape parameter is less than 1, the gamma distribution is also exponentially declining and asymptotic to the axes
- The gamma distribution has been used in **survival analysis** and modeling **duration-of-event data**.

$$f(x; k, \theta) = x^{k-1} \frac{e^{-x/\theta}}{\theta^k \Gamma(k)} \text{ for } x \geq 0 \text{ and } k, \theta > 0.$$



Generalized Linear Models

Multinomial



- This distribution is used when the dependent variable has a finite number of categories, such as text string values, or is ordinal.
- The distribution among categories, not shown, is arbitrary.

$$f(x_1, \dots, x_k; n, p_1, \dots, p_k) = \Pr(X_1 = x_1 \text{ and } \dots \text{ and } X_k = x_k)$$

$$= \begin{cases} \frac{n!}{x_1! \cdots x_k!} p_1^{x_1} \cdots p_k^{x_k}, & \text{when } \sum_{i=1}^k x_i = n \\ 0 & \text{otherwise,} \end{cases}$$

Generalized Linear Models

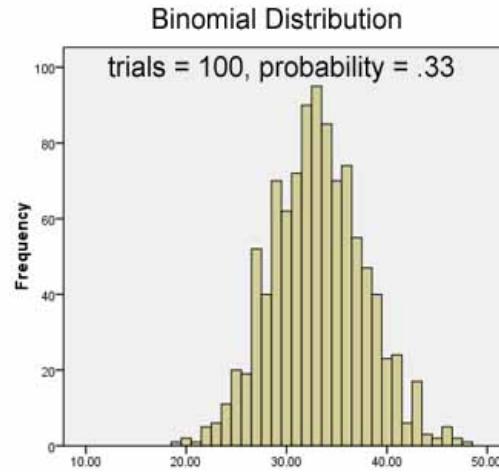
Binomial



- Used when the **dependent variable is binary**.
- The count of events in a fixed number of trials also has a binary distribution. Examples of binomial data are attributes “present/not present”, “innovation adopted/not adopted”, or “success/failure” data.
- It is assumed that the two values have a fixed rather than changing probability of occurrence (as in coin-flipping), even if that probability is not known

$$f(k; n, p) = \Pr(K = k) = \binom{n}{k} p^k (1 - p)^{n-k}$$

$$F(x; n, p) = \Pr(X \leq x) = \sum_{i=0}^{\lfloor x \rfloor} \binom{n}{i} p^i (1 - p)^{n-i}$$

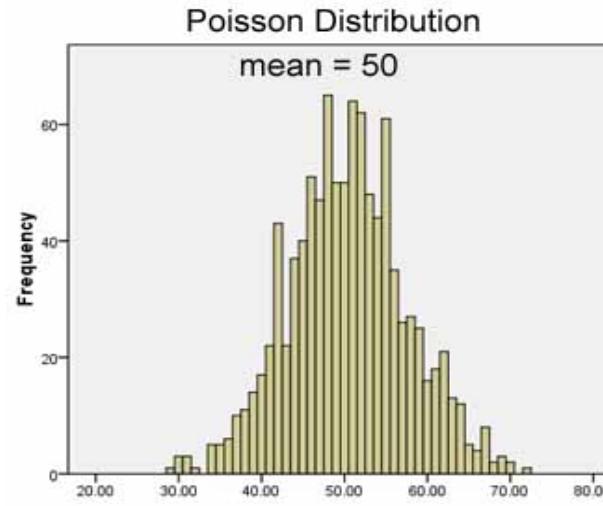


Generalized Linear Models

Poisson (COUNT DATA)

- The Poisson distribution is also used for count data and is preferred **when events are rare, as in modeling accidents, wars, or epidemics.**
- The binomial distribution is used when the dependent variable corresponds to **data counts of successes per given number of trials**
- **The Poisson distribution is used to count successes per given number of time units**

$$f(k; \lambda) = \frac{\lambda^k e^{-\lambda}}{k!},$$



Generalized Linear Models

Poisson (COUNT DATA)



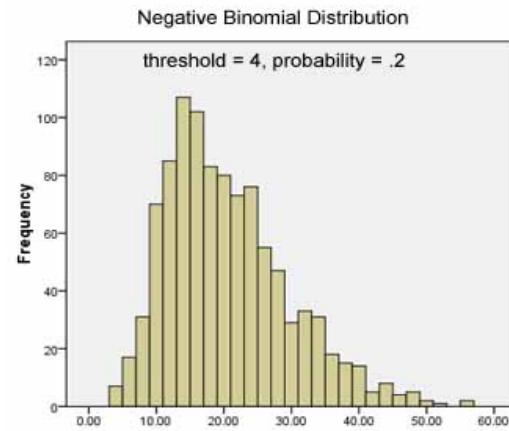
- A rule of thumb is to use a Poisson rather than binomial distribution when n is ≥ 100 , the probability of each event is below 0.05
- The Poisson distribution is also used when "events" can be counted but **non-occurrence of events cannot be counted**.
- In Poisson distributions, the mean equals the variance
 - Presence of homoscedasticity since variance doesn't change over data and λ is assumed constant
 - As such, there is no over-dispersion of data.
- All values are non-negative integers
 - **Thus count data, which cannot be negative, are better represented by Poisson than normal distributions**

Generalized Linear Models

Negative binomial

- It is similar to the Poisson distribution, also used for count data, but it is used when the variance is larger than the mean => over-dispersion of data.
- Typically this is characterized by "there being too many 0's."
 - As such, not all cases have an equal probability of experiencing the rare event, but instead, events may be clustered.
 - The negative binomial model is therefore sometimes called the "**over dispersed Poisson model**". Values must still be non-negative integers.
- The negative binomial is specified by an ancillary /dispersion parameter k (sometimes referred to as α or ψ).
 - When $k=0$, the negative binomial is equal to the Poisson distribution.

$$f(k) \equiv \Pr(X = k) = \binom{k + r - 1}{k} (1 - p)^r p^k \quad \text{for } k = 0, 1, 2, \dots$$



Generalized Linear Models

Systematic component



□ The linear predictor

- Quantity that incorporates the information about the independent variables into the model
- For a matrix of n observations and of p variables, the linear predictor η can be expressed as:

$$\eta_i = \sum_{j=1}^p x_{ij} \beta_j$$

□ Where

- ❖ each x_{ij} is the value of the j^{st} IV for the i^{st} observation
- ❖ β_j belong to a vector of unknown parameters to be estimated

Generalized Linear Models

Link function



- It provides the **relationship between the linear predictor and the mean of the distribution function**
 - The way the two previous components relate to each other
 - *In fact, the link function is a transformation of the response variable*
- It is a monotonous and differentiable function $g(\mu_i)$ that transforms μ_i in η_i where $g(\mu_i) = \eta_i$
- Inversely, μ_i can be obtained with (inverse function)

$$E[Y_i] = \mu_i = g^{-1}(\beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip}) + e_i$$

where μ_i is the expected value of Y_i
and X_{ij} are the predictors or explanatory variables

Generalized Linear Models

Link function



- It is used to **maintain a linear relationship** between the coefficients and predictors on the right hand side of the model equation and the dependent variable transformed by the link function on the left hand side of the equation

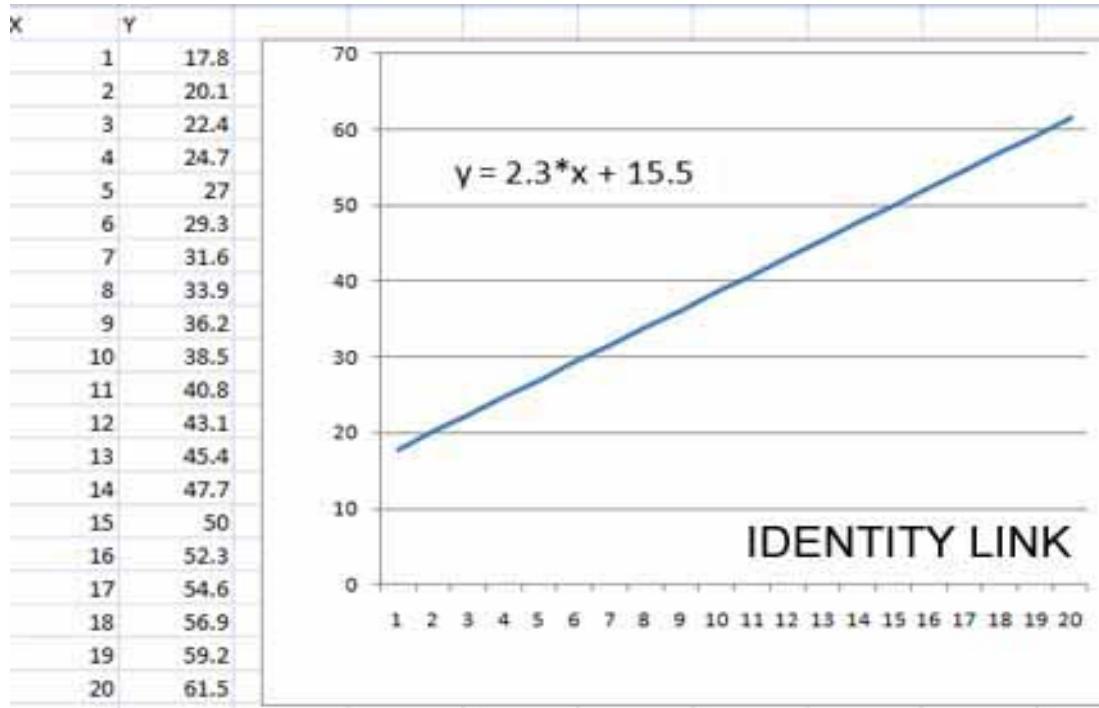
- **The choices of the link function depend on the natural parameter of the original distribution of the dependent variable Y_i**

Generalized Linear Models

Link function



- Normal distribution: Identity function



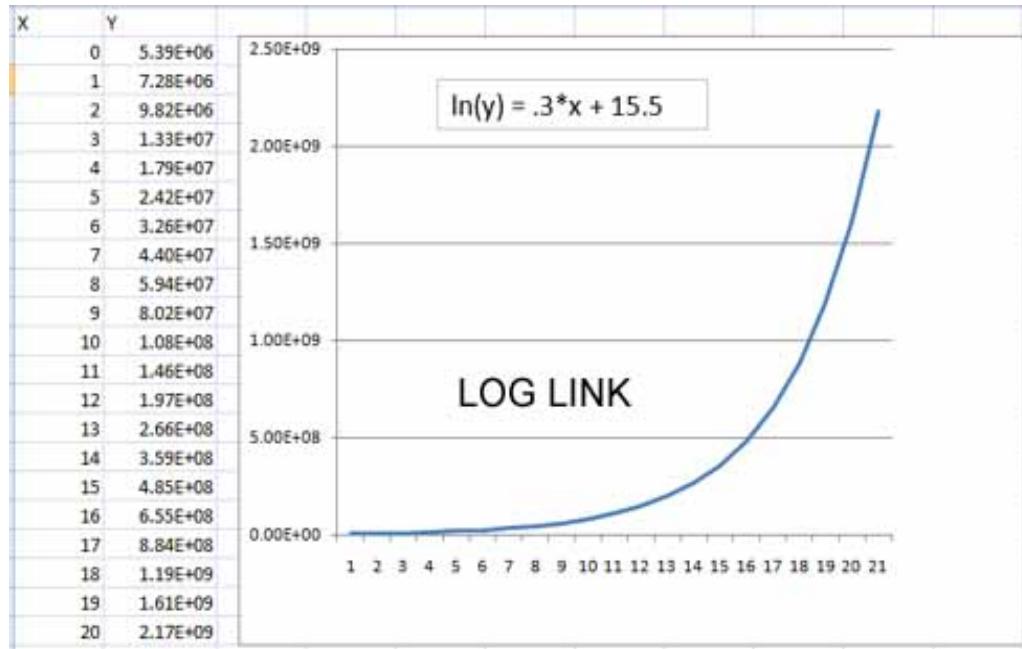
Generalized Linear Models

Link function



□ Poisson distribution: Log function

- **Loglinear models:** assume a Poisson distribution and use a log link function



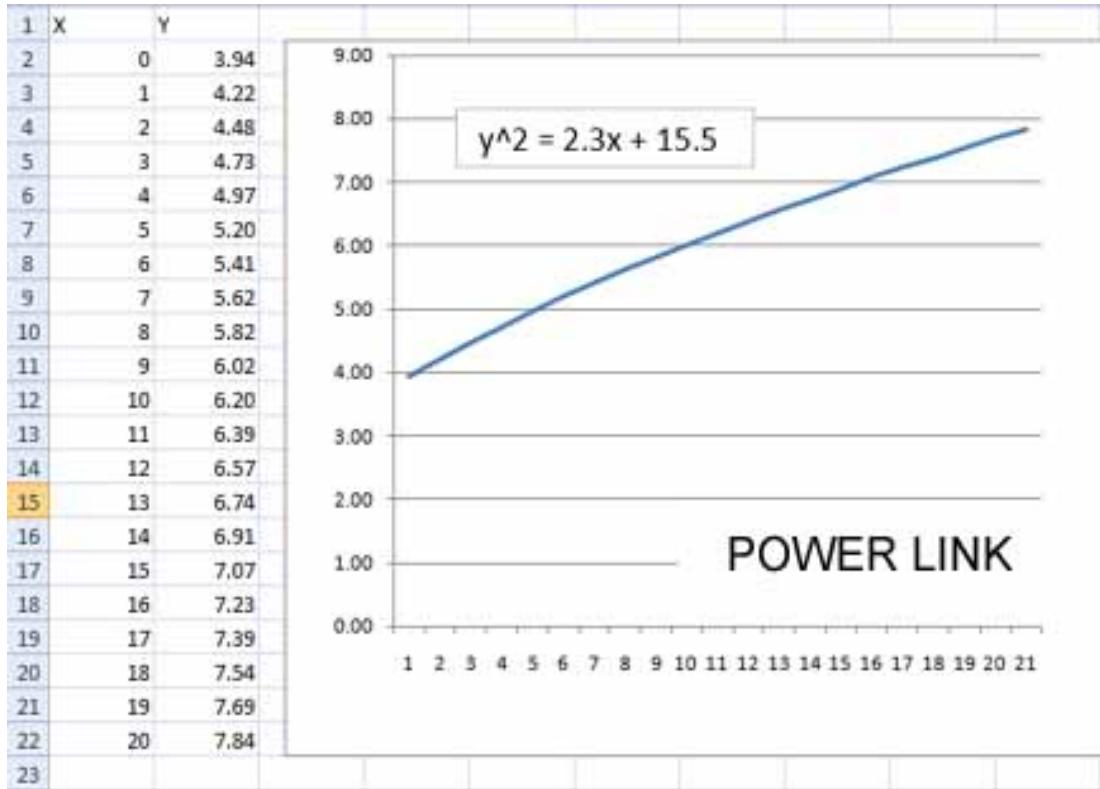
Generalized Linear Models

Link function



FEUP

- Many other distributions: Power functions



Generalized Linear Models

Link function



□ Common relations between distributions and link functions

Canonical Link Functions			
Distribution	Name	Link Function – $\eta_i = g(\mu_i)$	μ_i - Mean (Inverse) Function
Normal	Identity	$\mathbf{X}\beta = \mu$	$\mu = \mathbf{X}\beta$
Exponential	Inverse	$\mathbf{X}\beta = \mu^{-1}$	$\mu = (\mathbf{X}\beta)^{-1}$
Gamma			
Inverse Gaussian	Inverse squared	$\mathbf{X}\beta = \mu^{-2}$	$\mu = (\mathbf{X}\beta)^{-1/2}$
Poisson	Log	$\mathbf{X}\beta = \ln(\mu)$	$\mu = \exp(\mathbf{X}\beta)$
Binomial	Logit	$\mathbf{X}\beta = \ln\left(\frac{\mu}{1-\mu}\right)$	$\mu = \frac{\exp(\mathbf{X}\beta)}{1 + \exp(\mathbf{X}\beta)} = \frac{1}{1 + \exp(-\mathbf{X}\beta)}$
Multinomial			

GZLM for count of events

Poisson Regression Model

- The Poisson Distribution is commonly used to describe the count of events occurring at random in time or space
- The Poisson condition is that $E[Y_i] = VAR[Y_i] = \lambda_i$ or that
- Examples:
 - ❖ The number of cars passing through an intersection during a certain hour
 - ❖ The number of calls for emergency ambulance service during a tour of duty
 - ❖ The number of fires arising in a neighborhood
 - ❖ Number of vehicles waiting in a queue
 - ❖ Auto breakdowns in an express way in rush hour
 - ❖ Number of heart attack deaths per week in a county
 - ❖ Number of homes destroyed by a fire during the summer
 - ❖ Number of accidents in a road section or intersection



GZLM for count of events

Poisson Regression Model



- The most common relationship between the explanatory and the Poisson parameter is the log-linear model (because the logarithm of this function produces the linear combination of explanatory variables)

$$\lambda_i = e^{(\beta X_i)}$$

or

$$E[Y_i] = \lambda_i = \exp\left(\beta_0 + \sum_{j=1}^p x_{ij}\beta_j\right)$$

Inverse function

$$\ln(\lambda_i) = \beta X_i$$

- ❖ The expected number of accidents per period is given by

$$E[y_i] = \lambda_i = e^{(\beta X_i)}$$

GZLM for count of events

Poisson Regression Model



□ Poisson model

- For the case of count data (e.g., accidents), a variable Q_i is added and corresponds to the unit of exposure (e.g., vehicles per year)
- It is also referred to as the **offset value**

$$E[Y_i] = Q_i \times \exp \left(\beta_0 + \sum_{j=1}^p x_{ij} \beta_j \right)$$

GZLM for count of events

Poisson Regression Model - MLE

- Estimation by standard maximum likelihood methods, with the likelihood function given as

$$L(\beta) = \prod \frac{\text{EXP}[\text{EXP}(\beta X_i)][\text{EXP}(\beta X_i)]^{y_i}}{y_i!}$$

or

$$LL(\beta) = \sum_{i=1}^n [-\text{EXP}(\beta X_i) + y_i \beta X_i - \text{LN}(y_i!)]$$

- Maximum likelihood estimates produce Poisson parameters which are consistent, asymptotically normal, and asymptotically efficient

GZLM for count of events

Poisson Regression Model

- Elasticity is computed on the parameter estimation to evaluate the marginal effects of the independent variables
 - Effect of a 1% change in the variable on the expected frequency λ_i ,
 - Computed for each observation and then a single average is reported
 - **Continuous variables**

$$E_{xik}^{\lambda_i} = \frac{\delta\lambda_i}{\lambda_i} \times \frac{x_{ik}}{\delta x_{ik}} = \beta_k x_{ik}$$

- Count data

$$E_{xik}^{\lambda_i} = \frac{EXP(\beta_k) - 1}{EXP(\beta_k)}$$

GZLM for count of events

Negative Binomial Regression



- When the Poisson condition is violated (i.e., $E[Y_i] \neq VAR[Y_i]$), two situations can occur:
 - $E[Y_i] > VAR[Y_i]$ (dispersed)
 - $E[Y_i] < VAR[Y_i]$ (over dispersed)
- As such, the link function is rewritten
 - from $\lambda_i = EXP(\beta X_i)$ for each observation , with

$$\lambda_i = EXP(\beta X_i + \varepsilon_i)$$

➤ where ε_i is the dispersion term

GZLM for count of events

Negative Binomial Regression



- Therefore the variance differ from the mean through the addition of a quadratic term to the variance that represents over dispersion

$$\text{var}(Y_i) = \lambda_i + K(\lambda_i)^2$$

- The Poisson model is regarded as a limited model of the negative binomial as K approaches 0.
- This K parameter is called the **over dispersion** parameter.

GZLM for count of events

Negative Binomial Regression - MLE

□ Negative Binomial pdf

$$P(y_i) = \frac{\Gamma\left(y_i + \frac{1}{K}\right)}{y_i! \Gamma\left(\frac{1}{K}\right)} \left(\frac{K\lambda_i}{(1+K\lambda_i)}\right)^{y_i} \left(\frac{1}{1+K\lambda_i}\right)^{\frac{1}{K}}$$

➤ where Γ is a Gamma Function and K is an estimated parameter representative of dispersion

□ The corresponding likelihood function is:

$$L(\lambda_i) = \prod \frac{\Gamma\left(y_i + \frac{1}{K}\right)}{y_i! \Gamma\left(\frac{1}{K}\right)} \left(\frac{K\lambda_i}{(1+K\lambda_i)}\right)^{y_i} \left(\frac{1}{1+K\lambda_i}\right)^{\frac{1}{K}}$$

GZLM for count of events

Negative Binomial Regression



□ Negative Binomial with log link

- Specifies a negative binomial distribution (**with the ancillary K parameter = 1**) with a log link function
 - It is used to modeling count data that violates the Poisson assumption of equality of mean and variance.
 - Also, negative binomial regression is thought to be more stable than Poisson regression for **small datasets**.
 - An error term ξ of Gamma distribution and variance K_2 is added to the Poisson Regression

$$\log(\lambda_i) = \beta_0 + \sum_{j=1}^p x_{ij} \beta_j + \xi_i$$

GZLM for count of events

Negative Binomial Regression



- Tests on over dispersion - Lagrange Multiplier test (SPSS)
 - The Lagrange multiplier test may be used to **test if a negative binomial model is significantly different from a Poisson model**
 - Since the negative binomial model is the same as the Poisson model when the binomial model's ancillary (dispersion) parameter, $K=0$, **the Lagrange multiplier test analyses the null hypothesis that $K = 0$**
- A significant Lagrange test coefficient (i.e., **$p\text{-value} > 0,05$**) indicates that **K cannot be assumed to be different from 0**, and hence a Poisson model would be preferred over a negative binomial model (negative binomial models have one more parameter, k)

References

□ References

- **McCullagh, Peter; Nelder, John (1989). Generalized Linear Models, Second Edition. Boca Raton: Chapman and Hall/CRC. ISBN 0-412-31760-5.**
- J. B. S. Haldane, "On a Method of Estimating Frequencies", Biometrika, Vol. 33, No. 3 (Nov., 1945), pp. 222–225. JSTOR 2332299
- Hilbe, Joseph M., Negative Binomial Regression, Cambridge, UK: Cambridge University Press (2007) Negative Binomial Regression - Cambridge University Press
- **Washington, Simon P., Karlaftis, Mathew G. e Manning (2003) Statistical and Econometric Methods for Transportation Data Analysis, CRC**
- Lord, D., Washington, S. P., & Ivan, J. N. (2005). Poisson, Poisson-Gamma and zero-inflated regression models of motor vehicle crashes: balancing statistical fit and theory. Accident Analysis and Prevention , pp. 35-46.
- Fernandes, A. (2010) Programas de manutenção de características da superfície de pavimentos associados a critérios de segurança rodoviária. Tese de Doutoramento em Engenharia Civil. Instituto Superior Técnico, Universidade Técnica de Lisboa.