



---

## SHORT COURSE

# **Regression Models for Count Data: beyond Poisson model**

Wagner Hugo Bonat  
Walmes Marques Zeviani  
Eduardo Elias Ribeiro Jr

---



Laboratório de Estatística e Geoinformação  
Universidade Federal do Paraná  
<http://leg.ufpr.br>



# Regression Models for Count Data: beyond Poisson model

Wagner Hugo Bonat<sup>1 3</sup>

[www.leg.ufpr.br/~wagner](http://www.leg.ufpr.br/~wagner)

Walmes Marques Zeviani<sup>1 3</sup>

[www.leg.ufpr.br/~walmes](http://www.leg.ufpr.br/~walmes)

Eduardo Elias Ribeiro Jr<sup>2 3</sup>

[www.leg.ufpr.br/~eduardojr](http://www.leg.ufpr.br/~eduardojr)

<sup>1</sup>Departamento de Estatística (DEST)

Universidade Federal do Paraná (UFPR)

<sup>2</sup>Departamento de Ciências Exatas (LCE)

ESALQ - Universidade de São Paulo (ESALQ-USP)

<sup>3</sup>Laboratório de Estatística e Geoinformação (LEG)

<http://www.leg.ufpr.br>

Supplementary content: <http://www.leg.ufpr.br/rmcd>

Contact: [rmcd@leg.ufpr.br](mailto:rmcd@leg.ufpr.br)



# Contents

<b>Preface</b>	<b>3</b>
<b>1 Introduction</b>	<b>5</b>
<b>2 Count distributions: properties and regression models</b>	<b>9</b>
2.1 Poisson distribution . . . . .	9
2.2 Gamma-Count distribution . . . . .	10
2.3 Poisson-Tweedie distribution . . . . .	13
2.4 COM-Poisson distribution . . . . .	16
2.5 Comparing count distributions . . . . .	18
<b>3 The method of maximum likelihood</b>	<b>23</b>
<b>4 Models specified by second-moment assumptions</b>	<b>27</b>
4.1 Extended Poisson-Tweedie model . . . . .	27
4.2 Estimation and Inference . . . . .	29
<b>5 Applications</b>	<b>33</b>
5.1 Cotton Bolls . . . . .	33
5.2 Soybean pod and beans . . . . .	43
5.3 Number of vehicle claims . . . . .	60



# Preface

The main goal of this material is to provide a technical support for the students attending the course "Regression models for count data: beyond the Poisson model", given as part of the XV Brazilian School of Regression models - March/2017 in Goiânia, Goiás, Brazil.

The main goal of this course is to present a wide range of statistical models to deal with count data. We focus on parametric and second-moment specified models. We shall present the model specification along with strategies for model fitting and associated R(R Core Team, 2015) code. Furthermore, this book-course and supplementary materials, such as R code and data sets are available for the students on the web page <http://cursos.leg.ufpr.br/rmcd>.

We intend to keep the course in a level suitable for bachelor students who already attended a course on generalized linear models (Nelder and Wedderburn, 1972). However, since the course also covers updated topics, it can be of interest of postgraduate students and researches in general.

We designed the course for three hours of tuition. In the first part of the course, we shall present the analysis of count data based on fully parametric models. After a brief introduction and motivation on count data, we present the Poisson, Gamma-Count, Poisson-Tweedie and COM-Poisson distributions. We explore their properties through a consideration of dispersion, zero-inflated and heavy tail indexes. Furthermore, the estimation and inference for these models based on the likelihood paradigm is discussed along with the associated R code and worked examples.

In the second part of the course, we provide a brief introduction to the estimating function approach (Jørgensen and Knudsen, 2004, ; Bonat and Jørgensen, 2016) and discuss models based on second-moment assumptions in the style of Wedderburn (1974). In particular, we focus on the recently proposed Extended Poisson-Tweedie model (Bonat et al., 2016) and its special case the quasi-Poisson model. The estimating function approach adopted for

estimation and inference is presented along with R code and data examples. The use of the R package `mcglm` (Bonat, 2016) is discussed for fitting the extended Poisson-Tweedie model.

We acknowledge our gratitude to the scientific committee of XV Brazilian regression model school for this opportunity.

Department of Statistics, Paraná Federal University, Curitiba, PR, Brazil.

March 27, 2017.



# Chapter 1

## Introduction

The analysis of count data has received attention from the statistical community in the last four decades. Since the seminal paper published by Nelder and Wedderburn (1972), the class of generalized linear models (GLMs) have a prominent role for regression modelling of normal and non-normal data. GLMs are fitted by a simple and efficient Newton scoring algorithm relying only on second-moment assumptions for estimation and inference. Furthermore, the theoretical background for GLMs is well established in the class of dispersion models (Jørgensen, 1987, 1997) as a generalization of the exponential family of distributions.

In spite of the flexibility of the GLM class, the Poisson distribution is the only choice for the analysis of count data in this framework. Thus, in practice there is probably an over-emphasis on the use of the Poisson distribution for count data. A well known limitation of the Poisson distribution is its mean and variance relationship, which implies that the variance equals the mean, referred to as equidispersion. In practice, however, count data can present other features, namely underdispersion (mean  $>$  variance) and overdispersion (mean  $<$  variance). There are many different possible causes for departures from the equidispersion. Furthermore, in practical data analysis a number of these could be involved.

One possible cause of under/overdispersion is departure from the Poisson process. It is well known that the Poisson counts can be interpreted as the number of events in a given time interval where the arrival's times are exponential distributed. In the cases where this assumption is violated the resulting counts can be under or overdispersed (Zeviani et al., 2014). Another possibility and probably more frequent cause of overdispersion is unobserved heterogeneity of experimental units. It can be due, for exam-

ple, to correlation between individual responses, cluster sampling, omitted covariates and others.

In general, these departures from the Poisson distribution are manifested in the raw data as a zero-inflated or heavy-tailed count distribution. It is important to discuss the consequences of failing to take into account the under or overdispersion when analysing count data. In the case of overdispersion, the standard errors associated with the regression coefficients calculated under the Poisson assumption are too optimistic and associated hypothesis tests will tend to give false positive results by incorrectly rejecting null hypotheses. The opposite situation will appear in case of underdispersed data. In both cases, the Poisson model provides unreliable standard errors for the regression coefficients and hence potentially misleading inferences. However, the regression coefficients are still consistently estimated.

The strategies for constructing alternative count distributions are related with the causes of the non-equidispersion. When departures from the Poisson process are plausible the class of duration dependence models (Winkelmann, 2003) can be employed. This class of models changes the distribution of the time between events from the exponential to more general distributions, such as gamma and inverse Gaussian. In this course, we shall discuss one example of this approach, namely, the Gamma-Count distribution (Zeviani et al., 2014). This distribution assumes that the time between events is gamma distributed, thus it can deal with under, equi and overdispersed count data.

On the other hand, if unobserved heterogeneity is present its in general implies extra variability and consequently overdispersed count data. In this case, a Poisson mixtures is commonly applied. This approach consists of include random effects on the observation level, and thus take into account the unobserved heterogeneity. Probably, the most popular example of this approach is the negative binomial model, that corresponds to a Poisson-gamma mixtures. In this course, we shall present the Poisson-Tweedie family of distributions, which in turn corresponds to Poisson-Tweedie mixtures (Bonat et al., 2016; Jørgensen and Kokonendji, 2015). Finally, a third approach to deal with non-equidispersed count data consists of generalize the Poisson distribution by adding an extra parameter to model under and overdispersion. Such a generalization can be done using the class of weighted Poisson distributions (Del Castillo and Pérez-Casany, 1998). One popular example of this approach is the Conway–Maxwell–Poisson distribution (COM-Poisson) (Sellers and Shmueli, 2010). The COM-Poisson is a member of the exponential family, has the Poisson and geometric distributions as special cases and the Bernoulli distribution as a limiting case. It can deal with both under and overdispersed count data. Thus, given the nice properties of the COM-

Poisson distribution for handling count data, we chose to present this model as part of this course.

In this course, we shall highlight and compare the flexibility of these distributions to deal with count data through a consideration of dispersion, zero-inflated and heavy tail indexes. Furthermore, we specify regression models and illustrate their application with three worked examples.

In Chapter 2 we present the properties and regression models associated with the Poisson, Gamma-count, Poisson-Tweedie and COM-Poisson distributions. Furthermore, we compare these distributions using the dispersion, zero-inflated and heavy tail indexes. Estimation and inference for these models based on the likelihood paradigm are discussed in Chapter 3. In Chapter 4, we extend the Poisson-Tweedie model using only second-moment assumptions and present the fitting algorithm based on the estimating functions approach. Chapter 5 presents three worked examples.



## Chapter 2

# Count distributions: properties and regression models

In this chapter, we present the probability mass function and discuss the main properties of the Poisson, Gamma-Count, Poisson-Tweedie and COM-Poisson distributions.

### 2.1 Poisson distribution

The Poisson distribution is a notorious discrete distribution. It has a dual interpretation as a natural exponential family and as an exponential dispersion model. The Poisson distribution denoted by  $P(\mu)$  has probability mass function

$$\begin{aligned} p(y; \mu) &= \frac{\mu^y}{y!} \exp\{-\mu\} \\ &= \frac{1}{y!} \exp\{\phi y - \exp\{\phi\}\}, \quad y \in \mathbb{N}_0, \end{aligned} \tag{2.1}$$

where  $\phi = \log\{\mu\} \in \mathbb{R}$ . Hence the Poisson is a natural exponential family with cumulant generator  $\kappa(\phi) = \exp\{\phi\}$ . We have  $E(Y) = \kappa'(\phi) = \exp\{\phi\} =$

$\mu$  and  $\text{var}(Y) = \kappa''(\phi) = \exp\{\phi\} = \mu$ . The probability mass function (2.1) can be evaluated in R through the `dpois()` function.

In order to specify a regression model based on the Poisson distribution, we consider a cross-section dataset,  $(y_i, x_i)$ ,  $i = 1, \dots, n$ , where  $y_i$ 's are iid realizations of  $Y_i$  according to a Poisson distribution. The Poisson regression models is defined by

$$Y_i \sim P(\mu_i), \quad \text{with} \quad \mu_i = g^{-1}(x_i^\top \beta).$$

In this notation,  $x_i$  and  $\beta$  are  $(p \times 1)$  vectors of known covariates and unknown regression parameters, respectively. Moreover,  $g$  is a standard link function, for which we adopt the logarithm link function, but potentially any other suitable link function could be adopted.

## 2.2 Gamma-Count distribution

The Poisson distribution as presented in (2.1) follows directly from the natural exponential family and thus fits in the generalized linear models (GLMs) framework. Alternatively, the Poisson distribution can be derived by assuming independent and exponentially distributed times between events (Zeviani et al., 2014). This derivation allows for a flexible framework to specify more general models to deal with under and overdispersed count data.

As point out by Winkelmann (2003) the distributions of the arrival times determine the distribution of the number of events. Following Winkelmann (1995), let  $\tau_k$ ,  $k \in \mathbb{N}$  denote a sequence of waiting times between the  $(k-1)$ th and the  $k$ th events. Then, the arrival time of the  $y$ th event is given by  $v_y = \sum_{k=1}^y \tau_k$ , for  $y = 1, 2, \dots$ . Furthermore, denote  $Y$  the total number of events in the open interval between 0 and  $T$ . For fixed  $T$ ,  $Y$  is a count variable. Indeed, from the definitions of  $Y$  and  $v_y$  we have that  $Y < y$  iff  $v_y \geq T$ , which in turn implies  $P(Y < y) = P(v_y \geq T) = 1 - F_y(T)$ , where  $F_y(T)$  denotes the cumulative distribution function of  $v_y$ . Furthermore,

$$\begin{aligned} P(Y = y) &= P(Y < y+1) - P(Y < y) \\ &= F_y(T) - F_{y+1}(T). \end{aligned} \tag{2.2}$$

Equation (2.2) provides the fundamental relation between the distribution of arrival times and the distribution of counts. Moreover, this type of specification allows to derive a rich class of models for count data by choosing

a distribution for the arrival times. In this material, we shall explore the Gamma-Count distribution which is obtained by specifying the arrival times distribution as gamma distributed.

Let  $\tau_k$  be identically and independently gamma distributed, with density distribution (dropping the index  $k$ ) given by

$$f(\tau; \alpha, \gamma) = \frac{\gamma^\alpha}{\Gamma(\alpha)} \tau^{\alpha-1} \exp\{-\gamma\tau\}, \quad \alpha, \gamma \in \mathbb{R}^+. \quad (2.3)$$

In this parametrization  $E(\tau) = \alpha/\gamma$  and  $\text{var}(\tau) = \alpha/\gamma^2$ . Thus, by applying the convolution formula for gamma distributions, it is easy to show that the distribution of  $v_y$  is given by

$$f_y(v; \alpha, \gamma) = \frac{\gamma^{y\alpha}}{\Gamma(y\alpha)} v^{y\alpha-1} \exp\{-\gamma v\}. \quad (2.4)$$

To derive the new count distribution, we have to evaluate the cumulative distribution function, which after the change of variable  $u = \gamma\alpha$  can be written as

$$F_y(T) = \frac{1}{\Gamma(y\alpha)} \int_0^{\gamma T} u^{y\alpha-1} \exp\{-u\} du, \quad (2.5)$$

where the integral is the incomplete gamma function. We denote the right side of (2.5) as  $G(\alpha y, \gamma T)$ . Thus, the number of event occurrences during the time interval  $(0, T)$  has the two-parameter distribution function

$$P(Y = y) = G(\alpha y, \gamma T) - G(\alpha(y+1), \gamma T), \quad (2.6)$$

for  $y = 0, 1, \dots$ , where  $\alpha, \gamma \in \mathbb{R}^+$ . Winkelmann (1995) showed that for integer  $\alpha$  the probability mass function defined in (2.6) is given by

$$P(Y = y) = \exp\{-\gamma T\} \sum_{i=0}^{\alpha-1} \frac{(\gamma T)^{\alpha y + i}}{\alpha y + i}!. \quad (2.7)$$

For  $\alpha = 1$ ,  $f(\tau)$  is the exponential distribution and (2.6) clearly simplifies to the Poisson distribution. The following R function can be used to evaluate the probability mass function of the Gamma-Count distribution.

```

dgc <- function(y, gamma, alpha, log = FALSE) {
  p <- pgamma(q = 1,
             shape = y * alpha,
             rate = alpha * gamma) -
  pgamma(q = 1,
        shape = (y + 1) * alpha,
        rate = alpha * gamma)
  if(log == TRUE) {p <- log(p)}
  return(p)
}

```

Although, numerical evaluation of (2.6) can easily be done, the moments (mean and variance) cannot be obtained in closed form. Winkelmann (1995) showed for a random variable  $Y \sim GC(\alpha, \gamma)$ , where  $GC(\alpha, \gamma)$  denotes a Gamma-Count distribution with parameters  $\alpha$  and  $\gamma$ ,  $E(Y) = \sum_{i=1}^{\infty} G(\alpha i, \gamma T)$ .

Furthermore, for increasing  $T$  it holds that

$$Y(T) \overset{a}{\sim} N\left(\frac{\gamma T}{\alpha}, \frac{\gamma T}{\alpha^2}\right), \quad (2.8)$$

thus the limiting variance-mean ratio equals a constant  $1/\alpha$ . Consequently, the Gamma-Count distribution displays overdispersion for  $0 < \alpha < 1$  and underdispersion for  $\alpha > 1$ . Figure 2.1 presents the probability mass function for some Gamma-Count distributions. We fixed the parameter  $\gamma = 10$  and fit the parameter  $\alpha$  in order to have dispersion index ( $DI = \text{var}(Y)/E(Y)$ ) equalling to 0.5, 2, 5 and 20.

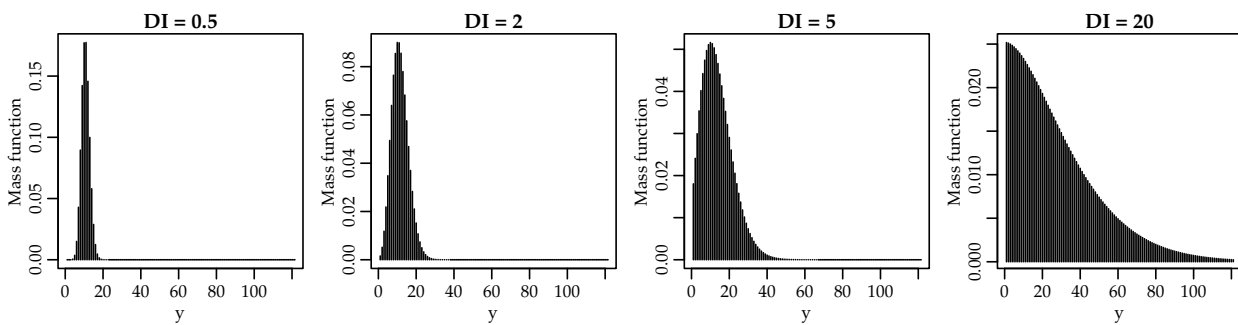


Figure 2.1: Gamma-Count probability mass function by values of the dispersion index (DI).

The Gamma-Count regression model assumes that the period at risk ( $T$ ) is identical for all observations, thus  $T$  may be set to unity without loss of generality. In the Gamma-count regression model, the parameters depend on a vector of individual covariates  $\mathbf{x}_i$ . Thus, the Gamma-Count regression model is defined by



$$E(\tau_i | \mathbf{x}_i) = \frac{\alpha}{\gamma} = g^{-1}(-\mathbf{x}_i^\top \boldsymbol{\beta}). \quad (2.9)$$

Consequently, the regression model is for the waiting times and not directly for the counts. Note that,  $E(N_i | \mathbf{x}_i) = E(\tau_i | \mathbf{x}_i)^{-1}$  iff  $\alpha = 1$ . Thus,  $\hat{\boldsymbol{\beta}}$  should be interpreted accordingly.  $-\boldsymbol{\beta}$  measures the percentage change in the expected waiting time caused by a unit increase in  $\mathbf{x}_i$ . The model parameters can be estimated using the maximum likelihood method as we shall discuss in Chapter 3.

## 2.3 Poisson-Tweedie distribution

The Poisson-Tweedie distribution (Bonat et al., 2016; Jørgensen and Kokonendji, 2015; El-Shaarawi et al., 2011) consists of include Tweedie distributed random effects on the observation level of Poisson random variables, and thus to take into account unobserved heterogeneity. The Poisson-Tweedie family is given by the following hierarchical specification

$$\begin{aligned} Y|Z &\sim \text{Poisson}(Z) \\ Z &\sim \text{Tw}_p(\mu, \phi), \end{aligned} \quad (2.10)$$

where  $\text{Tw}_p(\mu, \phi)$  denotes a Tweedie distribution (Jørgensen, 1987, ; Jørgensen, 1997) with probability function given by

$$f_Z(z; \mu, \phi, p) = a(z, \phi, p) \exp\{(z\psi - k_p(\psi))/\phi\}. \quad (2.11)$$

In this notation,  $\mu = k'_p(\psi)$  is the expectation,  $\phi > 0$  is the dispersion parameter,  $\psi$  is the canonical parameter and  $k_p(\psi)$  is the cumulant function. Furthermore,  $\text{var}(Z) = \phi V(\mu)$  where  $V(\mu) = k''_p(\psi)$  is the variance function. Tweedie densities are characterized by power variance functions of the form  $V(\mu) = \mu^p$ , where  $p \in (-\infty, 0] \cup [1, \infty)$  is an index determining the distribution. The support of the distribution depends on the value of the power parameter. For  $p \geq 2$ ,  $1 < p < 2$  and  $p = 0$  the support corresponds to the positive, non-negative and real values, respectively. In these cases  $\mu \in \Omega$ , where  $\Omega$  is the convex support (i.e. the interior of the closed convex hull of the corresponding distribution support). Finally, for  $p < 0$  the support corresponds to the real values, however the expectation  $\mu$  is positive. Here, we required  $p \geq 1$ , to make  $\text{Tw}_p(\mu, \phi)$  non-negative.

The function  $a(z, \phi, p)$  cannot be written in a closed form apart of the special cases corresponding to the Gaussian ( $p = 0$ ), Poisson ( $\phi = 1$  and  $p = 1$ ), non-central gamma ( $p = 3/2$ ), gamma ( $p = 2$ ) and inverse Gaussian ( $p = 3$ ) distributions (Jørgensen, 1997). The compound Poisson distribution is obtained when  $1 < p < 2$ . This distribution is suitable to deal with non-negative data with probability mass at zero and highly right-skewed (Andersen and Bonat, 2016).

The Poisson-Tweedie is an overdispersed factorial dispersion model (Jørgensen and Kokonendji, 2015) and its probability mass function for  $p > 1$  is given by

$$f(y; \mu, \phi, p) = \int_0^\infty \frac{z^y \exp -z}{y!} a(z, \phi, p) \exp\{(z\psi - k_p(\psi))/\phi\} dz. \quad (2.12)$$

The integral (2.12) has no closed-form apart of the special case corresponding to the negative binomial distribution, obtained when  $p = 2$ , i.e. a Poisson gamma mixture. In the case of  $p = 1$ , the integral (2.12) is replaced by a sum and we have the Neyman Type A distribution. Further special cases include the compound Poisson ( $1 < p < 2$ ), factorial discrete positive stable ( $p > 2$ ) and Poisson-inverse Gaussian ( $p = 3$ ) distributions (Jørgensen and Kokonendji, 2015, ; Kokonendji et al., 2004).

In spite of other approaches to compute the probability mass function of the Poisson-Tweedie distribution are available in the literature (Esnaola et al., 2013, ; Barabesi et al., 2016). In this material, we opted to compute it by numerical evaluation of the integral in (2.12) using the Monte Carlo method as implemented by the following functions.

```
# Integrand Poisson X Tweedie distributions
integrand <- function(x, y, mu, phi, power) {
  int = dpois(y, lambda = x)*dtweedie(x, mu = mu,
                                         phi = phi, power = power)
  return(int)
}

# Computing the pmf using Monte Carlo
dptw <- function(y, mu, phi, power, control_sample) {
  pts <- control_sample$pts
  norma <- control_sample$norma
  integral <- mean(integrand(pts, y = y, mu = mu, phi = phi,
                             power = power)/norma)
  return(integral)
}
dptw <- Vectorize(dptw, vectorize.args = "y")
```

When using the Monte Carlo method, we need to specify a proposal distribution, from which samples will be taken to compute the integral as an expectation. In the Poisson-Tweedie case is sensible to use the Tweedie distribution as proposal. Thus, in our function we use the argument `control_sample` to provide these values. The advantage of this approach is that we need to simulate values once and we can reuse them for all evaluations of the probability mass function, as shown in the following code.

```
require(tweedie)
set.seed(123)
pts <- rtweedie(n = 1000, mu = 10, phi = 1, power = 2)
norma <- dtweedie(pts, mu = 10, phi = 1, power = 2)
control_sample <- list("pts" = pts, "norma" = norma)
dptw(y = c(0, 5, 10, 15), mu = 10, phi = 1, power = 2,
      control_sample = control_sample)

## [1] 0.0937 0.0590 0.0354 0.0217

dnbinom(x = c(0, 5, 10, 15), mu = 10, size = 1)

## [1] 0.0909 0.0564 0.0350 0.0218
```

It is also possible to use the Gauss-Laguerre method to approximate the integral in (2.12). In the supplementary material Script2.R, we provide R functions using both Monte Carlo and Gauss-Laguerre methods to approximate the probability mass function of Poisson-Tweedie distribution.

Figure 2.2 presents the empirical probability mass function of some Poisson-Tweedie distributions computed based on a sample of size 100000 (gray). Furthermore, we present an approximation (black) for the probability mass function obtained by Monte Carlo integration. We considered different values of the Tweedie power parameter  $p = 1.1, 2$ , and  $3$  combined with different values of the dispersion index. In all scenarios the expectation  $\mu$  was fixed at  $10$ .

For all scenarios considered the Monte Carlo method provides a quite accurate approximation to the empirical probability mass function. For these examples, we used 5000 random samples from the proposal distribution.

Finally, the Poisson-Tweedie regression model is defined by

$$Y_i \sim \text{PTW}_p(\mu_i, \phi), \quad \text{with} \quad \mu_i = g^{-1}(\mathbf{x}_i^\top \boldsymbol{\beta}),$$

where  $\mathbf{x}_i$  and  $\boldsymbol{\beta}$  are  $(p \times 1)$  vectors of known covariates and unknown regression parameters. The estimation and inference of Poisson-Tweedie regression models based on the maximum likelihood method are challenged by the presence of an intractable integral in the probability mass function and non-trivial restrictions on the power parameter space. In Chapter 3,

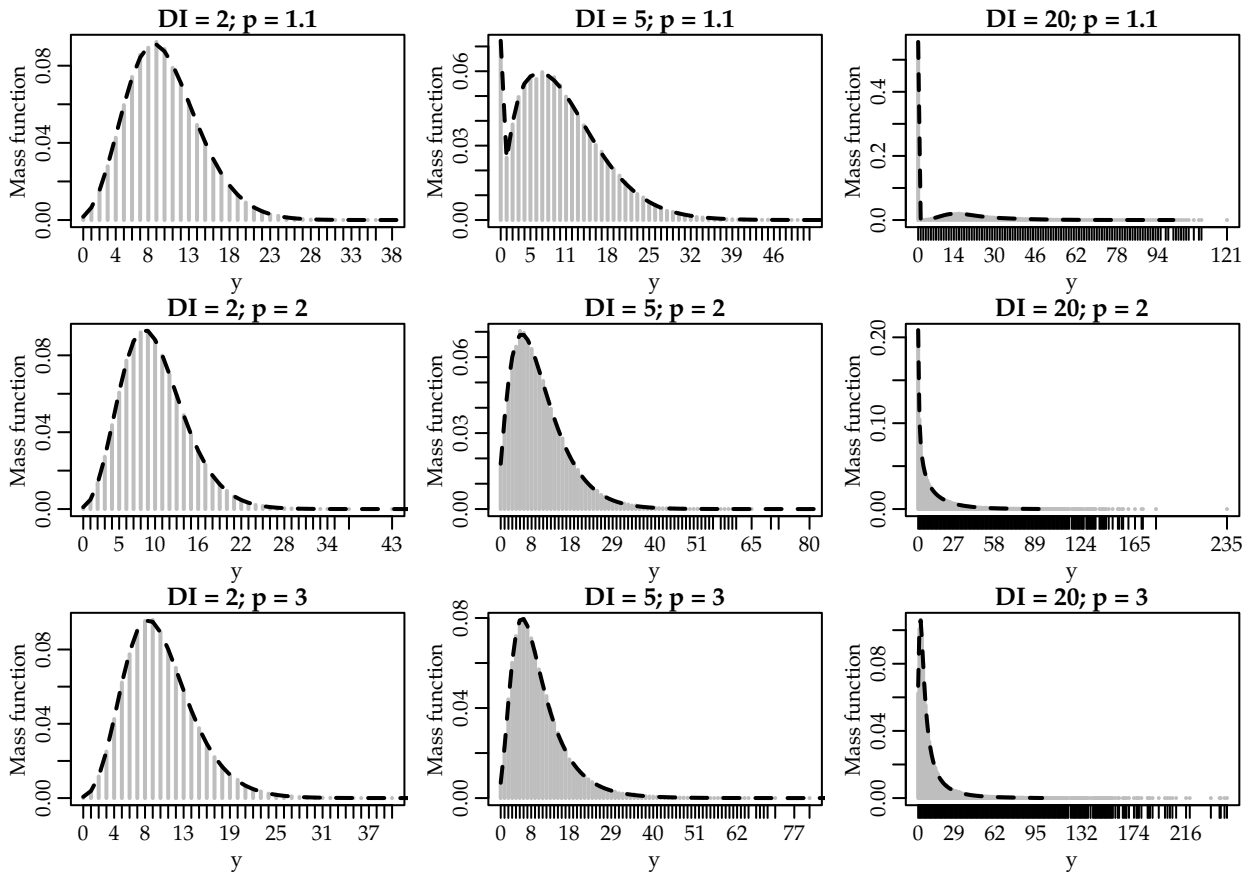


Figure 2.2: Empirical (gray) and approximated (black) Poisson-Tweedie probability mass function by values of the dispersion index (DI) and Tweedie power parameter.

we discuss maximum likelihood estimation for Poisson-Tweedie regression. Furthermore, in Chapter 4 we extended the Poisson-Tweedie model by using an estimating function approach in the style of Wedderburn (1974).

## 2.4 COM-Poisson distribution

The COM-Poisson distribution belongs to the family of weighted Poisson distributions. A random variable  $Y$  is a weighted Poisson distribution if its probability mass function can be written in the form

$$f(y; \lambda, \nu) = \frac{\exp\{-\lambda\} \lambda^y w_y}{W_y!}, \quad y = 0, 1, \dots,$$

where  $W = \sum_{i=0}^{\infty} \exp\{-\lambda\} \lambda^i w_i / i!$  is a normalizing constant (Sellers et al., 2012). The COM-Poisson is obtained when  $w_y = (y!)^{1-\nu}$  for  $\nu \geq 0$ . The series  $W$  for COM-Poisson distribution is denoted by  $Z(\lambda, \nu)$  and can be written as  $\sum_{i=0}^{\infty} \lambda^i / (i!)^\nu$ . Note that the series is theoretically divergent only

Table 2.1: Values for  $Z(\lambda, \nu)$  constant (calculated numerically) to combined values of  $\lambda$  (0.5 to 50) and  $\phi$  (0 to 1)

	0.5	1	5	10	30	50
0	2.00	Inf	Inf	Inf	Inf	Inf
0.1	1.92	7.64	Inf	Inf	Inf	Inf
0.2	1.86	5.25	3.17e+273	Inf	Inf	Inf
0.3	1.81	4.32	1.60e+29	2.54e+282	Inf	Inf
0.4	1.77	3.80	4.71e+10	1.33e+56	Inf	Inf
0.5	1.74	3.47	1.34e+06	3.67e+22	3.32e+196	Inf
0.6	1.72	3.23	2.05e+04	4.99e+12	1.73e+76	4.63e+177
0.7	1.70	3.06	2.37e+03	3.69e+08	4.93e+39	6.93e+81
0.8	1.68	2.92	6.49e+02	2.70e+06	5.09e+24	3.43e+46
0.9	1.66	2.81	2.74e+02	1.47e+05	1.80e+17	2.19e+30
1	1.65	2.72	1.48e+02	2.20e+04	1.07e+13	5.18e+21

when  $\nu = 0$  and  $\lambda \geq 1$ , but numerically for small values of  $\nu$  combined with large values of  $\lambda$ , the sum is so huge it causes overflow. The Table 2.1 shows the sums calculated with 1000 increments, in other words,  $\sum_{i=0}^{1000} \lambda^i / (i!)^\nu$  for different values of  $\nu$  (horizontal lines) and  $\lambda$  (vertical lines).

In general, the expectation and variance of the COM-Poisson distribution cannot be expressed in closed-form. However, they can be approximated by

$$E(Y) \approx \lambda^{1/\nu} - \frac{\nu - 1}{2\nu} \quad \text{and} \quad \text{var}(Y) \approx \frac{1}{\nu} \lambda^{1/\nu}.$$

These approximations are accurate when  $\nu \leq 1$  or  $\lambda > 10^\nu$ . The infinite sum involved in computing the probability mass function of the COM-Poisson distribution can be approximated to any level of precision. It can be evaluated in R using the function `dcom()` from the `compoisson` package (Dunn, 2012). Figure 2.3 presents some COM-Poisson probability mass functions. We tried to find parameters  $\lambda$  and  $\nu$  in order to have  $E(Y) = 10$  and dispersion index equals to  $DI = 0.5, 2, 5$  and  $20$ . However, we could not find any parameter combination to have  $DI = 20$ . Probably, it was due to overflow of the sum since  $\nu$  is inversely proportional to the dispersion index (see table 2.1).

Sellers and Shmueli (2010) proposed a regression model based on the COM-Poisson distribution where the parameter  $\lambda$  is described by the values of known covariates in a generalized linear models style. The COM-Poisson

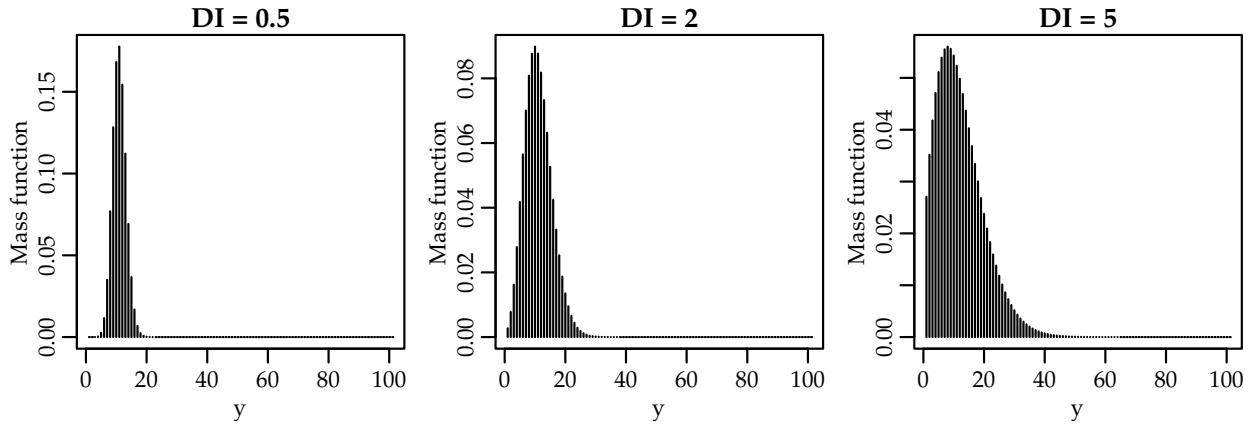


Figure 2.3: COM-Poisson probability mass function by values of the dispersion index (DI).

regression model is defined by

$$Y_i \sim CP(\lambda_i, \nu), \quad \text{with} \quad \lambda_i = g^{-1}(\mathbf{x}_i^\top \boldsymbol{\beta}).$$

In this notation, the parameter  $\nu$  is considered the dispersion parameter such that  $\nu > 1$  represents underdispersion and  $\nu < 1$  overdispersion. The Poisson model is obtained for  $\nu = 1$  and as usual we adopt the logarithm link function for  $g$ .

## 2.5 Comparing count distributions

Let  $Y$  be a count random variable and  $E(Y) = \mu$  and  $\text{var}(Y)$  denote its mean and variance, respectively. To explore and compare the flexibility of the models aforementioned, we introduce the dispersion (DI), zero-inflation (ZI) and heavy-tail (HT) indexes, which are respectively given by

$$DI = \frac{\text{var}(Y)}{E(Y)}, \quad ZI = 1 + \frac{\log P(Y = 0)}{E(Y)} \quad (2.13)$$

and

$$HT = \frac{P(Y = y + 1)}{P(Y = y)} \quad \text{for} \quad y \rightarrow \infty. \quad (2.14)$$

These indexes are defined in relation to the Poisson distribution. Thus, the dispersion index indicates underdispersion for  $DI < 1$ , equidispersion for  $DI = 1$  and overdispersion for  $DI > 1$ . Similarly, the zero-inflation index is

easily interpreted, since  $ZI < 0$  indicates zero-deflation,  $ZI = 0$  corresponds to no excess of zeroes and  $ZI > 0$  indicates zero-inflation. Finally,  $HT \rightarrow 1$  when  $y \rightarrow \infty$  indicates a heavy tail distribution.

For the Poisson distribution the dispersion index equals  $1 \forall \mu$ . In the Poisson case, it is easy to show that  $ZI = 0$  and  $HT \rightarrow 0$  when  $y \rightarrow \infty$ . Thus, it is quite clear that the Poisson model can deal only with equidispersed data and has no flexibility to deal with zero-inflation and/or heavy tail count data. In fact, the presented indexes were proposed in relation to the Poisson distribution in order to highlight its limitations.

Figure 2.4 presents the relationship between mean and variance, the dispersion and zero-inflation indexes as a function of the expected values  $\mu$  for different scenarios and count distributions. Scenario 1 corresponds to the case of underdispersion. Thus, we fixed the dispersion index at  $DI = 0.5$  when the mean equalling 10. Since the Poisson-Tweedie cannot deal with underdispersion, in this scenario we present only the Gamma-Count and COM-Poisson distributions. Similarly, scenarios 2 – 4 are obtained by fixing the dispersion index at  $DI = 2, 5$  and  $10$  when mean equalling 10. In the scenario 4 we could not find a parameter configuration in order to have a COM-Poisson distribution with dispersion index equals 20. Consequently, we present results only for the Gamma-Count and Poisson-Tweedie distributions. Furthermore, Figure 2.5 presents the heavy tail index for some extreme values of the random variable  $Y$ .

The indexes presented in Figures 2.4 and 2.5 show that for all considered scenarios the Gamma-Count and COM-Poisson distributions are quite similar. In general, for these distributions, the indexes slightly depend on the expected values and tend to stabilize for large values of  $\mu$ . Consequently, the mean and variance relationship is proportional to the dispersion parameter value. In the overdispersion case, the Gamma-Count and COM-Poisson distributions can handle with a limited amount of zero-inflation and are in general light tailed distributions, i.e.  $HT \rightarrow 0$  for  $y \rightarrow \infty$ .

Regarding the Poisson-Tweedie distributions the indexes show that for small values of the power parameter the Poisson-Tweedie distribution is suitable to deal with zero-inflated count data. In that case, the  $DI$  and  $ZI$  are almost not dependent on the values of the mean. Furthermore, the  $HT$  decreases as the mean increases. On the other hand, for large values of the power parameter the  $HT$  increases with increasing mean, showing that the model is specially suitable to deal with heavy-tailed count data. In this case, the  $DI$  and  $ZI$  increase quickly as the mean increases giving an extremely overdispersed model for large values of the mean. In general, the  $DI$  and  $ZI$  are larger than one and zero, respectively, which, of course, show that the corresponding Poisson-Tweedie distributions cannot deal with underdispersed and zero-

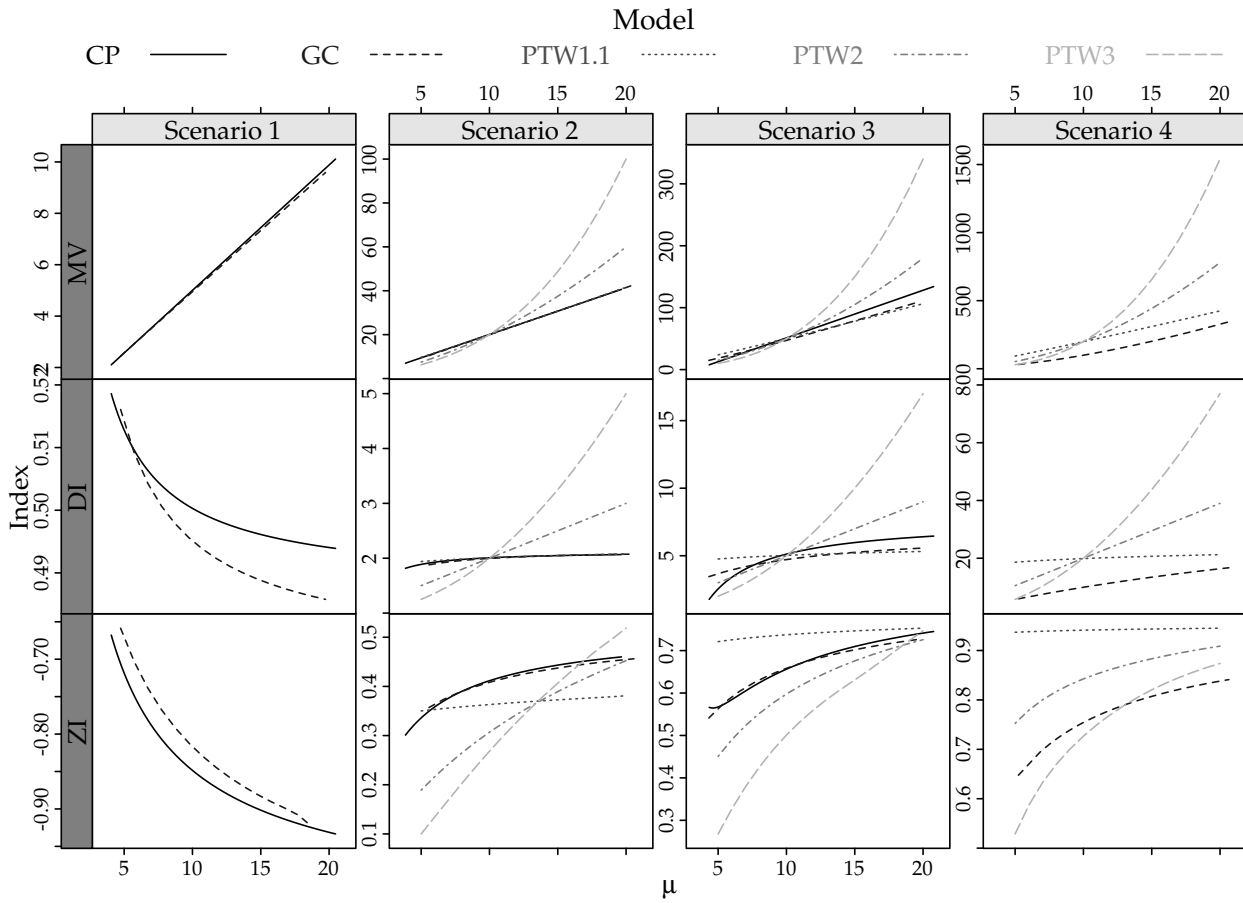


Figure 2.4: Mean and variance relationship (first line), dispersion (DI) and zero-inflation (ZI) indexes as a function of the expected values by simulation scenarios and count distributions.

deflated count data.

For multi-parameter probability function, a desirable feature is the orthogonality between parameters. This property leads a series of statistical implications as allows make inference for one parameter without worrying about the values of the other and computationally numerical methods for adjusting distributions with orthogonal parameters are more stable and fast.

The orthogonality is defined by second derivatives of the log-likelihood function, however for the Gamma-Count, Poisson-Tweedie and COM-Poisson and we cannot obtain such derivatives analytically. So we designed a simulation study to evaluate the properties about log-likelihood function for Gamma-Count and COM-Poisson. We used sample size  $n = 5000$  for simulate values of Gamma-Count and COM-Poisson following dispersion indexes  $DI = 0.5, 2, 5$  and  $20$  and plot the deviance contours around the maximum likelihood estimation. The Script6.R, in supplement material contains the codes for simulation study. The results are shows in Figure @ (fig:ortho).



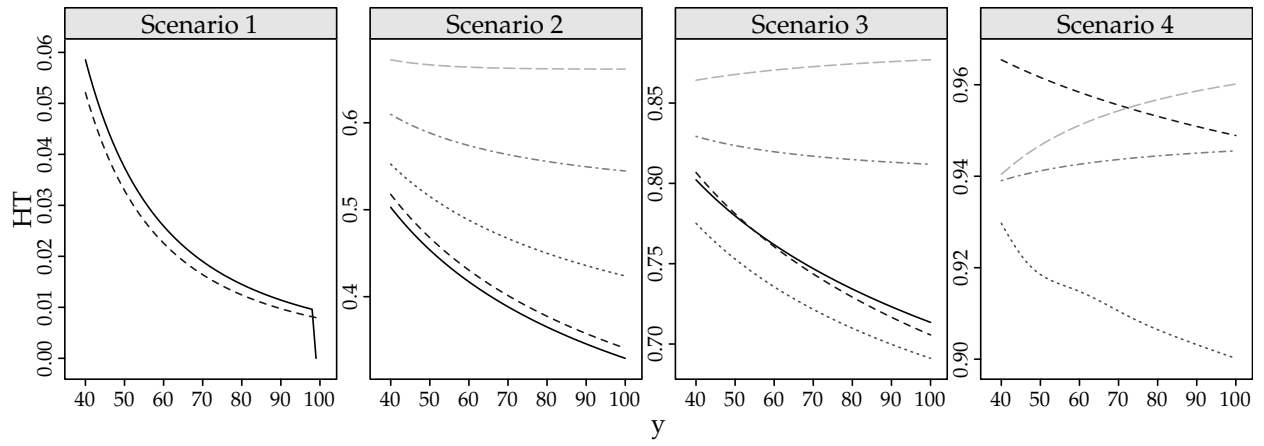


Figure 2.5: Heavy tail index for some extreme values of the random variable  $Y$  by simulation scenarios and count distributions.

This graphics presented in Figure 2.6 show that deviance contours are similar a quadratic function (blue dashed contours) for  $DI = 0.5, 2, 5$  in both distributions. However, for the Gamma-Count with  $DI = 20$  the quadratic aproximation not as good. For  $DI = 20$  and  $E[Y] = 10$  we could not find any parameter combination for COM-Poisson and for Gamma-Count the estimation is innacurate. With respect to orthogonality the Gamma-Count distribution is preferable, the deviance contours are well-behaved, while for COM-Poisson the contour are strongly flat in one direction making the estimation process difficult..

In terms of regression models, the Poisson and Poisson-Tweedie models are easy and convenient to interpret because the expected value is directly modelled as a function of known covariates in a generalized linear models manner. On the other hand, the Gamma-Count specifies the regression model for the expectation of the times between events and, thus requires careful interpretation. The COM-Poisson regression model is hard to interpret and compare with the traditional Poisson regression model, since it specifies the regression model for the parameter  $\lambda$  that has no easy interpretation in relation to the expectation of the count response variable.

Finally, in terms of computational implementation the simplicity of the Poisson regression model is unquestionable. The probability mass function of the Gamma-Count distribution requires the evaluation of the difference between two cumulative gamma distributions. For large values of the random variable  $Y$ , such a difference can be time consuming and inaccurately computed. Similarly, the COM-Poisson probability mass function involves the evaluation of a infinity sum, which can be computational expensive and inaccurate for large values of  $Y$ . Furthermore, for extreme values of  $\lambda$  combined with small values of  $\nu$  the infinite sum can numerically diverges

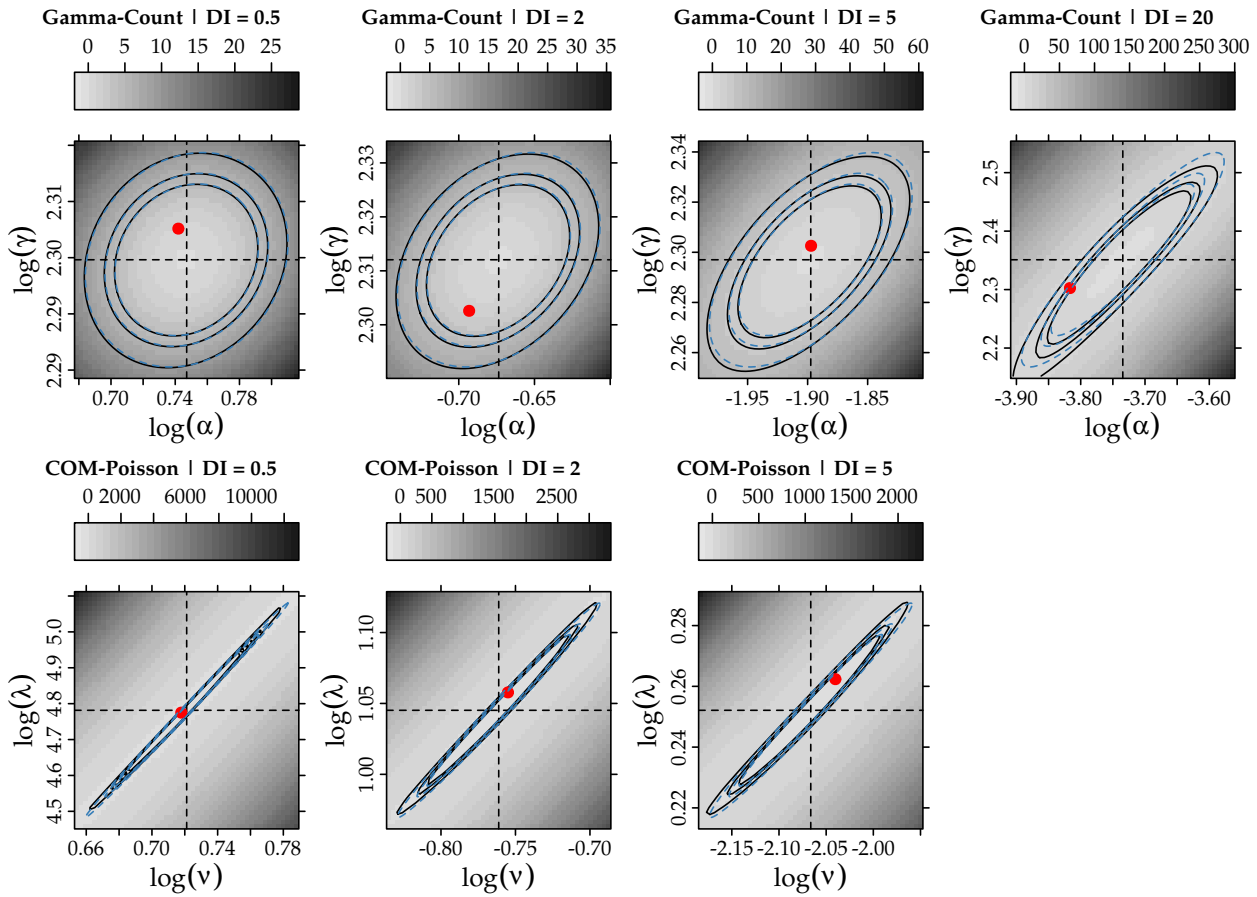


Figure 2.6: Deviance surfaces and quadratic approximation with confidence regions (90, 95 and, 99%) for the two parameters following dispersion indexes (DI) in the simulation study. Dashed lines represents the MLE estimation and red points the parameters used in simulation.

making impossible to evaluate the probability mass function. Finally, the Poisson-Tweedie probability mass function involves an intractable integral, which makes the estimation and inference based on likelihood methods computationally intensive.

## Chapter 3

# The method of maximum likelihood

The estimation and inference for the models discussed in Chapter 2 can be done by the method of maximum likelihood (Silvey, 1975). In this Chapter, we present the maximum likelihood method and its main properties along with some examples in R. The maximum likelihood method is applicable mainly in situations where the true distribution of the count random variable  $Y$  is known apart of the values of a finite number of unknown parameters. Let  $f(y; \theta)$  denote the true probability mass function of the count random variable  $Y$ . We assume that the family  $f(y; \theta)$  is labelled by a  $(p \times 1)$  parameter vector  $\theta$  taking values in  $\Theta$  a subset of  $\mathbb{R}^n$ . For a given observed value  $y$  of the a random variable  $Y$ , the likelihood function corresponding to the observation  $y$  is defined as  $L(\theta; y) = f(y; \theta)$ . It is important to highlight that  $f(y; \theta)$  is a probability mass function on the sample space. On the other hand,  $L(\theta; y) = f(y; \theta)$  is a function on the parameter space  $\Theta$ . The likelihood function expresses the plausibilities of different parameters after we have observed  $y$ , in the absence of any other information that we may have about these different values. In particular, for count random variables the likelihood function is the probability of the point  $y$  when  $\theta$  is the true parameter.

The method of maximum likelihood has a strong intuitive appeal and according to it, we estimate the true parameter  $\theta$  by any parameter which maximizes the likelihood function. In general, there is a unique maximizing parameter which is the most plausible and this is the maximum likelihood estimate (Silvey, 1975). In other words, a maximum likelihood estimate  $\hat{\theta}(y)$

is any element of  $\Theta$  such that  $L(\hat{\theta}(y); y) = \max_{\theta \in \Theta} L(\theta; y)$ . At this stage, we make the distinction between the estimate  $\hat{\theta}(y)$  and the estimator  $\hat{\theta}$ . However, we are not maintain this distinction and we shall use only  $\hat{\theta}$  leaving the context to make it clear whether we are thinking of  $\hat{\theta}$  as a function or as a particular value of a function.

Let  $Y_i$  be independent and identically distributed count random variables with probability mass function  $f(y; \theta)$ , whose observed values are denoted by  $y_i$  for  $i = 1, \dots, n$ . In this case, the likelihood function can be written as the product of the individuals probability mass distributions, i.e.

$$L(\theta; \mathbf{y}) = \prod_{i=1}^n L(\theta; y_i) = \prod_{i=1}^n f(y_i; \theta). \quad (3.1)$$

For convenience, in practical situations is advisable to work with the log-likelihood function obtained by taking the logarithm of Eq. (3.1). Thus, the maximum likelihood estimator (MLE) for the parameter vector  $\theta$  is obtained by maximizing the following log-likelihood function,

$$\ell(\theta) = \sum_{i=1}^n \log\{L(\theta; y_i)\}. \quad (3.2)$$

Often, it is not possible to find a relatively simple expression in closed form for the maximum likelihood estimates. However, it is usually possible to assume that maximum likelihood estimates emerge as a solution of the likelihood equations or also called score functions, i.e.

$$\mathcal{U}(\theta) = \left( \frac{\partial \ell(\theta)}{\partial \theta_1}, \dots, \frac{\partial \ell(\theta)}{\partial \theta_p} \right)^\top = \mathbf{0}. \quad (3.3)$$

The system of non-linear equations in (3.3) often have to be solved numerically. The entry  $(i, j)$  of the  $p \times p$  Fisher information matrix  $\mathcal{F}_\theta$  for the vector of parameter  $\theta$  is given by

$$\mathcal{F}_{\theta_{ij}} = -E \left\{ \frac{\partial^2 \ell(\theta)}{\partial \theta_i \partial \theta_j} \right\}. \quad (3.4)$$

In order to solve the system of equations  $\mathcal{U}(\theta) = \mathbf{0}$ , we employ the Newton scoring algorithm, defined by

$$\boldsymbol{\theta}^{(i+1)} = \boldsymbol{\theta}^{(i)} - \mathcal{F}_{\boldsymbol{\theta}}^{-1} \mathcal{U}(\boldsymbol{\theta}^{(i)}). \quad (3.5)$$

Finally, the well known distribution of the maximum likelihood estimator  $\hat{\boldsymbol{\theta}}$  is  $N(\boldsymbol{\theta}, \mathcal{F}_{\boldsymbol{\theta}}^{-1})$ . Thus, the maximum likelihood estimator is asymptotically consistent, unbiased and efficient.

A critical point of the approach described so far, is that we should be able to compute the first and second derivatives of the log-likelihood function. However, for the Gamma-Count where the log-likelihood function is given by the difference between two integrals, we cannot obtain such derivatives analytically. Similarly, for the COM-Poisson the log-likelihood function involves an infinite sum and consequently such derivatives cannot be obtained analytically. Finally, in the Poisson-Tweedie distribution the log-likelihood function is defined by an intractable integral, which implies that we cannot obtain a closed-form for the score function and Fisher information matrix.

Thus, an alternative approach is to maximize directly the log-likelihood function in equation (3.2) using a derivative-free algorithm as the Nelder-Mead method (Nelder and Mead, 1965) or some other numerical method for maximizing the log-likelihood function, examples include the BFGS, conjugate gradient and simulated annealing. All of them are implemented in R through the `optim()` function. The package `bbmle` (?) offers a suite of functions to work with numerical maximization of log-likelihood functions in R. As an example, consider the Gamma-Count distribution described in subsection 2.2. The log-likelihood function for the parameters  $\boldsymbol{\theta} = (\gamma, \alpha)$  in R is given by

```
ll_gc <- function(gamma, alpha, y) {
  ll <- sum(dgc(y = y, gamma = gamma, alpha = alpha, log = TRUE))
  return(-ll)
}
```

Thus, for a given vector of observed count values, we can numerically maximize the log-likelihood function above using the function `mle2()` from the `bbmle` package. It is important to highlight that by default the `mle2()` function requires the negative of the log-likelihood function instead of the log-likelihood itself. Thus, our function returns the negative value of the log-likelihood function.

```
require(bbmle)
y <- rpois(100, lambda = 10)
fit_gc <- mle2(ll_gc, start = list("gamma" = 10, "alpha" = 1),
               data = list("y" = y))
```

The great advantage of the `bbmle` package for maximum likelihood estimation in R, is that it already provides standard methods, such as `summary()`, `coef()`, `confint()`, `vcov()`, `profile()` and other for objects of `mle2` class.

```
summary(fit_gc)
```

```
## Maximum likelihood estimation
##
## Call:
## mle2(minuslogl = ll_gc, start = list(gamma = 10, alpha = 1),
##      data = list(y = y))
##
## Coefficients:
##      Estimate Std. Error z value Pr(z)
## gamma      9.842      0.335  29.38 < 2e-16 ***
## alpha      0.929      0.139   6.68 2.5e-11 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## -2 log L: 518
```

Similar functions can be done for the Poisson, Poisson-Tweedie and COM-Poisson distributions.

## Chapter 4

# Models specified by second-moment assumptions

In Chapter 2, we presented four statistical models to deal with count data and in the Chapter 3 the method of maximum likelihood was introduced to estimate the model's parameters. As discussed in Chapter 3 the method of maximum likelihood assumes that the true distribution of the count random variable  $Y$  is known apart of the values of a finite number of unknown parameters. In this Chapter, we shall present a different approach for model specification, estimation and inference based only on second-moment assumptions.

### 4.1 Extended Poisson-Tweedie model

The Poisson-Tweedie distribution as presented in subsection 2.3 provides a very flexible family of count distributions, however, such a family has two main drawbacks: it cannot deal with underdispersed count data and its probability mass function is given by an intractable integral, which implies that estimation based on the maximum likelihood method is computational demanding for practical data analysis.

In spite of these issues Jørgensen and Kokonendji (2015) showed using factorial cumulant function that for  $Y \sim \text{PTW}_p(\mu, \phi)$ ,  $E(Y) = \mu$  and  $\text{var}(Y) =$

$\mu + \phi\mu^p$ . This fact motivates Bonat et al. (2016) to specify a model by using only second-moment assumptions, i.e. mean and variance.

Thus, consider a cross-section dataset,  $(y_i, \mathbf{x}_i)$ ,  $i = 1, \dots, n$ , where  $y_i$ 's are i.i.d. realizations of  $Y_i$  according to an unspecified distribution, whose expectation and variance are given by

$$\begin{aligned} E(Y_i) &= \mu_i = g^{-1}(\mathbf{x}_i^\top \boldsymbol{\beta}) \\ \text{var}(Y_i) &= C_i = \mu_i + \phi\mu_i^p, \end{aligned} \quad (4.1)$$

where as before  $\mathbf{x}_i$  and  $\boldsymbol{\beta}$  are  $(p \times 1)$  vectors of known covariates and unknown regression parameters and  $g$  is the logarithm link function. The regression model specified in (4.1) is parametrized by  $\boldsymbol{\theta} = (\boldsymbol{\beta}^\top, \boldsymbol{\lambda}^\top)^\top$ , where  $\boldsymbol{\lambda} = (\phi, p)$ .

Note that, based on second-moment assumptions, the only restriction to have a proper model is that  $\text{var}(Y_i) > 0$ , thus

$$\phi > -\mu_i^{(1-p)},$$

which shows that at least at some extent negative values for the dispersion parameter are allowed. Consequently, the Poisson-Tweedie model can be extended to deal with underdispersed count data, however, in doing so the associated probability mass functions do not exist. However, in a regression modelling framework as discussed in this material, we are in general interested in the regression coefficient effects, thus such an issue does not imply any loss of interpretation and applicability. The formulation of the extended Poisson-Tweedie model is exactly the same of the quasi-binomial and quasi-Poisson models popular in the context of generalized linear models, see (Wedderburn, 1974, Nelder and Wedderburn (1972)) for details. Furthermore, note that for  $p = 1$  the extended Poisson-Tweedie regression model corresponds to a reparametrization of the popular quasi-Poisson regression model.

It is also important to highlight that in this case the relationship between mean and variance is proportional to the dispersion parameter  $\phi$  as in the Gamma-Count and COM-Poisson distributions. Thus, we expect for  $\phi < 0$  and  $p = 1$  the extended Poisson-Tweedie model presents results in terms of regression coefficients really similar the ones from the Gamma-Count and COM-Poisson regression models.



## 4.2 Estimation and Inference

Since the model presented in (4.1) is based only on second-moment assumptions the method of maximum likelihood cannot be employed. Bonat et al. (2016) based on ideas of Jørgensen and Knudsen (2004) and Bonat and Jørgensen (2016) proposed an estimating function approach for estimation and inference for the extended Poisson-Tweedie regression model. Bonat et al. (2016) combined the quasi-score and Pearson estimating functions for estimation of the regression and dispersion parameters respectively. Following Bonat et al. (2016) the quasi-score function for  $\beta$  has the following form,

$$\psi_{\beta}(\beta, \lambda) = \left( \sum_{i=1}^n \frac{\partial \mu_i}{\partial \beta_1} C_i^{-1} (Y_i - \mu_i), \dots, \sum_{i=1}^n \frac{\partial \mu_i}{\partial \beta_p} C_i^{-1} (Y_i - \mu_i) \right)^{\top},$$

where  $\partial \mu_i / \partial \beta_j = \mu_i x_{ij}$  for  $j = 1, \dots, p$ . The sensitivity matrix is defined as the expectation of the first derivative of the estimating function with respect to the model parameters. Thus, the entry  $(j, k)$  of the  $p \times p$  sensitivity matrix for  $\psi_{\beta}$  is given by

$$S_{\beta_{jk}} = E \left( \frac{\partial}{\partial \beta_k} \psi_{\beta_j}(\beta, \lambda) \right) = - \sum_{i=1}^n \mu_i x_{ij} C_i^{-1} x_{ik} \mu_i. \quad (4.2)$$

In a similar way, the variability matrix is defined as the variance of the estimating function. In particular, for the quasi-score function the entry  $(j, k)$  of the  $p \times p$  variability matrix is given by

$$V_{\beta_{jk}} = \text{Cov}(\psi_{\beta_j}(\beta, \lambda), \psi_{\beta_k}(\beta, \lambda)) = \sum_{i=1}^n \mu_i x_{ij} C_i^{-1} x_{ik} \mu_i.$$

The Pearson estimating function for the dispersion parameters has the following form,

$$\psi_{\lambda}(\lambda, \beta) = \left( - \sum_{i=1}^n \frac{\partial C_i^{-1}}{\partial \phi} [(Y_i - \mu_i)^2 - C_i], - \sum_{i=1}^n \frac{\partial C_i^{-1}}{\partial p} [(Y_i - \mu_i)^2 - C_i] \right)^{\top}.$$

Note that, the Pearson estimating functions are unbiased estimating functions for  $\lambda$  based on the squared residuals  $(Y_i - \mu_i)^2$  with expected value  $C_i$ .

The entry  $(j, k)$  of the  $2 \times 2$  sensitivity matrix for the dispersion parameters is given by

$$S_{\lambda_j \lambda_k} = E \left( \frac{\partial}{\partial \lambda_k} \psi_{\lambda_j}(\lambda, \beta) \right) = - \sum_{i=1}^n \frac{\partial C_i^{-1}}{\partial \lambda_j} C_i \frac{\partial C_i^{-1}}{\partial \lambda_k} C_i, \quad (4.3)$$

where  $\lambda_1$  and  $\lambda_2$  denote either  $\phi$  or  $p$ .

Similarly, the cross entries of the sensitivity matrix are given by

$$S_{\beta_j \lambda_k} = E \left( \frac{\partial}{\partial \lambda_k} \psi_{\beta_j}(\beta, \lambda) \right) = 0 \quad (4.4)$$

and

$$S_{\lambda_j \beta_k} = E \left( \frac{\partial}{\partial \beta_k} \psi_{\lambda_j}(\lambda, \beta) \right) = - \sum_{i=1}^n \frac{\partial C_i^{-1}}{\partial \lambda_j} C_i \frac{\partial C_i^{-1}}{\partial \beta_k} C_i. \quad (4.5)$$

Finally, the joint sensitivity matrix for the parameter vector  $\theta$  is given by

$$S_{\theta} = \begin{pmatrix} S_{\beta} & \mathbf{0} \\ S_{\lambda\beta} & S_{\lambda} \end{pmatrix},$$

whose entries are defined by equations (4.2), (4.3), (4.4) and (4.5).

We now calculate the asymptotic variance of the estimating function estimators denoted by  $\hat{\theta}$ , as obtained from the inverse Godambe information matrix, whose general form for a vector of parameter  $\theta$  is  $J_{\theta}^{-1} = S_{\theta}^{-1} V_{\theta} S_{\theta}^{-\top}$ , where  $-\top$  denotes inverse transpose. The variability matrix for  $\theta$  has the form

$$V_{\theta} = \begin{pmatrix} V_{\beta} & V_{\beta\lambda} \\ V_{\lambda\beta} & V_{\lambda} \end{pmatrix}, \quad (4.6)$$

where  $V_{\lambda\beta} = V_{\beta\lambda}^{\top}$  and  $V_{\lambda}$  depend on the third and fourth moments of  $Y_i$ , respectively. In order to avoid this dependence on higher-order moments, we use the empirical versions of  $V_{\lambda}$  and  $V_{\lambda\beta}$  as given by

$$\tilde{V}_{\lambda_j \lambda_k} = \sum_{i=1}^n \psi_{\lambda_j}(\lambda, \beta)_i \psi_{\lambda_k}(\lambda, \beta)_i \quad \text{and} \quad \tilde{V}_{\lambda_j \beta_k} = \sum_{i=1}^n \psi_{\lambda_j}(\lambda, \beta)_i \psi_{\beta_k}(\lambda, \beta)_i.$$

Finally, the well known asymptotic distribution of  $\hat{\theta}$  (Jørgensen and Knudsen, 2004) is given by

$$\hat{\theta} \sim N(\theta, J_{\theta}^{-1}), \quad \text{where} \quad J_{\theta}^{-1} = S_{\theta}^{-1} V_{\theta} S_{\theta}^{-\top}.$$

To solve the system of equations  $\psi_{\beta} = \mathbf{0}$  and  $\psi_{\lambda} = \mathbf{0}$  Jørgensen and Knudsen (2004) proposed the modified chaser algorithm, defined by

$$\begin{aligned} \beta^{(i+1)} &= \beta^{(i)} - S_{\beta}^{-1} \psi_{\beta}(\beta^{(i)}, \lambda^{(i)}) \\ \lambda^{(i+1)} &= \lambda^{(i)} - \alpha S_{\lambda}^{-1} \psi_{\lambda}(\beta^{(i+1)}, \lambda^{(i)}). \end{aligned}$$

The modified chaser algorithm uses the insensitivity property (4.4), which allows us to use two separate equations to update  $\beta$  and  $\lambda$ . We introduce the tuning constant,  $\alpha$ , to control the step-length. This algorithm is a special case of the flexible algorithm presented by Bonat and Jørgensen (2016) in the context of multivariate covariance generalized linear models. Hence, estimation for the extended Poisson-Tweedie model is easily implemented in R through the `mcglm` (Bonat, 2016) package.



# Chapter 5

## Applications

In this chapter, we will bring some applications based on real data sets to show how use R packages to analyse count data.

### 5.1 Cotton Bolls

Cotton production can be drastically reduced by attack of defoliating insects. Depending on the growth stage, the plant can recover from the caused damage and keeps production not affected or can have the production reduced by low intensity defoliation.

A greenhouse experiment with cotton plants (*Gossypium hirsutum*) was done under a completely randomized design with five replicates to assess the effects of five defoliation levels (0%, 25%, 50%, 75% and 100%) on the observed number of bolls produced by plants at five growth stages: vegetative, flower-bud, blossom, fig and cotton boll. The experimental unity was a pot with two plants (Silva et al., 2012, for more). The number of cotton bolls was recorded at the harvest of the experiment.

```
library(lattice)
library(latticeExtra)
library(gridExtra)
library(plyr)
library(car)
library(corrplot)
library(doBy)
library(multcomp)
library(mcglim)
```

```

library(MRDCr)
ls("package:MRDCr")

## [1] "apc"                "calc_mean_cmp"
## [3] "calc_mean_gcnt"     "calc_var_cmp"
## [5] "cambras"            "capdesfo"
## [7] "capmosca"           "cmp"
## [9] "conftemp"           "confterm"
## [11] "convergencez"       "dcmp"
## [13] "dgcnt"              "dpgnz"
## [15] "gcnt"               "led"
## [17] "llcmp"              "llgcnt"
## [19] "llpgnz"             "nematoide"
## [21] "ninfas"             "panel.beeswarm"
## [23] "panel.cbH"          "panel.groups.segplot"
## [25] "peixe"              "pgnz"
## [27] "postura"            "prepanel.cbH"
## [29] "seguros"            "soja"

# Documentation in Portuguese.
help(capdesfo, help_type = "html")

str(capdesfo)

## 'data.frame':   125 obs. of  4 variables:
## $ est : Factor w/ 5 levels "vegetativo","botão floral",...: 1 1 1 1 1 1 1 1 1 1 ...
## $ des : num  0 0 0 0 0 0.25 0.25 0.25 0.25 0.25 ...
## $ rept: int  1 2 3 4 5 1 2 3 4 5 ...
## $ ncap: int  10 9 8 8 10 11 9 10 10 10 ...

levels(capdesfo$est) <- c("vegetative",
                          "flower-bud",
                          "blossom",
                          "fig",
                          "cotton boll")

xtabs(~est + des, data = capdesfo)

##           des
## est      0 0.25 0.5 0.75 1
## vegetative 5    5    5    5 5
## flower-bud 5    5    5    5 5
## blossom    5    5    5    5 5
## fig        5    5    5    5 5
## cotton boll 5    5    5    5 5

```

Figure 5.1 (top) shows the beeswarm plot of number of cotton bolls recorded for each combination of defoliation level and growth stage. All the points in the sample means and variances dispersion diagram (bottom) are below the identity line, clearly suggesting data with underdispersion.

The exploratory data analysis, although simple, was able to detect departures from the Poisson equidispersion assumption. So, we have in advance few conditions met for the use of GLM Poisson as a regression model to analyse this experiment.

Poisson, as being a process derived from the memoryless waiting times Exponential distribution, implies that each boll is an independent event in the artificial subjacent domain, that can be thought was the natural resource domain that the plant has to allocate bolls. Its is easy to assume, based on plant fisiology, that the probability of a boll decreases with the number of previous bolls because the plant's resource to produce bolls is limited and it is a non memoryless process equivalent.

Based on the exploratory data analysis, a predictor with 2nd order effect of defoliation for each growth stage should be enough to model the number of bolls mean in a regression model. The analysis and assessment of the effects of the experimental factors are based on the Poisson, Gamma-count and Poisson Tweedie models.

```
m0 <- glm(ncap ~ est * (des + I(des^2)),
          data = capdesfo,
          family = poisson)
```

```
summary(m0)
```

```
##
## Call:
## glm(formula = ncap ~ est * (des + I(des^2)), family = poisson,
##      data = capdesfo)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.0771  -0.3098  -0.0228   0.2704   1.1665
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)      2.2142    0.1394  15.89  <2e-16 ***
## estflower-bud    -0.0800    0.2007  -0.40    0.69
## estblossom       -0.0272    0.2001  -0.14    0.89
## estfig           -0.1486    0.2051  -0.72    0.47
## estcotton boll    0.1129    0.1922   0.59    0.56
## des              0.3486    0.6805   0.51    0.61
## I(des^2)         -0.7384    0.6733  -1.10    0.27
## estflower-bud:des  0.1364    0.9644   0.14    0.89
## estblossom:des    -1.5819    1.0213  -1.55    0.12
## estfig:des        0.4755    1.0194   0.47    0.64
## estcotton boll:des -0.8210    0.9395  -0.87    0.38
```

```
## estflower-bud:I(des^2)    0.1044    0.9447    0.11    0.91
## estblossom:I(des^2)      1.4044    1.0191    1.38    0.17
## estfig:I(des^2)         -0.9294    1.0323   -0.90    0.37
## estcotton boll:I(des^2)  1.0757    0.9210    1.17    0.24
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for poisson family taken to be 1)
##
##      Null deviance: 75.514  on 124  degrees of freedom
## Residual deviance: 25.331  on 110  degrees of freedom
## AIC: 539.7
##
## Number of Fisher Scoring iterations: 4
```

**logLik(m0)**

```
## 'log Lik.' -255 (df=15)
```

We fit the GLM Poisson regression model using the standard `glm()` function in R. The fitted model summary shows the estimated parameters for the second order effect of defoliation crossed with growth stages levels. The residual deviance was 25.33 based on 110 degrees of freedom. The ratio  $25.33/110 = 0.23$  is a strong evidence against Poisson equidispersion assumption that uses a dispersion parameter equals 1.

**anova(m0, test = "Chisq")**

```
## Analysis of Deviance Table
##
## Model: poisson, link: log
##
## Response: ncap
##
## Terms added sequentially (first to last)
##
##
##              Df Deviance Resid. Df Resid. Dev Pr(>Chi)
## NULL              124      75.5
## est                4   19.96   120   55.6 0.00051 ***
## des                1   15.86   119   39.7 6.8e-05 ***
## I(des^2)           1    1.29   118   38.4 0.25557
## est:des            4    6.71   114   31.7 0.15212
## est:I(des^2)       4    6.36   110   25.3 0.17388
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

The analysis of deviance table did not stated effect of any interactions, neither second order effect of defoliation. Although, all these effects are noticeable



in Figure ??.

Figure ?? displays the four residual plots for the fitted model. Based on these plots, there is no concern about misspecifications regarding to the model predictor or influential observations. The only remarkable aspect is about the range of the standardized deviance residuals quite distant from the expected -3 to 3 from the normal distribution. Once more, these is another measure indicating a underdispersed count data.

The `gcnt()` is a function defined in the `MRDCr` package (Zeviani et al., 2016) to fit the Gamma-Count regression model. This function fits a GML-Poisson to use the estimates as initial values to optimize Gamma-Count likelihood using `optim()` through `bblme` package (Bolker and Team, 2016).

```
m1 <- gcnt(ncap ~ est * (des + I(des^2)),
           data = capdesfo)
summary(m1)

## Maximum likelihood estimation
##
## Call:
## bbmle::mle2(minuslogl = llgcnt, start = start, data = list(y = y,
##      X = X, offset = off), vecpar = TRUE)
##
## Coefficients:
##
##              Estimate Std. Error z value Pr(z)
## alpha              1.7110      0.1352  12.66 < 2e-16 ***
## (Intercept)         2.2580      0.0593  38.05 < 2e-16 ***
## estflower-bud       -0.0765      0.0854  -0.90 0.37025
## estblossom          -0.0253      0.0851  -0.30 0.76616
## estfig              -0.1398      0.0872  -1.60 0.10885
## estcotton boll       0.1084      0.0818   1.33 0.18476
## des                  0.3294      0.2896   1.14 0.25543
## I(des^2)            -0.6997      0.2866  -2.44 0.01464 *
## estflower-bud:des     0.1337      0.4105   0.33 0.74456
## estblossom:des       -1.5020      0.4345  -3.46 0.00055 ***
## estfig:des           0.4218      0.4336   0.97 0.33062
## estcotton boll:des   -0.7820      0.3998  -1.96 0.05046 .
## estflower-bud:I(des^2) 0.0943      0.4021   0.23 0.81452
## estblossom:I(des^2)   1.3382      0.4335   3.09 0.00202 **
## estfig:I(des^2)      -0.8333      0.4390  -1.90 0.05768 .
## estcotton boll:I(des^2) 1.0222      0.3920   2.61 0.00911 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## -2 log L: 408
```

During the optimization process for this dataset, `optim()` has found NaN

when evaluating the likelihood. This occurs due little numerical precision to calculate the difference of Gamma CDFs on tails or for extreme values, resulting in numerical zeros and corresponding  $-\text{Inf}$  log-likelihood. This is a numerical problem that can narrow, or make things difficult, the use of Gamma-Count regression model.

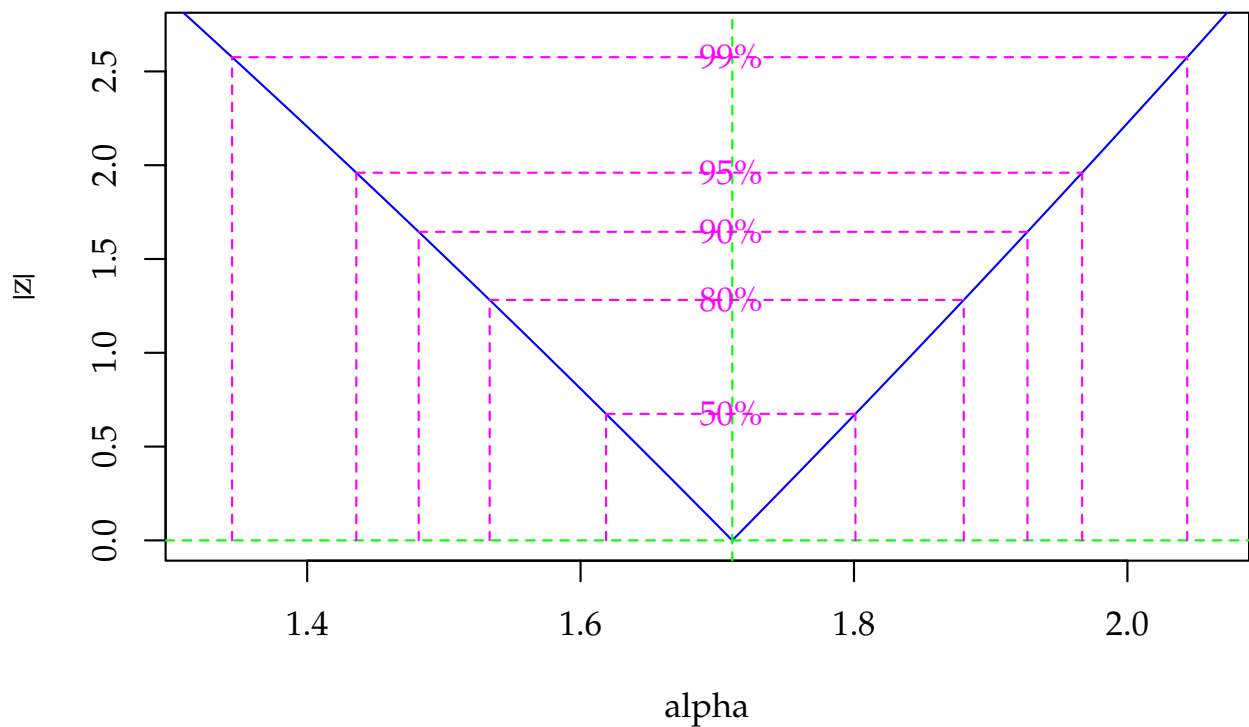
The dispersion parameter is the first position in the parameter vector. The optimization was carried out on the log scale to avoid problems regarding to bounded parameter spaces. As the dispersion parameter is in fact interpreted as a precision coefficient, the positive estimate indicates an underdispersed count. According to the  $z$  statistic,  $\hat{\alpha}$  is significantly different from zero (Poisson case). Poisson is special case of Gamma-Count when  $\alpha = 0$ , so we can perform a likelihood ratio test to the hypothesis  $H_0 : \alpha = 0$ .

```
# Likelihood ratio test.
chi <- 2 * (logLik(m1) - logLik(m0))
pval <- 2 * pchisq(chi, df = 1, lower.tail = FALSE)
cat("Likelihood Ratio Test\n",
    "Chisq:\t\t ", chi, "\n",
    "Pr(>Chisq):\t ", pval, "\n",
    sep = "")
```

```
## Likelihood Ratio Test
## Chisq:          102
## Pr(>Chisq):     1.01e-23
```

```
# Log-likelihood profile for alpha.
plot(profile(m1, which = "alpha"))
```

## Likelihood profile: alpha



```
cbind(c(0, coef(m0)), coef(m1))
```

```
##           [,1]    [,2]
##           0.0000  1.7110
## (Intercept)  2.2142  2.2580
## estflower-bud -0.0800 -0.0765
## estblossom    -0.0272 -0.0253
## estfig        -0.1486 -0.1398
## estcotton boll  0.1129  0.1084
## des           0.3486  0.3294
## I(des^2)      -0.7384 -0.6997
## estflower-bud:des  0.1364  0.1337
## estblossom:des   -1.5819 -1.5020
## estfig:des       0.4755  0.4218
## estcotton boll:des -0.8210 -0.7820
## estflower-bud:I(des^2)  0.1044  0.0943
## estblossom:I(des^2)    1.4044  1.3382
## estfig:I(des^2)      -0.9294 -0.8333
## estcotton boll:I(des^2)  1.0757  1.0222
```

```
rstd <- summary(m1)@coef[-1, 2]/summary(m0)$coeff[, 2]
plyr::each(mean, range)(rstd)
```

```
##   mean range1 range2
##  0.426  0.425  0.426
```

The estimates for the location parameters were very close. The ratio between

Gamma-Count parameters standard error and Poisson ones, on the other hand, were 0.426 for all estimates, for 3 decimals of precision. This leads to the conclusion that TODO relação linear no parâmetro de dispersão.

```
# Wald test for the interaction.
a <- c(0, attr(model.matrix(m0), "assign"))
ai <- a == max(a)
L <- t(replicate(sum(ai), rbind(coef(m1) * 0), simplify = "matrix"))
L[, ai] <- diag(sum(ai))
```

```
linearHypothesis(model = m0, # m0 is not being used here.
                  hypothesis.matrix = L,
                  vcov. = vcov(m1),
                  coef. = coef(m1))
```

```
## Linear hypothesis test
##
## Hypothesis:
## estflower - bud:I(des^2) = 0
## estblossom:I(des^2) = 0
## estfig:I(des^2) = 0
## estcotton boll:I(des^2) = 0
##
## Model 1: restricted model
## Model 2: ncap ~ est * (des + I(des^2))
##
## Note: Coefficient covariance matrix supplied.
##
##   Res.Df Df Chisq Pr(>Chisq)
## 1      114
## 2      110  4  30.5    3.8e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

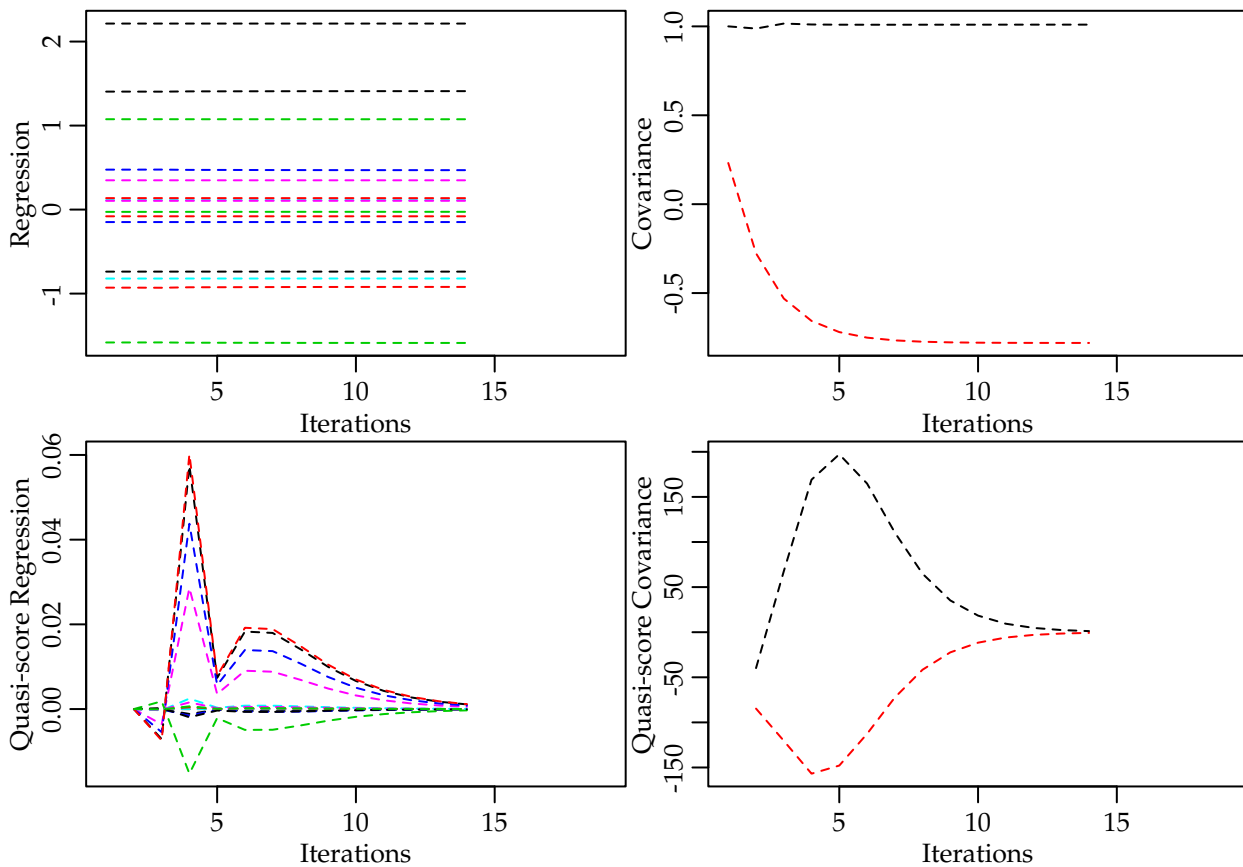
```
# Fitting Poisson-Tweedie model.
m2 <- mcglm(linear_pred = c(ncap ~ est * (des + I(des^2))),
            matrix_pred = list(mc_id(data = capdesfo)),
            link = "log",
            variance = "poisson_tweedie",
            power_fixed = FALSE,
            data = capdesfo,
            control_algorithm = list(verbose = FALSE,
                                     max_iter = 100,
                                     tuning = 0.5,
                                     correct = FALSE))
```

```
## Automatic initial values selected.
```

```
# Parameter estimates.
summary(m2)
```

```
## Call: ncap ~ est * (des + I(des^2))
##
## Link function: log
## Variance function: poisson_tweedie
## Covariance function: identity
## Regression:
##
##               Estimates Std.error Z value
## (Intercept)      2.2143    0.0627  35.308
## estflower-bud    -0.0800    0.0904  -0.886
## estblossom       -0.0265    0.0900  -0.294
## estfig           -0.1485    0.0924  -1.607
## estcotton boll    0.1128    0.0864   1.306
## des              0.3486    0.3065   1.137
## I(des^2)         -0.7385    0.3037  -2.432
## estflower-bud:des  0.1361    0.4345   0.313
## estblossom:des    -1.5875    0.4612  -3.442
## estfig:des        0.4693    0.4602   1.020
## estcotton boll:des -0.8204    0.4228  -1.940
## estflower-bud:I(des^2) 0.1048    0.4259   0.246
## estblossom:I(des^2)  1.4098    0.4609   3.059
## estfig:I(des^2)    -0.9205    0.4671  -1.971
## estcotton boll:I(des^2) 1.0752    0.4149   2.592
##
## Power:
##      Estimates Std.error Z value
## 1      1.01      0.136   7.43
##
## Dispersion:
##      Estimates Std.error Z value
## 1     -0.781     0.217  -3.59
##
## Algorithm: chaser
## Correction: FALSE
## Number iterations: 14
```

```
plot(m2, type = "algorithm")
```



```
# Wald test for fixed effects.
```

```
anova(m2)
```

```
## Wald test for fixed effects
```

```
## Call: anova(m2) ~ est * (des + I(des^2))
```

```
##
```

	Covariate	Chi.Square	Df	p.value
## 1	estflower-bud	9.54	4	0.0489
## 2	des	1.29	1	0.2555
## 3	I(des^2)	5.91	1	0.0150
## 4	estflower-bud:des	25.00	4	0.0001
## 5	estflower-bud:I(des^2)	30.72	4	0.0000

```
# New data values for prediction.
```

```
pred <- with(capdesfo,
  expand.grid(est = levels(est),
    des = seq(0, 1, length.out = 30)))
```

```
# Corresponding model matrix.
```

```
X <- model.matrix(formula(m0)[-2], data = pred)
```

```
#-----
```

```
# Poisson.
```

```
yp <- predict(m0, newdata = pred, se.fit = TRUE)
```

```
em <- outer(yp$se.fit,
```

```

      c(lwrP = -1, fitP = 0, uprP = 1) * qnorm(0.975),
      FUN = "*"")
ci <- sweep(em, MARGIN = 1, STATS = yp$fit, FUN = "+")
ci <- m0$family$linkinv(ci)

pred <- cbind(pred, as.data.frame(ci))
str(pred)

## 'data.frame':    150 obs. of  5 variables:
## $ est : Factor w/ 5 levels "vegetative","flower-bud",...: 1 2 3 4 5 1 2 3 4 5 ...
## $ des : num  0 0 0 0 0 ...
## $ lwrP: num  6.97 6.37 6.72 5.88 7.91 ...
## $ fitP: num  9.15 8.45 8.91 7.89 10.25 ...
## $ uprP: num  12 11.2 11.8 10.6 13.3 ...

# TODO predito pelo 3 modelos.

```

TODO: pares de estimativa e erros padrões.

## 5.2 Soybean pod and beans

The tropical soils, usually poor in potassium (K), demand potassium fertilization when cultivated with soybean (*Glycine max* L.) to obtain satisfactory yields. Soybean production is affected by long exposition to water deficit. As potassium is a nutrient involved in the water balance in plant, by photosynthesis, a good supply of potassium avoids lose production.

The aim of this experiment was to evaluate the effects of K doses and soil humidity levels on soybean production. The experiment was carried out in a greenhouse, in pots with two plants, containing 5 dm<sup>3</sup> of soil. The experimental design was completely randomized block with treatments in a 5 x 3 factorial arrangement. The K doses were 0, 30, 60, 120 and 180 mg dm<sup>-3</sup>, and the soil humidity ranged from 35 to 40, 47.5 to 52.5, and 60 to 65% of the total porosity (Serafim et al., 2012, for more details).

Two count variables were recorded in this experiment: the total number of pods per plot and the total number of grains per plot. The ratio, grains/pod, can also be analysed, since the fisiological response can change it.

There is an outlier in the dataset at position 74 that must be removed. Potassium amount (K) will be converted to a categorical factor, despite it is a numerical one, to prevent concerns with lack of fit, that is not the main scope of the following analysis.

```

data(soja)
str(soja)

```

```
## 'data.frame':   75 obs. of  5 variables:
## $ K    : int   0 30 60 120 180 0 30 60 120 180 ...
## $ umid: Factor w/ 3 levels "37,5","50","62,5": 1 1 1 1 1 2 2 2 2 2 ...
## $ bloc: Factor w/ 5 levels "I","II","III",...: 1 1 1 1 1 1 1 1 1 1 ...
## $ ngra: int  136 159 156 171 190 140 193 200 208 237 ...
## $ nvag: int   56 62 66 68 82 63 86 94 86 97 ...

# Removing an outlier.
soja <- soja[-74, ]
soja <- transform(soja, K = factor(K))
```

### 5.2.1 Number of pods

The pod is a dehiscent fruit of a leguminous plant such as the beans and soybeans. The pod is the basic unit of production in soybean, so factors that reduces or increases the number of pods have impact in crop production.

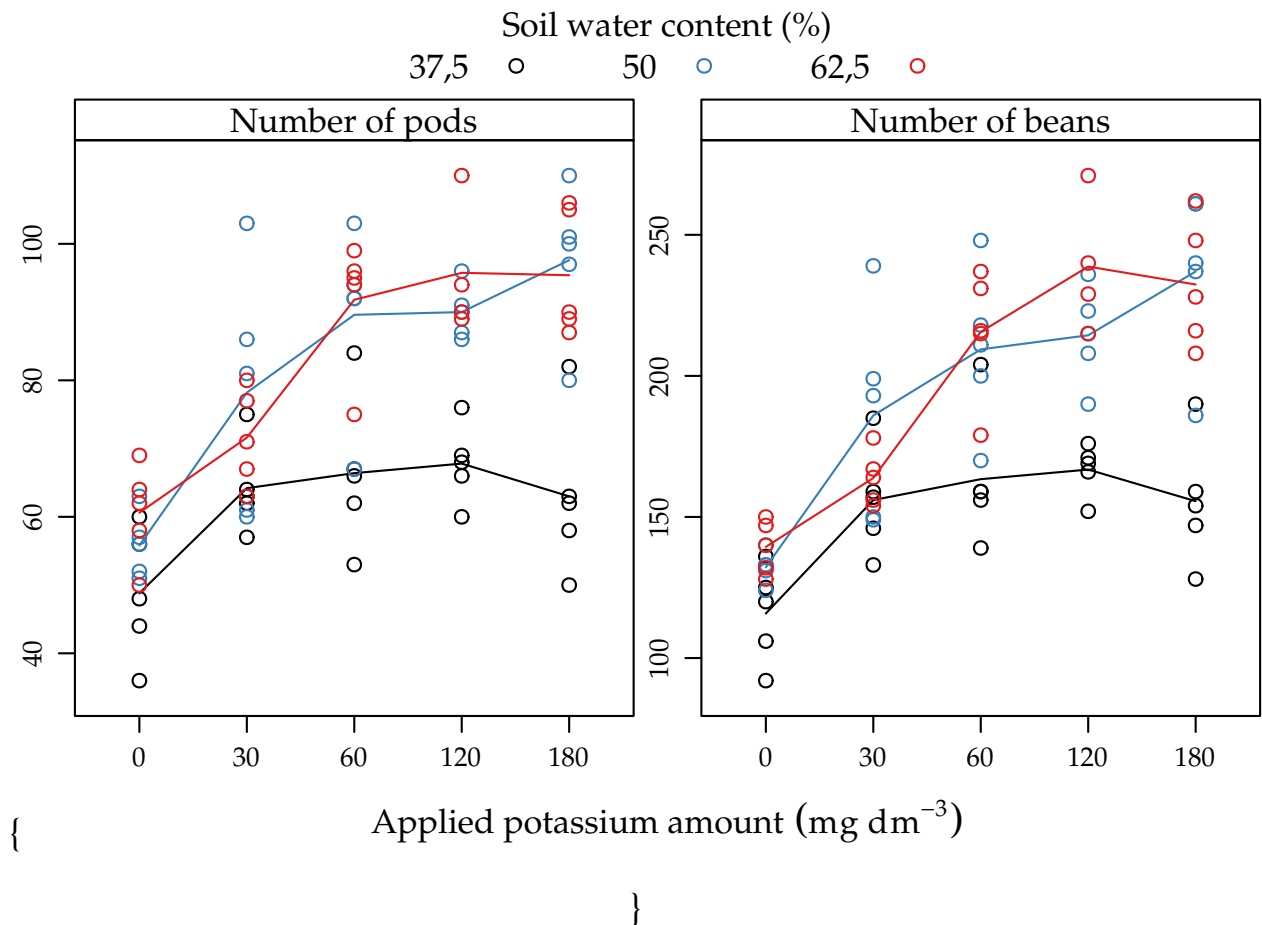
The potassium amount (K) increased the number of pods and number of beans in all soil water content levels (A) (5.2.1). For lowest soil water level, the mean was fairly low than the other levels that perform very similar. The pattern driven by the mean lines suggests interaction between applied potassium amount and soil water content for both variables.

We could obtain the sample variance and sample mean for each cell combination for checking dispersion level but it would ignore the block effect, that can enlarges the variance.

```
xyplot(nvag + ngra ~ K,
       groups = umid,
       outer = TRUE,
       data = soja,
       type = c("p", "a"),
       scales = "free",
       ylab = NULL,
       xlab = expression("Applied potassium amount" ~ (mg ~ dm^{-3})),
       auto.key = list(title = "Soil water content (%)",
                       cex.title = 1,
                       columns = 3),
       strip = strip.custom(
         factor.levels = c("Number of pods",
                           "Number of beans")))
```

`\begin{figure}[h]`





\caption{Number of pods and beans as function of potassium amount (K) for each soil water content level (%). Lines passes on the average of points. The mean response pattern is the same on the two variables.} \end{figure}

The following analysis will be carried out in parallel, that is, results will be showed together in each step for the sake of comparison.

We fit Poisson (P), Gamma-Count (GC) and Poisson-Tweedie (TW) regression models. The first two were fit by maximum likelihood and the last by moments specification.

```
#-----
# Poisson.

m0 <- glm(nvag ~ bloc + umid * K,
          data = soja,
          family = poisson)

#-----
# Gamma-Count.

m1 <- gcnt(formula(m0), data = soja)

#-----
```

```
# Tweedie.
```

```
m2 <- mcglm(linear_pred = c(nvag ~ bloc + umid * K),
            matrix_pred = list(mc_id(data = soja)),
            link = "log",
            variance = "poisson_tweedie",
            power_fixed = TRUE,
            data = soja,
            control_algorithm = list(verbose = FALSE,
                                    max_iter = 100,
                                    tuning = 0.5,
                                    correct = FALSE))
```

```
## Automatic initial values selected.
```

To fit Poisson and Gamma-Count, the correspond function were used in the default setting. To fit Poisson Tweedie, we `power_fixed = TRUE` option, since the number of pods is a close to equidispersed count variable. The Poisson-Tweedie loses identifiability in the equidispersion zone because the variance is  $\text{var}(Y) = \mu + \phi\mu^p$ , then  $p$  can be any value if  $\phi$  goes to zero. We fixed  $p = 1$ .

Figure 5.3 shows the 4 plots based on residuals. This residuals didn't show any departure pattern. On the contray, the axes of the qq-norm plot shows the same range, indicating a equidispersed count variable.

The maximised log-likelihood were very close for Poisson and Gamma-Count models. The profile log-likelihood for Gamma-Count dispersion parameter contains 0 inside (Figure 5.4), so indicating a close to Poisson case.

```
#-----
# Comparing models.
```

```
# Log-likelihood.
c(P = logLik(m0), GC = logLik(m1), TW = NA)
```

```
##      P      GC      TW
## -260 -259      NA
```

```
cap <-
```

```
  "Profile log-likelihood for the Gamma-Count dispersion parameter. The confidence i
# Likelihood profile for Gamma-Count dispersion parameter.
plot(profile(m1, which = "alpha"))
abline(v = 0, lty = 2)
```

The estimates and standard erros also were close on the (location) regression parameters for all models. They differ only in the dispersion parameter by construction.

```

c0 <- summary(m0)$coefficients[, 1:2]
c1 <- summary(m1)$coef[, 1:2]
c2 <- rbind(summary(m2)[[1]]$tau[, 1:2],
             summary(m2)[[1]]$Regression[, 1:2])

# Parameter estimates according to each model.
c4 <- cbind("P" = rbind(NA, c0),
            "GC" = c1,
            "TW" = c2)
colnames(c4) <- substr(colnames(c4), 1, 6)
round(c4, digits = 4)

##           P.Esti P.Std.  GC.Est GC.Std  TW.Est TW.Std
##           NA      NA   0.1288 0.1655 -0.1088 0.1458
## (Intercept)  3.9537 0.0689  3.9549 0.0646  3.9537 0.0651
## blocII      -0.0293 0.0409 -0.0293 0.0383 -0.0293 0.0386
## blocIII     -0.0727 0.0414 -0.0726 0.0388 -0.0727 0.0390
## blocIV      -0.1254 0.0419 -0.1253 0.0393 -0.1254 0.0396
## blocV       -0.1079 0.0430 -0.1079 0.0403 -0.1079 0.0406
## umid50       0.1340 0.0877  0.1339 0.0822  0.1340 0.0827
## umid62,5     0.2166 0.0860  0.2163 0.0806  0.2166 0.0812
## K30          0.2743 0.0849  0.2740 0.0796  0.2743 0.0802
## K60          0.3080 0.0843  0.3076 0.0790  0.3080 0.0796
## K120         0.3288 0.0840  0.3285 0.0787  0.3288 0.0793
## K180         0.2554 0.0853  0.2551 0.0799  0.2554 0.0805
## umid50:K30   0.0632 0.1156  0.0632 0.1083  0.0632 0.1091
## umid62,5:K30 -0.1075 0.1154 -0.1073 0.1081 -0.1075 0.1089
## umid50:K60   0.1656 0.1137  0.1655 0.1066  0.1656 0.1073
## umid62,5:K60 0.1074 0.1122  0.1073 0.1052  0.1074 0.1059
## umid50:K120  0.1492 0.1134  0.1491 0.1063  0.1492 0.1070
## umid62,5:K120 0.1184 0.1140  0.1184 0.1069  0.1184 0.1077
## umid50:K180  0.3037 0.1136  0.3035 0.1065  0.3037 0.1072
## umid62,5:K180 0.1984 0.1126  0.1983 0.1055  0.1984 0.1063

```

Until now, all models perform very close suggesting to keep the Poisson by parsimony. But for testing purposes, Gamma-Count and Poisson-Tweedie got a more significant p-value for the potassium amount  $\times$  soil water content interaction. If a 5% significance is adopted, models lead to different practical conclusions.

```

# Analysis of deviance table.
anova(m0, test = "Chisq")

## Analysis of Deviance Table
##
## Model: poisson, link: log
##
## Response: nvag

```

```
##
## Terms added sequentially (first to last)
##
##
##           Df Deviance Resid. Df Resid. Dev Pr(>Chi)
## NULL                73        323
## bloc      4      14.3      69        308  0.0064 **
## umid      2      92.9      67        215 <2e-16 ***
## K         4     136.1      63         79 <2e-16 ***
## umid:K    8      14.2      55         65  0.0779 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

# Wald test for interaction.
a <- c(0, attr(model.matrix(m0), "assign"))
ai <- a == max(a)
L <- t(replicate(sum(ai), rbind(coef(m1) * 0), simplify = "matrix"))
L[, ai] <- diag(sum(ai))
linearHypothesis(model = m0, # m0 is not being used here.
                 hypothesis.matrix = L,
                 vcov. = vcov(m1),
                 coef. = coef(m1))

## Linear hypothesis test
##
## Hypothesis:
## umid50:K30 = 0
## umid62,5:K30 = 0
## umid50:K60 = 0
## umid62,5:K60 = 0
## umid50:K120 = 0
## umid62,5:K120 = 0
## umid50:K180 = 0
## umid62,5:K180 = 0
##
## Model 1: restricted model
## Model 2: nvag ~ bloc + umid * K
##
## Note: Coefficient covariance matrix supplied.
##
##   Res.Df Df Chisq Pr(>Chisq)
## 1      63
## 2     55  8    16    0.043 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

# Wald test for fixed effects.
anova(m2)
```

```
## Wald test for fixed effects
## Call: nvag ~ bloc + umid * K
##
##      Covariate Chi.Square Df p.value
## 1      blocII      13.92  4 0.0076
## 2      umid50       7.15  2 0.0280
## 3         K30      20.93  4 0.0003
## 4 umid50:K30      15.76  8 0.0459
```

out that more flexible models are powerful in detecting effects.

## 5.2.2 Number of grains

For the analysis of number of beans, the same steps will be carried out, just adaptating the code when needed. The first adaptation we need make is to overcome a numerical problem in the Gamma-Count implementation.

The Gamma-Count mass function, and also the likelihood, computes the difference of Gamma CDF. This difference can be numerically zero by lack of precision and this leads to a -Inf in the log-likelihood. To overcome this, we can use an artifitial offset that can prevent those zeros. The code below shows the effect of offset on the probabilities.

dgcnt

```
## function (y, lambda, alpha)
## {
##     p <- pgamma(q = 1, shape = y * alpha, rate = alpha * lambda) -
##         pgamma(q = 1, shape = (y + 1) * alpha, rate = alpha *
##             lambda)
##     return(p)
## }
## <environment: namespace:MRDCr>

y <- 0:30
lambda <- 30
alpha <- 5
off <- 1
pgamma(q = off,
      shape = y * alpha,
      rate = alpha * lambda) -
pgamma(q = off,
      shape = (y + 1) * alpha,
      rate = alpha * lambda)

## [1] 0.00e+00 0.00e+00 0.00e+00 0.00e+00 0.00e+00 0.00e+00 0.00e+00
## [8] 0.00e+00 0.00e+00 0.00e+00 0.00e+00 0.00e+00 1.78e-15 1.07e-13
## [15] 4.43e-12 1.32e-10 2.86e-09 4.62e-08 5.65e-07 5.31e-06 3.89e-05
## [22] 2.24e-04 1.03e-03 3.81e-03 1.14e-02 2.80e-02 5.67e-02 9.54e-02
## [29] 1.34e-01 1.59e-01 1.59e-01

off <- off * 0.5
pgamma(q = off,
      shape = y * alpha,
      rate = alpha * lambda) -
pgamma(q = off,
```

```

        shape = (y + 1) * alpha,
        rate = alpha * lambda)

## [1] 0.00e+00 0.00e+00 0.00e+00 1.24e-14 6.30e-12 1.15e-09 9.09e-08
## [8] 3.48e-06 7.10e-05 8.29e-04 5.87e-03 2.63e-02 7.76e-02 1.56e-01
## [15] 2.18e-01 2.19e-01 1.60e-01 8.68e-02 3.56e-02 1.12e-02 2.74e-03
## [22] 5.25e-04 8.00e-05 9.77e-06 9.67e-07 7.80e-08 5.18e-09 2.85e-10
## [29] 1.31e-11 5.03e-13 1.64e-14

```

The practical implication of an artificial offset is changing the size of the unit were the counts were observed, so the estimates are multiples of the those with offset 1. To be consistent, we will use the same artificial offset, 10, in all models. This offset, artificially, assumes that the count were observed as a sum of 10 pots.

```

soja$off <- 10
fivenum(with(soja, ngra/off))

## [1] 9.2 14.7 17.1 21.6 27.1

#-----
# Poisson.

m0 <- glm(ngra ~ offset(log(off)) + bloc + umid * K,
          data = soja,
          family = poisson)

#-----
# Gamma-Count.

m1 <- gcnt(formula(m0), data = soja)

#-----
# Tweedie.

m2 <- mcglm(linear_pred = c(ngra ~ bloc + umid * K),
            matrix_pred = list(mc_id(data = soja)),
            link = "log",
            offset = list(log(soja$off)),
            variance = "poisson_tweedie",
            power_fixed = FALSE,
            data = soja,
            control_algorithm = list(verbose = FALSE,
                                     max_iter = 100,
                                     tuning = 0.2,
                                     correct = FALSE))

## Automatic initial values selected.

```

To fit Poisson Tweedie for number of beans, we set `power_fixed = FALSE` option, because the number of beans is a more dispersed variable than number of pods. We used 0.2 for tuning to get convergence.

Figure 5.5 shows the 4 plots based on residuals. These residuals didn't show any departure pattern regarding misspecification of the predictor (lack of fit, for example). The y axis of the qq-norm plot has range on -3 to 4, indicating an overdispersed count variable.

The maximised log-likelihood were different between Poisson and Gamma-Count models. The profile log-likelihood for Gamma-Count dispersion parameter does not contain 0 inside (Figure 5.6), so indicating an overdispersed case. The profile shows a symmetric shape with almost linear or "V" shape that indicates a quadratic profile likelihood function.

```
#-----
# Comparing models.

# Log-likelihood.
c(P = logLik(m0), GC = logLik(m1), TW = NA)

##      P      GC      TW
## -322 -316      NA

cap <-
  "Profile log-likelihood for the Gamma-Count dispersion parameter. The confidence interval is
# Likelihood profile for Gamma-Count dispersion parameter.
plot(profile(m1, which = "alpha"))
```

The estimates for location parameters were the same for all three models.

The standard error for Poisson estimates were smaller than those for Gamma-Count and Poisson-Tweedie. The ratio of standard errors between Gamma-Count and Poisson were 1.3 (mean) and for Poisson-Tweedie 1.26 (mean).

```
c0 <- summary(m0)$coefficients[, 1:2]
c1 <- summary(m1)$coef[, 1:2]
c2 <- rbind(summary(m2)[[1]]$tau[, 1:2],
            summary(m2)[[1]]$Regression[, 1:2])

# Parameter estimates according to each model.
c4 <- cbind("P" = rbind(NA, c0),
           "GC" = c1,
           "TW" = c2)
colnames(c4) <- substr(colnames(c4), 1, 6)
round(c4, digits = 4)

##              P.Esti P.Std.  GC.Est GC.Std  TW.Est TW.Std
##              NA      NA -0.5231 0.1651  0.0053 0.0314
```



```
## (Intercept)      2.4994 0.0448  2.4965 0.0582  2.4995 0.0549
## blocII           -0.0194 0.0266 -0.0194 0.0345 -0.0216 0.0352
## blocIII          -0.0366 0.0267 -0.0367 0.0347 -0.0376 0.0353
## blocIV           -0.1056 0.0272 -0.1058 0.0353 -0.1003 0.0356
## blocV            -0.0931 0.0279 -0.0933 0.0362 -0.0974 0.0365
## umid50            0.1325 0.0569  0.1328 0.0740  0.1333 0.0696
## umid62,5          0.1855 0.0562  0.1860 0.0731  0.1871 0.0690
## K30               0.2980 0.0549  0.2988 0.0713  0.2992 0.0678
## K60               0.3443 0.0543  0.3452 0.0706  0.3447 0.0674
## K120              0.3649 0.0541  0.3659 0.0703  0.3662 0.0672
## K180              0.2954 0.0549  0.2962 0.0714  0.2950 0.0679
## umid50:K30        0.0434 0.0748  0.0434 0.0973  0.0401 0.0936
## umid62,5:K30     -0.1367 0.0753 -0.1371 0.0979 -0.1392 0.0940
## umid50:K60        0.1156 0.0736  0.1157 0.0957  0.1146 0.0926
## umid62,5:K60      0.0917 0.0729  0.0917 0.0948  0.0896 0.0920
## umid50:K120       0.1186 0.0733  0.1187 0.0953  0.1184 0.0923
## umid62,5:K120     0.1627 0.0737  0.1628 0.0958  0.1591 0.0937
## umid50:K180       0.2883 0.0733  0.2887 0.0953  0.2884 0.0923
## umid62,5:K180     0.2157 0.0729  0.2159 0.0947  0.2142 0.0920
```

```
# Ratios between standard errors.
```

```
cbind(GC = summary(c4[-1, 4]/c4[-1, 2]),
      TW = summary(c4[-1, 6]/c4[-1, 2]))
```

```
##           GC    TW
## Min.      1.3 1.22
## 1st Qu.   1.3 1.24
## Median    1.3 1.26
## Mean      1.3 1.26
## 3rd Qu.   1.3 1.27
## Max.      1.3 1.32
```

The Poisson model gave the higher statistic to the rejection of the null hypothesis than Gamma-Count and Poisson-Tweedie because it is assuming a dispersion of 1 that is not the case. Gamma-Count and Poisson-Tweedie perform very similar to test the interaction.

```
# Analysis of deviance table.
```

```
anova(m0, test = "Chisq")
```

```
## Analysis of Deviance Table
```

```
##
```

```
## Model: poisson, link: log
```

```
##
```

```
## Response: ngra
```

```
##
```

```
## Terms added sequentially (first to last)
```

```
##
```

```
##
##           Df Deviance Resid. Df Resid. Dev Pr(>Chi)
## NULL                73          761
## bloc      4          28          69          733 1.1e-05 ***
## umid      2         185          67          548 < 2e-16 ***
## K         4         380          63          168 < 2e-16 ***
## umid:K    8          43          55          125 8.8e-07 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

# Wald test for interaction.
a <- c(0, attr(model.matrix(m0), "assign"))
ai <- a == max(a)
L <- t(replicate(sum(ai), rbind(coef(m1) * 0), simplify = "matrix"))
L[, ai] <- diag(sum(ai))
linearHypothesis(model = m0, # m0 is not being used here.
                  hypothesis.matrix = L,
                  vcov. = vcov(m1),
                  coef. = coef(m1))

## Linear hypothesis test
##
## Hypothesis:
## umid50:K30 = 0
## umid62,5:K30 = 0
## umid50:K60 = 0
## umid62,5:K60 = 0
## umid50:K120 = 0
## umid62,5:K120 = 0
## umid50:K180 = 0
## umid62,5:K180 = 0
##
## Model 1: restricted model
## Model 2: ngra ~ offset(log(off)) + bloc + umid * K
##
## Note: Coefficient covariance matrix supplied.
##
##   Res.Df Df Chisq Pr(>Chisq)
## 1      63
## 2      55  8  25.3    0.0014 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

# Wald test for fixed effects.
anova(m2)

## Wald test for fixed effects
## Call: ngra ~ bloc + umid * K
```



### 5.2.3 Number of grains per pod

Analyse the number of beans per plot is better than the total number of beans, because the leter can be a side effect of the number of pods. We will analyse the number of beans using the number of pods as an offset.

```
#-----
# Poisson.

m0 <- glm(ngra ~ offset(log(nvag)) + bloc + umid * K,
          data = soja,
          family = poisson)

#-----
# Gamma-Count.

m1 <- gcnt(formula(m0), data = soja)

#-----
# Tweedie.

m2 <- mcglm(linear_pred = c(ngra ~ bloc + umid * K),
            matrix_pred = list(mc_id(data = soja)),
            link = "log",
            offset = list(log(soja$nvag)),
            variance = "poisson_tweedie",
            power_fixed = TRUE,
            data = soja,
            control_algorithm = list(verbose = FALSE,
                                     max_iter = 100,
                                     tuning = 0.5,
                                     correct = FALSE))

## Automatic initial values selected.
```

To accomplish the fitting, we set `power_fixed = TRUE` otherwise convergence wasn't met.

Figure 5.7 shows the 4 plots based on residuals. The y axis of the qq-norm plot has range on -1.5 to 1.5, indicating an underdispersed count variable.

The maximised log-likelihood were different between Poisson and Gamma-Count models. The profile log-likelihood for Gamma-Count dispersion parameter does not contain 0 inside (Figure 5.8), so indicating an underdispersed count. The profile likelihood is "V" shape that represents a quadratic profile log-likelihood function.

```
#-----
```

```
# Comparing models.
```

```
# Log-likelihood.
```

```
c(P = logLik(m0), GC = logLik(m1), TW = NA)
```

```
##      P      GC      TW
```

```
## -271 -255      NA
```

```
cap <-
```

```
  "Profile log-likelihood for the Gamma-Count dispersion parameter. The confide
```

```
# Likelihood profile for Gamma-Count dispersion parameter.
```

```
plot(profile(m1, which = "alpha"))
```

The point estimates for location parameters were, once more, the same for all three models. The standard error for Poisson estimates were greater than those for Gamma-Count and Poisson-Tweedie.

```
c0 <- summary(m0)$coefficients[, 1:2]
```

```
c1 <- summary(m1)$coef[, 1:2]
```

```
c2 <- rbind(summary(m2)[[1]]$tau[, 1:2],
            summary(m2)[[1]]$Regression[, 1:2])
```

```
# Parameter estimates according to each model.
```

```
c4 <- cbind("P" = rbind(NA, c0),
```

```
          "GC" = c1,
```

```
          "TW" = c2)
```

```
colnames(c4) <- substr(colnames(c4), 1, 6)
```

```
round(c4, digits = 4)
```

```
##          P.Esti P.Std.  GC.Est GC.Std  TW.Est TW.Std
##          NA      NA   1.1225 0.1648 -0.6764 0.0446
## (Intercept)  0.8482 0.0447  0.8510 0.0255  0.8482 0.0254
## blocII      0.0113 0.0266  0.0114 0.0152  0.0113 0.0151
## blocIII     0.0363 0.0267  0.0363 0.0152  0.0363 0.0152
## blocIV      0.0194 0.0272  0.0195 0.0155  0.0194 0.0155
## blocV       0.0161 0.0280  0.0163 0.0159  0.0161 0.0159
## umid50      -0.0019 0.0569 -0.0023 0.0325 -0.0019 0.0324
## umid62,5    -0.0312 0.0563 -0.0317 0.0321 -0.0312 0.0320
## K30         0.0229 0.0549  0.0221 0.0313  0.0229 0.0312
## K60         0.0345 0.0544  0.0337 0.0310  0.0345 0.0309
## K120        0.0353 0.0541  0.0344 0.0309  0.0353 0.0308
## K180        0.0406 0.0549  0.0399 0.0313  0.0406 0.0312
## umid50:K30  -0.0187 0.0748 -0.0187 0.0427 -0.0187 0.0426
## umid62,5:K30 -0.0285 0.0753 -0.0281 0.0429 -0.0285 0.0428
## umid50:K60  -0.0482 0.0736 -0.0483 0.0420 -0.0482 0.0419
## umid62,5:K60 -0.0140 0.0729 -0.0140 0.0416 -0.0140 0.0415
## umid50:K120 -0.0304 0.0733 -0.0305 0.0418 -0.0304 0.0417
## umid62,5:K120 0.0455 0.0738  0.0454 0.0421  0.0455 0.0420
## umid50:K180 -0.0163 0.0733 -0.0167 0.0418 -0.0163 0.0417
```

```
## umid62,5:K180 0.0169 0.0729 0.0167 0.0416 0.0169 0.0415
```

```
# Ratios between standard errors.
```

```
cbind(GC = summary(c4[-1, 4]/c4[-1, 2]),
      TW = summary(c4[-1, 6]/c4[-1, 2]))
```

```
##          GC      TW
## Min.    0.57 0.569
## 1st Qu. 0.57 0.569
## Median  0.57 0.569
## Mean    0.57 0.569
## 3rd Qu. 0.57 0.569
## Max.    0.57 0.569
```

None of the experimental factors had effect on the number of beans per pod. Although, Gamma-Count and Poisson-Tweedie showed more favorable statistics to the rejection of the null hypothesis than Poisson.

```
# Analysis of deviance table.
```

```
anova(m0, test = "Chisq")
```

```
## Analysis of Deviance Table
```

```
##
```

```
## Model: poisson, link: log
```

```
##
```

```
## Response: ngra
```

```
##
```

```
## Terms added sequentially (first to last)
```

```
##
```

```
##
```

	Df	Deviance	Resid.	Df	Resid.	Dev	Pr(>Chi)
## NULL			73		34.2		
## bloc	4	2.03	69		32.2	0.73	
## umid	2	1.75	67		30.5	0.42	
## K	4	3.83	63		26.6	0.43	
## umid:K	8	2.65	55		24.0	0.95	

```
# Wald test for interaction.
```

```
a <- c(0, attr(model.matrix(m0), "assign"))
```

```
ai <- a == max(a)
```

```
L <- t(replicate(sum(ai), rbind(coef(m1) * 0), simplify = "matrix"))
```

```
L[, ai] <- diag(sum(ai))
```

```
linearHypothesis(model = m0, # m0 is not being used here.
```

```
  hypothesis.matrix = L,
```

```
  vcov. = vcov(m1),
```

```
  coef. = coef(m1))
```

```
## Linear hypothesis test
```

```
##
```

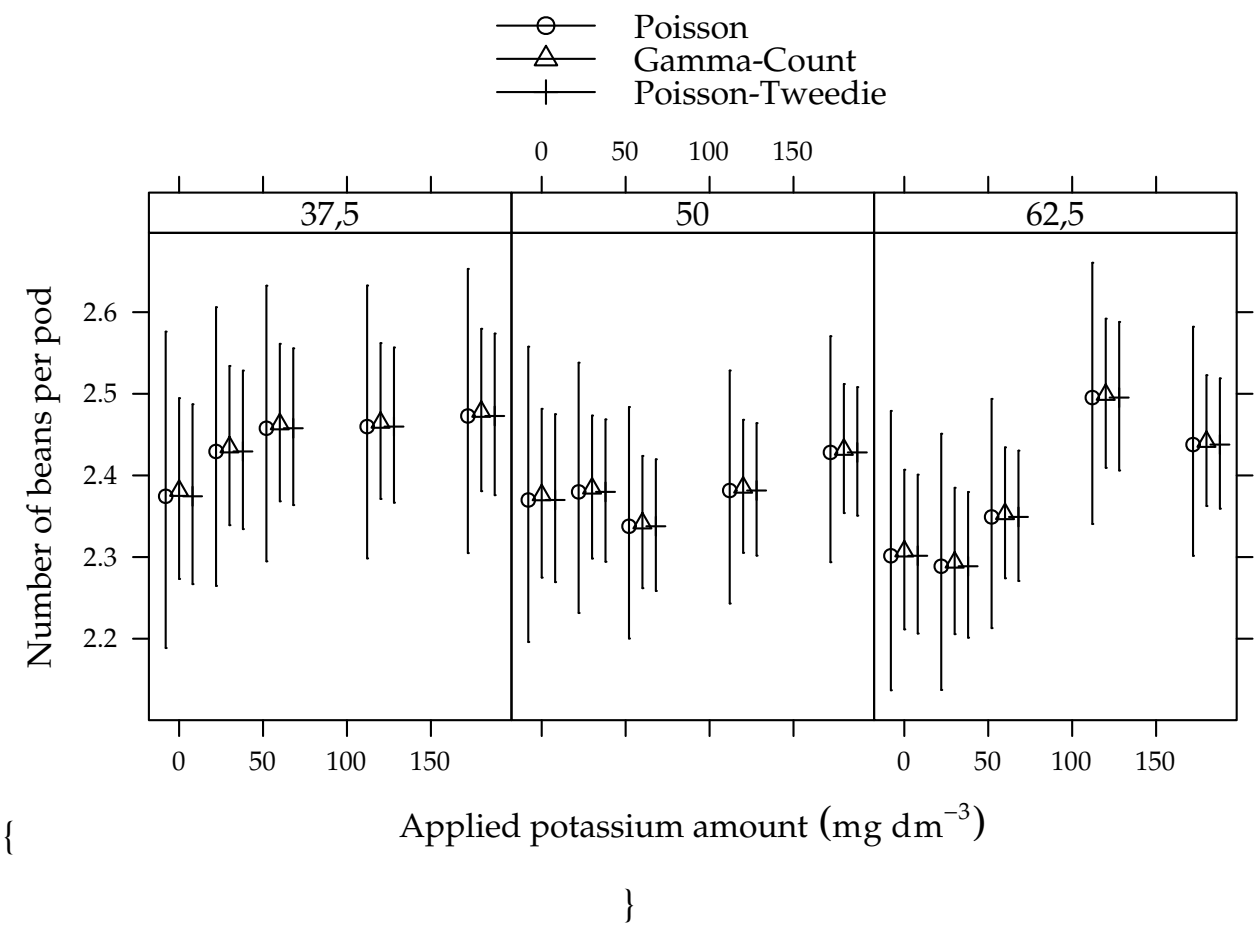
```
## Hypothesis:
## umid50:K30 = 0
## umid62,5:K30 = 0
## umid50:K60 = 0
## umid62,5:K60 = 0
## umid50:K120 = 0
## umid62,5:K120 = 0
## umid50:K180 = 0
## umid62,5:K180 = 0
##
## Model 1: restricted model
## Model 2: ngra ~ offset(log(nvag)) + bloc + umid * K
##
## Note: Coefficient covariance matrix supplied.
##
##   Res.Df Df Chisq Pr(>Chisq)
## 1      63
## 2      55  8   8.1      0.42

# Wald test for fixed effects.
anova(m2)
```

```
## Wald test for fixed effects
## Call: ngra ~ bloc + umid * K
##
##   Covariate Chi.Square Df p.value
## 1      blocII      6.03  4  0.197
## 2      umid50      1.26  2  0.534
## 3         K30      2.09  4  0.719
## 4 umid50:K30      8.20  8  0.414
```

Figure 5.2.3 shows the estimated cells means with 95% confidence intervals. All estimated means are equals along models in each cell. Gamma-count and Poisson-Tweedie have more shorter confidence intervals than Poisson.

\begin{figure}[h]



\caption{Estimated cell means based on Poisson, Gamma-Count and Poisson-Tweedie regression models. Segments are 95% individual coverage confidence intervals.} \end{figure}

We analysed the number of beans por pod in a experiment with soybeans. This variable showed underdispersion so Gamma-Count and Poisson-Tweedie perform better than Poisson.

### 5.3 Number of vehicle claims



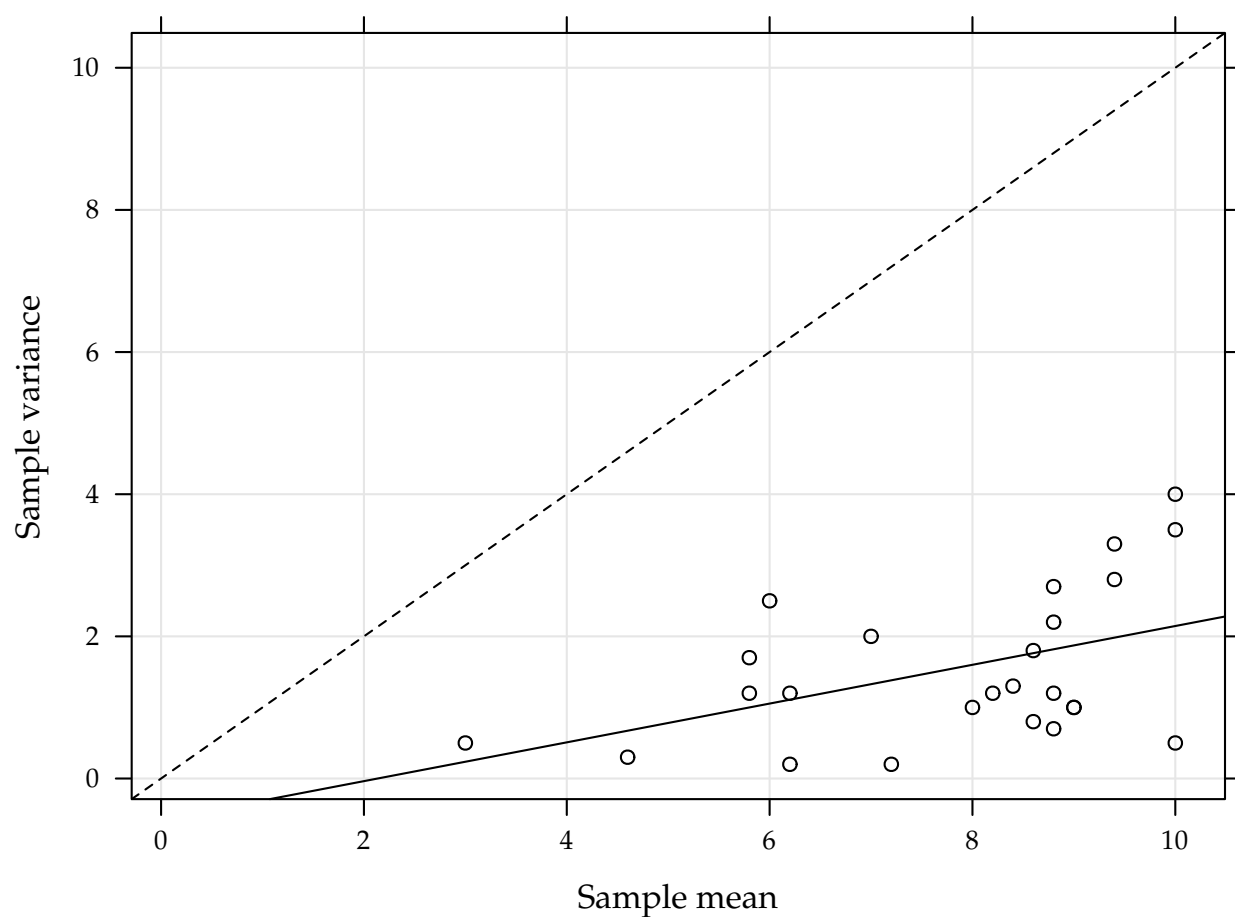
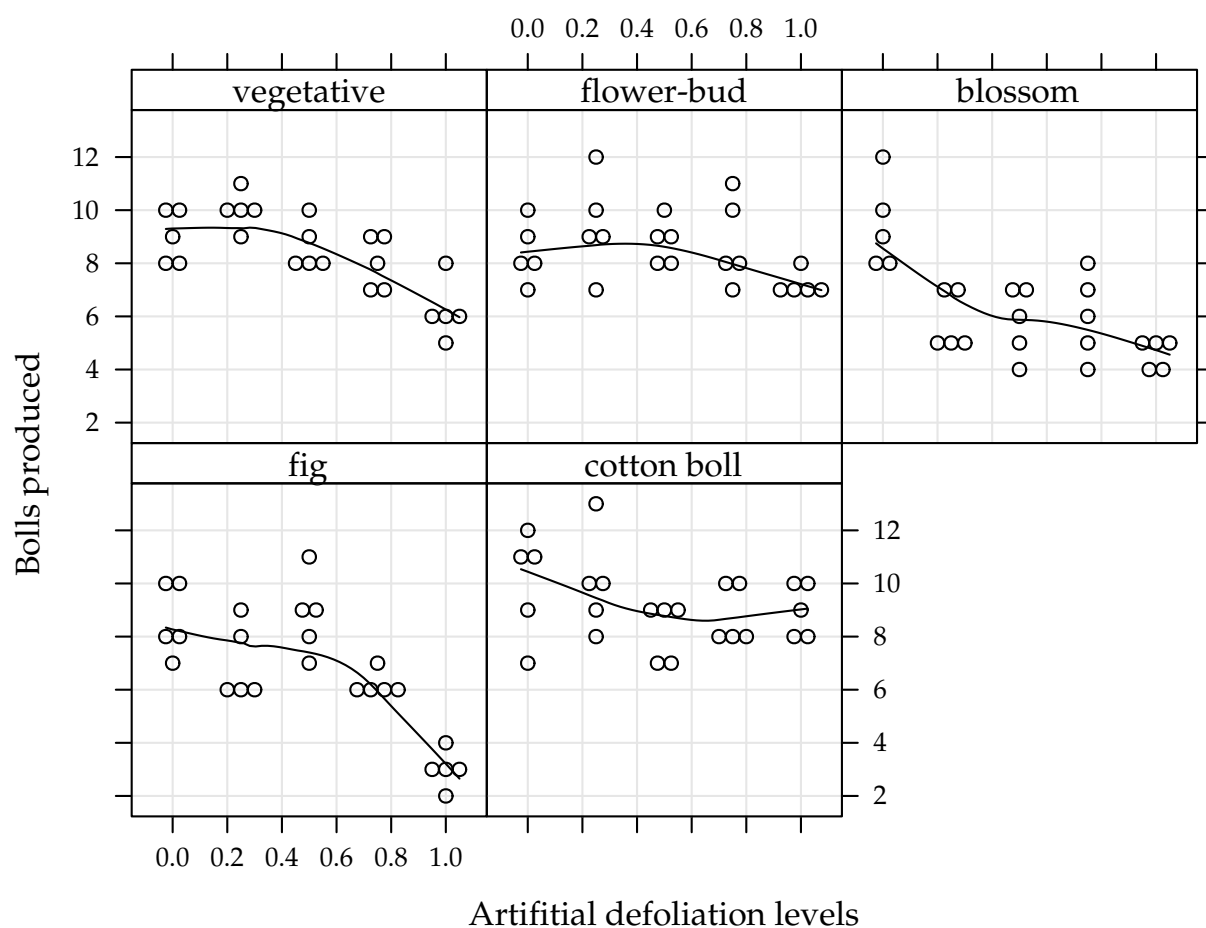


Figure 5.1: (top) Number of bolls produced for each artificial defoliation level and each growth stage. (bottom) Sample variance against the sample mean of the five replicates for each combination of defoliation level and

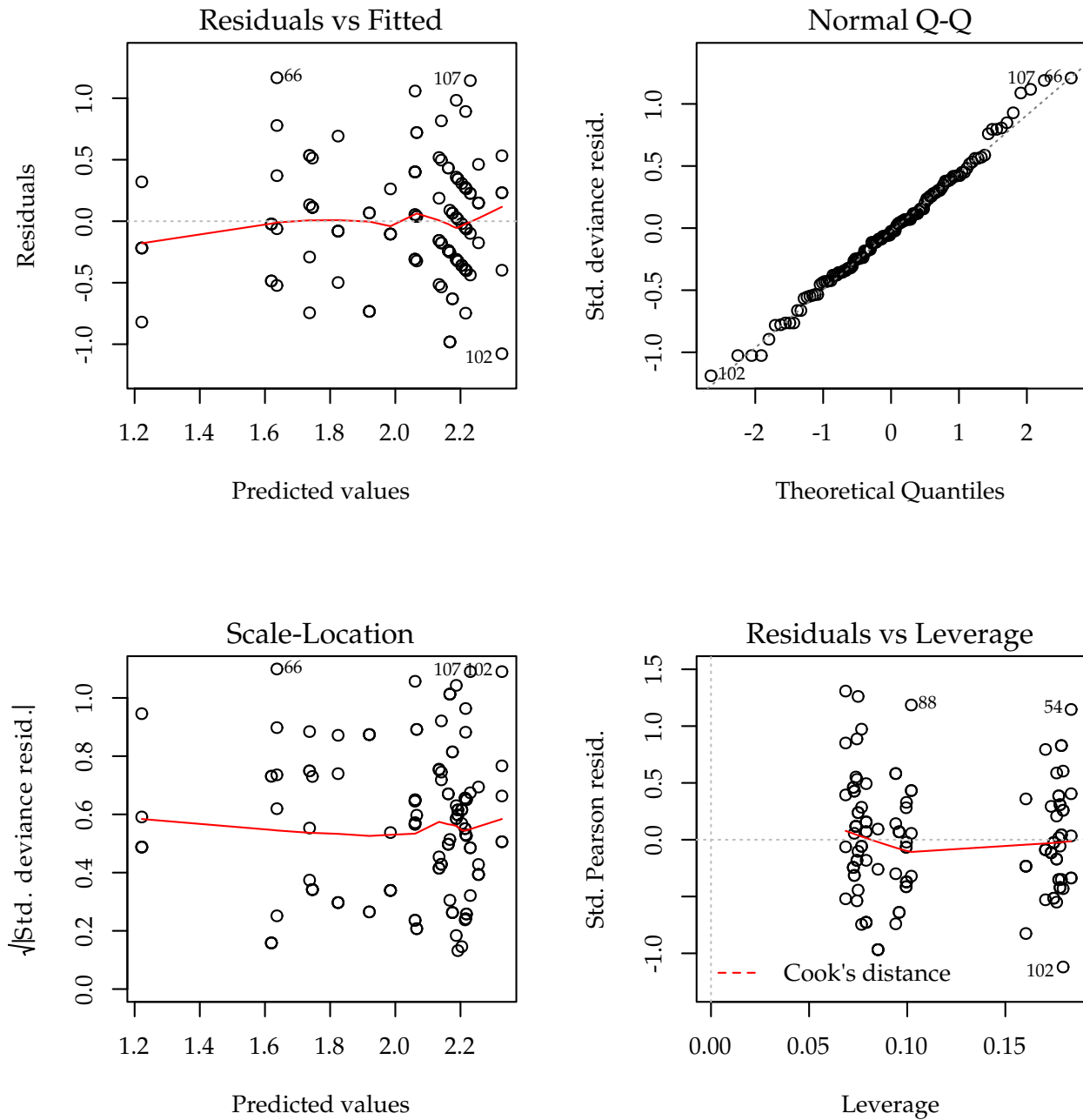


Figure 5.2: The 4 plots for checking departures of assumptions in the GLM-Poisson regression model for the number of cotton bolls.

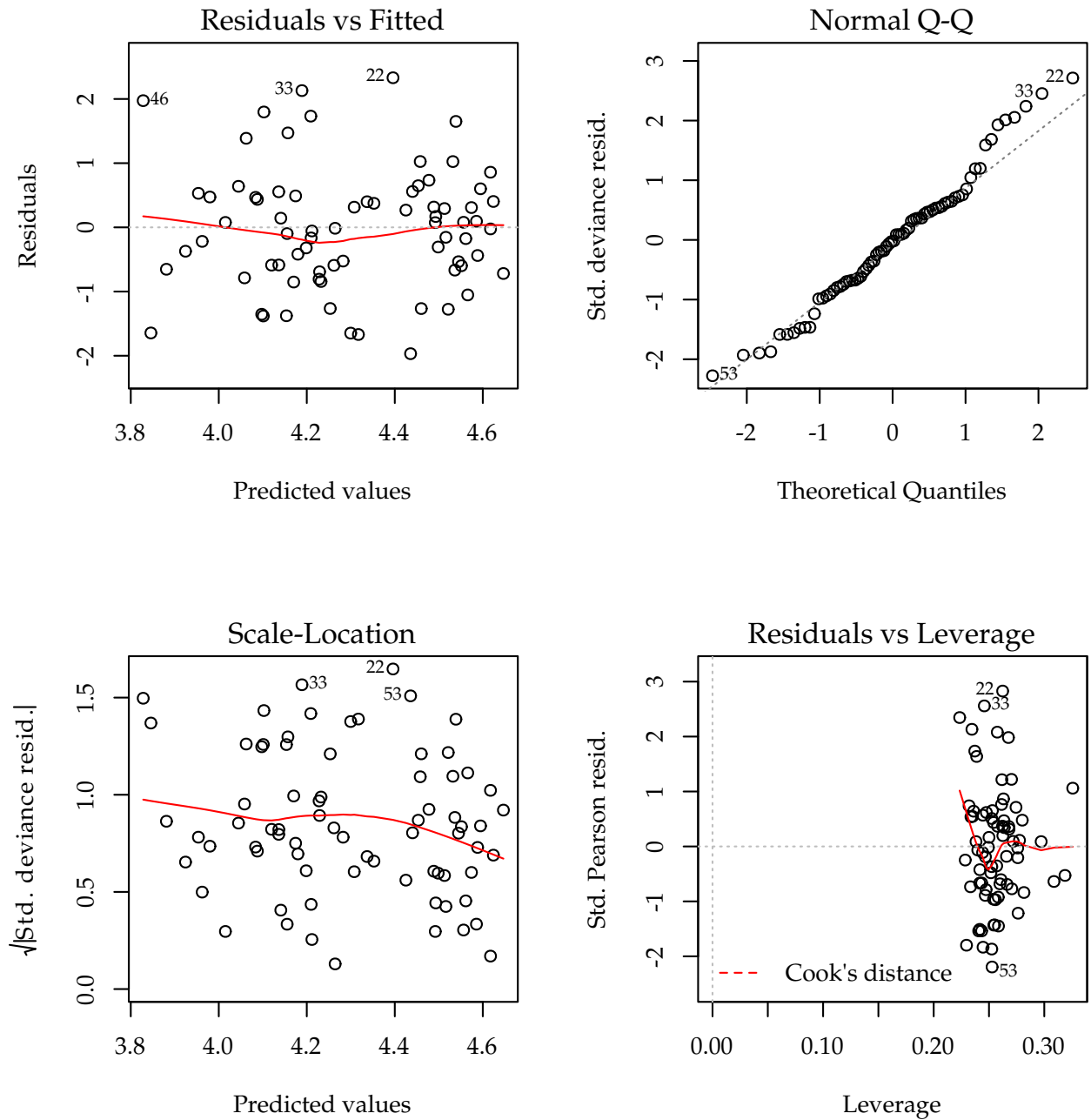


Figure 5.3: The 4 plots for checking departures of assumptions in the GLM-Poisson regression model for number of soybean pods.

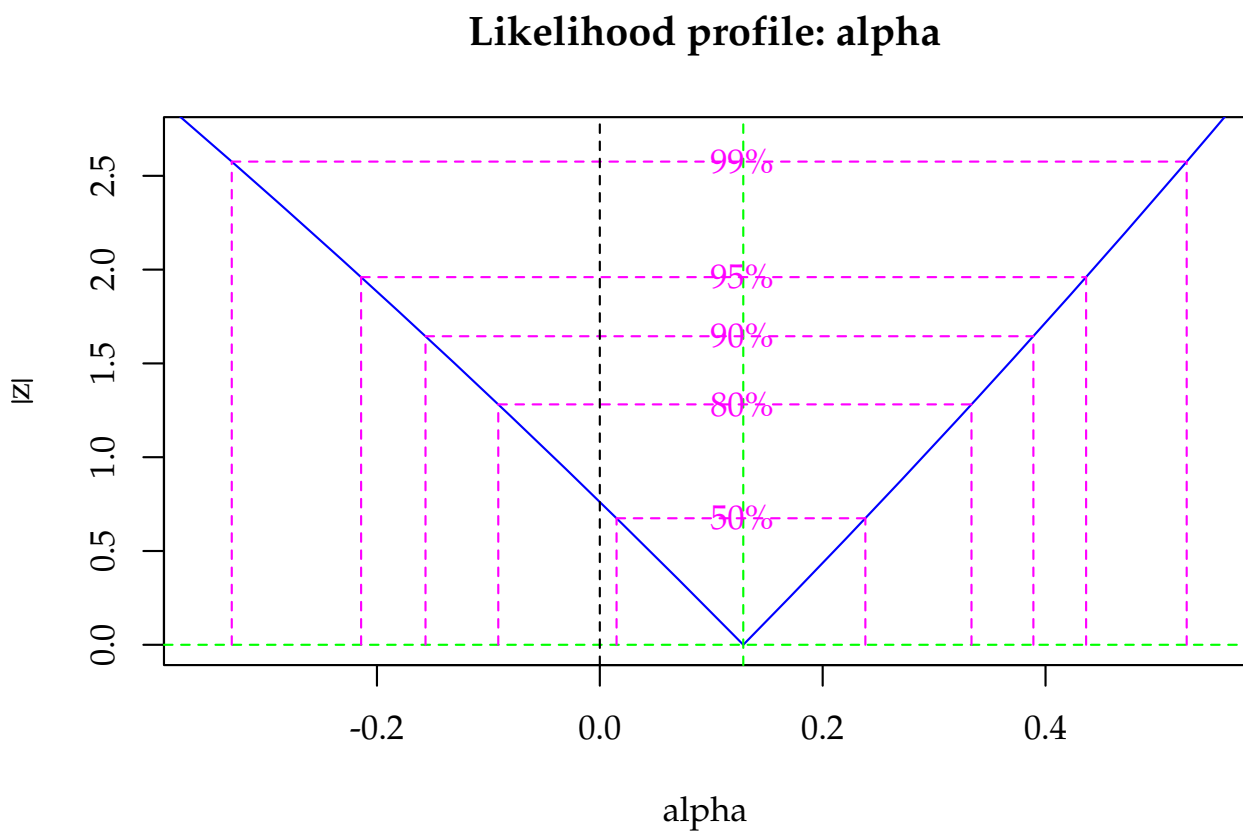


Figure 5.4: Profile log-likelihood for the Gamma-Count dispersion parameter. The confidence interval based on profile likelihood contains 0 inside as indicated by the solid vertical line.

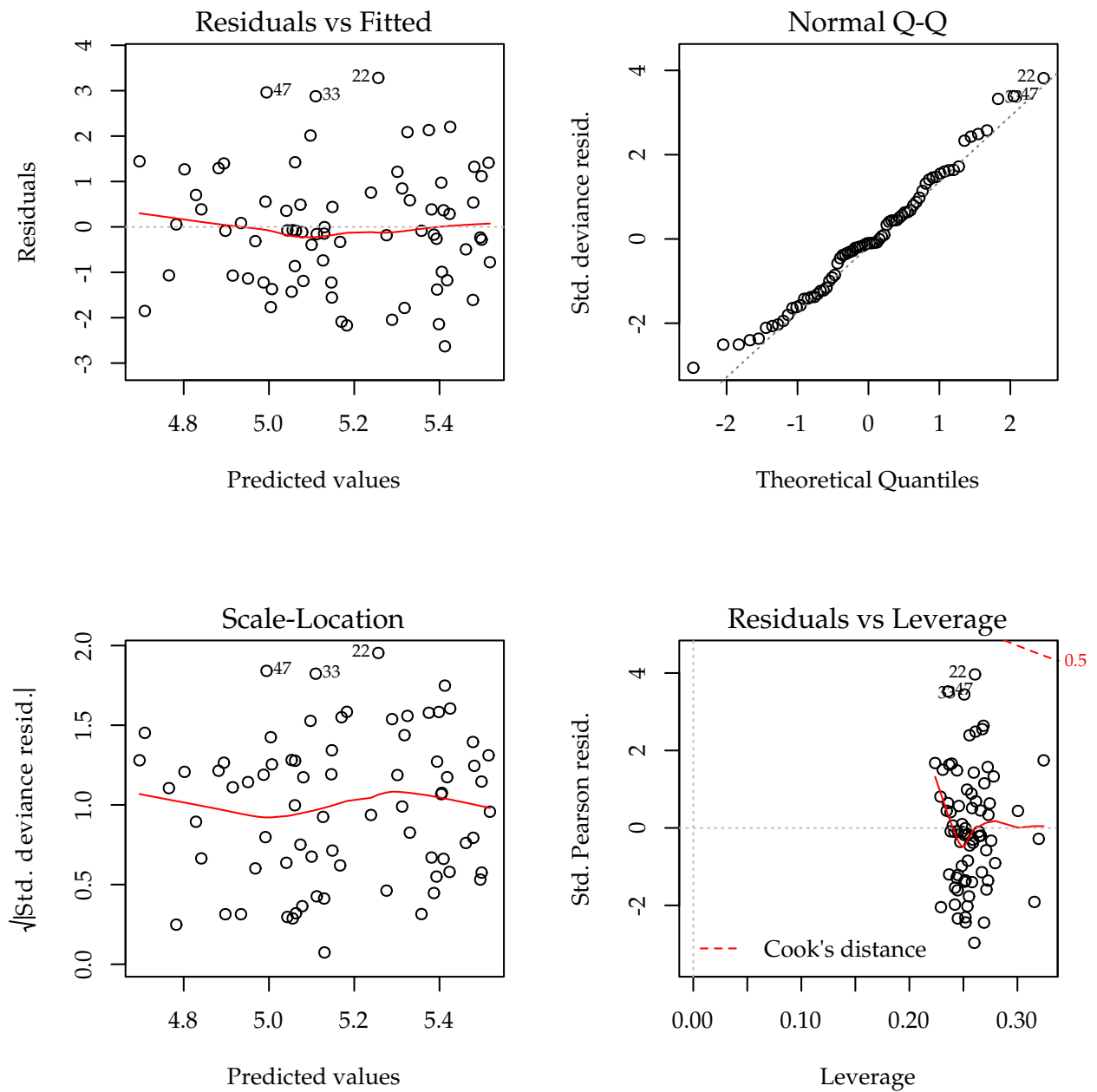


Figure 5.5: The 4 plots for checking departures of assumptions in the GLM-Poisson regression model for number of soybean pods.

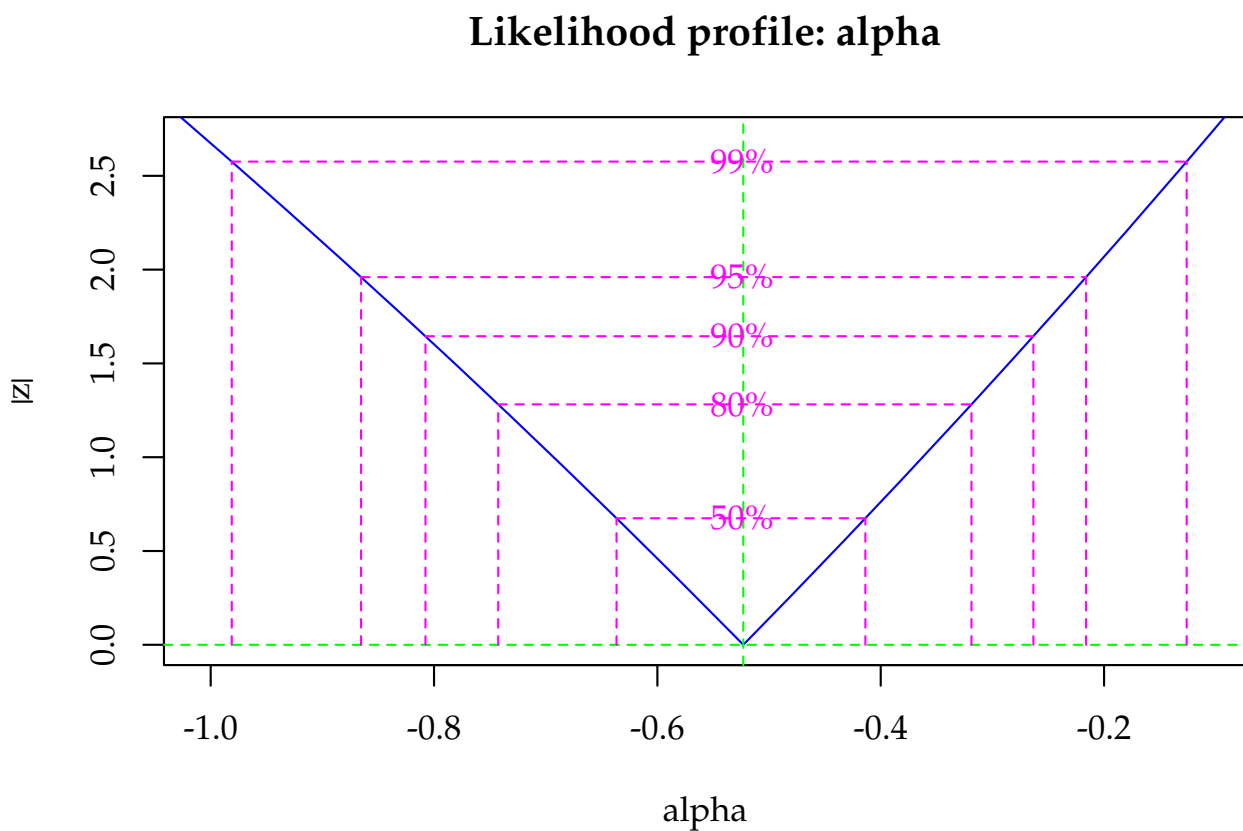


Figure 5.6: Profile log-likelihood for the Gamma-Count dispersion parameter. The confidence interval based on profile likelihood contains 0 inside as indicated by the solid vertical line.

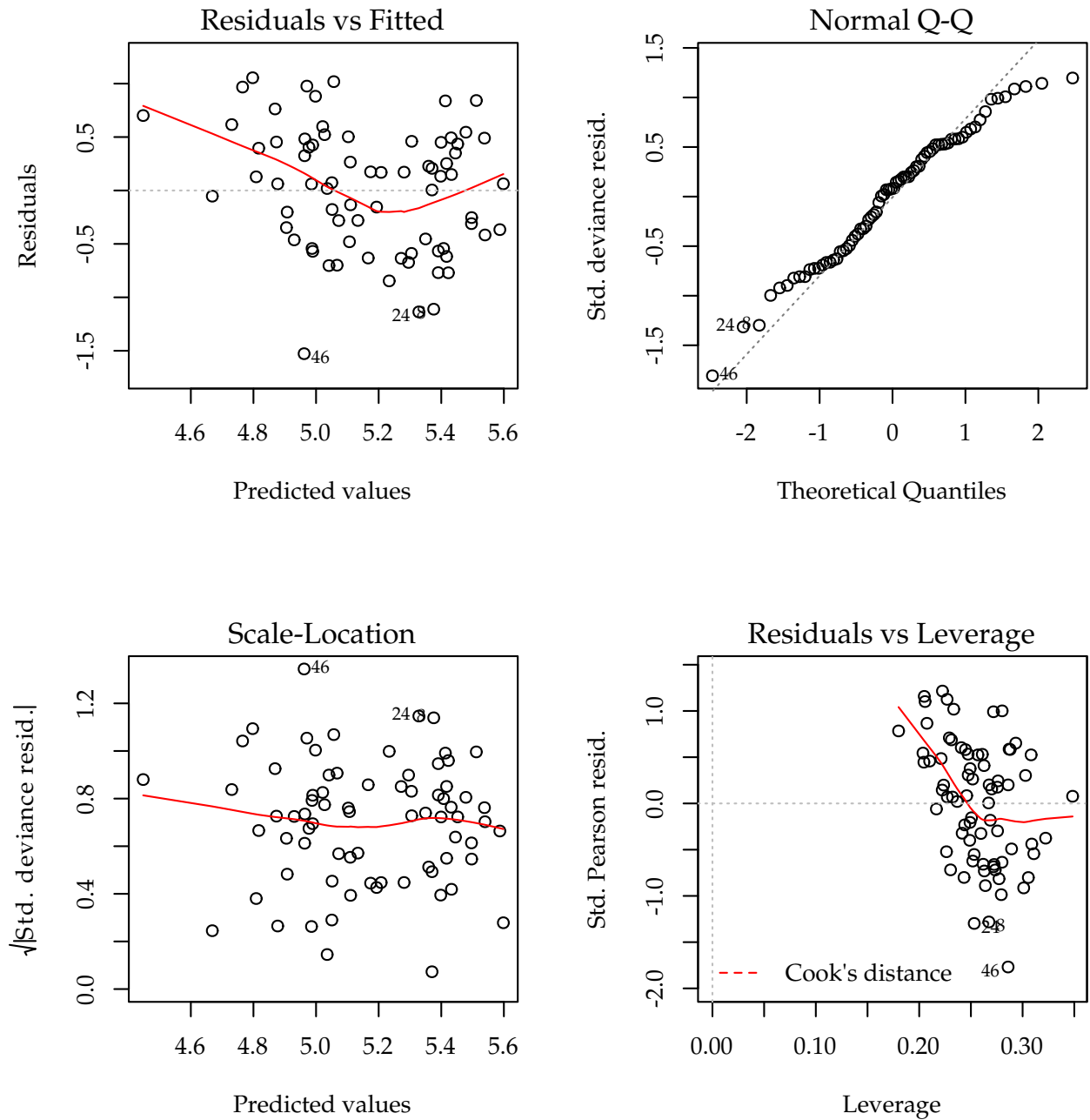


Figure 5.7: The 4 plots for checking departures of assumptions in the GLM-Poisson regression model for number of soybean pods.

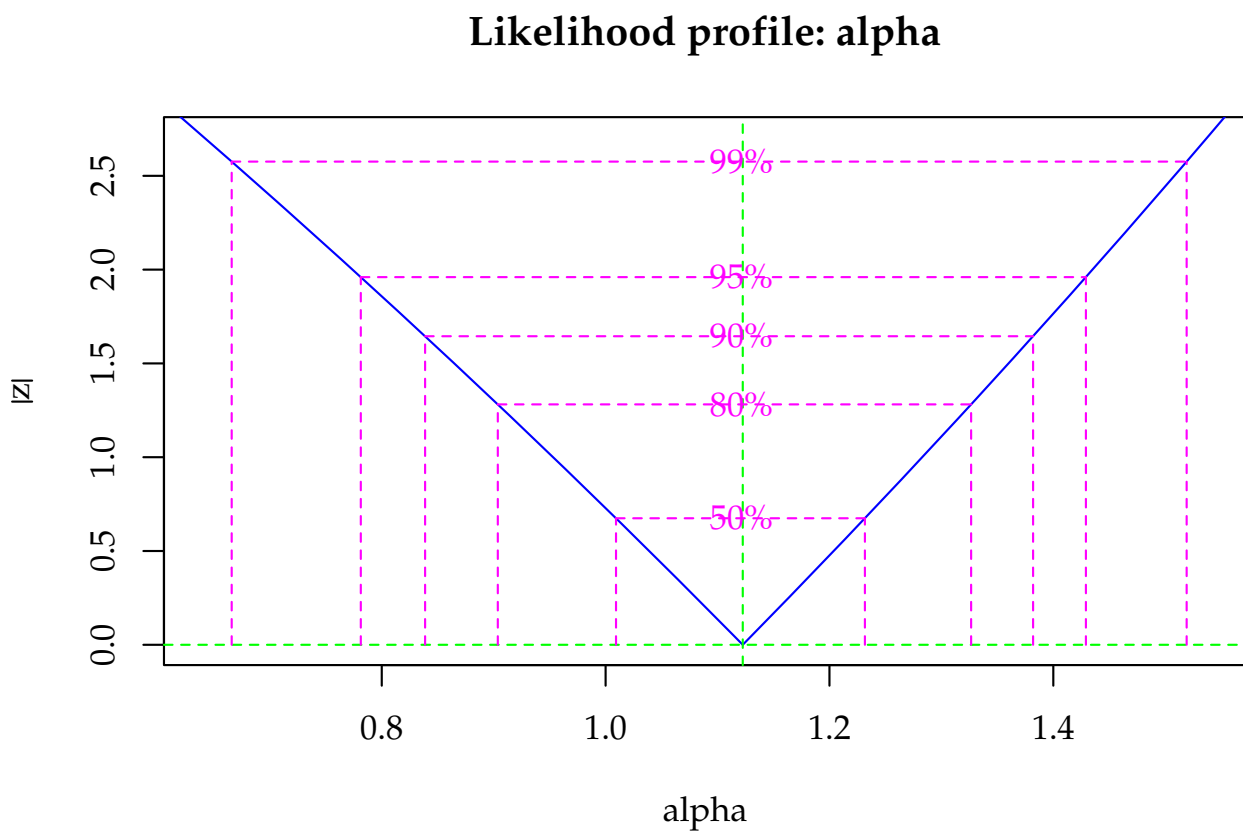


Figure 5.8: Profile log-likelihood for the Gamma-Count dispersion parameter. The confidence interval based on profile likelihood contains 0 inside as indicated by the solid vertical line.



# Bibliography

- Andersen, D. A. and Bonat, W. H. (2016). Double generalized linear compound poisson models to insurance claims data. *ArXiv*. to appear.
- Barabesi, L., Becatti, C., and Marcheselli, M. (2016). The tempered discrete Linnik distribution. *ArXiv e-prints*.
- Bolker, B. and Team, R. D. C. (2016). *bbmle: Tools for General Maximum Likelihood Estimation*. R package version 1.0.18.
- Bonat, W. H. (2016). mcglm: Multivariate covariance generalized linear models. <http://git.leg.ufpr.br/wbonat/mcglm>. R package version 0.3.0.
- Bonat, W. H. and Jørgensen, B. (2016). Multivariate covariance generalized linear models. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 65:649–675.
- Bonat, W. H., Jørgensen, B., Kokonendji, C. C. andHinde, J., and Démetrio, C. G. B. (2016). Extended poisson-tweedie: properties and regression model for count data. *Arxiv*.
- Del Castillo, J. and Pérez-Casany, M. (1998). Weighted poisson distributions for overdispersion and underdispersion situations. *Annals of the Institute of Statistical Mathematics*, 50(3):567–585.
- Dunn, J. (2012). *compoisson: Conway-Maxwell-Poisson Distribution*. R package version 0.3.
- El-Shaarawi, A. H., Zhu, R., and Joe, H. (2011). Modelling species abundance using the Poisson-Tweedie family. *Environmetrics*, 22(2):152–164.
- Esnaola, M., Puig, P., Gonzalez, D., Castelo, R., and Gonzalez, J. R. (2013). A flexible count data model to fit the wide diversity of expression profiles arising from extensively replicated rna-seq experiments. *BMC Bioinformatics*, 14(1):254–276.

- Jørgensen, B. (1987). Exponential dispersion models. *Journal of the Royal Statistical Society. Series B (Methodological)*, 49(2):127–162.
- Jørgensen, B. (1997). *The Theory of Dispersion Models*. Chapman & Hall, London.
- Jørgensen, B. and Knudsen, S. J. (2004). Parameter orthogonality and bias adjustment for estimating functions. *Scandinavian Journal of Statistics*, 31(1):93–114.
- Jørgensen, B. and Kokonendji, C. (2015). Discrete dispersion models and their Tweedie asymptotics. *AStA Advances in Statistical Analysis*. in press.
- Kokonendji, C. C., Dossou-Gbété, S., and Demétrio, C. G. B. (2004). Some discrete exponential dispersion models: Poisson-Tweedie and Hinde-Demétrio classes. *Statistics and Operations Research Transactions*, 28(2):201–214.
- Nelder, J. A. and Mead, R. (1965). A Simplex method for function minimization. *The Computer Journal*, 7:308–313.
- Nelder, J. A. and Wedderburn, R. W. M. (1972). Generalized linear models. *Journal of the Royal Statistical Society. Series A*, 135(3):370–384.
- R Core Team (2015). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. ISBN 3-900051-07-0.
- Sellers, K. F., Borle, S., and Shmueli, G. (2012). The com-poisson model for count data: a survey of methods and applications. *Applied Stochastic Models in Business and Industry*, 28(2):104–116.
- Sellers, K. F. and Shmueli, G. (2010). A flexible regression model for count data. *Ann. Appl. Stat.*, 4(2):943–961.
- Serafim, M. E., Ono, F. B., Zeviani, W. M., Novelino, J. O., and Silva, J. V. (2012). Umidade do solo e doses de potássio na cultura da soja. *Revista Ciência Agronômica*, 43(2):222–227.
- Silva, A. M., Degrande, P. E., Suekane, R., Fernandes, M. G., and Zeviani, W. M. (2012). Impacto de diferentes níveis de desfolha artificial nos estádios fenológicos do algodoeiro. *Revista de Ciências Agrárias*, 35(1):163–172.
- Silvey, S. (1975). *Statistical Inference*. Chapman and Hall/CRC Monographs on Statistics and Applied Probability. Taylor and Francis.
- Wedderburn, R. W. M. (1974). Quasi-likelihood functions, generalized linear models, and the gauss—newton method. *Biometrika*, 61(3):439–447.

- Winkelmann, R. (1995). Duration dependence and dispersion in count-data models. *Journal of Business & Economic Statistics*, 13(4):467–474.
- Winkelmann, R. (2003). *Econometric Analysis of Count Data*. Springer.
- Zeviani, W. M., Junior, E. E. R., and Taconeli, C. A. (2016). *MRDCr: Modelos de Regressão para Dados de Contagem*. R package version 0.0-2.
- Zeviani, W. M., Ribeiro Jr, P. J., Bonat, W. H., Shimakura, S. E., and Muniz, J. A. (2014). The gamma-count distribution in the analysis of experimental underdispersed data. *Journal of Applied Statistics*, 41(12):2616–2626.