



HAL
open science

Closed form Maximum Likelihood Estimator for Generalized Linear Models in the case of categorical explanatory variables: Application to insurance loss modelling

Alexandre Brouste, Christophe Dutang, Tom Rohmer

► To cite this version:

Alexandre Brouste, Christophe Dutang, Tom Rohmer. Closed form Maximum Likelihood Estimator for Generalized Linear Models in the case of categorical explanatory variables: Application to insurance loss modelling. 2018. hal-01781504v1

HAL Id: hal-01781504

<https://hal.science/hal-01781504v1>

Preprint submitted on 30 Apr 2018 (v1), last revised 25 Aug 2019 (v3)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Closed form Maximum Likelihood Estimator for Generalized Linear Models in the case of categorical explanatory variables: Application to insurance loss modelling

Alexandre Brouste^{*}, Christophe Dutang^{**} & Tom Rohmer^{*,1}

^{*}Institut du Risque de l'Assurance & Laboratoire Manceau de Mathématiques
Le Mans Université, Avenue Olivier Messiaen, 72085 LE MANS, France

^{**}CEREMADE, CNRS, Univ. Paris-Dauphine,
Place du Maréchal de Lattre de Tassigny, 75016 PARIS, France

Abstract

Generalized Linear Models with categorical explanatory variables are considered and parameters of the model are estimated with an original exact maximum likelihood method. The existence of a sequence of maximum likelihood estimators is discussed and considerations on possible link functions are proposed. A focus is then given on two particular positive distributions: the Pareto 1 distribution and the shifted log-normal distributions. Finally, the approach is illustrated on an actuarial dataset to model insurance losses.

Keywords: Generalized linear model, heavy-tailed distributions, insurance claim modeling

1. Introduction

The assumption of identical distributions for random variables in an observation sample is relaxed for regression models by considering explanatory variables. Generalized Linear Models (GLMs) were introduced by Nelder and Wedderburn (1972) and popularized in McCullagh and Nelder (1989). GLMs rely on probability distribution functions of exponential type for the response variable which include most of the light and medium tailed distributions (such as normal, gamma or inverse Gaussian). Asymptotic properties of sequences of maximum likelihood estimators for GLMs were studied by Fahrmeir and Kaufmann (1985).

Regression models for heavy-tailed distributions have been mainly studied through the point-of-view of extreme value analysis, see Beirlant et al. (2004) for a review. A regression model for the generalized Pareto distribution (GPD) where the scale parameter depends on covariates are described in Davison and Smith (1990) with a least square estimation procedure and a model checking method. Beirlant et al. (1998) propose a Burr model by regressing the shape parameter with an exponential link on explanatory variables. In the aforementioned article, a simulation study with one explanatory variable is detailed as well

¹Corresponding author: tom.rohmer@univ-lemans.fr

as an application to fire insurance. Residual plots and asymptotic convergence towards the normal distribution are also discussed. Similarly, Ozkok et al. (2012) propose a regression model for Burr distribution where the scale parameter depends on covariates.

An estimation of the extremal tail index (used in generalized extreme value (GEV) distributions and GPD) by considering a class of distribution function with an exponential link on explanatory variables is also described in Beirlant and Goegebeur (2003). Using generalized residuals of explanatory variable makes possible the estimation of the tail index. Still by the extreme value theory approach, Chavez-Demoulin et al. (2015) and Hambuckers et al. (2016) both propose a semi-parametric regression model for GEV and GPD where the explanatory variables are time or known factor levels. They assume that all parameters depend on covariates and also use exponentially distributed residuals.

Few papers seem to study this topic outside the extreme value theory framework. Mainly, Rigby and Stasinopoulos (2005) propose a general regression framework where all parameters are modeled by explanatory variable and the distribution is not restricted to exponential family. The only restriction that the authors impose is the twice differentiability of the density function w.r.t. parameters. However, there is no clear convergence result of the proposed estimators. Among the proposed distributions, authors use 1-parameter Pareto, log-logistic (a special case of the Pareto 3 distribution) and GEV distributions.

In this paper, we propose closed-form estimators for generalized linear models in the case of categorical variables. The expression is valid for any distribution belonging to the one-parameter exponential family. Then, the paper will continue by the application of such formulas not on classical distributions, but on distributions such as the log-transformed variable has a distribution in the exponential family. Therefore, the choice of probability distributions of this paper is led by two aspects: distributions with positive values and distributions as the log-transformed variable belongs to the exponential family. The considered distributions have heavier tails than the exponential distribution. We choose to study two distributions: the Pareto 1 distribution and the shifted lognormal distribution with fixed threshold parameters. We could have considered log-logistic and GEV distributions being also appropriate in many situations but these distributions do not belong to the exponential family.

Applications of this distribution can be found in various disciplines such as finance, insurance, reliability theory, etc. Here, we are interested with an application to large insurance loss modeling. Indeed, pricing non-life insurance relies on estimating the claim frequency and the claim severity distributions. The former is generally estimated by a regression model such as Poisson or zero-inflated models. However for modeling claim severity, we commonly split the claim dataset between attritional and atypical claims. A threshold μ is chosen either from the extreme value theory or by expert judgments. A classical GLM such as gamma or inverse-Gaussian is fitted on attritional claim amounts below μ , see e.g. Ohlsson and Johansson (2010). Atypical claim amounts above μ are not necessarily modeled at all. An empirical rule of the insurance pricing is used to mutualize atypical claims over the portfolio, i.e. the aggregate sum of atypical claims is shared equally among all policies. We aim at providing a regression model for those claims above μ in order to refine this empirical rule.

The threshold μ can also be interpreted in another insurance context. Generally in non-life insurance, contracts are underwritten with a deductible. This has two consequences: the policyholder will retain the risk of claims below the deductible; and the insurer will only

know and be interested in claims above the deductible. In the numerical section, we consider only the example of large claim modeling.

The paper is organized as follows. In Section 2, we present the Generalized Linear Models. Section 3 provide exact formulas for maximum likelihood estimators in the case of categorical explanatory variables. Section 4 is dedicated to the Pareto 1 GLM, while Section 5 is dedicated to the shifted lognormal GLM. Finally, an application to an actuarial dataset is carried out in Section 6, before Section 7 concludes.

2. Preliminaries on Generalized Linear Models

In this section, we consider the estimation problem in GLMs. We consider deterministic exogenous variables $\mathbf{y}_1, \dots, \mathbf{y}_n$, with $\mathbf{y}_i = (y_i^{(1)}, \dots, y_i^{(p)}) \in \mathbb{R}^p$ for $i = 1, \dots, n$. In the following, for the sake of clarity, bold notations are reserved for vector of \mathbb{R}^p and bold notations with an underline are reserved for vector of \mathbb{R}^n .

The index $i \in I = \{1, \dots, n\}$ is reserved for the observations, while the indexes j, k, l are used for the explanatory variables.

In this setting, the sample $\underline{\mathbf{X}} = (X_1, \dots, X_n)$ is composed of real-valued independent random variables; each one belongs to a family of probability measures of one-parameter exponential type with respective parameters $\lambda_1, \dots, \lambda_n$ valued in $\Lambda \subset \mathbb{R}$.

Precisely, the likelihood L associated to the statistical experiment generated by $X_i, i \in I$ verifies

$$\log L(\boldsymbol{\vartheta} | x_i) = \frac{\lambda_i(\boldsymbol{\vartheta})x_i - b(\lambda_i(\boldsymbol{\vartheta}))}{a(\phi)} + c(x_i, \phi), \quad x_i \in \mathbb{X} \subset \mathbb{R}, \quad (1)$$

and $-\infty$ if $x_i \notin \mathbb{X}$, where $a : \mathbb{R} \rightarrow \mathbb{R}$, $b : \Lambda \rightarrow \mathbb{R}$ and $c : \mathbb{X} \times \mathbb{R} \rightarrow \mathbb{R}$ are known real-valued measurable functions and ϕ is the dispersion parameter, e.g. McCullagh and Nelder (1989, Section 2.2).

In Equation (1), the parameters $\lambda_1, \dots, \lambda_n$ depend on a finite-dimensional parameter $\boldsymbol{\vartheta} \in \Theta \subset \mathbb{R}^p$. Direct computations lead to

$$b'(\lambda_i(\boldsymbol{\vartheta})) = \mathbf{E}_{\boldsymbol{\vartheta}} X_i \quad \text{and} \quad b''(\lambda_i(\boldsymbol{\vartheta}))a(\phi) = \mathbf{Var}_{\boldsymbol{\vartheta}} X_i. \quad (2)$$

Using a twice continuously differentiable and bijective function g from $b'(\Lambda)$ to \mathbb{R} , the GLM are defined by assuming the following relation between the expectation and the predictor

$$g(b'(\lambda_i(\boldsymbol{\vartheta}))) = \langle \mathbf{y}_i, \boldsymbol{\vartheta} \rangle = \eta_i, \quad \text{for all } \boldsymbol{\vartheta} \in \Theta,$$

where η_i are the linear predictors and $\langle \cdot, \cdot \rangle$ denotes the scalar product. In other words, the bijective function $\ell = (b')^{-1} \circ g^{-1}$ is setted; then we have

$$\lambda_i(\boldsymbol{\vartheta}) = \ell(\eta_i). \quad (3)$$

We summarize with the following relations

$$Y \times \Theta \xrightarrow{\langle \cdot, \cdot \rangle} D \xrightleftharpoons[\ell]{\ell^{-1}} \Lambda$$

where D is the space of linear predictor and Y the possible set of value of \mathbf{y}_i for $i \in I$. Here ℓ is chosen and, consecutively Θ , Λ and D must be set.

The parameter $\boldsymbol{\vartheta} \in \Theta \subset \mathbb{R}^p$ is to be estimated and g is called the link function in the regression framework. We talk of canonical link function, when ℓ is the identity function.

Let us compute the log-likelihood of $\underline{\mathbf{x}} = (x_1, \dots, x_n)$:

$$\log L(\boldsymbol{\vartheta} | \underline{\mathbf{x}}) = \sum_{i=1}^n \frac{1}{a(\phi)} (x_i \ell(\eta_i) - b(\ell(\eta_i))) + \sum_{i=1}^n c(x_i, \phi), \quad (4)$$

with b , h and ℓ being respectively defined in (1) and (3). Here, the vector of the parameters $\boldsymbol{\vartheta}$ is unknown. If the model is identifiable, it can be shown that the sequence of maximum likelihood estimators $(\widehat{\boldsymbol{\vartheta}}_n)_{n \geq 1}$ defined by $\widehat{\boldsymbol{\vartheta}}_n = \arg \max_{\boldsymbol{\vartheta} \in \Theta} L(\boldsymbol{\vartheta} | \underline{\mathbf{x}})$ asymptotically exists and is consistent (for example Fahrmeir and Kaufmann, 1985, Theorem 2, 4).

The maximum likelihood estimator (MLE) $\widehat{\boldsymbol{\vartheta}}_n$, if it exists, is the solution of the non linear system

$$S_j(\boldsymbol{\vartheta}) = 0, \quad j = 1, \dots, p, \quad (5)$$

with $S_j(\boldsymbol{\vartheta})$ are the component of the Score vector defined by

$$S_j(\boldsymbol{\vartheta}) = \frac{1}{a(\phi)} \sum_{i=1}^n y_i^{(j)} \ell'(\eta_i) (x_i - b'(\ell(\eta_i))).$$

Note that the MLE $\widehat{\boldsymbol{\vartheta}}_n$ does not depend on the value of the dispersion parameter ϕ . Indeed, the dispersion parameter is estimated in a second step using the sum of square residuals or the log-likelihood, see e.g. (McCullagh and Nelder, 1989, Chap. 9).

In a general setting, the system (5) does not have a closed-form solution and generalized linear models are generally fitted using a Newton-type method, such that an iteratively re-weighted least square (IWLS) algorithm also refereed to Fisher Scoring algorithm, see e.g. McCullagh and Nelder (1989).

Nevertheless, the choice of the initial value for the Newton-type method sometimes is problematic. A misspecification of this initial value can lead to divergence of the algorithm. Moreover, for a small data set (small n) or large number of explanatory variables, the (non-asymptotic) existence of the MLE is not guaranteed.

In the case of categorical explanatory variables described later on, the non-asymptotic existence of the MLE depends on the conditional distribution and the chosen link function (see Examples 1, 2 and 3 on Section 4.2).

3. A closed form MLE for categorical explanatory variables

In any regression model, categorical or nominal explanatory variables have to be coded since their value is a name or a category. When the possible values are unordered, it is common to use a binary incidence matrix or dummy variables where each row has a single unity in the column of the class to which it belongs. In the case of ordered values, a contrast matrix has to be used, see e.g. Venables and Ripley (2002).

3.1. A single explanatory variable

Let us first consider the case of a single categorical explanatory variable. That is $p = 2$ and for all $i = 1, \dots, n$, $y_i^{(1)} = 1$ is the intercept and $y_i^{(2)}$ takes values in a set of d modalities $\{v_1, \dots, v_d\}$. We define the incidence matrix $(y_i^{(2),j})_{i,j}$ where $y_i^{(2),j} = \mathbf{1}_{y_i^{(2)}=v_j}$ is the binary dummy of the j th category for $i \in I = \{1, \dots, n\}$ and $j \in J = \{1, \dots, d\}$. From this incidence matrix, we compute the number of appearance $m_j > 0$ of the j th category and $\bar{x}_n^{(j)}$ the mean value of \underline{x} taking over the j th category by

$$m_j = \sum_{i=1}^n y_i^{(2),j}, \quad j \in J \quad \text{and} \quad \bar{x}_n^{(j)} = \frac{1}{m_j} \sum_{i=1}^n x_i y_i^{(2),j}, \quad j \in J.$$

By construction, this incidence matrix has rows that sum to 1. Therefore if we use the combination of the incidence matrix with a 1-column for the intercept $(y_i^{(1)}, y_i^{(2),j})_{i,j}$: a redundancy appears. We must choose either to use *no intercept*, to *drop one column* for a particular modality of $y_i^{(2)}$, or to use a *zero-sum condition* on the parameters. We investigate below these three options in a single framework.

Consider the following GLM for the explanatory variables $y_i^{(1)}, y_i^{(2),1}, \dots, y_i^{(2),d}$

$$g(\mathbf{E}X_i) = \vartheta_{(1)} + \sum_{j=1}^d y_i^{(2),j} \vartheta_{(2),j}, \quad i = 1, \dots, n, \quad (6)$$

where $\boldsymbol{\vartheta} = (\vartheta_{(1)}, \vartheta_{(2),1}, \dots, \vartheta_{(2),d})$ is the unknown vector parameters. The model being not identifiable, we impose exactly one linear equation on $\boldsymbol{\vartheta}$

$$\langle \mathbf{R}, \boldsymbol{\vartheta} \rangle = 0, \quad (7)$$

with $\mathbf{R} = (r_{(1)}, r_{(2),1}, \dots, r_{(2),d})$ any real vector of size $d + 1$. A theorem and a corollary are given below and corresponding proofs are postponed to Appendix A.

Theorem 3.1. *Suppose that for all $i \in \{1, \dots, n\}$, X_i takes values in $b'(\Lambda)$. If the vector \mathbf{R} is such that $\sum_{j=1}^d r_{(2),j} - r_{(1)} \neq 0$, then there exists a unique, consistent and explicit MLE $\hat{\boldsymbol{\vartheta}}_n = (\hat{\boldsymbol{\vartheta}}_{n,(1)}, \hat{\boldsymbol{\vartheta}}_{n,(2),1}, \dots, \hat{\boldsymbol{\vartheta}}_{n,(2),d})$ of $\boldsymbol{\vartheta}$ given by*

$$\hat{\boldsymbol{\vartheta}}_{n,(1)} = \frac{\sum_{j=1}^d r_{(2),j} g(\bar{X}_n^{(j)})}{\sum_{j=1}^d r_{(2),j} - r_{(1)}}, \quad \hat{\boldsymbol{\vartheta}}_{n,(2),j} = g(\bar{X}_n^{(j)}) - \frac{\sum_{j=1}^d r_{(2),j} g(\bar{X}_n^{(j)})}{\sum_{j=1}^d r_{(2),j} - r_{(1)}}, \quad j = 1, \dots, d. \quad (8)$$

Note that if $\bar{X}_n^{(j)}$ does not belong to $b'(\Lambda)$, $g(\bar{X}_n^{(j)})$ and hence $\hat{\boldsymbol{\vartheta}}_{n,(l),j}$ are not defined.

We give below the three most common examples of linear constraint, some details of these calculus are given in Appendix A.

Example 3.1. *No-intercept model*

The no-intercept model is obtained with $\mathbf{R} = (1, 0, \dots, 0)$ leading to $\vartheta_{(1)} = 0$. Therefore the unique, consistent and explicit MLE $\hat{\boldsymbol{\vartheta}}_n$ of $\boldsymbol{\vartheta}$ is

$$\hat{\boldsymbol{\vartheta}}_{n,(1)} = 0, \quad \hat{\boldsymbol{\vartheta}}_{n,(2),j} = g(\bar{X}_n^{(j)}), \quad j \in J. \quad (9)$$

Example 3.2. *Model without first modality*

The model without first modality is obtained with $R = (0, 1, \dots, 0)$ leading to $\vartheta_{(2),1} = 0$. Therefore, the unique, consistent and explicit MLE $\widehat{\boldsymbol{\vartheta}}_n$ of $\boldsymbol{\vartheta}$ is

$$\widehat{\vartheta}_{n,(1)} = g\left(\overline{X}_n^{(1)}\right), \quad \widehat{\vartheta}_{n,(2),1} = 0, \quad \widehat{\vartheta}_{n,(2),j} = g\left(\overline{X}_n^{(j)}\right) - \widehat{\vartheta}_{n,1}, \quad j \in J \setminus \{1\}.$$

Example 3.3. *Zero-sum condition* The zero-sum model is obtained with $R = (0, 1, \dots, 1)$ leading to $\sum_{j=1}^d \vartheta_{(2),j} = 0$. Therefore, the unique, consistent and explicit MLE $\widehat{\boldsymbol{\vartheta}}_n$ of $\boldsymbol{\vartheta}$ is

$$\widehat{\vartheta}_{n,(1)} = \frac{1}{d} \sum_{j=1}^d g\left(\overline{X}_n^{(j)}\right), \quad \widehat{\vartheta}_{n,(2),j} = g\left(\overline{X}_n^{(j)}\right) - \widehat{\vartheta}_{n,1}, \quad j \in J.$$

Remark 3.1. In Theorem 3.1, it is worth noting that the value of $\widehat{\boldsymbol{\vartheta}}_n$ does not depend on the distribution of the X_i .

Remark 3.2. The three different parametrizations (Examples 3.1, 3.2 and 3.3) depends on the type of application and on the modeler choice. In statistical software, there is a default choice: for instance in the statistical software R, the model without the first modality is the default parametrization (see functions `lm()`, `glm()` by R Core Team (2017)). The first option without intercept may be justified when no group can be chosen as the reference group.

Remark 3.3. When g is the identity function, the third option with a zero-sum condition can be interpreted as a generalized analysis of variance (ANOVA) for Z_i with respect to groups defined by the explanatory variable $\mathbf{y}^{(2)}$. Even for non-Gaussian random variables, some applications may justify this option.

Theorem 3.1 has two interesting corollaries which give some clues on the choice of the link function g . This corollary tempers the importance of the link function since it will not affect the predicted moments in the case of a single explanatory variable.

Corollary 3.1. The value of the log-likelihood defined in (4) taken on the exact MLE $\widehat{\boldsymbol{\vartheta}}_n$ (if it exists) given by (8), under constraint (7), does not depend on the link function g . More precisely, we have $\forall i \in I, \quad \ell(\widehat{\eta}_i) = (b')^{-1}(\overline{x}_n^{(j)})$ for $j \in J$ such that $y_i^{(2),j} = 1$ and

$$\log L(\widehat{\boldsymbol{\vartheta}}_n | \mathbf{x}) = \frac{1}{a(\phi)} \sum_{j=1}^d \sum_{i, y_i^{(2),j}=1} \left(x_i (b')^{-1}(\overline{x}_n^{(j)}) - b\left((b')^{-1}(\overline{x}_n^{(j)})\right) \right) + \sum_{i=1}^n c(x_i, \phi).$$

The estimator of ϕ is obtained by maximizing $\log L(\widehat{\boldsymbol{\vartheta}}_n | \mathbf{x})$ with respect to ϕ given a, b, c functions.

Corollary 3.2. The predicted mean and predicted variance for the i th individual is estimated by $\widehat{\mathbf{E}}X_i = b'(\ell(\widehat{\eta}_i))$ and $\widehat{\mathbf{Var}}X_i = a(\widehat{\phi})b''(\ell(\widehat{\eta}_i))$ respectively using (2). Both estimates do not depend on the link function g and the predicted mean does not depend on the function b . More precisely, when v_j is the modality of the i th individual (i.e. $y_i^{(2),j} = 1$), the predicted mean and predicted variance are given by

$$\widehat{\mathbf{E}}X_i = \overline{x}_n^{(j)}, \quad \widehat{\mathbf{Var}}X_i = a(\widehat{\phi})b'' \circ (b')^{-1}(\overline{x}_n^{(j)}).$$

Corollary 3.2 may be surprising because the predicted mean does not depend on the conditional distribution of X_1, \dots, X_n . The predicted mean is just the mean of the response variable taken over the class j i.e. observations x_i such that the covariate $y_i^{(2)}$ takes the modality v_j .

The formula (8) of the MLE $\hat{\boldsymbol{\vartheta}}_n$ can be reformulated as

$$\hat{\boldsymbol{\vartheta}}_n = \begin{pmatrix} Q \\ \mathbf{R} \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{g}(\bar{\mathbf{X}}) \\ 0 \end{pmatrix}$$

with Q is the $d \times (1 + d)$ matrix defined by $Q = (A_0, A_1)$, with A_0 is the ones vector of size d , A_1 the identity matrix of size d , and $\mathbf{g}(\bar{\mathbf{X}})$ the vector $(g(\bar{X}_n^{(1)}), \dots, g(\bar{X}_n^{(d)}))$. We have

$$\sum_{j=1}^d r_{(2),j} - r_{(1)} \neq 0 \Leftrightarrow \text{rank} \begin{pmatrix} Q \\ \mathbf{R} \end{pmatrix} = d + 1.$$

With this formulation, the estimator of $\boldsymbol{\vartheta}$ is $\hat{\boldsymbol{\vartheta}}_n = (Q'Q + R'R)^{-1}Q'\mathbf{g}(\bar{\mathbf{X}})$. This general form is particularly useful in the case of two categorical variables of the next subsection.

3.2. Two explanatory variables

Now, we consider the case of two explanatory categorical variables. That is $p = 3$ and for all $i = 1, \dots, n$, $y_i^{(1)} = 1$ is the intercept and $y_i^{(2)}, y_i^{(3)}$ take values in $\{v_{j1}, \dots, v_{jd_j}\}$ with d_2, d_3 modalities respectively. We define by $y_i^{(2),k}$ and $y_i^{(3),l}$, $k \in K = \{1, \dots, d_2\}$ and $l \in L = \{1, \dots, d_3\}$ the binary dummies of the k th and l th resp. categories, $m_k^{(j)} > 0$ the number of appearance of the k th modality of the j th variable, $j = 1, 2$, m_{kl} the number of appearance of the k th and l th category simultaneously and $\bar{x}_n^{(k,l)}$ the mean value of $\underline{\mathbf{x}}$ taking over the k th and l th categories. That is

dummy	frequency	mean
$y_i^{(2),k} = 1(y_i^{(2)} = v_{2k})$	$m_k^{(2)} = \sum_{i=1}^n y_i^{(2),k} \quad k \in K$	$\bar{x}_n^{(2),k} = \frac{1}{m_k^{(2)}} \sum_{i=1}^n x_i y_i^{(2),k} \quad k \in K$
$y_i^{(3),l} = 1(y_i^{(3)} = v_{3l})$	$m_l^{(3)} = \sum_{i=1}^n y_i^{(3),l} \quad l \in L$	$\bar{x}_n^{(3),l} = \frac{1}{m_l^{(3)}} \sum_{i=1}^n x_i y_i^{(3),l} \quad l \in L$
$y_i^{(k,l)} = y_i^{(2),k} y_i^{(3),l}$	$m_{kl} = \sum_{i=1}^n y_i^{(k,l)} \quad (k, l) \in K \times L$	$\bar{x}_n^{(k,l)} = \frac{1}{m_{kl}} \sum_{i=1}^n y_i^{(k,l)} x_i \quad (k, l) \in KL^*$

with $KL^* = (K \times L) \setminus \{(k, l) \in K \times L; m_{kl} = 0\}$. Set $d_{2,3}^* = \#KL^*$, for $l \in L$, $K_l^* = \{(k, l) \in K^* \times \{l\}; m_{kl} > 0\}$, $d_{(3),l}^* = \#K_l^*$ and for $k \in K^*$, $L_k^* = \{(k, l) \in \{k\} \times L^*; m_{kl} > 0\}$, $d_{(2),k}^* = \#L_k^*$.

Note that $(m_{kl})_{kl}$ are absolute frequencies of the contingency table resulting from cross-classifying factors and can be computed very easily. Be careful that KL^* is not equal to $K \times L$ but $\bigcup_{l \in L} K_l^* = \bigcup_{k \in K} L_k^* = KL^*$, and $d_{2,3}^* = d_2 d_3 - r$, where $r = \#\{(k, l) \in K \times L; m_{kl} = 0\}$.

Let Q be the $d_{2,3}^* \times (1 + d_2 + d_3 + d_{2,3}^*)$ real matrix defined by $Q = (A_0, A_1, A_2, A_{12})$ with $A_0 = \mathbf{1}_{d_{2,3}^*}$ the $d_{2,3}^* \times 1$ ones matrix; $A_1 = (\text{diag}(\mathbf{1}_{d_{(2),k}^*}))_{k \in K}$, the $d_{2,3}^* \times K$ diagonal block matrix of ones vector of size $d_{(2),k}^*$; $A_2 = (I_{d_3}^{*,k})_{k \in K}$, the $d_{2,3}^* \times L$ matrix where $I_{d_3}^{*,k}$ is the identity matrix of size d_3 without rows l for which $m_{kl} = 0$; $A_{12} = I_{d_{2,3}^*}$ the $d_{2,3}^* \times d_{2,3}^*$ identity matrix.

Consider the following GLM for explanatory variables $y_i^{(1)}, y_i^{(2),j}, y_i^{(3),j}$

$$g(\mathbf{E}X_i) = \vartheta_1 + \sum_{k=1}^{d_2} y_i^{(2),k} \vartheta_{(2),k} + \sum_{l=1}^{d_3} y_i^{(3),l} \vartheta_{(3),l} + \sum_{(k,l) \in KL^*} y_i^{(k,l)} \vartheta_{kl}, \quad (10)$$

where $\vartheta_{(1)}, (\vartheta_{(2),k})_{k \in K}, (\vartheta_{(3),l})_{l \in L}, (\vartheta_{kl})_{(k,l) \in KL^*}$ are the $d_2 + d_3 + d_{2,3}^* + 1$ unknown parameters. Again at this stage, the model is not identifiable because of the redundancy on the vectors $(y_1^{(2),k}, \dots, y_n^{(2),k}), k \in K$, the vectors $(y_1^{(3),l}, \dots, y_n^{(3),l}), l \in L$ and the ones vector. As previously, We need to impose $q \geq 1 + d_2 + d_3$ linear constraints on the vector parameters $\boldsymbol{\vartheta}$

$$R\boldsymbol{\vartheta} = \mathbf{0}_q, \quad (11)$$

where R is a $q \times (1 + d_2 + d_3 + d_{2,3}^*)$ real matrix of linear contrasts, with $\text{rank}(R) = 1 + d_2 + d_3$ and $\mathbf{0}_q$ the zeros vector of size q .

Theorem 3.2. *Suppose that for all $i \in \{1, \dots, n\}$, X_i takes values in $b'(\Lambda)$. Under constraint (11) and if R such that $(Q'R')$ is of rank $d_{2,3}^*$, there exists a unique, consistent and explicit MLE $\hat{\boldsymbol{\vartheta}}_n$ of $\boldsymbol{\vartheta}$ given by*

$$\hat{\boldsymbol{\vartheta}}_n = (Q'Q + R'R)^{-1}Q'g(\bar{\mathbf{X}}), \quad (12)$$

where the vector $g(\bar{\mathbf{X}})$ is $((g(\bar{X}_n^{(k,l)}))_{l \in L_k^*})_{k \in K}$.

Example 3.4. *No intercept and no single-variable dummy*

The model with no intercept and no single-variable dummy is $\vartheta_1 = 0$ and $\vartheta_{(2),k} = \vartheta_{(3),l} = 0 \forall k \in K \forall l \in L$. Therefore, the unique, consistent and explicit MLE $\hat{\boldsymbol{\vartheta}}_n$ of $\boldsymbol{\vartheta}$ is

$$\hat{\vartheta}_{n,kl} = g\left(\bar{X}_n^{(k,l)}\right), \quad (k,l) \in KL^*.$$

Example 3.5. *Zero-sum conditions* The model with zero-sum conditions assumes

$$\sum_{k \in K} m_k^{(2)} \vartheta_{(2),k} = \sum_{l \in L} m_l^{(3)} \vartheta_{(3),l} = 0, \quad \forall l \in L, \sum_{k \in K_l^*} m_{kl} \vartheta_{kl} = 0, \quad \forall k \in K, \sum_{l \in L_k^*} m_{kl} \vartheta_{kl} = 0.$$

Therefore, the unique, consistent and explicit MLE $\hat{\boldsymbol{\vartheta}}_n$ of $\boldsymbol{\vartheta}$ is

$$\left\{ \begin{array}{l} \hat{\vartheta}_{n,(1)} = \frac{1}{n} \sum_{(k,l) \in KL^*} m_{kl} g\left(\bar{X}_n^{(k,l)}\right) \\ \hat{\vartheta}_{n,(2),k} = \frac{1}{m_k^{(2)}} \sum_{l \in L_k^*} m_{kl} g\left(\bar{X}_n^{(k,l)}\right) - \hat{\vartheta}_{n,1}, \quad k \in K \\ \hat{\vartheta}_{n,(3),l} = \frac{1}{m_l^{(3)}} \sum_{k \in K_l^*} m_{kl} g\left(\bar{X}_n^{(k,l)}\right) - \hat{\vartheta}_{n,1}, \quad l \in L \\ \hat{\vartheta}_{n,kl} = g\left(\bar{X}_n^{(k,l)}\right) - \hat{\vartheta}_{n,(2),k} - \hat{\vartheta}_{n,(3),l} - \hat{\vartheta}_{n,1}, \quad (k,l) \in KL^*. \end{array} \right.$$

For simplicity, we consider only the cases of one and two explanatory categorical variables. With a higher number of explanatory variables, we can perform a similar analysis to obtain an explicit solution of the MLE.

As for one explanatory variable, Theorem 3.2 has two interesting corollaries on the value of the log-likelihood and the predicted moments.

Corollary 3.3. *The value of log-likelihood defined in (4) taken on the exact MLE $\widehat{\boldsymbol{\vartheta}}_n$ (if it exists) given by (12), under constraint (11), does not depend on the link function g . More precisely, we have $\forall i \in I, \ell(\widehat{\eta}_i) = (b')^{-1}(\bar{x}_n^{(k,l)})$ for $l \in L$ and $k \in K$ such that $y_i^{(2),j} = 1$ and $y_i^{(3),k} = 1$ and*

$$\log L(\widehat{\boldsymbol{\vartheta}}_n | \underline{\mathbf{x}}) = \frac{1}{a(\widehat{\phi})} \sum_{(k,l) \in KL^*} \sum_{i | y_i^{(2),k} = y_i^{(3),l} = 1} (x_i(b')^{-1}(\bar{x}_n^{(k,l)}) - b((b')^{-1}(\bar{x}_n^{(k,l)}))) + \sum_{i=1}^n c(x_i, \phi).$$

The estimator of ϕ is obtained by maximizing $\log L(\widehat{\boldsymbol{\vartheta}}_n | \underline{\mathbf{x}})$ with respect to ϕ given a, b, c functions.

Corollary 3.4. *The predicted mean and predicted variance for the i th individual is estimated by $\widehat{\mathbf{E}}X_i = b'(\ell(\widehat{\eta}_i))$ and $\widehat{\mathbf{Var}}X_i = a(\widehat{\phi})b''(\ell(\widehat{\eta}_i))$ respectively using (2). Both estimates do not depend on the link function g and the predicted mean does not depend on the function b . Let v_{2k} and v_{3l} be the modalities of the i th individual of the two explanatory variables, i.e. $y_i^{(2),k} = 1$ and $y_i^{(3),l} = 1$. The predicted mean and variance are given by*

$$\widehat{\mathbf{E}}X_i = \bar{x}_n^{(k,l)}, \quad \widehat{\mathbf{Var}}X_i = a(\widehat{\phi})b'' \circ (b')^{-1}(\bar{x}_n^{(k,l)}).$$

In the next two sections, we apply previous theorems and corollaries to two particular distributions: Pareto 1 and lognormal distribution. Our results do not only apply to continuous distributions but also for discrete distributions. But we choose these distributions in order to model insurance losses.

4. GLM for Pareto I distribution with categorical explanatory variables

4.1. Characterization

Consider the sample $\underline{\mathbf{X}} = (X_1, \dots, X_n)$ composed of independent Pareto Type 1. Precisely, we assume that the independent random variables X_1, \dots, X_n are Pareto with known threshold parameter μ and respective shape parameter (depending on the unknown parameter $\boldsymbol{\vartheta}$) $\lambda_1(\boldsymbol{\vartheta}), \dots, \lambda_n(\boldsymbol{\vartheta}) \in \Lambda = (0, \infty)$. The density f of Pareto distribution with scale and shape parameter μ and $\lambda_i(\boldsymbol{\vartheta})$, $i \in I$ is

$$f(x) = \lambda_i(\boldsymbol{\vartheta}) \frac{\mu^{\lambda_i(\boldsymbol{\vartheta})}}{x^{\lambda_i(\boldsymbol{\vartheta})+1}}, \quad x \in \mathbb{X} = [\mu, \infty), \quad (13)$$

and 0 if $x < \mu$.

We recall that for the Pareto Type 1 distribution

$$\mathbf{E}X_i = \frac{\lambda_i(\boldsymbol{\vartheta})\mu}{\lambda_i(\boldsymbol{\vartheta}) - 1} < +\infty \quad \text{iff } \lambda_i(\boldsymbol{\vartheta}) > 1 \quad \text{and} \quad \mathbf{E}X_i^2 = \frac{\lambda_i(\boldsymbol{\vartheta})\mu^2}{\lambda_i(\boldsymbol{\vartheta}) - 2} < +\infty \quad \text{iff } \lambda_i(\boldsymbol{\vartheta}) > 2.$$

Unlike the known parameter μ , the parameter $\boldsymbol{\vartheta}$ is to be estimated. These closed-form formulas are particularly useful in an insurance context since the expectation and the variance

are used in most premium computation. For instance, $\mathbf{E}X$ is the pure premium and for $\gamma > 0$, $\mathbf{E}X + \gamma \mathbf{Var}X$ is the variance principle (see Bühlmann and Gisler, 2006, Section 1.2.2).

In the following, instead of $\underline{\mathbf{X}}$ we consider the sample $\underline{\mathbf{z}} = (T(X_1), \dots, T(X_n))$. With the re-parametrization $z_i = T(x_i) = -\log(x_i/\mu)$, $i \in I$, this distribution belongs to the exponential family as defined in (1), with

$$a(\phi) = 1, \quad b(\lambda) = -\log(\lambda), \quad \text{and} \quad c(z, \phi) = 0, \quad z \in T(\mathbb{X}) = \mathbb{R}^-, \quad \lambda \in \Lambda. \quad (14)$$

In particular, for the Pareto I distribution, there is no dispersion parameter. It is also worth mentioning that $-Z_i$ is exponential with parameter $\lambda_i(\boldsymbol{\vartheta})$. Consecutively, all moments of Z_i exist and are given by $\mathbf{E}Z_i^m = (-1)^m m! / \lambda_i(\boldsymbol{\vartheta})^m$, $m \in \mathbb{N}^*$.

Consider the regression model with a link function g , a response variable X_i Pareto I distributed where $Z_i = -\log(X_i/\mu)$ and

$$g(\mathbf{E}Z_i) = \vartheta_{(1)} + y_i^{(2),1} \vartheta_{(2),1} + \dots + y_i^{(2),d} \vartheta_{(2),d} = \langle \mathbf{y}_i, \boldsymbol{\vartheta} \rangle, \quad i \in I \quad (15)$$

with for $i \in I$, $\mathbf{y}_i = (1, y_i^{(2),1}, \dots, y_i^{(2),d})^T$ are the covariate vectors and $\boldsymbol{\vartheta} = (\vartheta_{(1)}, \vartheta_{(2),1}, \dots, \vartheta_{(2),d})^T$ is the unknown parameter vector.

The choice of the link function g appearing in (15) is a crucial point. Let us start with the canonical link. That is, the chosen function g so that $\ell = (b')^{-1} \circ g^{-1}$ is the identity function. For our Pareto model, $g(t) = -\frac{1}{t}$ since $b'(\lambda) = -\frac{1}{\lambda}$. From (4), the choice of the canonical link function imposes constraints on the linear predictor space $D \subset \Lambda \subset (0, +\infty)$ in that case. Since D results from the scalar product of $\boldsymbol{\vartheta}$ parameters with explanatory variables \mathbf{y}_i , some negative values might be produced when the covariables take negative values. This make the choice of the canonical link inappropriate.

In order to remedy this issue, we can choose a link function such that the values of the function ℓ falls in $\Lambda \subset (0, \infty)$. A natural choice is $\ell(\eta) = \exp(\eta)$ which guarantees a positive parameter. The choice $\ell(\eta) = \exp(\eta) + 1$ guarantee a finite expectation for the random variables X_i , $i = 1, \dots, n$. We summarize in Table 1 the tested ℓ functions in our application in Section 6, and in Table 2, the spaces given a link function.

Table 1: Table of typical link functions for Pareto I

Names	$\ell(\eta_i)$	$g^{-1}(t)$	$g(t)$
canonical	η_i	$-\frac{1}{t}$	$-\frac{1}{t}$
log-inv	e^{η_i}	$-e^{-t}$	$\log(-\frac{1}{t})$
shifted log-inv	$e^{\eta_i} + 1$	$-\frac{1}{e^t + 1}$	$\log(-\frac{1}{t} - 1)$

4.2. Estimation for categorical exogenous variables

Consider the case of one categorical exogenous variable. We expose the case of the re-parametrization without intercept ($\langle \mathbf{R}, \boldsymbol{\vartheta} \rangle = 0$ with $\mathbf{R} = (1, 0, \dots, 0)$).

Table 2: Summary of spaces for Pareto I

Link name	Parameter–covariable space $\boldsymbol{\vartheta} \times Y_i$	Linear predictor space D	Parameter space Λ	$b'(\Lambda)$	
unspecified	$\boldsymbol{\vartheta} \times Y_i \subset \mathbb{R}^p \times \mathbb{R}^p$	$\langle \cdots \rangle$	$D \subset \mathbb{R}$	$\frac{\ell^{-1}}{\ell}$	$\Lambda \subset (0, +\infty)$
canonical	$\{(\boldsymbol{\vartheta}, \mathbf{y}_i) \in \mathbb{R}^p \times \mathbb{R}^p, \langle \boldsymbol{\vartheta}, \mathbf{y}_i \rangle \geq 0\}$	$\langle \cdots \rangle$	$(0, +\infty)$	$\frac{\text{id}}{\text{id}}$	$(0, +\infty)$ $(-\infty, 0)$
log-inv	$\mathbb{R}^p \times \mathbb{R}^p$	$\langle \cdots \rangle$	\mathbb{R}	$\frac{\log}{\exp}$	$(0, +\infty)$ $(-\infty, 0)$
shifted log-inv	$\mathbb{R}^p \times \mathbb{R}^p$	$\langle \cdots \rangle$	\mathbb{R}	$\frac{\log(x-1)}{\exp(x)+1}$	$(1, +\infty)$ $(-1, 0)$

Let $\widehat{\boldsymbol{\vartheta}}_n$ the MLE defined in (9), if it exists, of $\boldsymbol{\vartheta}$. Using $\sum_{j=1}^d \sum_{i, y_i^{(2),j}=1} z_i (b')^{-1}(\bar{z}_n^{(j)}) = n$, the log likelihood evaluated on $\widehat{\boldsymbol{\vartheta}}_n$ for both the transformed sample $\underline{\mathbf{z}}$ and the original sample $\underline{\mathbf{x}}$ with one categorical exogenous variable (Corollary 3.1) is

$$\log L(\widehat{\boldsymbol{\vartheta}}_n | \underline{\mathbf{z}}) = n - \sum_{j=1}^d m_j \log(-\bar{z}_n^{(j)}), \quad \log L(\widehat{\boldsymbol{\vartheta}}_n | \underline{\mathbf{x}}) = n - \sum_{j=1}^d m_j \log(-\bar{z}_n^{(j)}) - \sum_{i=1}^n \log(x_i). \quad (16)$$

The log-likelihood $\log L(\widehat{\boldsymbol{\vartheta}}_n | \underline{\mathbf{x}})$ is detailed on Appendix B. The second remark is that $-Z_i$ are exponentially distributed $\mathcal{E}(\ell(\eta_i))$. Hence for $j \in J$, the estimators $\widehat{\vartheta}_{n,(2),j}$ of $\vartheta_{(2),j}$ are known transforms of a gamma random variable $\mathcal{G}a(m_j, m_j \ell(\vartheta_{2,j}))$. Below we analyze the choice of the link functions considered in Table 1 in Examples 4.1, 4.2, 4.3 and plotted in Figure C.3a.

Example 4.1. *canonical link*

In the special case of canonical Pareto model, because $z_i < 0$ for all $i \in I$, we have $\bar{z}_n^{(j)} \in b'(\Lambda) = (-\infty, 0)$ for all $j \in J$ (g and Λ are respectively defined in Tables 1 and 2). With no-intercept using Equation (9), the MLE is

$$\widehat{\vartheta}_{n,(2),j} = -m_j \left(\sum_{i=1}^n y_i^{(2),j} Z_i \right)^{-1} = -\frac{1}{\bar{Z}_n^{(j)}}, \quad j \in J.$$

Hence, for $j \in J$, $\widehat{\vartheta}_{n,(2),j}$ follows an Inverse Gamma distribution with shape parameter m_j and rate parameter $m_j \vartheta_{(2),j}$, see e.g. (Johnson et al., 2000, Ch. 17). We can compute the moments of the Inverse Gamma distribution, for $m_j > 2$,

$$\mathbf{E} \widehat{\vartheta}_{n,(2),j} = \frac{m_j}{m_j - 1} \vartheta_{(2),j}, \quad \text{and} \quad \mathbf{Var} \widehat{\vartheta}_{n,(2),j} = \frac{m_j^2}{(m_j - 1)^2 (m_j - 2)} \vartheta_{(2),j}^2, \quad j \in J.$$

An unbiased estimator of $\vartheta_{(2),j}$ is then $\widehat{\vartheta}_{n,(2),j}^* = \frac{m_j - 1}{m_j} \widehat{\vartheta}_{n,(2),j}$ which has a lower variance

$$\mathbf{Var} \widehat{\vartheta}_{n,(2),j}^* = \frac{\vartheta_{(2),j}^2}{m_j - 2} \leq \mathbf{Var} \widehat{\vartheta}_{n,(2),j}, \quad j \in J.$$

A similar bias is also obtained by Bühlmann and Gisler (2006) in a credibility context. Of course this unbiased estimator is also applicable for two exogenous variables with the first parametrization of the Theorem 3.2. When some modalities (or couple of modalities) aren't much represented, it can be relevant to use this unbiased estimator.

Example 4.2. *log-inverse link*

In the special case of the log-inv Pareto model, we also have $\bar{z}_n^{(j)} \in b'(\Lambda) = (-\infty, 0)$ for all $j \in J$ (g and Λ are respectively defined in Tables 1 and 2). With no-intercept using Equation (9), the MLE is

$$\hat{\vartheta}_{n,(2),j} = -\log \left(\frac{1}{-m_j} \sum_{i=1}^n y_i^{(2),j} Z_i \right) = -\log \left(-\bar{Z}_n^{(j)} \right), \quad j \in J.$$

Here, for $j \in J$, the distribution of $-\hat{\vartheta}_{n,(2),j}$ is the distribution of the log of the gamma distribution with shape m_j and rate $m_j \exp(\vartheta_{(2),j})$. We can derive moments of this distribution which should not be confused with the log-gamma distribution studied e.g. in Hogg and Klugman (1984).

Let $L = \log(G)$ when G is gamma distributed with shape parameter $a > 0$ and rate parameter $\lambda > 0$. We have by elementary manipulations the moment generating function of L :

$$M_L(t) = \mathbf{E}e^{tL} = \frac{\Gamma(a+t)}{\Gamma(a)} \lambda^{-t}, \quad t > -a,$$

where Γ denotes the usual gamma function. Therefore by differentiating and evaluating at 0, we deduce that the expectation and the variance of L are $\mathbf{E}L = \psi(a) - \log \lambda$ and $\mathbf{Var}L = \psi'(a)$, where the functions ψ and ψ' are the digamma and trigamma function, see e.g. Olver et al. (2010). Getting back to our example, we deduce that

$$\mathbf{E}\hat{\vartheta}_{n,(2),j} = \vartheta_{(2),j} + \log m_j - \psi(m_j) \quad \text{and} \quad \mathbf{Var}\hat{\vartheta}_{n,(2),j} = \psi'(m_j), \quad j \in J.$$

From Olver et al. (2010), we know that $\log(m_j) - \psi(m_j)$ tends to 0 as m_j tend to infinity. Hence $\hat{\vartheta}_{n,(2),j}$ is asymptotically unbiased, and an unbiased estimator of $\vartheta_{(2),j}$ is

$$\hat{\vartheta}_{n,(2),j}^* = \hat{\vartheta}_{n,(2),j} - (\log(m_j) - \psi(m_j)), \quad j \in J.$$

Example 4.3. *shifted log-inverse link*

In the special case of the shifted log-inv Pareto model, $\bar{z}_n^{(j)}$ is not necessarily in $b'(\Lambda) = (-1, 0)$ for all $j \in J$ (g and Λ are respectively defined in Tables 1 and 2). If there is an index j such as $\bar{z}_n^{(j)} \leq -1$, the MLE is not defined and we couldn't use the shifted log-inv link with the same incidence matrix.

Nevertheless, for sufficiently large n , for j such that $y_i^{(2)} = v_j$, $\bar{Z}_n^{(j)} \rightarrow \mathbf{E}Z_i$ almost surely, where $\mathbf{E}Z_i = -1/(\exp(\vartheta_{2,j}) + 1) > -1$. Hence for sufficiently large n , the conditions of Theorem 3.1 are satisfied. With no-intercept using Equation (9), the MLE (provided it exists) is

$$\hat{\vartheta}_{n,(2),j} = \log \left(\frac{m_j}{-\sum_{i=1}^n y_i^{(2),j} Z_i} - 1 \right) = \log \left(-1/\bar{Z}_n^{(j)} - 1 \right), \quad j \in J.$$

The expectation of $\widehat{\vartheta}_{n,(2),j}$ is complex and should be done numerically. However by the strong law of large numbers and the continuity of the link function, $\widehat{\vartheta}_{n,(2),j} = -\log\left(-1/\overline{Z}_n^{(j)} - 1\right)$ converge almost surely to $\log((\exp(\vartheta_{(2),j}) + 1) - 1) = \vartheta_{(2),j}$.

Remark 4.1. In Theorem 3.1, the condition X_i takes values in $b'(\Lambda)$ might seem too restrictive. In fact the condition $\overline{x}_n^{(j)} \in b'(\Lambda)$ for all $j \in J$ is sufficient to define a vector value $\widehat{\boldsymbol{\vartheta}}_n$ which maximise the likelihood. But $\widehat{\boldsymbol{\vartheta}}_n$ fails to be a MLE estimator because the random variable $g(\overline{X}_n^{(j)})$ can to be not defined. Nevertheless, when m_j tends to infinity for any $j \in J$, the random variables X_i 's defined by (6) such that $y_i^{(2),j} = 1$ are i.i.d. (not only independent) and the strong law of large numbers implies that $\overline{X}_n^{(j)}$ converges almost surely to $\mathbf{E}X_i = b'(\ell(\eta_i)) \in b'(\Lambda)$. Hence, the probability $P(\overline{X}_n^{(j)} \notin b'(\Lambda))$ tends to zero which guarantees the asymptotically existence of the MLE estimator.

4.3. Model diagnostic

In this paragraph, we propose residuals adapted at the case of Pareto problem. First note that X_i is Pareto I with shape $\ell(\eta_i)$ and threshold μ and for the parametrization (14) $-Z_i = \log(X_i/\mu) \sim \mathcal{E}(\ell(\eta_i))$. Let define the residuals

$$R_i = -\ell(\eta_i)Z_i, \quad i \in I.$$

Hence R_1, \dots, R_n are i.i.d. and have an exponential distribution $\mathcal{E}(1)$. The consistency of the MLE makes it possible to say that the estimated residuals $\widehat{R}_{n,i} = -\ell(\widehat{\eta}_i)Z_i$, $i \in I$, with $\widehat{\eta}_i = \langle \mathbf{y}_i, \widehat{\boldsymbol{\vartheta}}_n \rangle$ are asymptotically i.i.d..

It is also possible to verify the assumption of the Pareto distribution for X_i conditionally to \mathbf{y}_i with graphical diagnostic as an exponential Quantile-Quantile plot on the residuals $\widehat{R}_{n,i}$.

In the case of a single explanatory variable, for $i \in I$, the residuals $\widehat{R}_{n,i}$ do not depend on ℓ function. Their explicit forms are

$$\widehat{R}_{n,i} = \frac{Z_i}{\overline{Z}_n^{(j)}} \quad j \text{ such that } y_i^{(2)} = v_j, \quad i \in I. \quad (17)$$

Furthermore, the summation of $\widehat{R}_{n,1}, \dots, \widehat{R}_{n,n}$ has the surprising property to be deterministic and is exactly equal to n .

5. GLM for shifted log-normal distribution with categorical explanatory variables

5.1. Characterization

Secondly, consider the sample $\underline{\mathbf{X}} = (X_1, \dots, X_n)$ composed of independent shifted log-normal variables respectively with mean $\lambda_1(\boldsymbol{\vartheta}), \dots, \lambda_n(\boldsymbol{\vartheta})$, dispersion $\phi = \sigma^2$ and a known threshold μ . The shifted log-normal is also known as the 3-parameter log-normal. Precisely, the density of X_i is

$$f(x) = \frac{1}{(x - \mu)\sqrt{2\pi\phi}} \exp\left(-\frac{(\log(x - \mu) - \lambda_i(\boldsymbol{\vartheta}))^2}{2\phi}\right), \quad x \in \mathbb{X} = (\mu, \infty), \quad (18)$$

and 0 for $x \leq \mu$. It is well known that the lognormal distribution has finite moment, see e.g. Johnson et al. (2000). In particular, the expectation and the variance are given by

$$\mathbf{E}X_i = \mu + \exp(\lambda_i(\boldsymbol{\vartheta}) + \phi/2), \quad \mathbf{Var}X_i = (\exp(\phi) - 1) \exp(2\lambda_i(\boldsymbol{\vartheta}) + \phi).$$

The transformed sample $\underline{\mathbf{Z}} = T(\underline{\mathbf{X}}) = (T(X_1), \dots, T(X_n))$ with $T(x) = \log(x - \mu)$ is belongs the exponential family with

$$a(\phi) = \phi, \quad b(\lambda) = \lambda^2/2, \quad c(z, \phi) = -\frac{1}{2} \left(\frac{z^2}{\phi} + \log(2\pi\phi) \right), \quad z \in \mathbb{R}+, \lambda \in \mathbb{R}.$$

It is worth mentioning that Z_i are normally distributed with mean $\lambda_i(\boldsymbol{\vartheta})$ and variance ϕ . As a consequence, all moments of Z_i exists and $\mathbf{E}(Z_i - \lambda_i)^m = (2m)! \phi^m / (2^m m!)$ for m even and 0 for m odd. Consider the regression model with a link function g , a response variable X_i lognormally distributed where $Z_i = \log(X_i - \mu)$ and

$$g(\mathbf{E}Z_i) = \vartheta_{(1)} + \vartheta_{(2),1} y_i^{(2),1} + \dots + \vartheta_{(2),d} y_i^{(2),d} = \langle \mathbf{y}_i, \boldsymbol{\vartheta} \rangle, \quad i \in I \quad (19)$$

with for $i \in I$, $\mathbf{y}_i = (1, y_i^{(2),1}, \dots, y_i^{(2),d})^T$ are the covariate vectors and $\boldsymbol{\vartheta} = (\vartheta_{(1)}, \vartheta_{(2),1}, \dots, \vartheta_{(2),d})^T \in \mathbb{R}^d$ is the unknown parameter vector.

The choice of the link function g for Equation (19) is less restrictive than for the Pareto case. Any differentiable invertible function from \mathbb{R} to \mathbb{R} will work. Since $b'(x) = x$, the canonical link function is obtained by choosing g such that $\ell = \text{id} \circ g^{-1} = g^{-1}$ is the identity function. In other words, the canonical link function is the identity function.

Unlike the previous section, any moment of X_i exist and there is no particular link needed to guarantee their existence. In the numerical application, we will also consider another link function: a real version of the logarithm.

5.2. Estimation for categorical exogenous variables

Again we consider the case of categorical variables and without intercep model, that is with a predictor $\eta_i = y_i^{(2),1} \vartheta_{(2),1} + \dots + y_i^{(2),d} \vartheta_{(2),d}$. In the case of the lognormal dispersion, there is a dispersion to be estimated. The log-likelihood is given by

$$\log L(\boldsymbol{\vartheta} | \underline{\mathbf{z}}) = -\frac{1}{2\phi} \sum_{i=1}^n (z_i - \lambda_i(\boldsymbol{\vartheta}))^2 - \frac{n \log(2\pi\phi)}{2}.$$

The components of the MLE of $\boldsymbol{\vartheta}$ are given by $\hat{\boldsymbol{\vartheta}}_{n,(2),j} = g(\bar{z}_n^{(j)})$, $j \in J$, and the estimated log likelihood for the transformed sample $\underline{\mathbf{z}}$ is

$$\log L(\hat{\boldsymbol{\vartheta}}_n | \underline{\mathbf{z}}) = -\frac{1}{2\phi} \sum_{j \in J} \sum_{i, y_i^{(2),j}=1} (z_i - \bar{z}_n^{(j)})^2 - \frac{n \log(2\pi\phi)}{2}$$

Maximizing over ϕ the log-likelihood leads to the empirical variance

$$\hat{\phi} = \frac{1}{n} \sum_{j \in J} \sum_{i, y_i^{(2),j}=1} (z_i - \bar{z}_n^{(j)})^2 \quad (20)$$

Hence the estimator of ϕ is the intra-class variance. This closed-form estimate $\hat{\phi}$ differs from what classical statistical softwares carry out, where the dispersion parameter is estimated by a quasi-likelihood approach.

Using the previous equations, we compute the log likelihood evaluated on $\hat{\phi}$ and on $\hat{\vartheta}_n$ for both the transformed sample \underline{z} and the original sample \underline{x} with one categorical exogenous variable (Corollary 3.1) is

$$\log L(\hat{\vartheta}_n | \underline{z}) = -\frac{n}{2}(1 + \log(2\pi\hat{\phi})), \quad \log L(\hat{\vartheta}_n | \underline{x}) = -\frac{n}{2}(1 + \log(2\pi\hat{\phi})) - \sum_{i=1}^n z_i. \quad (21)$$

The log-likelihood $\log L(\hat{\vartheta}_n | \underline{x})$ is detailed on Appendix B. Below we analyze the choice of the link functions considered in Table 3 in Examples 5.1, 5.2 and plotted in Figure C.3b.

Table 3: Table of typical link functions for lognormal

Names	$\ell(\eta_i)$	$g^{-1}(t)$	$g(t)$
canonical	η_i	t	t
sym. log	$e^{\eta_i} 1_{\eta_i \geq 0} + (2 - e^{-\eta_i}) 1_{\eta_i < 0}$	$e^t 1_{t \geq 0} + (2 - e^{-t}) 1_{t < 0}$	$\log(t) 1_{t \geq 1} - \log(2 - t) 1_{t < 1}$

Table 4: Summary of spaces for lognormal

Link name	Parameter-covariable space $\vartheta \times Y_i$	Linear predictor space D	Parameter space Λ
canonical	$\mathbb{R}^p \times \mathbb{R}^p$	$\langle \cdot, \cdot \rangle$	\mathbb{R}
sym. log	$\mathbb{R}^p \times \mathbb{R}^p$	$\langle \cdot, \cdot \rangle$	\mathbb{R}

Example 5.1. *canonical link*

With the canonical link function, there is no issue between the parameter space and the linear predictor space since $D = \Lambda = \mathbb{R}$. With no-intercept using Equation (9), the MLE is

$$\hat{\vartheta}_{n,(2),j} = \frac{1}{m_j} \sum_{i=1}^n y_i^{(2,j)} Z_i = \bar{Z}_n^{(j)}, \quad j \in J.$$

Hence, the distribution $\hat{\vartheta}_{n,(2),j}$ is simply a normal distribution with mean $\vartheta_{(2),j}$ and variance ϕ/m_j . Therefore, the MLE is unbiased and converges in almost surely to $\vartheta_{(2),j}$.

Example 5.2. *symmetrical log link*

We consider a central symmetry of the logarithm function given in Table 3 leading to $l_g(x) = e^x 1_{x \geq 0} + (2 - e^{-x}) 1_{x < 0}$. With this symmetrical log link function, there is no issue

between the parameter space and the linear predictor space since again $D = \Lambda = \mathbb{R}$. With no-intercept using Equation (9), the MLE is

$$\widehat{\vartheta}_{n,(2),j} = \log \left(\bar{z}_n^{(j)} 1_{\bar{z}_n^{(j)} \geq 1} \right) - \log \left(2 - \bar{z}_n^{(j)} 1_{\bar{z}_n^{(j)} < 1} \right), \quad j \in J.$$

The expectation of $\widehat{\vartheta}_{n,(2),j}$ is complex and should be done numerically. However by the strong law of large numbers and the continuity of the link function, $\widehat{\vartheta}_{n,(2),j} = l_g(\bar{Z}_n^{(j)})$ converge almost surely to $l_g(\mathbf{E}\bar{Z}_n^{(j)}) = \vartheta_{(2),j}$.

5.3. Model diagnostic

In this paragraph, we give some details about residuals in the lognormal case. As already said, the transformed variables $Z_i = \log(X_i - \mu)$ is normally distributed with mean $\ell(\eta_i)$ and variance ϕ . Let define the residuals

$$R_i = \frac{Z_i - \ell(\eta_i)}{\sqrt{\phi}}, \quad i \in I.$$

Hence R_1, \dots, R_n are i.i.d. and have a normal distribution $\mathcal{N}(0, 1)$. The consistency of the MLE makes it possible to say that the $\widehat{R}_{n,i} = \frac{Z_i - \ell(\widehat{\eta}_i)}{\sqrt{\widehat{\phi}}}$, $i \in I$, with $\widehat{\eta}_i = \langle \mathbf{y}_i, \widehat{\boldsymbol{\vartheta}}_n \rangle$ are asymptotically i.i.d. Furthermore, the summation of $\widehat{R}_{n,1}, \dots, \widehat{R}_{n,n}$ is exactly equal to 0.

It is also possible to verify the assumption of the lognormal distribution for X_i conditionally to \mathbf{y}_i with graphical diagnostic as a normal Quantile-Quantile plot on the residuals $\widehat{R}_{n,i}$.

6. Application to large claim modeling

This section is devoted to the numerical illustration: all computations are carried out thanks to the R statistical software R Core Team (2018). In our application, we focus on modeling non-life insurance losses of corporate business lines. Our data set comes from an anonymous private insurer: for privacy reason, amounts have been randomly scaled, dates randomly rearranged, variable modalities renamed. The data set consists of 211,739 claims which occurred between 2000 and 2010. In addition to the claim amount level, various explanatory variables are available.

We provide in Table C.9 in Appendix C a short descriptive analysis of the two most important variables (risk class and guarantee type with respectively 5 and 7 modalities). Due to the very high value of skewness and kurtosis, we observe that claim amount is particularly heavy tailed.

In the sequel, we consider only large claims which are in our context claims above $\mu = 340,000$ (in euros). The threshold value has been chosen by expert opinion of practitioners. We refer to e.g. Reiss and Thomas (2007) for advanced selection methods based on extreme value theory.

6.1. A single explanatory variable

Firstly, we consider both Pareto 1 GLM and Shifted log-normal GLM with only one explanatory variable: the guarantee type. We choose Guarantee 1 as the reference level implying that $\vartheta_{(2),1} = 0$. So, $\vartheta_{(1)}$ representing the effect of the reference category and $(\vartheta_{(2),j})_j$ representing the differential effect of categories j relative to the reference category will be estimated through (15) and (19). Observations x_1, \dots, x_n are observed claim amounts either from Pareto 1 (13) or shifted log-normal (18).

For these two models, we have many possible choices for the link function g . Naturally, we choose link functions appearing in Tables 1 and 3 respectively. In accordance to Corollaries 3.1, the choice of g does not impact the values respectively given on (16) and (21) of the log-likelihoods applied on the MLE of ϑ .

Furthermore, for Pareto GLM, the choice of shifted log-inverse link function seems attractive because it guarantees the existence of $\mathbf{E}X_i$. Nevertheless, alternative link functions (canonical or log-inv) allow to construct an unbiased estimator (see Section 4). For shifted log-normal model, the choice of canonical link function is more attractive because it leads to an unbiased and simpler MLE estimator (see Section 5).

Coefficients are estimated by explicit formulas given in Sections 4 and 5. We compare the fitted coefficients with the result of the IWLS algorithm described in McCullagh and Nelder (1989): we found exactly the same value for the MLE (within the numeric tolerance). In Table 5, the estimated coefficients are given in the five considered situations. Positive values of $\vartheta_{(2),j}$ in the Pareto GLM increase the shape parameter of the Pareto 1 distribution leading to a decrease in heavy-tailedness. Regarding the log-normal model, positive values of $\vartheta_{(2),j}$ increase the scale parameter of the log-normal distribution leading to a shrink of the distribution.

Irrespectively of the considered link function, the sign of the fitted coefficients are same except for intercept (Table 5) given a distribution. This convinces us that different model assumptions (i.e. link) do not lead to opposite conclusions on the claim severity. Furthermore from Table C.9, we retrieve the fact that all guarantees except Guarantee 2 have heavier tails than the reference Guarantee 1.

Table 5: Coefficients for the guarantee variable

Model Variable	Pareto 1			Shifted log normal	
	canonical	loginv	shifted.loginv	canonical	symlog
Intercept	1.89	0.64	-0.11	11.75	2.46
Guarantee 2	0.04	0.02	0.04	0.10	0.01
Guarantee 3	-0.67	-0.43	-1.36	0.75	0.06
Guarantee 4	-0.86	-0.60	-3.13	1.04	0.08
Guarantee 5	-0.71	-0.47	-1.55	0.72	0.06
Guarantee 6	-0.42	-0.25	-0.63	0.42	0.04
Guarantee 7	-0.48	-0.29	-0.78	0.59	0.05
log likelihood	-14507.53	-14507.53	-14507.53	-14517.37	-14517.37

Whatever the considered link function g , the residuals defined in Section 4.3 by $\widehat{R}_{n,i} = -\ell(\widehat{\eta}_i)Z_i$, $i \in I$, do not depend on ℓ and are given by Equation (17). We show on Figure 2 (left) the quantile/quantile plots of residuals described on Section 4.3 against the standard

exponential distribution. The assumption of Pareto 1 for X_1, \dots, X_n does not seem to be contradictory. However the assumption of log-normal distribution is more questionable, see Figure 1. Comparing the value of the log-likelihood in Table 5, Pareto 1 distribution is also the best choice. In the following, we focus only on the Pareto 1 distribution.

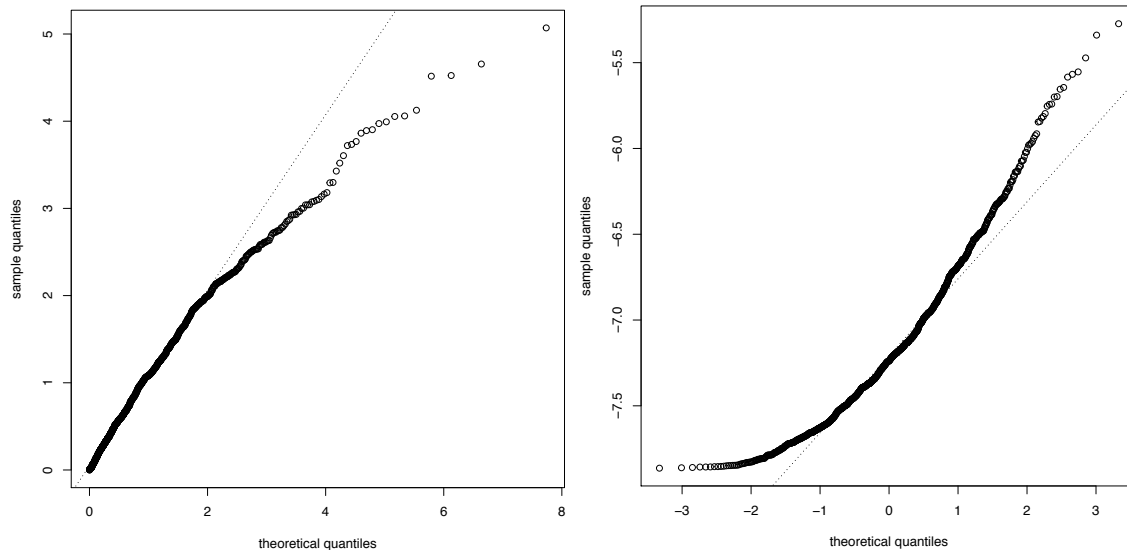


Figure 1: Quantile-quantile plots of the residuals defined in Section 4.3: (left) for `pareto1`, (right) for `shifted log-normal`

For all coefficient, let us compute the p-values of statistical tests with the null hypothesis $\vartheta_{(1)} = 0$ (Intercept null), and for $j \in J \setminus \{1\}$, $\vartheta_{(2),j} = 0$ (no differential effect of the j th Guarantee). Table 6 reports the value of the coefficient, its standard error, the student statistics and the associated p-value. We observe that some modalities of the guarantee variable are statistically significant at the 5% level. Except for Guarantee 2 and Guarantee 6, other p-values are relatively small showing the Pareto 1 distribution with explanatory variables is relevant in this context.

Table 6: Statistics and p.values for the tests $\vartheta_{(1)} = 0$ and $\vartheta_j = 0$, $j \in J$ in the Pareto GLM model for the log-inverse link.

	Estimate	Std. Error	z value	$\Pr(> z)$
Intercept	0.6391	0.1644	3.8877	0.0001
Guarantee 2	0.0214	0.2219	0.0966	0.9230
Guarantee 3	-0.4332	0.1950	-2.2217	0.0263
Guarantee 4	-0.6009	0.1708	-3.5180	0.0004
Guarantee 5	-0.4660	0.1805	-2.5817	0.0098
Guarantee 6	-0.2485	0.2295	-1.0827	0.2789
Guarantee 7	-0.2949	0.1834	-1.6075	0.1080

6.2. Two explanatory variables

Secondly, we consider the Pareto GLM models and Shifted log-normal GLM models with the two explanatory variables (guarantee and risk class) without intercept nor single-variable

(c.f. model (10) and example 3.4), that are

$$g(\mathbf{E}Z_i) = \sum_{(k,l) \in KL^*} \vartheta_{kl} y_i^{(k,l)}, \quad i \in I \quad (22)$$

with $Z_i = -\log(X_i/\mu)$ for the Pareto 1 modeling $Z_i = \log(X_i - \mu)$ for the shifted log normal modeling and where for $(k,l) \in KL^*$ the unknown parameters ϑ_{kl} represent the effect of the couple of the modalities k and l for the first and the second variable. In these examples, as it describes in Table 7, we have $K = \{1, \dots, 7\}$, $L = \{1, \dots, 5\}$ but $KL^* = \{1, \dots, 7\} \times \{1, \dots, 5\} \setminus \{(1,2), (6,5)\}$.

Consider the estimation procedures in 22. Computing claim numbers according Guarantee and Risk is done in Table 7. This claim number per class might be too short to ensure the existence of the MLE with the shifted log-inv link. We arbitrary choose the simple case of the canonical link and an unbiased estimator is relevant in this context.

The coefficients of the model are estimated using the exact method described in Section 3 and then unbiased in the same way of Example 1. The fitted coefficients are not shown but are available upon request to the authors. Furthermore, we compute the p-values of the statistical test $\vartheta_{kl} = 0$ in Table 8. We observe that most computed p-values are small: either less than 10^{-6} or less than 1%. Only 5 on the 33 p-values are above the usual 5% level, corresponding to the couples Guarantee/Risk class (1,5), (2,2), (2,3), (2,5) and (7,5) (claim number of 1,2 or 3). In the two variables setting, the Pareto 1 GLM is thus still relevant.

Table 7: Number of claim par couple Guarantee/Risk class.

Claim number	Risk class 1	Risk class 2	Risk class 3	Risk class 4	Risk class 5
Guarantee 1	39	0	4	6	1
Guarantee 2	26	2	3	16	3
Guarantee 3	48	7	11	29	4
Guarantee 4	232	40	75	147	20
Guarantee 5	68	18	36	72	6
Guarantee 6	24	7	4	11	0
Guarantee 7	94	9	22	57	3

Table 8: p-values for the tests $\vartheta_{kl} = 0$, $(k,l) \in KL^*$ in 22 for the canonical link.

P.values	Risk class 1	Risk class 2	Risk class 3	Risk class 4	Risk class 5
Guarantee 1	$< 10^{-6}$	-	0.04550	0.01431	0.31731
Guarantee 2	$< 10^{-6}$	0.15730	0.08301	0.00006	0.08326
Guarantee 3	$< 10^{-6}$	0.00815	0.00091	$< 10^{-6}$	0.04550
Guarantee 4	$< 10^{-6}$	$< 10^{-6}$	$< 10^{-6}$	$< 10^{-6}$	0.00001
Guarantee 5	$< 10^{-6}$	0.00002	$< 10^{-6}$	$< 10^{-6}$	0.01431
Guarantee 6	$< 10^{-6}$	0.00815	0.04550	0.00091	-
Guarantee 7	$< 10^{-6}$	0.00270	$< 10^{-6}$	$< 10^{-6}$	0.08326

7. Conclusion

In this paper, we deal with regression models where the response variable belongs to the general formulation of the exponential family, the so-called GLM. We focus on the estimation

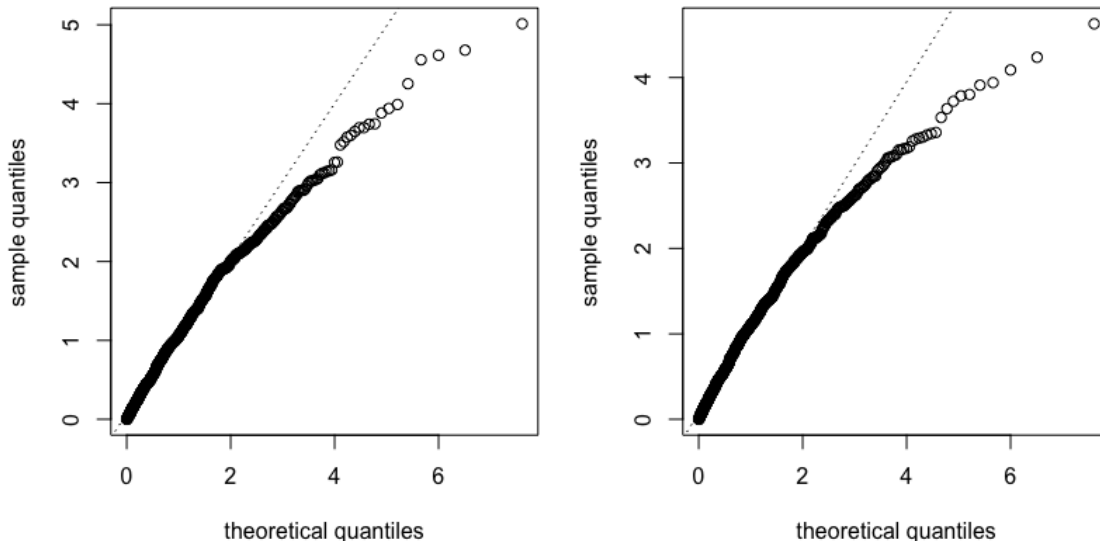


Figure 2: Quantile-quantile plots of the residuals defined in Section 4.3: (left) for one explanatory variable, (right) for two explanatory variables

of parameters of GLMs and derive explicit formulas for MLE in the case of categorical explanatory variables. In this case, the closed-form estimators do not require any use of numerical algorithms, in particular the well-known IWLS algorithm. This is logical, because in the special setting of categorical variables, a regression model is equivalent to fitting the same distribution on subgroups defined with respect to explanatory variables. Hence, we get back to usual explicit solutions for the exponential family in the i.i.d. case.

Yet we work with one or two explanatory variables for the two derived theorems, this is not a limit because the general setting of d categorical variables can be rewritten as a single categorical variable defined as the observed combination of the d variables.

The explicit formulas are exemplified on two particular positive distributions particularly useful in an insurance context: the Pareto 1 distribution and the shifted log-normal distribution. In both cases, we present typical link functions and derive in most cases the distribution of the MLE. In relevant cases, we also give an unbiased estimator. Finally, we illustrate the estimation process for both distributions on an actuarial data set.

For future research, a natural extension is to propose regression models for distribution outside the exponential family. Typically we could consider generalized Pareto distribution based on the peak over thresholds approach. A second natural extension could also be to jointly estimate the threshold μ and the parameters of the distribution.

Acknowledgments

The authors thank Vanessa Desert for her active support during the writing of this paper. This work is supported by the research project “PANORisk” and Région Pays de la Loire.

A. Proofs of Section 3

Proof of Theorem 3.1. We have to solve the system

$$\begin{cases} S(\boldsymbol{\vartheta}) = 0 \\ \mathbf{R}\boldsymbol{\vartheta} = 0. \end{cases} \quad (\text{A.1})$$

The system $S(\boldsymbol{\vartheta}) = 0$ is

$$\begin{cases} \sum_{i=1}^n \ell'(\eta_i) (x_i - b' \circ \ell(\eta_i)) = 0 \\ \sum_{i=1}^n y_i^{(2),j} \ell'(\eta_i) (x_i - b' \circ \ell(\eta_i)) = 0, \quad \forall j \in J. \end{cases}$$

that is

$$\begin{cases} \sum_{j \in J} \ell'(\vartheta_{(1)} + \vartheta_{(2),j}) \left(\sum_{i=1}^n y_i^{(2),j} x_i - m_j b' \circ \ell(\vartheta_{(1)} + \vartheta_{(2),j}) \right) = 0 \\ \ell'(\vartheta_{(1)} + \vartheta_{(2),j}) \left(\sum_{i=1}^n y_i^{(2),j} x_i - m_j b' \circ \ell(\vartheta_{(1)} + \vartheta_{(2),j}) \right) = 0, \quad \forall j \in J. \end{cases}$$

The first equation in the previous system is redundancy, and

$$S(\boldsymbol{\vartheta}) = 0 \Leftrightarrow \ell'(\vartheta_{(1)} + \vartheta_{(2),j}) \left(\sum_{i=1}^n y_i^{(2),j} x_i - m_j b' \circ \ell(\vartheta_{(1)} + \vartheta_{(2),j}) \right) = 0, \quad \forall j \in J.$$

Hence if X_i takes values in $\mathcal{X} \subset b'(\Lambda)$, and ℓ injective, we have

$$\vartheta_{(1)} + \vartheta_{(j)} = g(\bar{X}_n^{(j)}) \quad \forall j \in J.$$

The system (A.1) is

$$\begin{cases} Q\boldsymbol{\vartheta} = \mathbf{g}(\bar{\mathbf{X}}) \\ \mathbf{R}\boldsymbol{\vartheta} = 0. \end{cases} \Leftrightarrow \begin{pmatrix} Q \\ \mathbf{R} \end{pmatrix} \boldsymbol{\vartheta} = \begin{pmatrix} \mathbf{g}(\bar{\mathbf{X}}) \\ 0 \end{pmatrix}. \quad (\text{A.2})$$

Let us compute the determinant of the matrix $M_d = \begin{pmatrix} Q \\ \mathbf{R} \end{pmatrix}$. Consider $\mathbf{R} = (r_0, r_1, \dots, r_d)$.

We have

$$M_d = \begin{pmatrix} \mathbf{1}_d & I_d \\ r_0 & \mathbf{r} \end{pmatrix} = \begin{pmatrix} 1 & 1 & 0 & \dots & \\ 1 & 0 & 1 & 0 & \\ \vdots & \vdots & \ddots & \ddots & \ddots \\ 1 & 0 & \dots & 0 & 1 \\ r_0 & r_1 & & \dots & r_d \end{pmatrix}, \text{ with } \mathbf{r} = (r_1 \ \dots \ r_d), \mathbf{1}_d = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}.$$

The determinant can be computed recursively

$$\det(M_d) = r_d \begin{vmatrix} 1 & 1 & 0 & \dots \\ 1 & 0 & \ddots & 0 \\ \vdots & \vdots & \ddots & 1 \\ 1 & 0 & \dots & 0 \end{vmatrix} - \begin{vmatrix} 1 & 1 & 0 & \dots \\ 1 & 0 & \ddots & 0 \\ \vdots & \vdots & \ddots & 1 \\ r_0 & r_1 & \dots & r_{d-1} \end{vmatrix} = (-1)^{d+1} r_d - \det(M_{d-1})$$

Since $\det(M_1) = -r_0 + r_1$ and $\det(M_2) = -r_2 - (-r_0 + r_1) = r_0 - r_1 - r_2$, we get $\det(M_d) = (-1)^d r_0 + (-1)^{d+1}(r_1 + \dots + r_d) = (-1)^d(r_0 - r_1 - \dots - r_d)$. This determinant is non zero as long as $r_0 \neq \sum_{j=1}^d r_j$.

Now we compute the inverse of matrix M_d by a direct inversion.

$$\begin{pmatrix} \mathbf{1}_d & I_d \\ r_0 & \mathbf{r} \end{pmatrix} \begin{pmatrix} \mathbf{a}' & b \\ C & \mathbf{d} \end{pmatrix} = \begin{pmatrix} I_d & \mathbf{0} \\ \mathbf{0}' & 1 \end{pmatrix} \Leftrightarrow \begin{cases} \mathbf{1}_d \mathbf{a}' + I_d C = I_d \\ b \mathbf{1}_d + I_d \mathbf{d} = \mathbf{0} \\ r_0 \mathbf{a}' + \mathbf{r} C = \mathbf{0}' \\ b r_0 + \mathbf{r} \mathbf{d} = 1 \end{cases} \Leftrightarrow \begin{cases} C = I_d - \frac{1}{-r_0 + \mathbf{r} \mathbf{1}_d} \mathbf{1}_d \mathbf{r}' \\ \mathbf{d} = \frac{1}{-r_0 + \mathbf{r} \mathbf{1}_d} \mathbf{1}_d \\ \mathbf{a}' = \frac{\mathbf{r}}{-r_0 + \mathbf{r} \mathbf{1}_d} \\ b = \frac{-1}{-r_0 + \mathbf{r} \mathbf{1}_d} \end{cases}$$

Let us check the inverse of M_d

$$\begin{pmatrix} \mathbf{1}_d & I_d \\ r_0 & \mathbf{r} \end{pmatrix} \begin{pmatrix} \frac{\mathbf{r}}{-r_0 + \mathbf{r} \mathbf{1}_d} & \frac{-1}{-r_0 + \mathbf{r} \mathbf{1}_d} \\ I_d - \frac{\mathbf{1}_d \mathbf{r}}{-r_0 + \mathbf{r} \mathbf{1}_d} & \frac{\mathbf{1}_d}{-r_0 + \mathbf{r} \mathbf{1}_d} \end{pmatrix} = \begin{pmatrix} \frac{\mathbf{1}_d \mathbf{r}}{-r_0 + \mathbf{r} \mathbf{1}_d} + I_d - \frac{\mathbf{1}_d \mathbf{r}}{-r_0 + \mathbf{r} \mathbf{1}_d} & \frac{-\mathbf{1}_d}{-r_0 + \mathbf{r} \mathbf{1}_d} + \frac{\mathbf{1}_d}{-r_0 + \mathbf{r} \mathbf{1}_d} \\ r_0 \frac{\mathbf{r}}{-r_0 + \mathbf{r} \mathbf{1}_d} + \mathbf{r} - \frac{\mathbf{r} \mathbf{1}_d \mathbf{r}}{-r_0 + \mathbf{r} \mathbf{1}_d} & \frac{-r_0}{-r_0 + \mathbf{r} \mathbf{1}_d} + \frac{\mathbf{r} \mathbf{1}_d}{-r_0 + \mathbf{r} \mathbf{1}_d} \end{pmatrix} = \begin{pmatrix} I_d & \mathbf{0} \\ \mathbf{0} & 1 \end{pmatrix}$$

So as long as $r_0 \neq \sum_{j=1}^d r_j$

$$\widehat{\boldsymbol{\vartheta}}_n = \begin{pmatrix} \frac{\mathbf{r}}{-r_0 + \mathbf{r} \mathbf{1}_d} & \frac{-1}{-r_0 + \mathbf{r} \mathbf{1}_d} \\ I_d - \frac{\mathbf{1}_d \mathbf{r}}{-r_0 + \mathbf{r} \mathbf{1}_d} & \frac{\mathbf{1}_d}{-r_0 + \mathbf{r} \mathbf{1}_d} \end{pmatrix} \begin{pmatrix} \mathbf{g}(\bar{\mathbf{X}}) \\ 0 \end{pmatrix} = \begin{pmatrix} \frac{\mathbf{r} \mathbf{g}(\bar{\mathbf{X}})}{-r_0 + \mathbf{r} \mathbf{1}_d} \\ \mathbf{g}(\bar{\mathbf{X}}) - \mathbf{1}_d \frac{\mathbf{r} \mathbf{g}(\bar{\mathbf{X}})}{-r_0 + \mathbf{r} \mathbf{1}_d} \end{pmatrix}.$$

In an other way, the system (A.2) is equivalent to

$$\begin{pmatrix} Q' & \mathbf{R}' \end{pmatrix} \begin{pmatrix} Q \\ \mathbf{R} \end{pmatrix} \boldsymbol{\vartheta} = Q' \mathbf{g}(\bar{\mathbf{X}}),$$

and for $(Q \ \mathbf{R})$ of full rank, the matrix $(Q'Q + \mathbf{R}'\mathbf{R})$ is invertible and $\boldsymbol{\vartheta} = (Q'Q + \mathbf{R}'\mathbf{R})^{-1} Q' \mathbf{g}(\bar{\mathbf{X}})$. \square

Examples - Choice of the contrast vector R

1. Taking $r_0 = 1, \mathbf{r} = \mathbf{0}$ leads to

$$-r_0 + \mathbf{r} \mathbf{1}_d = -1 \Rightarrow \widehat{\boldsymbol{\vartheta}}_n = \begin{pmatrix} 0 \\ \mathbf{g}(\bar{\mathbf{X}}) \end{pmatrix}$$

2. Taking $r_0 = 0, \mathbf{r} = (1, \mathbf{0})$ leads to

$$-r_0 + \mathbf{r} \mathbf{1}_d = 1 \Rightarrow \widehat{\boldsymbol{\vartheta}}_n = \begin{pmatrix} \mathbf{g}(\bar{\mathbf{X}}_n^{(1)}) \\ 0 \\ \mathbf{g}(\bar{\mathbf{X}}_n^{(2)}) - \mathbf{g}(\bar{\mathbf{X}}_n^{(1)}) \\ \vdots \\ \mathbf{g}(\bar{\mathbf{X}}_n^{(d)}) - \mathbf{g}(\bar{\mathbf{X}}_n^{(1)}) \end{pmatrix}.$$

3. Taking $r_0 = 0$, $\mathbf{r} = \mathbf{1}$ leads to

$$-r_0 + \mathbf{r}\mathbf{1}_d = d \Rightarrow \widehat{\boldsymbol{\vartheta}}_n = \begin{pmatrix} \overline{\mathbf{g}(\bar{\mathbf{X}})} \\ g(\bar{X}_n^{(1)}) - \overline{\mathbf{g}(\bar{\mathbf{X}})} \\ \dots \\ g(\bar{X}_n^{(d)}) - \overline{\mathbf{g}(\bar{\mathbf{X}})} \end{pmatrix}$$

$$\text{with } \overline{\mathbf{g}(\bar{\mathbf{X}})} = \frac{1}{d} \sum_{j=1}^d g(\bar{X}_n^{(j)}).$$

Proof of Corollaries 3.1. The log likelihood of $\widehat{\boldsymbol{\vartheta}}_n$ is defined by

$$\log L(\widehat{\boldsymbol{\vartheta}}_n | \mathbf{x}) = \frac{1}{a(\phi)} \sum_{i=1}^n (x_i \ell(\widehat{\eta}_i) - b(\ell(\widehat{\eta}_i))) + \sum_{i=1}^n c(x_i, \phi).$$

In fact, we must be verified that $\ell(\widehat{\eta}_i)$ does not depend on g function. If we consider $\widehat{\boldsymbol{\vartheta}}_n$ defined by (8), we have $Q\widehat{\boldsymbol{\vartheta}}_n = \mathbf{g}(\bar{\mathbf{x}})$, since $\widehat{\boldsymbol{\vartheta}}_n$ is solution of the system (A.1), i.e. $Q(Q'Q + R'R)^{-1}Q' = I$ Using $\widehat{\eta}_i = (Q\widehat{\boldsymbol{\vartheta}}_n)_j$ for i such that $y_i^{(2),j} = 1$ we obtain

$$\ell(\widehat{\eta}_i) = \sum_{j=1}^d \ell \circ g(\bar{x}_n^{(j)}) y_i^{(2),j} = \sum_{j=1}^d \ell \circ \ell^{-1} \circ (b')^{-1}(\bar{x}_n^{(j)}) y_i^{(2),j} = \sum_{j=1}^d (b')^{-1}(\bar{x}_n^{(j)}) y_i^{(2),j},$$

and

$$\log L(\widehat{\boldsymbol{\vartheta}}_n | \mathbf{x}) = \frac{1}{a(\phi)} \sum_{j=1}^d \sum_{i, y_i^{(2),j} = 1} \left(x_i (b')^{-1}(\bar{x}_n^{(j)}) - b\left((b')^{-1}(\bar{x}_n^{(j)})\right) \right) + \sum_{i=1}^n c(x_i, \phi).$$

In the same way,

$$\widehat{\mathbf{E}}X_i = b'(\ell(\widehat{\eta}_i)) = \sum_{j=1}^d \bar{x}_n^{(j)} y_i^{(2),j} \quad \widehat{\mathbf{Var}}X_i = a(\phi) b''(\ell(\widehat{\eta}_i)) = a(\phi) \sum_{j=1}^d b'' \circ (b')^{-1}(\bar{x}_n^{(j)}) y_i^{(2),j}$$

□

Proof of Theorem 3.2. The system $S(\boldsymbol{\vartheta}) = 0$ is

$$\left\{ \begin{array}{l} \sum_{i=1}^n \ell'(\eta_i) (x_i - b' \circ \ell(\eta_i)) = 0 \\ \sum_{i=1}^n y_i^{(3),l} \ell'(\eta_i) (x_i - b' \circ \ell(\eta_i)) = 0, \quad \forall l \in L \\ \sum_{i=1}^n y_i^{(2),k} \ell'(\eta_i) (x_i - b' \circ \ell(\eta_i)) = 0, \quad \forall k \in K \\ \sum_{i=1}^n y_i^{kl} \ell'(\eta_i) (x_i - b' \circ \ell(\eta_i)) = 0, \quad \forall (k, l) \in KL^*. \end{array} \right.$$

that is

$$\left\{ \begin{array}{l} \sum_{(k,l) \in KL^*} \ell'(\vartheta_{(1)} + \vartheta_{(2),k} + \vartheta_{(3),l} + \vartheta_{kl}) \left(\sum_{i=1}^n y_i^{(k,l)} x_i - m_{kl} b' \circ \ell(\vartheta_{(1)} + \vartheta_{(2),k} + \vartheta_{(3),l} + \vartheta_{kl}) \right) = 0 \\ \sum_{k \in K_l^*} \ell'(\vartheta_{(1)} + \vartheta_{(2),k} + \vartheta_{(3),l} + \vartheta_{kl}) \left(\sum_{i=1}^n y_i^{(k,l)} x_i - m_{kl} b' \circ \ell(\vartheta_{(1)} + \vartheta_{(2),k} + \vartheta_{(3),l} + \vartheta_{kl}) \right) = 0 \quad \forall l \in L \\ \sum_{l \in L_k^*} \ell'(\vartheta_{(1)} + \vartheta_{(2),k} + \vartheta_{(3),l} + \vartheta_{kl}) \left(\sum_{i=1}^n y_i^{(k,l)} x_i - m_{kl} b' \circ \ell(\vartheta_{(1)} + \vartheta_{(2),k} + \vartheta_{(3),l} + \vartheta_{kl}) \right) = 0 \quad \forall k \in K \\ \ell'(\vartheta_{(1)} + \vartheta_{(2),k} + \vartheta_{(3),l} + \vartheta_{kl}) \left(\sum_{i=1}^n y_i^{(k,l)} x_i - m_{kl} b' \circ \ell(\vartheta_{(1)} + \vartheta_{(2),k} + \vartheta_{(3),l} + \vartheta_{kl}) \right) = 0 \quad \forall (k,l) \in KL^* \end{array} \right.$$

The system have exactly $1 + d_2 + d_3$ redundancies, and $S(\boldsymbol{\vartheta}) = 0$ is

$$\ell'(\vartheta_{(1)} + \vartheta_{(2),k} + \vartheta_{(3),l} + \vartheta_{kl}) \left(\sum_{i=1}^n y_i^{(k,l)} x_i - m_{kl} b' \circ \ell(\vartheta_{(1)} + \vartheta_{(2),k} + \vartheta_{(3),l} + \vartheta_{kl}) \right) = 0 \quad \forall (k,l) \in KL^*.$$

Hence the system has rank KL^* and if X_i takes values in $\mathcal{X} \subset b'(\Lambda)$, and ℓ injective, we have

$$\vartheta_{(1)} + \vartheta_{(2),k} + \vartheta_{(3),l} + \vartheta_{kl} = g(\bar{X}_n^{(k,l)}) \quad \forall (k,l) \in KL^*.$$

In the same way of proof of theorem 1, we have to solve

$$\begin{cases} Q\boldsymbol{\vartheta} = \mathbf{g}(\bar{\mathbf{X}}) \\ R\boldsymbol{\vartheta} = \mathbf{0}. \end{cases}$$

that is, because $QQ' + RR'$ is full rank, in the same way of proof of Theorem 1

$$\boldsymbol{\vartheta} = (Q'Q + R'R)^{-1} Q' \mathbf{g}(\bar{\mathbf{X}}).$$

1. Under linear contrasts (\tilde{C}_0), the model (10) is equivalent to model (6) with $J = KL^*$ modalities. Hence the solution is evident.
2. Under linear contrasts (\tilde{C}_Σ), the system

$$\vartheta_{(1)} + \vartheta_{(2),k} + \vartheta_{(3),l} + \vartheta_{kl} = g(\bar{X}_n^{(k,l)}) \quad \forall (k,l) \in KL^*$$

implies that

$$\sum_{(k,l) \in KL^*} m_{kl} (\vartheta_{(1)} + \vartheta_{(2),k} + \vartheta_{(3),l} + \vartheta_{kl}) = \sum_{(k,l) \in KL^*} m_{kl} g(\bar{X}_n^{(k,l)}).$$

Using

$$\sum_{(k,l) \in KL^*} m_{kl} = n, \quad \sum_{(k,l) \in KL^*} m_{kl} \vartheta_{(2),k} = \sum_{k \in K} \sum_{l \in L_k^*} m_{kl} \vartheta_{(2),k} = \sum_{k \in K} m_k^{(2)} \vartheta_{(2),k} = 0,$$

$$\sum_{(k,l) \in KL^*} m_{kl} \vartheta_{(3),l} = \sum_{l \in L} \sum_{k \in K_l^*} m_{kl} \vartheta_{(3),l} = \sum_{l \in L} m_l^{(3)} \vartheta_{(3),l} = 0, \quad \sum_{(k,l) \in KL^*} m_{kl} \vartheta_{kl} = 0,$$

we have $\vartheta_{(1)} = \frac{1}{n} \sum_{(k,l) \in KL^*} m_{kl} g(\bar{X}_n^{(k,l)})$. In the same way, taking summation on K_l^* for $l \in L$ and on L_k^* for $k \in K$, we found $\vartheta_{(2),k}$ and $\vartheta_{(3),l}$, and then ϑ_{kl} .

□

B. Calculus of the Log-likelihoods appearing in Section 4 and 5

Consider the Pareto GLM described on (13) and (15). The b function is $b(\lambda) = -\log(\lambda)$, using corollary 3.1 we have $\ell(\hat{\eta}_i) = (b')^{-1}(\bar{z}_n^{(j)}) = -(\bar{z}_n^{(j)})^{-1}$ for j such that $y_i^{(2),j} = 1$ and

$$\begin{aligned} \log L(\hat{\boldsymbol{\vartheta}}_n | \underline{\mathbf{z}}) &= \sum_{j=1}^d \sum_{i, y_i^{(2),j}=1} (z_i / \bar{z}_n^{(j)} - \log(-\bar{z}_n^{(j)})) \\ &= n - \sum_{j=1}^d m_j \log(-\bar{z}_n^{(j)}) \end{aligned}$$

Compute the original log likelihood of Pareto 1 distribution:

$$\log L(\boldsymbol{\vartheta} | \underline{\mathbf{x}}) = \sum_{i=1}^n (\log \ell(\eta_i) + \ell(\eta_i) \log \mu - (\ell(\eta_i) + 1) \log x_i).$$

Hence with $z_i = -\log(x_i/\mu)$,

$$\begin{aligned} \log L(\hat{\boldsymbol{\vartheta}}_n | \underline{\mathbf{x}}) &= \sum_{j=1}^d \sum_{i, y_i^{(2),j}=1} \left(-\log(-\bar{z}_n^{(j)}) - \frac{\log \mu}{\bar{z}_n^{(j)}} + \frac{\log(x_i)}{\bar{z}_n^{(j)}} - \log x_i \right) \\ &= n - \sum_{j=1}^d m_j \log(-\bar{z}_n^{(j)}) - \sum_{i=1}^n \log x_i \end{aligned}$$

Now consider the shifted log-normal GLM described on (18) and (19). Here, the b function is $b(\lambda) = \lambda^2/2$, hence using Corollary 3.1, we have $\ell(\hat{\eta}_i) = (b')^{-1}(\bar{z}_n^{(j)}) = \bar{z}_n^{(j)}$ for j such that $y_i^{(2),j} = 1$ and equation (21) holds.

Let us compute the original log likelihood of the shifted log normal distribution:

$$\begin{aligned} \log L(\boldsymbol{\vartheta} | \underline{\mathbf{x}}) &= \sum_{i=1}^n \left(-\log(x_i - \mu) - \log(\sqrt{2\pi\phi}) - \frac{(\log(x_i - \mu) - \ell(\eta_i))^2}{2\phi} \right) \\ &= -\sum_{i=1}^n z_i - n \log(\sqrt{2\pi\phi}) - \sum_{i=1}^n \frac{(z_i - \ell(\eta_i))^2}{2\phi}, \end{aligned}$$

with $z_i = \log(x_i - \mu)$. Hence

$$\log L(\hat{\boldsymbol{\vartheta}} | \mathbf{x}) = -\sum_{i=1}^n z_i - n \log(\sqrt{2\pi\hat{\phi}}) - \frac{1}{2\hat{\phi}} \sum_{j=1}^d \sum_{i, y_i^{(2), j}=1} (z_i - \bar{z}_n^{(j)})^2$$

Using $\hat{\phi} = \frac{1}{n} \sum_{j \in J} \sum_{i, y_i^{(2), j}=1} (z_i - \bar{z}_n^{(j)})^2$ we have

$$\log L(\hat{\boldsymbol{\vartheta}} | \mathbf{x}) = -\sum_{i=1}^n z_i - \frac{n}{2} \log(2\pi\hat{\phi}) - \frac{n}{2}$$

C. Link functions and descriptive statistics

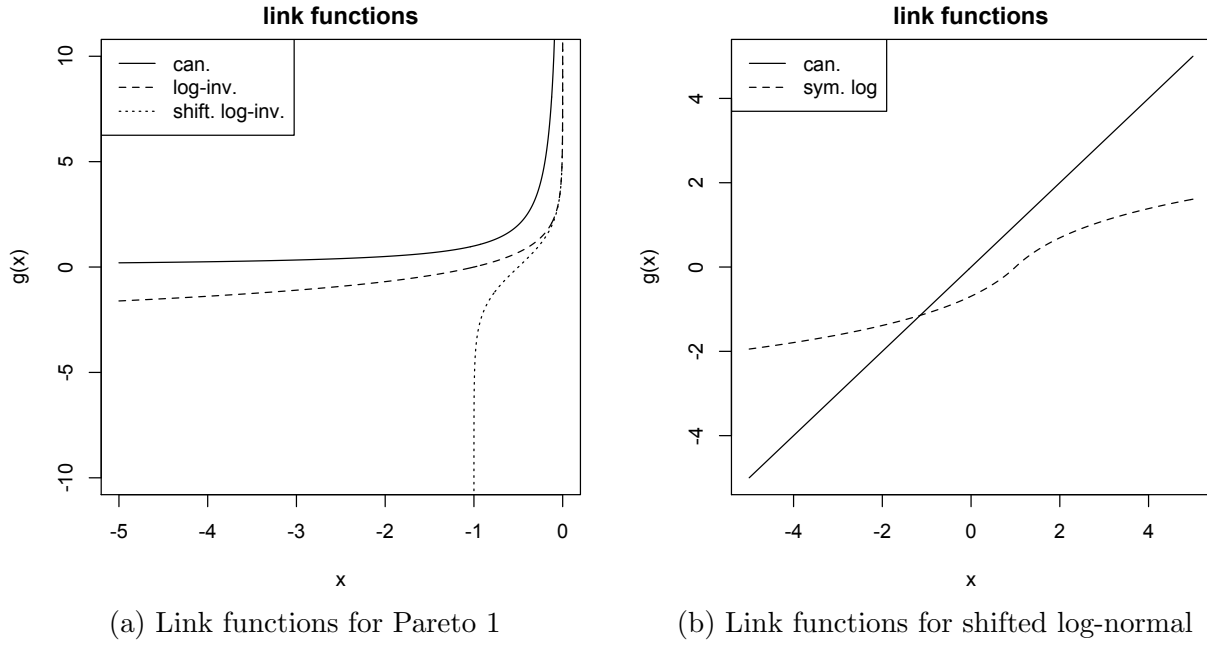


Figure C.3: Graphs of link functions

Table C.9: Empirical quantiles and moments (in euros)

	Amount	Risk class 1	Risk class 2	Risk class 3	Risk class 4	Risk class 5
Min.	0	0	0	0	0	0
1st Qu.	51	140	39	130	150	27
Median	761	1120	652	1073	1015	737
3rd Qu.	3,003	4,169	2,474	4,486	4,155	3,113
Max.	15,688,300	15,315,173	15,688,300	11,916,121	6,078,593	10,833,825
Mean	10,745	14,508	7,265	28,082	18,193	11,179
Std dev.	128,146	148,380	98,141	275,175	140,607	125,004
Skewness	54	48	96	24	24	38
Kurtosis	4,473	3,933	12,753	751	796	2,124

	Guarantee 1	Guarantee 2	Guarantee 3	Guarantee 4	Guarantee 5	Guarantee 6	Guarantee 7
Min.	0	0	0	0	0	0	0
1st Qu.	123	155	235	128	2	1	2
Median	1,253	814	1,955	893	2,977	2	564
3rd Qu.	4,994	2,664	8,246	3,726	39,647	1,560	2,097
Max.	3,882,524	4,529,249	15,315,173	14,272,522	15,688,300	4,888,656	4,670,686
Mean	7,022	4,055	28,429	32,328	110,056	8,388	7,157
Std dev.	39,581	24,620	280,500	273,958	534,337	74,969	60,916
Skewness	49	85	38	22	16	42	35
Kurtosis	3,955	13,620	1,833	738	366	2,399	1,927

References

- J. Beirlant and Y. Goegebeur. Regression with response distributions of pareto-type. *Computational statistics & data analysis*, 42(4):595–619, 2003. 2
- J. Beirlant, Y. Goegebeur, R. Verlaak, and P. Vynckier. Burr regression and portfolio segmentation. *Insurance: Mathematics and Economics*, 23(3):231–250, 1998. 1
- J. Beirlant, Y. Goegebeur, J. Teugels, and J. Segers. *Statistics of extremes: Theory and applications*. Wiley & Sons, 2004. 1
- H. Bühlmann and A. Gisler. *A course in credibility theory and its applications*. Springer Science & Business Media, 2006. 10, 12
- V. Chavez-Demoulin, P. Embrechts, and M. Hofert. An extreme value approach for modeling operational risk losses depending on covariates. *Journal of Risk and Insurance*, 2015. 2
- A. Davison and R. Smith. Models for exceedances over high thresholds. *Journal of the Royal Statistical Society. Series B*, 52(3):393–442, 1990. 1
- L. Fahrmeir and H. Kaufmann. Consistency and asymptotic normality of the maximum likelihood estimator in generalized linear models. *The Annals of Statistics*, pages 342–368, 1985. 1, 4
- J. Hambuckers, C. Heuchenne, and O. Lopez. A semiparametric model for generalized pareto regression based on a dimension reduction assumption. HAL, 2016. URL <https://hal.archives-ouvertes.fr/hal-01362314/>. 2
- R. V. Hogg and S. A. Klugman. *Loss distributions*. John Wiley & Sons, 1984. 12
- N. Johnson, S. Kotz, and N. Balakrishnan. *Continuous Univariate Distributions*, volume 1. Wiley, 2nd edition, 2000. 11, 14

- P. McCullagh and J. A. Nelder. Generalized linear models, volume 37. CRC press, 1989. 1, 3, 4, 17
- J. A. Nelder and R. W. M. Wedderburn. Generalized linear models. Journal of the Royal Statistical Society. Series A, 135(3):370–384, 1972. 1
- E. Ohlsson and B. Johansson. Non-Life Insurance Pricing with Generalized Linear Models. Springer, 2010. 2
- F. W. J. Olver, D. W. Lozier, R. F. Boisvert, and C. W. Clark, editors. NIST Handbook of Mathematical Functions. Cambridge University Press, 2010. URL <http://dlmf.nist.gov/>. 12
- E. Ozkok, G. Streftaris, H. R. Waters, and A. D. Wilkie. Bayesian modelling of the time delay between diagnosis and settlement for critical illness insurance using a burr generalised-linear-type model. Insurance: Mathematics and Economics, 50(2):266 – 279, 2012. ISSN 0167-6687. doi: <https://doi.org/10.1016/j.insmatheco.2011.12.001>. URL <http://www.sciencedirect.com/science/article/pii/S0167668711001326>. 2
- R Core Team. R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing, Vienna, Austria, 2017. URL <https://www.R-project.org/>. 6
- R Core Team. R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing, Vienna, Austria, 2018. URL <https://www.R-project.org/>. 16
- R. Reiss and M. Thomas. Statistical Analysis of Extreme Values. Basel: Birkhauser, 3rd edition, 2007. 16
- R. Rigby and D. Stasinopoulos. Generalized additive models for location, scale and shape. Applied Statistics, 54(3):507–554, 2005. 2
- W. Venables and B. Ripley. Modern Applied Statistics with S. Springer, 2002. 4