

Nonlinear Regression Models with Applications in Insurance

Rastislav Potocký and Milan Stehlík*

Department of Applied Mathematics and Statistics, Faculty of Mathematics, Physics and Informatics, Comenius University, Mlynská dolina, 842 48 Bratislava 4, Slovak Republic

Abstract: Two possible applications of nonlinear regression models in insurance are discussed. The first part deals with modelling IBNR reserves when a cubic approximation to the solution locus is used instead of linear or quadratic ones. A formula is given for construction of improved confidence regions for parameters in such models. Using this approach IBNR reserves for a data set are computed. In the second part a method is proposed of how to measure the influence of additive perturbations on nonlinear regression model parameters. An example is given which shows how this method can be used to preserve privacy of sensitive data in insurance business.

Keywords: Nonlinear regression models, confidence regions, additive perturbations, IBNR reserves, privacy of data.

1. INTRODUCTION

Due to the latest developments in insurance and reinsurance industries nonlinear regression models became more popular with steadily increasing importance. Let us mention two up-to date applications: nonlinear credibility estimation and nonlinear modelling of IBNR reserves (i.e. reserves for incurred but not reported losses).

De Vylder [1] extended Hachemeister's linear regression credibility model to a nonlinear regression model by assuming that observations are an arbitrary function $f(\beta(\theta))$ of the unknown vector $\beta(\theta)$. This model lacks of robustness of credibility estimators. Pitselis [2] applied robust statistics to De Vylder's nonlinear credibility estimation and presented an application to Hachemeister's data.

The process of IBNR modelling and estimation has been studied by actuaries for many years since both RBNS (reported but not settled) and IBNR reserves are the largest liabilities of insurance companies. Recall that incurred but not reported loss or IBNR is the difference between ultimate loss and reported loss. Thus quantifying the uncertainty in estimation of IBNR plays the important role in insurance business. The classical approach makes use of run-off triangles (e.g. the chain-ladder method and its modifications). There exist also direct methods for computing or modelling of IBNR, for instance those using copulas and indirect methods based on estimation of loss development factors. For instance, Stelltjes [3] presents a model for predicting losses as a function of exposures, calendar period and development age. Typically, then a nonlinear regression model is used for estimating the 95% confidence interval of IBNR loss for an accident period. However using the quadratic approximation

of the model function may lead to inaccurate confidence regions for parameters of the model.

Consider the usual nonlinear regression model

$$y_a = f(x_a, \theta) + e_a \quad a = 1, \dots, n \quad (1)$$

$Y = (y_1, \dots, y_n)'$ denotes the vector of observations, the function $f(\theta) = (f(x_1, \theta), \dots, f(x_n, \theta))$ has a known form dependent on p unknown parameters $\theta = (\theta_1, \dots, \theta_p)'$, x_a are known vector-valued variables, e_a are independent and normally distributed random errors with zero mean and variance σ^2 .

The problem of finding acceptable confidence regions for θ has been discussed by many authors, see e.g. [4-6]. Cook and Goldberg [7] showed examples of models for which the Bates-Watts methodology based on quadratic approximations did not work. On the other hand Clarke [8] presented a method of constructing regions with higher precision. However, his investigations deal with a single parameter θ_i , not with the whole vector θ . The aim of our paper is to construct confidence regions based on cubic approximation which are more accurate than those currently used.

In what follows pre- and post-, square bracket and \otimes -multiplications of a three-dimensional array $U_{ij}^a = (U_{ij}^a)$ or a four-dimensional array $U_{ijk}^a = (U_{ijk}^a)$ by a matrix E mean summation over the indexes i, j, a, k , respectively. The reader is recommended to consult Table 1 for better understanding of these operations. Recall that if $U_{ij}^a = (U_{ij}^a)$, $a=1, \dots, n$, $i=1, \dots, p$, $j=1, \dots, q$ is an $n \times p \times q$ array, then its a -th face is $p \times q$ matrix (U_{ij}^a) and its ij -th column $(U_{ij}^1, \dots, U_{ij}^n)'$ is n -vector.

*Address correspondence to this author at the Institut für angewandte Statistik, Johannes Kepler University in Linz, Freistädter Strasse 315, 2. Stock, A-4040 Linz a. D., Austria; Tel: +43 732 2468 5881; Fax: +43 732 2468 9846; E-mails: mlnstehlik@gmail.com, Milan.Stehlik@jku.at

Table 1. Rules for Multiplication

$U_{..}^* - n \times p \times q ; E - s \times p$	$EU_{..}^* - n \times s \times q$	$(EU_{..}^*)_{ij}^a = \sum_t E_{it} U_{tj}^a$
$U_{..}^* - n \times p \times q ; E - q \times s$	$U_{..}^* E - n \times p \times s$	$(U_{..}^* E)_{ij}^a = \sum_t U_{it}^a E_{tj}$
$U_{..}^* - n \times p \times q ; E - s \times n$	$[E] [U_{..}^*] - s \times p \times q$	$([E][U_{..}^*])_{ij}^a = \sum_t E_{at} U_{tj}^a$
$U_{...}^* - n \times r \times p \times q ; E - s \times r$	$E \otimes U_{...}^* - n \times s \times p \times q$	$(E \otimes U_{...}^*)_{ijk}^a = \sum_t E_{it} U_{tjk}^a$

Put

$$V = \left(\frac{\partial f_a}{\partial \theta_i} \right)_{\theta = \hat{\theta}}, \quad V_{..}^* = \left(\frac{\partial^2 f_a}{\partial \theta_i \partial \theta_j} \right)_{\theta = \hat{\theta}} \quad \text{and} \quad V_{...}^* = \left(\frac{\partial^3 f_a}{\partial \theta_i \partial \theta_j \partial \theta_k} \right)_{\theta = \hat{\theta}}$$

where $a=1, \dots, n$, $i, j, k=1, \dots, p$ and $\hat{\theta}$ is the least-squares estimate of θ . Let $V = UR$ be the unique QR-decomposition of V where R is an upper triangular matrix and the columns of U form an orthogonal basis for the column space of V .

Recall that by the Bates-Watts parameter-effects array we understand the $p \times p \times p$ array

$$A_{..}^* = [U^T] \left[(R^{-1})^T V_{..}^* R^{-1} \right] \quad \text{and} \quad \text{similarly} \quad A_{...}^* = [U^T] \left[R^{-1} \otimes \left((R^{-1})^T V_{...}^* R^{-1} \right) \right]$$

is the four dimensional parameter-effects array (for details see [9]). Bates and Watts have used $A_{..}^*$ to compute the maximum parameter-effects curvature which, together with the maximum intrinsic curvature, give answer to the question whether the model may be considered sufficiently linear or not for the set of values of parameters we are interested in. As mentioned above the quadratic approximation of the model function does not give satisfactory results in some cases, so a cubic approximation is needed..

Consider a partition $\theta = (\theta_{(1)}^T, \theta_{(2)}^T)^T$, $\theta_{(1)} = (\theta_1, \dots, \theta_q)^T$, $\theta_{(2)} = (\theta_{q+1}, \dots, \theta_p)^T$, where $\theta_{(2)}$ is the parameter of interest, $\theta_{(1)}$ being a nuisance parameter. Let $S(\theta)$ be the usual sum of squares for the model (1) and $g(\theta_{(2)}) = (g_1(\theta_{(2)}), \dots, g_q(\theta_{(2)}))^T$ be the value of $\theta_{(1)}$ which minimizes $S(\theta_{(1)}) = S(\theta_{(1)}, \theta_{(2)})$ for given $\theta_{(2)}$. Put $\alpha(\theta_{(2)}) = (g(\theta_{(2)}), \theta_{(2)})$ and $h(\theta_{(2)}) = f(\alpha(\theta_{(2)}))$.

The approximate confidence region based on the likelihood ratio (the likelihood region for short) is such a set of values $\theta_{(2)}$ for which

$$\left(Y - h(\theta_{(2)}) \right)^T \left(Y - h(\theta_{(2)}) \right) - S(\hat{\theta}) \leq c^2 \rho^2, \tag{2}$$

where $\rho = \sigma$ and $c^2 = \chi^2(p - q, \delta)$ if σ is known or $\rho^2 = (p - q) s^2$ and $c^2 = F(p - q, n - p, \delta)$ if σ^2 is estimated by $s^2 = S(\hat{\theta}) / (n - p)$.

In the first part the solution of the above-mentioned problem for parameter vector is given which is then compared to a result of Clarke [8]. Next we discuss the effect of perturbation of the response vector on the values of parameters in a nonlinear model. In the final part our results are applied to the problem of finding acceptable reserves for claims in insurance business and the problem of preserving privacy of stored data.

2. APPROXIMATE CONFIDENCE REGIONS

We assume that

- 1) $f(\theta)$ is a continuous function of θ with finite derivatives up to and including degree 4
- 2) the vector $\frac{\partial S(\theta)}{\partial \theta}$ vanishes at one point $\hat{\theta}$ where S is the corresponding sum of squares
- 3) the Bates-Watts intrinsic curvature as well as $\left(\hat{e} \right)^T (\varphi)^T \otimes \varphi^T V_{...}^* \varphi$ and $\hat{e}^T (\varphi^T \otimes (\varphi^T V_{...}^* \varphi))$ can be neglected where $\varphi = \theta - \hat{\theta}$ and $\hat{e} = Y - f(\hat{\theta})$.

Following the argument in [9] and using the cubic approximation of $h(\theta_{(2)})$ we rewrite (2) as

$$\begin{aligned} & \eta_2^T \eta_2 + \eta_2^T (\eta_2^T A_{22}^2 \eta_2) + \\ & \frac{1}{3} \eta_2^T (\eta_2^T \otimes (\eta_2^T A_{222}^2 \eta_2)) - \eta_2^T [\eta_2^T] [A_{21}^2] (\eta_2^T A_{22}^1 \eta_2) - \\ & \| [\eta_2^T] [A_{12}^2] \eta_2 \|^2 + \frac{1}{4} \| \eta_2^T A_{22}^2 \eta_2 \|^2 + \dots \leq c^2 \rho^2, \end{aligned} \tag{3}$$

where $\eta_2 = R_{22}(\theta_{(2)} - \hat{\theta}_{(2)})$ and the neglected terms are of degree 4 and higher in η .

Put $\Psi_2 = \lambda b, \Psi_2^* = \lambda^* b$, where b is a unit vector of R^q , $\lambda, \lambda^* \geq 0$, Ψ_2 and Ψ_2^* are boundary points of the region (3) and the L -region $\eta_2^T \eta_2 \leq c^2 \rho^2$, respectively. Inserting $\Psi_2 = \lambda b$ into (3)

and expressing λ as a power series in $c\rho$ we obtain

Theorem 1. Under the above assumptions

$$\lambda(b) = \lambda^*(b) \left\{ 1 - \frac{c\rho}{2} \Gamma(b) + \frac{c^2 \rho^2}{2} B(b) + \dots \right\} \quad (4)$$

where b is a unit vector of R^q and

$$\begin{aligned} \Gamma(b) &= b^T [b^T] [A_{22}^2] b, \\ B(b) &= \frac{5}{4} (b^T [b^T] [A_{22}^2] b)^2 - \frac{1}{4} \|b^T A_{22}^2 b\|^2 + \\ &\| [b^T] [A_{12}^2] b \|^2 + b^T [b^T] [A_{21}^2] (b^T A_{22}^1 b) - \\ &\frac{1}{3} b^T \otimes (b^T [b^T] [A_{222}^2] b) \\ \lambda^*(b) &= c\rho \| R_{22}^{-1} b \| \end{aligned}$$

The expression (4) provides a method for construction of the likelihood region (2): by substituting different values for b we obtain the bounds of this region.

If there are no nuisance parameters, we obtain the likelihood region for the whole parameter vector θ .

Theorem 2. Let d be a unit vector in R^p . Then

$$\lambda(d) = \lambda^*(d) \left\{ 1 - (c\rho/2) \Gamma(d) + (c^2 \rho^2 / 2) \beta(d) + \dots \right\} \quad (5)$$

where the neglected terms are of degree 4 and higher in $c\rho$ and

$$\begin{aligned} \lambda^*(d) &= c\rho \| R^{-1} d \| \\ \Gamma(d) &= d^T [d^T] [A_{..}] d \\ \beta(d) &= \frac{5}{4} (d^T [d^T] [A_{..}] d)^2 - \frac{1}{4} \|d^T A_{..} d\|^2 - \\ &\frac{1}{3} d^T \otimes (d^T [d^T] [A_{...}] d) \end{aligned} \quad (6)$$

For every array $U, U_{..}, U_{...}$ let us denote $U_r = (U_{1r}, \dots, U_{nr})^T, U_r = (U_{r1}, \dots, U_{rp})$, $U_r^a =$

$$\begin{aligned} (U_{1r}^a, \dots, U_{pr}^a)^T, U_{rs} = (U_{rs}^1, \dots, U_{rs}^n)^T, \text{ and so on. Put} \\ G = V^T V = (g_{ij}), G^{-1} = (V^T V)^{-1} = (g^{ij}). \end{aligned}$$

For one-dimensional parameter vector (2) becomes the likelihood interval with the bound points defined from the equation

$$\theta_p - \hat{\theta}_p = (g^{pp} \rho^2)^{1/2} \bar{c} \left(1 - \left(\bar{c} \rho / 2 \right) \Gamma + \left(\bar{c} \rho \right)^2 B / 2 \right) \quad (7)$$

where

$$\bar{c} = \pm c$$

$$\Gamma = A_{pp}^p$$

$$B = -\left(A_{pp}^p \right)^2 + \| A_{p.p}^p \|^2 + A_{p.p}^p A_{pp}^p - 1 / 3 A_{ppp}^p \quad (8)$$

It follows from (7) that there will be one value of θ_p for $+c$ and another for $-c$ (c is always taken positive). Two values of θ_p so determined will not be symmetric with respect to $\hat{\theta}_p$.

Clarke [8] proved the similar expression

$$\theta_p - \hat{\theta}_p = (g^{pp} \rho^2)^{1/2} \bar{c} \left(1 - \left(c \rho / 2 \right) \Gamma + \left(c^2 \rho^2 / 2 \right) \Psi + \dots \right) \quad (9)$$

where

$$\Psi = -\Gamma^2 + g^{pp} \left(\left(\Gamma^s \Gamma^s \right) + \left(\Gamma^s \Gamma_s \right) - \frac{1}{3} k \right) \quad (10)$$

with

$$\begin{aligned} \Gamma^s \Gamma^s &= r_{pp}^6 \left(G^{-1} [G^{-1} V^T] [V_{..}] G^{-1} \right)_p \left([G^{-1} V^T] [V_{..}] G^{-1} \right)_p \\ \Gamma^s \Gamma_s &= r_{pp}^6 \left(G^{-1} [G^{-1} V^T] [V_{..}] \right)_p \left(G^{-1} [G^{-1} V^T] [V_{..}] G^{-1} \right)_p \\ k &= r_{pp}^6 \left(G^{-1} \otimes \left(G^{-1} [G^{-1}] V^T [V_{..}] G^{-1} \right) \right)_{ppp} \end{aligned} \quad (11)$$

where r_{pp} means the (p,p) -element of the matrix R .

It is shown in [10] that $B = \Psi$ i.e. (7) and (9) are the same.

3. THE PROBLEM OF PERTURBATIONS

In the second part the problem of additive perturbations of the observation vector \mathbf{Y} is discussed. Suppose modifying \mathbf{Y} by addition of a vector \mathbf{u} . Let $L(\theta)$ denote the log-likelihood corresponding to the postulated model and $L(\theta / \mathbf{u})$ the log-likelihood corresponding to the perturbed model, i.e.

$$L(\theta / \mathbf{u}) = -n/2 \ln(2\pi\sigma^2) - 1/2\sigma^2 | \mathbf{Y} + \mathbf{u} - \mathbf{f}(\theta) |^2$$

St. Laurent and Cook [11, 12] have shown that the so-called generalized leverage vector and Jacobian leverage vector can be used to assess the influence of perturbations on

fitted values in nonlinear regression. Our aim is to give a tool for assessment of the influence of additive perturbations on the values of parameters themselves. Then the relationship among the new measure, leverages introduced by Laurent and Cook and the intrinsic curvature is explored.

Let θ_u^* denote the least squares estimator of θ under $L(\theta/u)$. We have $\theta^* = \theta_0^*$ and $L(\theta) = L(\theta/\theta)$. Let V and W be the first and second derivatives of $f(\theta)$ with elements

$$V_j^a = (\partial f_a / \partial \theta_j)_{\theta^*} \quad \text{and} \quad W_{jk}^a = (\partial^2 f_a / \partial \theta_j \partial \theta_k)_{\theta^*}$$

respectively. For each $a=1, \dots, n$ denote by T_a the $p \times p \times p$ array with elements

$$(\partial^3 f_a / \partial \theta_j \partial \theta_k \partial \theta_l)_{\theta^*}. \text{ Put } e^* = Y - f(\theta^*).$$

Consider a perturbation of Y in the form $u = c b$, where $c \in R$ and $\|b\| = 1$.

We define the parameter leverage vector due to a perturbation of the data by c in the direction b as

$$P(c, b) = (\theta_{cb}^* - \theta^*) / c$$

and the Jacobian parameter leverage vector as

$$P(b) = \lim_{c \rightarrow 0} P(c, b)$$

Finally we define the Jacobian parameter leverage matrix as

$$P = (V'V - [e^*][W])^{-1} V'$$

where the square bracket multiplication is as above.

Theorem 3. *The parameter leverage vector due to a perturbation of the data by c in the direction b is, up to terms of order c [13],*

$$P(c, b) = A^{-1} V' b + 1/2 b' (P_1 + P_2 + P_3 + P_4) b c \quad (12)$$

where

$$\begin{aligned} A &= V'V - [e^*][W] \\ P_1 &= - [A^{-1}][VA^{-1}WA^{-1}V'] \\ (P_2)_{uw}^s &= ([I - VA^{-1}V'] [A^{-1}WA^{-1}V'])_{su}^t \\ (P_3)_{uw}^s &= ([I - VA^{-1}V'] [A^{-1}WA^{-1}V'])_{st}^u \\ P_4 &= [A^{-1}][VA^{-1}(\sum e_i T_i) A^{-1}V'] \end{aligned}$$

Corollary 1. *The Jacobian parameter leverage vector due to a perturbation of the data by c in the direction b is $P(b) = P b$.*

In order to get deeper insight into the matrix P we express it in a standard form. Let the QR decomposition of V be given by $V = UR$. Consider now the reparameterization

$$\tau = U' (f(\theta) - f(\theta^*))$$

where τ is called the normal parameter. It can be shown that the Jacobian parameter leverage matrix with respect to the normal parameter is

$$P = (I - B)^{-1} U' \quad (13)$$

where

$$B = [e^*][L'WL], \quad L = R^{-1}$$

The relationship between the Jacobian parameter leverage matrix and the intrinsic curvature of Bates and Watts [4] follows from the fact that

$$B = [e^*N][C]$$

where $N'e^*$ is the rotated residual vector and C is the intrinsic curvature array.

The norm $\|Pb\|$ measures the influence of a perturbation of data in the direction b on values of parameters. In the case of the derived linear model $f(\theta^*) + V(\theta - \theta^*)$, the Jacobian parameter leverage matrix for the normal parameter is $H = U'$. The matrix P takes into account the normal curvatures of the expectation surface at $f(\theta^*)$, while H does not.

For the normal parameter

$$(P)_{\max} = \max \| (I - B)^{-1} U' b \| = (1 - \lambda_p)^{-1} \quad (14)$$

$$\|b\| = 1$$

where λ_p is the largest eigenvalue of B . Denoting by b_i the i -th standard basis vector in R^n ,

we have

$$P_i = \|P b_i\| = \| (I - B)^{-1} (U')_i \| \quad (15)$$

where $(U')_i$ means the i -th column of U' . Evidently

$$P_i^2 = b_i' U (I - B)^{-2} U' b_i \quad (16)$$

We recall that the Jacobian leverage matrix introduced by St. Laurent and Cook [11, 12] is $J = U (I - B)^{-1} U'$. Its i -th diagonal element j_{ii} measures the rate of change in the fitted value y_i^* with respect to the rate of change in y_i . It follows from (16) that if $(I - B)^{-1}$ differs from $(I - B)^{-2}$, the case maximizing j_{ii} may be different from that with the largest Jacobian parameter leverage. If $(I - B)^{-1} = (I - B)^{-2}$, then $B = O$. Two special cases in which this occurs are when the model is intrinsically linear (the elements of $L'WL$ are zeros) and when the model provides an exact fit to the data (the residual vector e^* is zero).

4. APPLICATIONS

4.1. Application of Approximate Confidence Regions to IBNR Reserves

Our results will be applied to the case of computing IBNR reserves by the method described above which is completely different from chain-ladder methods. The results obtained show that asymmetric (with respect to the least squares estimators of parameters) confidence intervals are more accurate than intervals based on an approximation of lower order.

The model used by Stelltjes [3] $f(x, \theta) = \alpha \exp(\beta x) + \gamma \exp(\omega x)$, i.e. 4-parameter bi-exponential model is by computation ill-conditioned. More precisely, the practical im-

plementation in computational software S-plus leads to the error message "singular gradient matrix" (which is more or less expectable if one exponential fits data better) and "step factor reduced below minimum". In the present paper we do not publish the original data of [3] since paper is online and insurance data are considered as privacy data.

If correlations among incremental pure premiums are not negligible, what is also case of the data given by [3], usage of the correct model produces different values. For example, the PROC NLIN based method in [3] gave resulting estimators 3.1994, -0.0754, 29.4446, -0.5480 and 95% confidence intervals for estimators (2.0596, 4.3392), (-0.0942, -0.0566), (18.5334, 40.3557) and (-0.6986,-0.3974) of parameters α , β , γ , ω , respectively. However, using the correlation, as addressed in [3], we obtain the estimators 2.36488, -0.07767, 21.61184, -0.566532 and 95% confidence intervals for estimators (1.2746, 3.4551), (-0.0959, -0.05937), (10.664, 32.5596) and (-0.72473,-0.4083) of parameters α , β , γ , ω , respectively.

One of the main practical and theoretical problems of the IBNR confidence intervals implementation is the efficient estimation of the covariance structure and, consequently, sums of squares $S(\hat{\theta})$. This problem evidently goes behind the framework of this contribution and will be addressed in some future work. For the related topic see designs of experiments literature addressing the covariance estimation, for recent results in the parametrized covariance see e.g. [3].

Applying the results of the second part of this paper reveals the important fact, namely, that we can neglect the effect of Bin (7). Here we have the following quantities needed for calculation of confidence intervals by means of (7):

$$g^{11} = 0.1129 \times 10^{-6}, g^{22} = 0.31 \times 10^{-10}, g^{33} = 0.1033 \times 10^{-4},$$

$$g^{44} = 0.1968 \times 10^{-8},$$

$$A_{11}^1 = 0, A_{22}^2 = 0.2134 \times 10^{-3}, A_{33}^3 = 0.1231 \times 10^{-3}, A_{44}^4 = 0.3603 \times 10^{-3},$$

$$c = \sqrt{F(1, 36, 0.95)} \approx \sqrt{4.13}.$$

Using them we obtain

$$\alpha : \text{upper bound } 3.1994 + 0.5807=3.7801 ; \text{ lower bound } 3.1994 - 0.5807= 2.6187$$

$$\beta \text{ upper bound } -0.0754 +0.013= -0.0624; \text{ lower bound } -0.0754 -0.027= -0.0781$$

$$\gamma \text{ upper bound } 29.4446 +8.85= 38.2946 \text{ lower bound } 29.4446-13.73= 15.7146$$

$$\omega \text{ upper bound } -0.5480 +0.057=-0.4903 \text{ lower bound } -0.5480-0.157= -0.7050.$$

Comparing them to the original results of Stelltjes shows that our method gives narrower intervals than the classical approach (the only exception being γ). Recall that the

intervals are not symmetric with respect to the least-squares estimates thus reflecting the nonlinearity of the model.

Admitting for a moment that the hypothesis $\beta = 0$ holds, we get the famous Mitcherlitz model $f(x, \theta)=\alpha +\gamma \exp(\omega x)$. A thorough analysis of this model can be found in [7, 8, 14] to illustrate that Bates- Watts method is not always reliable . It follows that for this model our approach outperforms the classical one based on quadratic approximation as it produces confidence regions whose true coverage is closer to the nominal level than for the classical ones.

The outlook of the research in this direction brings a possibility to compute 2-parameter confidence interval, for instance for (α, γ) or (γ, ω) . Notice that Stelltjes [3] has computed $S(\hat{\theta})$ as the sum of weighted least squares and not as reported in (2.3.1), p. 359, it has been a typo. Thus, we are using the same approach as Stehlík [15] and employ a weight function that is inversely proportional to the variance of the data.

Probably, the main objective of the actuarial work is to construct the confidence interval for IBNR reserves for various accident quarters. A classical approach how to construct approximate confidence intervals is given by Bates and Watts [4], p. 58. For instance, 95% confidence interval for total IBNR for data in [3] is (27459851, 32751218) and (25254267, 40083031) when the parameter variance - covariance matrix is used (see [3]). The outlook is to construct a confidence interval based on cubic approximation of the model function (1). However, one should take care about the fact, that traditional nonlinear regression assumes that the error terms are normal which is a symmetric distribution with a range of whole real line. Incremental pure premium data may actually be skewed and can hardly ever be highly negative, therefore, using the normal distribution is approximation at best.

4.2. Application of Additive Perturbation Models for Privacy

Preserving Data Mining in Insurance

In insurance companies, many data sets have to be stored. A large fraction of them uses randomized data distortion techniques to mask the data for preserving the privacy of sensitive data. The additive perturbation attempts to hide the sensitive data by randomly modifying the data values often using additive noise. Random matrices have ‘predictable’ structures in the spectral domain and it develops a random matrix-based ‘Spectral Filtering Technique’ (SPF) to retrieve original data from the dataset distorted by adding random values (see [16]). In the present example we consider the same model as in section 4.1. As mentioned above the PROC NLIN based method gave estimators 3.1994, -0.0754, 29.4446, -0.5480 of model parameters. Then for $x=10$ we have $f=1.628027$.

Assuming that the error has a Gaussian $N(0,1)$ form we get an observed value of 3.97282 and additive perturbation by $N(0,0.001)$ will give us -0.0001432867, while the additive

perturbation by $N(0,0.01)$ will give 3.949242. Thus tuning of the standard deviation can tune the necessary level of privacy of the data. The control over this procedure is guaranteed by the theoretical results in the section 3, and may be properly analysed for every insurance model.

ACKNOWLEDGEMENTS

Authors are thankful to the referee comments which improved the quality of the paper. Work was partially supported by AKTION 54p13 and Vega grant 1/0077/09.

REFERENCES

- [1] F. De Vylder, "Non-linear regression in credibility theory", *Insurance: Mathematics and Economics*, vol. 4, No. 3, pp. 163-172, 1985.
- [2] G. Pitselis, "De Vylder's robust nonlinear regression credibility", *Belgian Actuarial Bulletin*, vol. 4, no. 1, pp. 1-4, 2004.
- [3] S. Stelljes, "Nonlinear Regression Model of Incurred but Not Reported Losses", *Casualty Actuarial Society Forum Casualty Actuarial Society – Arlington: Virginia, Fall, 2006*, pp. 353-377.
- [4] D. M. Bates, and D. G. Watts, *Nonlinear Regression Analysis and its Applications*, Wiley Series in Probability and Statistics: New York, 2007.
- [5] D. Hamilton, Confidence regions for parameter subsets in nonlinear regression, *Biometrika*, vol. 73, pp. 57-64, 1986.
- [6] L. M. Haines, T. E. O'Brien, and G. P. Y. Clarke, "Kurtosis and curvature measures for nonlinear regression models", *Statistica Sinica*, vol. 14, pp. 547-570, 2004.
- [7] R. D. Cook, and M. L. Goldberg, "Curvatures for parameter subsets in nonlinear regression", *Annals Statistics*, vol. 14, pp. 1399-1418, 1986.
- [8] G. P. Y. Clarke, "Marginal curvatures and their usefulness in the analysis of nonlinear regression models", *Journal of the American Statistical Association*, vol. 82, pp. 844-850, 1987.
- [9] R. Potocký, and T. Van Ban, "Confidence regions in nonlinear regression models", *Application of Mathematics*, vol. 37, pp. 29-39, 1992.
- [10] R. Potocký, and M. Stehlík, "Improved non-linear regression modelling of parameters of IBNR reserves", *Communications in Dependability and Quality Managements*, vol. 11, no. 4, pp. 54-60, 2008.
- [11] R. J. St. Laurent, and R. D. Cook, "Leverage and superleverage in nonlinear regression", *Journal of the American Statistical Association*, vol. 87, pp. 985-990, 1992.
- [12] R. J. St. Laurent, and R. D. Cook, "Leverage, local influence and curvature in nonlinear regression", *Biometrika*, vol. 80, pp. 99-106, 1993.
- [13] R. Potocký, and T. Van Ban, "On additive perturbations in nonlinear regression models", *Tatra Mountains Mathematical Publications*, vol. 17, pp. 65-71, 1999.
- [14] R. Potocký, "Approximate confidence regions for parameter subsets in a non-linear model", *Tatra Mountains Mathematical Publications*, vol. 26, pp. 227-236, 2003.
- [15] M. Stehlík, "Covariance related properties of D-optimal correlated designs", in the *5th St. Petersburg Workshop on Simulation*, 2005, pp. 645-652.
- [16] S. Datta, "On Random Additive Perturbation for Privacy Preserving Data Mining", M.S. thesis, University of Maryland, MD, USA, 2004.

Received: June 18, 2010

Revised: August 17, 2010

Accepted: August 29, 2010

© Potocký and Stehlík; Licensee *Bentham Open*.

This is an open access article licensed under the terms of the Creative Commons Attribution Non-Commercial License (<http://creativecommons.org/licenses/by-nc/3.0/>) which permits unrestricted, non-commercial use, distribution and reproduction in any medium, provided the work is properly cited.