

# Insurance Premium Prediction via Gradient Tree-Boosted Tweedie Compound Poisson Models

YI YANG\*, WEI QIAN<sup>†</sup> AND HUI ZOU<sup>‡</sup>

April 22, 2016

## Abstract

The Tweedie GLM is a widely used method for predicting insurance premiums. However, the structure of the logarithmic mean is restricted to a linear form in the Tweedie GLM, which can be too rigid for many applications. As a better alternative, we propose a gradient tree-boosting algorithm and apply it to Tweedie compound Poisson models for pure premiums. We use a profile likelihood approach to estimate the index and dispersion parameters. Our method is capable of fitting a flexible nonlinear Tweedie model and capturing complex interactions among predictors. A simulation study confirms the excellent prediction performance of our method. As an application, we apply our method to an auto insurance claim data and show that the new method is superior to the existing methods in the sense that it generates more accurate premium predictions, thus helping solve the adverse selection issue. We have implemented our method in a user-friendly R package that also includes a nice visualization tool for interpreting the fitted model.

---

\*McGill University

<sup>†</sup>Rochester Institute of Technology

<sup>‡</sup>Corresponding author, zoux019@umn.edu, University of Minnesota

# 1 Introduction

One of the most important problems in insurance business is to set the premium for the customers (policyholders). In a competitive market, it is advantageous for the insurer to charge a fair premium according to the expected loss of the policyholder. In personal car insurance, for instance, if an insurance company charges too much for old drivers and charges too little for young drivers, then the old drivers will switch to its competitors, and the remaining policies for the young drivers would be underpriced. This results in the *adverse selection* issue (Dionne et al., 2001): the insurer loses profitable policies and is left with bad risks, resulting in economic loss both ways.

To appropriately set the premiums for the insurer’s customers, one crucial task is to predict the size of actual (currently unforeseeable) claims. In this paper, we will focus on modeling claim loss, although other ingredients such as safety loadings, administrative costs, cost of capital, and profit are also important factors for setting the premium. One difficulty in modeling the claims is that the distribution is usually highly right-skewed, mixed with a point mass at zero. Such type of data cannot be transformed to normality by power transformation, and special treatment on zero claims is often required. As an example, Figure 1 shows the histogram of an auto insurance claim data (Yip and Yau, 2005), in which there are 6,290 policy records with zero claims and 4,006 policy records with positive losses.

The need for predictive models emerges from the fact that the expected loss is highly dependent on the characteristics of an individual policy such as age and motor vehicle record points of the policyholder, population density of the policyholder’s residential area, and age and model of the vehicle. Traditional methods used generalized linear models (GLM; Nelder and Wedderburn, 1972) for modeling the claim size (e.g. Renshaw, 1994; Haberman and Renshaw, 1996). However, the authors of the above papers performed their analyses on a subset of the policies, which have at least one claim. Alternative approaches have employed Tobit models by treating zero outcomes as censored below some cutoff points (Van de Ven and van Praag, 1981; Showers and Shotick, 1994), but these approaches rely on a normality assumption of the latent response. Alternatively, Jørgensen and de Souza (1994) and Smyth and Jørgensen (2002) used GLMs with a Tweedie distributed outcome to simultaneously model frequency and severity of insurance claims. They assume Poisson arrival of claims and gamma distributed amount for individual claims so that the size of the total claim amount

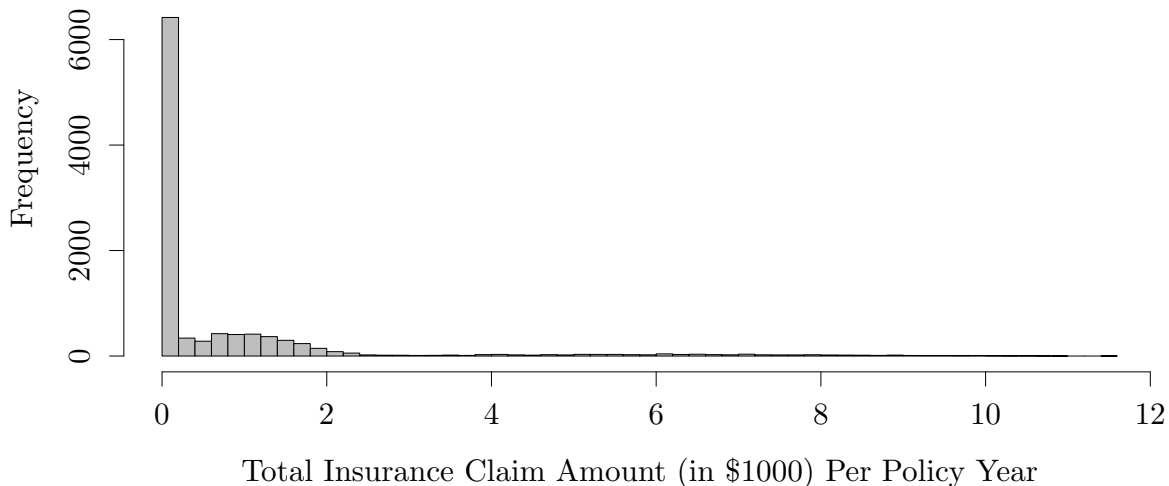


Figure 1: Histogram of the auto insurance claim data as analyzed in Yip and Yau (2005). It shows that there are 6290 policy records with zero total claims per policy year, while the remaining 4006 policy records have positive losses.

follows a Tweedie compound Poisson distribution. Due to its ability to simultaneously model the zeros and the continuous positive outcomes, the Tweedie GLM has been a widely used method in actuarial studies (Mildenhall, 1999; Murphy et al., 2000; Peters et al., 2008).

Despite of the popularity of the Tweedie GLM, a major limitation is that the structure of the logarithmic mean is restricted to a linear form, which can be too rigid for real applications. In auto insurance, for example, it is known that the risk does not monotonically decrease as age increases (Anstey et al., 2005). Although nonlinearity may be modeled by adding splines (Zhang, 2011), low-degree splines are often inadequate to capture the non-linearity in the data, while high-degree splines often result in the over-fitting issue that produces unstable estimates. Generalized additive models (GAM; Hastie and Tibshirani, 1990; Wood, 2006) overcome the restrictive linear assumption of GLMs, and can model the continuous variables by smooth functions estimated from data. The structure of the model, however, has to be determined *a priori*. That is, one has to specify the main effects and interaction effects to be used in the model. As a result, misspecification of non-ignorable effects is likely to adversely affect prediction accuracy.

In this paper, we aim to model the insurance claim size by a nonparametric Tweedie compound Poisson model, and propose a gradient tree-boosting algorithm (TDboost henceforth) to fit this model. To our knowledge, before this work, there is no existing nonparametric Tweedie method available. Additionally, we also implemented the proposed method as an easy-to-use R package, which is publicly available.

Gradient boosting is one of the most successful machine learning algorithms for nonparametric regression and classification. Boosting adaptively combines a large number of relatively simple prediction models called *base learners* into an ensemble learner to achieve high prediction performance. The seminal work on the boosting algorithm called *AdaBoost* (Freund and Schapire, 1997) was originally proposed for classification problems. Later Breiman (1998) and Breiman (1999) pointed out an important connection between the AdaBoost algorithm and a functional gradient descent algorithm. Friedman et al. (2000) and Hastie et al. (2009) developed a statistical view of boosting and proposed gradient boosting methods for both classification and regression. There is a large body of literature on boosting. We refer interested readers to Bühlmann and Hothorn (2007) for a comprehensive review of boosting algorithms.

The TDboost model is motivated by the proven success of boosting in machine learning for classification and regression problems (Friedman, 2001, 2002; Hastie et al., 2009). Its advantages are threefold. First, the model structure of TDboost is learned from data and not predetermined, thereby avoiding an explicit model specification. Non-linearities, discontinuities, complex and higher order interactions are naturally incorporated into the model to reduce the potential modeling bias and to produce high predictive performance, which enables TDboost to serve as a benchmark model in scoring insurance policies, guiding pricing practice, and facilitating marketing efforts. Feature selection is performed as an integral part of the procedure. In addition, TDboost handles the predictor and response variables of any type without the need for transformation, and it is highly robust to outliers. Missing values in the predictors are managed almost without loss of information (Elith et al., 2008). All these properties make TDboost a more attractive tool for insurance premium modeling. On the other hand, we acknowledge that its results are not as straightforward as those from the Tweedie GLM model. Nevertheless, TDboost does not have to be regarded as a black box. It can provide interpretable results, by means of the partial dependence plots, and relative importance of the predictors.

The remainder of this paper is organized as follows. We briefly review the gradient boosting algorithm and the Tweedie compound Poisson model in Section 2 and Section 3, respectively. We present the main methodological development with implementation details in Section 4. In Section 5, we use simulation to show the high predictive accuracy of TDboost. As an application, we apply TDboost to analyze an auto insurance claim data in Section 6.

## 2 Gradient Boosting

Gradient boosting (Friedman, 2001) is a recursive, nonparametric machine learning algorithm that has been successfully used in many areas. It shows remarkable flexibility in solving different loss functions. By combining a large number of base learners, it can handle higher order interactions and produce highly complex functional forms. It provides high prediction accuracy and often outperforms many competing methods, such as linear regression/classification, bagging (Breiman, 1996), splines and CART (Breiman et al., 1984).

To keep the paper self-contained, we briefly explain the general procedures for the gradient boosting. Let  $\mathbf{x} = (x_1, \dots, x_p)^\top$  be a  $p$ -dimensional column vector for the predictor variables and  $y$  be the one-dimensional response variable. The goal is to estimate the optimal prediction function  $\tilde{F}(\cdot)$  that maps  $\mathbf{x}$  to  $y$  by minimizing the expected value of a loss function  $\Psi(\cdot, \cdot)$  over the function class  $\mathcal{F}$ :

$$\tilde{F}(\cdot) = \arg \min_{F(\cdot) \in \mathcal{F}} E_{y, \mathbf{x}}[\Psi(y, F(\mathbf{x}))],$$

where  $\Psi$  is assumed to be differentiable with respect to  $F$ . Given the observed data  $\{y_i, \mathbf{x}_i\}_{i=1}^n$ , where  $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})^\top$ , estimation of  $\tilde{F}(\cdot)$  can be done by minimizing the empirical risk function

$$\min_{F(\cdot) \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \Psi(y_i, F(\mathbf{x}_i)). \quad (1)$$

For the gradient boosting, each candidate function  $F \in \mathcal{F}$  is assumed to be an ensemble of  $M$  base learners

$$F(\mathbf{x}) = F^{[0]} + \sum_{m=1}^M \beta^{[m]} h(\mathbf{x}; \boldsymbol{\xi}^{[m]}), \quad (2)$$

where  $h(\mathbf{x}; \boldsymbol{\xi}^{[m]})$  usually belongs to a class of some simple functions of  $\mathbf{x}$  called base learners (e.g., regression/decision tree) with the parameter  $\boldsymbol{\xi}^{[m]}$  ( $m = 1, 2, \dots, M$ ).  $F^{[0]}$  is a constant

scalar and  $\beta^{[m]}$  is the expansion coefficient. Note that differing from the usual structure of an additive model, there is no restriction on the number of predictors to be included in each  $h(\cdot)$ , and consequently, high-order interactions can be easily considered using this setting.

A forward stagewise algorithm is adopted to approximate the minimizer of (1), which builds up the components  $\beta^{[m]}h(\mathbf{x}; \boldsymbol{\xi}^{[m]})$  ( $m = 1, 2, \dots, M$ ) sequentially through a gradient-descent-like approach. At each iteration stage  $m$  ( $m = 1, 2, \dots$ ), suppose that the current estimate for  $\tilde{F}(\cdot)$  is  $\hat{F}^{[m-1]}(\cdot)$ . To update the estimate from  $\hat{F}^{[m-1]}(\cdot)$  to  $\hat{F}^{[m]}(\cdot)$ , the gradient boosting fits a negative gradient vector (as the working response) to the predictors using a base learner  $h(\mathbf{x}; \boldsymbol{\xi}^{[m]})$ . This fitted  $h(\mathbf{x}; \boldsymbol{\xi}^{[m]})$  can be viewed as an approximation of the negative gradient. Subsequently, the expansion coefficient  $\beta^{[m]}$  can then be determined by a line search minimization with the empirical risk function, and the estimation of  $\tilde{F}(\mathbf{x})$  for the next stage becomes

$$\hat{F}^{[m]}(\mathbf{x}) := \hat{F}^{[m-1]}(\mathbf{x}) + \nu\beta^{[m]}h(\mathbf{x}; \boldsymbol{\xi}^{[m]}), \quad (3)$$

where  $0 < \nu \leq 1$  is the shrinkage factor (Friedman, 2001) that controls the update step size. A small  $\nu$  imposes more shrinkage while  $\nu = 1$  gives complete negative gradient steps. Friedman (2001) has found that the shrinkage factor reduces over-fitting and improves the predictive accuracy.

### 3 Compound Poisson Distribution and Tweedie Model

In insurance premium prediction problems, the total claim amount for a covered risk usually has a continuous distribution on positive values, except for the possibility of being exact zero when the claim does not occur. One standard approach in actuarial science in modeling such data is using Tweedie compound Poisson models, which we briefly introduce in this section.

Let  $N$  be a Poisson random variable denoted by  $\text{Pois}(\lambda)$ , and let  $\tilde{Z}_d$ 's ( $d = 0, 1, \dots, N$ ) be i.i.d. gamma random variables denoted by  $\text{Gamma}(\alpha, \gamma)$  with mean  $\alpha\gamma$  and variance  $\alpha\gamma^2$ . Assume  $N$  is independent of  $\tilde{Z}_d$ 's. Define a random variable  $Z$  by

$$Z = \begin{cases} 0 & \text{if } N = 0 \\ \tilde{Z}_1 + \tilde{Z}_2 + \dots + \tilde{Z}_N & \text{if } N = 1, 2, \dots \end{cases}. \quad (4)$$

Thus  $Z$  is the Poisson sum of independent Gamma random variables. In insurance applications, one can view  $Z$  as the total claim amount,  $N$  as the number of reported claims and  $\tilde{Z}_d$ 's as the insurance payment for the  $d$ th claim. The resulting distribution of  $Z$  is referred to as the compound Poisson distribution (Jørgensen and de Souza, 1994; Smyth and Jørgensen, 2002), which is known to be closely connected to exponential dispersion models (EDM) (Jørgensen, 1987). Note that the distribution of  $Z$  has a probability mass at zero:  $Pr(Z = 0) = \exp(-\lambda)$ . Then based on that  $Z$  conditional on  $N = j$  is  $\text{Gamma}(j\alpha, \gamma)$ , the distribution function of  $Z$  can be written as

$$\begin{aligned} f_Z(z|\lambda, \alpha, \gamma) &= Pr(N = 0)d_0(z) + \sum_{j=1}^{\infty} Pr(N = j)f_{Z|N=j}(z) \\ &= \exp(-\lambda)d_0(z) + \sum_{j=1}^{\infty} \frac{\lambda^j e^{-\lambda}}{j!} \frac{z^{j\alpha-1} e^{-z/\gamma}}{\gamma^{j\alpha} \Gamma(j\alpha)}, \end{aligned}$$

where  $d_0$  is the Dirac delta function at zero and  $f_{Z|N=j}$  is the conditional density of  $Z$  given  $N = j$ . Smyth (1996) pointed out that the compound Poisson distribution belongs to a special class of EDMs known as Tweedie models (Tweedie, 1984), which are defined by the form

$$f_Z(z|\theta, \phi) = a(z, \phi) \exp \left\{ \frac{z\theta - \kappa(\theta)}{\phi} \right\}, \quad (5)$$

where  $a(\cdot)$  is a normalizing function,  $\kappa(\cdot)$  is called the cumulant function, and both  $a(\cdot)$  and  $\kappa(\cdot)$  are known. The parameter  $\theta$  is in  $\mathbb{R}$  and the dispersion parameter  $\phi$  is in  $\mathbb{R}^+$ . For Tweedie models the mean  $E(Z) \equiv \mu = \dot{\kappa}(\theta)$  and the variance  $\text{Var}(Z) = \phi \ddot{\kappa}(\theta)$ , where  $\dot{\kappa}(\theta)$  and  $\ddot{\kappa}(\theta)$  are the first and second derivatives of  $\kappa(\theta)$ , respectively. Tweedie models have the power mean-variance relationship  $\text{Var}(Z) = \phi \mu^\rho$  for some index parameter  $\rho$ . Such mean-variance relation gives

$$\theta = \begin{cases} \frac{\mu^{1-\rho}}{1-\rho}, & \rho \neq 1 \\ \log \mu, & \rho = 1 \end{cases}, \quad \kappa(\theta) = \begin{cases} \frac{\mu^{2-\rho}}{2-\rho}, & \rho \neq 2 \\ \log \mu, & \rho = 2 \end{cases}. \quad (6)$$

One can show that the compound Poisson distribution belongs to the class of Tweedie models. Indeed, if we reparametrize  $(\lambda, \alpha, \gamma)$  by

$$\lambda = \frac{1}{\phi} \frac{\mu^{2-\rho}}{2-\rho}, \quad \alpha = \frac{2-\rho}{\rho-1}, \quad \gamma = \phi(\rho-1)\mu^{\rho-1}, \quad (7)$$

the compound Poisson model will have the form of a Tweedie model with  $1 < \rho < 2$  and  $\mu > 0$ . As a result, for the rest of this paper, we only consider the model (4), and simply refer to (4) as the Tweedie model (or Tweedie compound Poisson model), denoted by  $\text{Tw}(\mu, \phi, \rho)$ , where  $1 < \rho < 2$  and  $\mu > 0$ .

It is straightforward to show that the log-likelihood of the Tweedie model is

$$\log f_Z(z|\mu, \phi, \rho) = \frac{1}{\phi} \left( z \frac{\mu^{1-\rho}}{1-\rho} - \frac{\mu^{2-\rho}}{2-\rho} \right) + \log a(z, \phi, \rho), \quad (8)$$

where the normalizing function  $a(\cdot)$  can be written as

$$a(z, \phi, \rho) = \begin{cases} \frac{1}{z} \sum_{t=1}^{\infty} W_t(z, \phi, \rho) = \frac{1}{z} \sum_{t=1}^{\infty} \frac{z^{t\alpha}}{(\rho-1)^{t\alpha} \phi^{t(1+\alpha)} (2-\rho)^t t! \Gamma(t\alpha)} & \text{for } z > 0 \\ 1 & \text{for } z = 0 \end{cases},$$

and  $\alpha = (2 - \rho)/(\rho - 1)$  and  $\sum_{t=1}^{\infty} W_t$  is an example of Wright's generalized Bessel function (Tweedie, 1984).

## 4 Our Proposal

In this section, we propose to integrate the Tweedie model to the tree-based gradient boosting algorithm to predict insurance claim size. Specifically, our discussion focuses on modeling the personal car insurance as an illustrating example (see Section 6 for a real data analysis), since our modeling strategy is easily extended to other lines of non-life insurance business.

Given an auto insurance policy  $i$ , let  $N_i$  be the number of claims (known as the claim frequency) and  $\tilde{Z}_{d_i}$  be the size of each claim observed for  $d_i = 1, \dots, N_i$ . Let  $w_i$  be the policy duration, that is, the length of time that the policy remains in force. Then  $Z_i = \sum_{d_i=1}^{N_i} \tilde{Z}_{d_i}$  is the total claim amount. In the following, we are interested in modeling the ratio between the total claim and the duration  $Y_i = Z_i/w_i$ , a key quantity known as the pure premium (Ohlsson and Johansson, 2010).

Following the settings of the compound Poisson model, we assume  $N_i$  is Poisson distributed, and its mean  $\lambda_i w_i$  has a multiplicative relation with the duration  $w_i$ , where  $\lambda_i$  is a policy-specific parameter representing the expected claim frequency under unit duration. Conditional on  $N_i$ , assume  $Z_{d_i}$ 's ( $d_i = 1, \dots, N_i$ ) are i.i.d.  $\text{Gamma}(\alpha, \gamma_i)$ , where  $\gamma_i$  is a



policy-specific parameter that determines claim severity, and  $\alpha$  is a constant. Furthermore, we assume that under unit duration (i.e.,  $w_i = 1$ ), the mean-variance relation of a policy satisfies  $\text{Var}(Y_i^*) = \phi[E(Y_i^*)]^\rho$  for all policies, where  $Y_i^*$  is the pure premium under unit duration,  $\phi$  is a constant, and  $\rho = (\alpha + 2)/(\alpha + 1)$ . Then, it is known that  $Y_i \sim \text{Tw}(\mu_i, \phi/w_i, \rho)$ , the details of which are provided in Appendix Part A.

Then, we consider a portfolio of policies  $\{(y_i, \mathbf{x}_i, w_i)\}_{i=1}^n$  from  $n$  independent insurance contracts, where for the  $i$ th contract,  $y_i$  is the policy pure premium,  $\mathbf{x}_i$  is a vector of explanatory variables that characterize the policyholder and the risk being insured (e.g. house, vehicle), and  $w_i$  is the duration. Assume that the expected pure premium  $\mu_i$  is determined by a predictor function  $F : \mathbb{R}^p \rightarrow \mathbb{R}$  of  $\mathbf{x}_i$ :

$$\log\{\mu_i\} = \log\{E(Y_i|\mathbf{x}_i)\} = F(\mathbf{x}_i). \quad (9)$$

In this paper, we do not impose a linear or other parametric form restriction on  $F(\cdot)$ . Given the flexibility of  $F(\cdot)$ , we call such setting as the boosted Tweedie model (as opposed to the Tweedie GLM). Given  $\{(y_i, \mathbf{x}_i, w_i)\}_{i=1}^n$ , the log-likelihood function can be written as

$$\begin{aligned} \ell(F(\cdot), \phi, \rho | \{y_i, \mathbf{x}_i, w_i\}_{i=1}^n) &= \sum_{i=1}^n \log f_Y(y_i | \mu_i, \phi/w_i, \rho), \\ &= \sum_{i=1}^n \frac{w_i}{\phi} \left( y_i \frac{\mu_i^{1-\rho}}{1-\rho} - \frac{\mu_i^{2-\rho}}{2-\rho} \right) + \log a(y_i, \phi/w_i, \rho). \end{aligned} \quad (10)$$

## 4.1 Estimating $F(\cdot)$ via TDboost

We estimate the predictor function  $F(\cdot)$  by integrating the boosted Tweedie model into the tree-based gradient boosting algorithm. To develop the idea, we assume that  $\phi$  and  $\rho$  are given for the time being. The joint estimation of  $F(\cdot)$ ,  $\phi$  and  $\rho$  will be studied in Section 4.2.

Given  $\rho$  and  $\phi$ , we replace the general objective function in (1) by the negative log-likelihood derived in (10), and target the minimizer function  $F^*(\cdot)$  over a class  $\mathcal{F}$  of base learner functions in the form of (2). That is, we intend to estimate

$$F^*(\mathbf{x}) = \underset{F \in \mathcal{F}}{\operatorname{argmin}} \left\{ -\ell(F(\cdot), \phi, \rho | \{y_i, \mathbf{x}_i, w_i\}_{i=1}^n) \right\} = \underset{F \in \mathcal{F}}{\operatorname{argmin}} \sum_{i=1}^n \Psi(y_i, F(\mathbf{x}_i) | \rho), \quad (11)$$

where

$$\Psi(y_i, F(\mathbf{x}_i) | \rho) = w_i \left\{ -\frac{y_i \exp[(1 - \rho)F(\mathbf{x}_i)]}{1 - \rho} + \frac{\exp[(2 - \rho)F(\mathbf{x}_i)]}{2 - \rho} \right\}.$$

Note that in contrast to (11), the function class targeted by Tweedie GLM (Smyth, 1996) is restricted to a collection of linear functions of  $\mathbf{x}$ .

We propose to apply the forward stagewise algorithm described in Section 2 for solving (11). The initial estimate of  $F^*(\cdot)$  is chosen as a constant function that minimizes the negative log-likelihood:

$$\begin{aligned} \hat{F}^{[0]} &= \underset{\eta}{\operatorname{argmin}} \sum_{i=1}^n \Psi(y_i, \eta | \rho) \\ &= \log \left( \frac{\sum_{i=1}^n w_i y_i}{\sum_{i=1}^n w_i} \right). \end{aligned}$$

This corresponds to the best estimate of  $F$  without any covariates. Let  $\hat{F}^{[m-1]}$  be the current estimate before the  $m$ th iteration. At the  $m$ th step, we fit a base learner  $h(\mathbf{x}; \boldsymbol{\xi}^{[m]})$  via

$$\hat{\boldsymbol{\xi}}^{[m]} = \underset{\boldsymbol{\xi}^{[m]}}{\operatorname{argmin}} \sum_{i=1}^n [u_i^{[m]} - h(\mathbf{x}_i; \boldsymbol{\xi}^{[m]})]^2, \quad (12)$$

where  $(u_1^{[m]}, \dots, u_n^{[m]})^\top$  is the current negative gradient of  $\Psi(\cdot | \rho)$ , i.e.,

$$u_i^{[m]} = - \left. \frac{\partial \Psi(y_i, F(\mathbf{x}_i) | \rho)}{\partial F(\mathbf{x}_i)} \right|_{F(\mathbf{x}_i) = \hat{F}^{[m-1]}(\mathbf{x}_i)} \quad (13)$$

$$= w_i \left\{ -y_i \exp[(1 - \rho)\hat{F}^{[m-1]}(\mathbf{x}_i)] + \exp[(2 - \rho)\hat{F}^{[m-1]}(\mathbf{x}_i)] \right\}, \quad (14)$$

and use an  $L$ -terminal node regression tree

$$h(\mathbf{x}; \boldsymbol{\xi}^{[m]}) = \sum_{l=1}^L u_l^{[m]} I(\mathbf{x} \in R_l^{[m]}) \quad (15)$$

with parameters  $\boldsymbol{\xi}^{[m]} = \{R_l^{[m]}, u_l^{[m]}\}_{l=1}^L$  as the base learner. To find  $R_l^{[m]}$  and  $u_l^{[m]}$ , we use a fast top-down “best-fit” algorithm with a least squares splitting criterion (Friedman et al., 2000) to find the splitting variables and corresponding split locations that determine the fitted terminal regions  $\{\hat{R}_l^{[m]}\}_{l=1}^L$ . Note that estimating the  $R_l^{[m]}$  entails estimating the  $u_l^{[m]}$

as the mean falling in each region:

$$\bar{u}_l^{[m]} = \text{mean}_{i:\mathbf{x}_i \in \hat{R}_l^{[m]}}(u_i^{[m]}) \quad l = 1, \dots, L.$$

Once the base learner  $h(\mathbf{x}; \boldsymbol{\xi}^{[m]})$  has been estimated, the optimal value of the expansion coefficient  $\beta^{[m]}$  is determined by a line search

$$\begin{aligned} \beta^{[m]} &= \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n \Psi(y_i, \hat{F}^{[m-1]}(\mathbf{x}_i) + \beta h(\mathbf{x}_i; \hat{\boldsymbol{\xi}}^{[m]}) \mid \rho) \\ &= \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n \Psi(y_i, \hat{F}^{[m-1]}(\mathbf{x}_i) + \beta \sum_{l=1}^L \bar{u}_l^{[m]} I(\mathbf{x}_i \in \hat{R}_l^{[m]}) \mid \rho). \end{aligned} \quad (16)$$

The regression tree (15) predicts a constant value  $\bar{u}_l^{[m]}$  within each region  $\hat{R}_l^{[m]}$ , so we can solve (16) by a separate line search performed within each respective region  $\hat{R}_l^{[m]}$ . The problem (16) reduces to finding a best constant  $\eta_l^{[m]}$  to improve the current estimate in each region  $\hat{R}_l^{[m]}$  based on the following criterion:

$$\hat{\eta}_l^{[m]} = \underset{\eta}{\operatorname{argmin}} \sum_{i:\mathbf{x}_i \in \hat{R}_l^{[m]}} \Psi(y_i, \hat{F}^{[m-1]}(\mathbf{x}_i) + \eta \mid \rho), \quad l = 1, \dots, L, \quad (17)$$

where the solution is given by

$$\hat{\eta}_l^{[m]} = \log \left\{ \frac{\sum_{i:\mathbf{x}_i \in \hat{R}_l^{[m]}} w_i y_i \exp[(1 - \rho) \hat{F}^{[m-1]}(\mathbf{x}_i)]}{\sum_{i:\mathbf{x}_i \in \hat{R}_l^{[m]}} w_i \exp[(2 - \rho) \hat{F}^{[m-1]}(\mathbf{x}_i)]} \right\}, \quad l = 1, \dots, L. \quad (18)$$

Having found the parameters  $\{\hat{\eta}_l^{[m]}\}_{l=1}^L$ , we then update the current estimate  $\hat{F}^{[m-1]}(\mathbf{x})$  in each corresponding region

$$\hat{F}^{[m]}(\mathbf{x}) = \hat{F}^{[m-1]}(\mathbf{x}) + \nu \hat{\eta}_l^{[m]} I(\mathbf{x} \in \hat{R}_l^{[m]}), \quad l = 1, \dots, L, \quad (19)$$

where  $0 < \nu \leq 1$  is the shrinkage factor. Following (Friedman, 2001), we set  $\nu = 0.005$  in our implementation. More discussions on the choice of tuning parameters are in Section 4.4.

In summary, the complete TDboost algorithm is shown in Algorithm 1. The boosting step is repeated  $M$  times and we report  $\hat{F}^{[M]}(\mathbf{x})$  as the final estimate.

---

**Algorithm 1** TDboost

---

1. Initialize  $\hat{F}^{[0]}$

$$\hat{F}^{[0]} = \log \left( \frac{\sum_{i=1}^n w_i y_i}{\sum_{i=1}^n w_i} \right).$$

2. For  $m = 1, \dots, M$  repeatedly do steps 2.(a)–2.(d)

- 2.(a) Compute the negative gradient  $(u_1^{[m]}, \dots, u_n^{[m]})^\top$

$$u_i^{[m]} = w_i \{ -y_i \exp[(1 - \rho)\hat{F}^{[m-1]}(\mathbf{x}_i)] + \exp[(2 - \rho)\hat{F}^{[m-1]}(\mathbf{x}_i)] \} \quad i = 1, \dots, n.$$

- 2.(b) Fit the negative gradient vector  $(u_1^{[m]}, \dots, u_n^{[m]})^\top$  to  $(\mathbf{x}_1, \dots, \mathbf{x}_n)^\top$  by an  $L$ -terminal node regression tree, where  $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})^\top$  for  $i = 1, \dots, n$ , giving us the partitions  $\{\hat{R}_l^{[m]}\}_{l=1}^L$ .

- 2.(c) Compute the optimal terminal node predictions  $\eta_l^{[m]}$  for each region  $\hat{R}_l^{[m]}$ ,  $l = 1, 2, \dots, L$

$$\hat{\eta}_l^{[m]} = \log \left\{ \frac{\sum_{i: \mathbf{x}_i \in \hat{R}_l^{[m]}} w_i y_i \exp[(1 - \rho)\hat{F}^{[m-1]}(\mathbf{x}_i)]}{\sum_{i: \mathbf{x}_i \in \hat{R}_l^{[m]}} w_i \exp[(2 - \rho)\hat{F}^{[m-1]}(\mathbf{x}_i)]} \right\}.$$

- 2.(d) Update  $\hat{F}^{[m]}(\mathbf{x})$  for each region  $\hat{R}_l^{[m]}$ ,  $l = 1, 2, \dots, L$

$$\hat{F}^{[m]}(\mathbf{x}) = \hat{F}^{[m-1]}(\mathbf{x}) + \nu \hat{\eta}_l^{[m]} I(\mathbf{x} \in \hat{R}_l^{[m]}) \quad l = 1, 2, \dots, L.$$

3. Report  $\hat{F}^{[M]}(\mathbf{x})$  as the final estimate.
-

## 4.2 Estimating $(\rho, \phi)$ via profile likelihood

Following Dunn and Smyth (2005), we use the profile likelihood to estimate the dispersion  $\phi$  and the index parameter  $\rho$ , which jointly determine the mean-variance relation  $Var(Y_i) = \phi \mu_i^\rho / w_i$  of the pure premium. We exploit the fact that in Tweedie models the estimation of  $\mu$  depends only on  $\rho$ : given a fixed  $\rho$ , the mean estimate  $\mu^*(\rho)$  can be solved in (11) without knowing  $\phi$ . Then conditional on this  $\rho$  and the corresponding  $\mu^*(\rho)$ , we maximize the log-likelihood function with respect to  $\phi$  by

$$\phi^*(\rho) = \underset{\phi}{\operatorname{argmax}} \{ \ell(\mu^*(\rho), \phi, \rho) \}, \quad (20)$$

which is a univariate optimization problem that can be solved using a combination of golden section search and successive parabolic interpolation (Brent, 2013). In such a way, we have determined the corresponding  $(\mu^*(\rho), \phi^*(\rho))$  for each fixed  $\rho$ . Then we acquire the estimate of  $\rho$  by maximizing the profile likelihood with respect to 50 equally spaced values  $\{\rho_1, \dots, \rho_{50}\}$  on  $(0, 1)$ : should be (1,2)

$$\rho^* = \underset{\rho \in \{\rho_1, \dots, \rho_{50}\}}{\operatorname{argmax}} \{ \ell(\mu^*(\rho), \phi^*(\rho), \rho) \}. \quad (21)$$

Finally, we apply  $\rho^*$  in (11) and (20) to obtain the corresponding estimates  $\mu^*(\rho^*)$  and  $\phi^*(\rho^*)$ . Some additional computational issues for evaluating the log-likelihood functions in (20) and (21) are discussed in Appendix Part B.

## 4.3 Model interpretation

Compared to other nonparametric statistical learning methods such as neural networks and kernel machines, our new estimator provides interpretable results. In this section, we discuss some ways for model interpretation after fitting the boosted Tweedie model.

### 4.3.1 Marginal effects of predictors

The main effects and interaction effects of the variables in the boosted Tweedie model can be extracted easily. In our estimate we can control the order of interactions by choosing the tree size  $L$  (the number of terminal nodes) and the number  $p$  of predictors. A tree with  $L$  terminal nodes produces a function approximation of  $p$  predictors with interaction order of at most  $\min(L - 1, p)$ . For example, a stump ( $L = 2$ ) produces an additive TDboost model

with only the main effects of the predictors, since it is a function based on a single splitting variable in each tree. Setting  $L = 3$  allows both main effects and second order interactions.

Following Friedman (2001) we use the so-called partial dependence plots to visualize the main effects and interaction effects. Given the training data  $\{y_i, \mathbf{x}_i\}_{i=1}^n$ , with a  $p$ -dimensional input vector  $\mathbf{x} = (x_1, x_2, \dots, x_p)^\top$ , let  $\mathbf{z}_s$  be a subset of size  $s$ , such that  $\mathbf{z}_s = \{z_1, \dots, z_s\} \subset \{x_1, \dots, x_p\}$ . For example, to study the main effect of the variable  $j$ , we set the subset  $\mathbf{z}_s = \{z_j\}$ , and to study the second order interaction of variables  $i$  and  $j$ , we set  $\mathbf{z}_s = \{z_i, z_j\}$ . Let  $\mathbf{z}_{\setminus s}$  be the complement set of  $\mathbf{z}_s$ , such that  $\mathbf{z}_{\setminus s} \cup \mathbf{z}_s = \{x_1, \dots, x_p\}$ . Let the prediction  $\hat{F}(\mathbf{z}_s | \mathbf{z}_{\setminus s})$  be a function of the subset  $\mathbf{z}_s$  conditioned on specific values of  $\mathbf{z}_{\setminus s}$ . The partial dependence of  $\hat{F}(\mathbf{x})$  on  $\mathbf{z}_s$  then can be formulated as  $\hat{F}(\mathbf{z}_s | \mathbf{z}_{\setminus s})$  averaged over the marginal density of the complement subset  $\mathbf{z}_{\setminus s}$

$$\hat{F}_s(\mathbf{z}_s) = \int \hat{F}(\mathbf{z}_s | \mathbf{z}_{\setminus s}) p_{\setminus s}(\mathbf{z}_{\setminus s}) d\mathbf{z}_{\setminus s}, \quad (22)$$

where  $p_{\setminus s}(\mathbf{z}_{\setminus s}) = \int p(\mathbf{x}) d\mathbf{z}_s$  is the marginal density of  $\mathbf{z}_{\setminus s}$ . We estimate (22) by

$$\bar{F}_s(\mathbf{z}_s) = \frac{1}{n} \sum_{i=1}^n \hat{F}(\mathbf{z}_s | \mathbf{z}_{\setminus s, i}), \quad (23)$$

where  $\{\mathbf{z}_{\setminus s, i}\}_{i=1}^n$  are evaluated at the training data. We then plot  $\bar{F}_s(\mathbf{z}_s)$  against  $\mathbf{z}_s$ . We have included the partial dependence plot function in our R package ‘‘TDboost’’. We will demonstrate this functionality in Section 6.

### 4.3.2 Variable importance

In many applications identifying relevant predictors of the model in the context of tree-based ensemble methods is of interest. The TDboost model defines a variable importance measure for each candidate predictor  $X_j$  in the set  $X = \{X_1, \dots, X_p\}$  in terms of prediction/explanation of the response  $Y$ . The major advantage of this variable selection procedure, as compared to univariate screening methods, is that the approach considers the impact of each individual predictor as well as multivariate interactions among predictors simultaneously.

We start by defining the variable importance (VI henceforth) measure in the context of a single tree. First introduced by Breiman et al. (1984), the VI measure  $\mathcal{I}_{X_j}(T_m)$  of the

variable  $X_j$  in a single tree  $T_m$  is defined as the total heterogeneity reduction of the response variable  $Y$  produced by  $X_j$ , which can be estimated by adding up all the decreases in the squared error reductions  $\hat{\delta}_l$  obtained in all  $L - 1$  internal nodes when  $X_j$  is chosen as the splitting variable. Denote  $v(X_j) = l$  the event that  $X_j$  is selected as the splitting variable in the internal node  $l$ , and let  $I_{jl} = I(v(X_j) = l)$ . Then

$$\mathcal{I}_{X_j}(T_m) = \sum_{l=1}^{L-1} \hat{\delta}_l I_{jl}, \quad (24)$$

where  $\hat{\delta}_l$  is defined as the squared error difference between the constant fit and the two sub-region fits (the sub-region fits are achieved by splitting the region associated with the internal node  $l$  into the left and right regions). Friedman (2001) extended the VI measure  $\mathcal{I}_{X_j}$  for the boosting model with a combination of  $M$  regression trees, by averaging (24) over  $\{T_1, \dots, T_M\}$ :

$$\mathcal{I}_{X_j} = \frac{1}{M} \sum_{m=1}^M \mathcal{I}_{X_j}(T_m). \quad (25)$$

Despite of the wide use of the VI measure, Breiman et al. (1984) and White and Liu (1994) among others have pointed out that the VI measures (24) and (25) are biased: even if  $X_j$  is a non-informative variable to  $Y$  (not correlated to  $Y$ ),  $X_j$  may still be selected as a splitting variable, hence the VI measure of  $X_j$  is non-zero by Equation (25). Following Sandri and Zuccolotto (2008) and Sandri and Zuccolotto (2010) to avoid the variable selection bias, in this paper we compute an adjusted VI measure for each explanatory variable by permutating each  $X_j$ , the computational details are provided in Appendix Part C.

## 4.4 Implementation

We have implemented our proposed method in an R package “TDboost”, which is publicly available from the Comprehensive R Archive Network at <http://cran.r-project.org/web/packages/TDboost/index.html>. Here, we discuss the choice of three meta parameters in Algorithm 1:  $L$  (the size of the trees),  $\nu$  (the shrinkage factor) and  $M$  (the number of boosting steps).

To avoid over-fitting and improve out-of-sample predictions, the boosting procedure can be regularized by limiting the number of boosting iterations  $M$  (early stopping; Zhang and

Yu, 2005) and the shrinkage factor  $\nu$ . Empirical evidence (Friedman, 2001; Bühlmann and Hothorn, 2007; Ridgeway, 2007) showed that the predictive accuracy is almost always better with a smaller shrinkage factor at the cost of more computing time. However, smaller values of  $\nu$  usually requires a larger number of boosting iterations  $M$  and hence induces more computing time (Friedman, 2001). We choose a “sufficiently small”  $\nu = 0.005$  throughout and determine  $M$  by the data.

The value  $L$  should reflect the true interaction order in the underlying model, but we almost never have such prior knowledge. Therefore we choose the optimal  $M$  and  $L$  using  $K$ -fold cross validation, starting with a fixed value of  $L$ . The data are split into  $K$  roughly equal-sized folds. Let an index function  $\pi(i) : \{1, \dots, n\} \mapsto \{1, \dots, K\}$  indicate the fold to which observation  $i$  is allocated. Each time, we remove the  $k$ th fold of the data ( $k = 1, 2, \dots, K$ ), and train the model using the remaining  $K - 1$  folds. Denoting by  $\hat{F}_{-k}^{[M]}(\mathbf{x})$  the resulting model, we compute the validation loss by predicting on each  $k$ th fold of the data removed:

$$\text{CV}(M, L) = \frac{1}{n} \sum_{i=1}^n \Psi(y_i, \hat{F}_{-\pi(i)}^{[M]}(\mathbf{x}_i; L) \mid \rho). \quad (26)$$

We select the optimal  $M$  at which the minimum validation loss is reached

$$\widehat{M}_L = \underset{M}{\operatorname{argmin}} \text{CV}(M, L).$$

If we need to select  $L$  too, then we repeat the whole process for several  $L$  (e.g.  $L = 2, 3, 4, 5$ ) and choose the one with the smallest minimum generalization error

$$\widehat{L} = \underset{L}{\operatorname{argmin}} \text{CV}(L, \widehat{M}_L).$$

For a given  $\nu$ , fitting trees with higher  $L$  leads to smaller  $M$  being required to reach the minimum error.

## 5 Simulation Studies

In this section, we compare TDboost with the Tweedie GLM model (TGLM: Jørgensen and de Souza, 1994) and the Tweedie GAM model in terms of the function estimation performance. The Tweedie GAM model is proposed by Wood (2001), which is based on a



penalized regression spline approach with automatic smoothness selection. There is an R package “MGCV” accompanying the work, available at <http://cran.r-project.org/web/packages/mgcv/index.html>. In all numerical examples below using the TDboost model, five-fold cross validation is adopted for selecting the optimal  $(M, L)$  pair, while the shrinkage factor  $\nu$  is set to its default value of 0.005.

## 5.1 Case I

In this simulation study, we demonstrate that TDboost is well suited to fit target functions that are non-linear or involve complex interactions. We consider two true target functions:

- **Model 1** (Discontinuous function): The target function is discontinuous as defined by  $F(x) = 0.5I(x > 0.5)$ . We assume  $x \sim \text{Unif}(0, 1)$ , and  $y \sim \text{Tw}(\mu, \phi, \rho)$  with  $\rho = 1.5$  and  $\phi = 0.5$ .
- **Model 2** (Complex interaction): The target function has two hills and two valleys.

$$F(x_1, x_2) = e^{-5(1-x_1)^2+x_2^2} + e^{-5x_1^2+(1-x_2)^2},$$

which corresponds to a common scenario where the effect of one variable changes depending on the effect of another. We assume  $x_1, x_2 \sim \text{Unif}(0, 1)$ , and  $y \sim \text{Tw}(\mu, \phi, \rho)$  with  $\rho = 1.5$  and  $\phi = 0.5$ .

We generate  $n = 1000$  observations for training and  $n' = 1000$  for testing, and fit the training data using TDboost, MGCV, and TGLM. Since the true target functions are known, we consider the mean absolute deviation (MAD) as performance criteria,

$$\text{MAD} = \frac{1}{n'} \sum_{i=1}^{n'} |F(\mathbf{x}_i) - \hat{F}(\mathbf{x}_i)|,$$

where both the true predictor function  $F(\mathbf{x}_i)$  and the predicted function  $\hat{F}(\mathbf{x}_i)$  are evaluated on the test set. The resulting MADs on the testing data are reported in Table 1, which are averaged over 100 independent replications. The fitted functions from Model 2 are plotted in Figure 2. In both cases, we find that TDboost outperforms TGLM and MGCV in terms of the ability to recover the true functions and gives the smallest prediction errors.

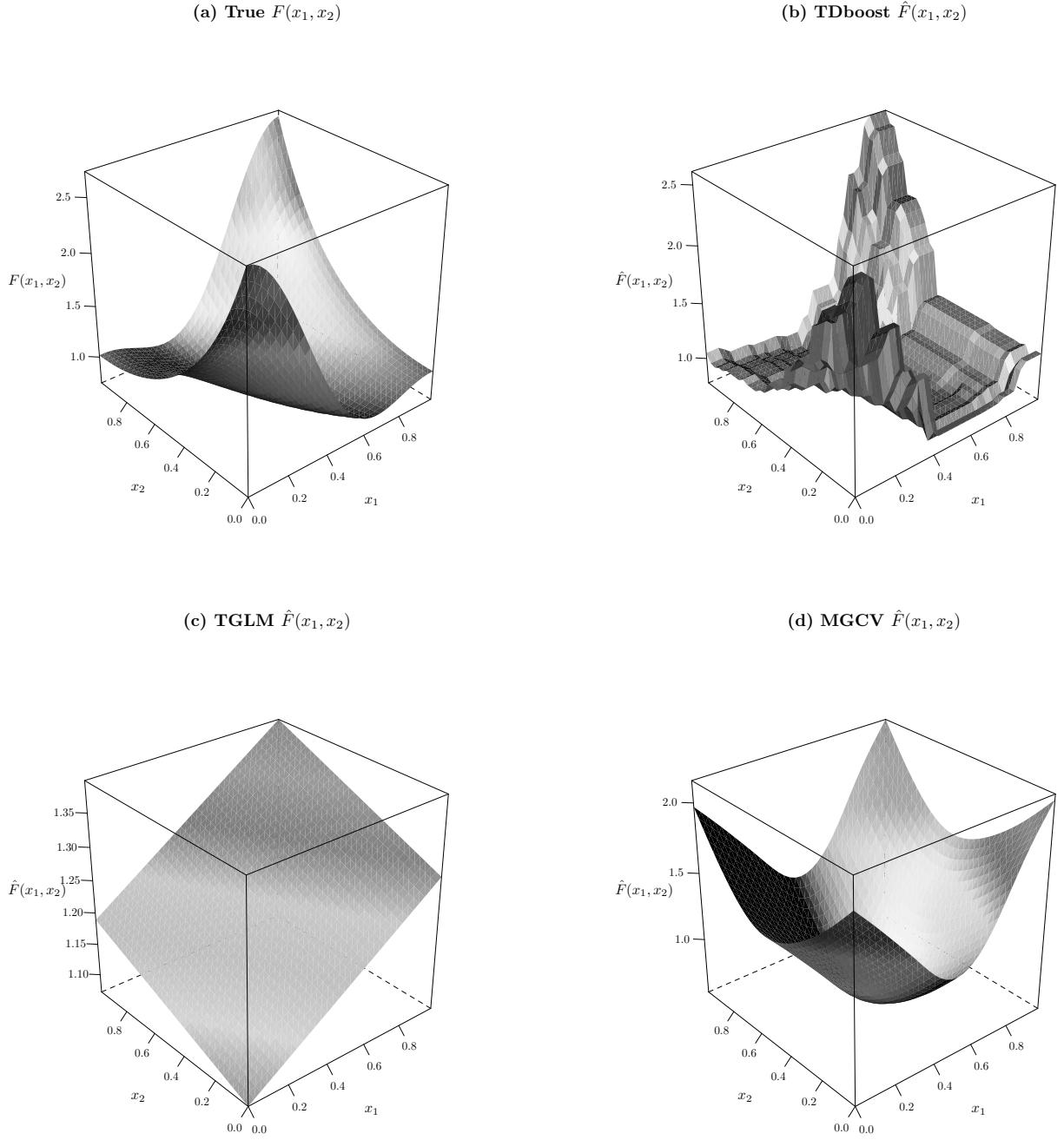


Figure 2: Fitted curves that recover the target function defined in Model 2. The top left figure shows the true target function. The top right, bottom left, and bottom right figures show the predictions on the testing data from TDboost, TGLM, and MGCV, respectively.

Model	TGLM	MGCV	TDboost
1	0.1102 (0.0006)	0.0752 (0.0016)	0.0595 (0.0021)
2	0.3516 (0.0009)	0.2511 (0.0004)	0.1034 (0.0008)

Table 1: The averaged MADs and the corresponding standard errors based on 100 independent replications.

## 5.2 Case II

The idea is to see the performance of the TDboost estimator and MGCV estimator on a variety of very complicated, randomly generated predictor functions, and study how the size of the training set, distribution settings and other characteristics of problems affect final performance of the two methods. We use the “random function generator” (RFG) model by Friedman (2001) in our simulation. The true target function  $F$  is randomly generated as a linear expansion of functions  $\{g_k\}_{k=1}^{20}$ :

$$F(\mathbf{x}) = \sum_{k=1}^{20} b_k g_k(\mathbf{z}_k). \quad (27)$$

Here each coefficient  $b_k$  is a uniform random variable from  $\text{Unif}[-1, 1]$ . Each  $g_k(\mathbf{z}_k)$  is a function of  $\mathbf{z}_k$ , where  $\mathbf{z}_k$  is defined as a  $p_k$ -sized subset of the ten-dimensional variable  $\mathbf{x}$  in the form

$$\mathbf{z}_k = \{x_{\psi_k(j)}\}_{j=1}^{p_k}, \quad (28)$$

where each  $\psi_k$  is an independent permutation of the integers  $\{1, \dots, p\}$ . The size  $p_k$  is randomly selected by  $\min(\lfloor 2.5 + r_k \rfloor, p)$ , where  $r_k$  is generated from an exponential distribution with mean 2. Hence the expected order of interactions presented in each  $g_k(\mathbf{z}_k)$  is between four and five. Each function  $g_k(\mathbf{z}_k)$  is a  $p_k$ -dimensional Gaussian function:

$$g_k(\mathbf{z}_k) = \exp \left\{ -\frac{1}{2}(\mathbf{z}_k - \mathbf{u}_k)^\top \mathbf{V}_k (\mathbf{z}_k - \mathbf{u}_k) \right\}, \quad (29)$$

where each mean vector  $\mathbf{u}_k$  is randomly generated from  $N(0, \mathbf{I}_{p_k})$ . The  $p_k \times p_k$  covariance matrix  $\mathbf{V}_k$  is defined by

$$\mathbf{V}_k = \mathbf{U}_k \mathbf{D}_k \mathbf{U}_k^\top, \quad (30)$$

where  $\mathbf{U}_k$  is a random orthonormal matrix,  $\mathbf{D}_k = \text{diag}\{d_k[1], \dots, d_k[p_k]\}$ , and the square root of each diagonal element  $\sqrt{d_k[j]}$  is a uniform random variable from  $\text{Unif}[0.1, 2.0]$ . We

generate data  $\{y_i, \mathbf{x}_i\}_{i=1}^n$  according to

$$y_i \sim \text{Tw}(\mu_i, \phi, \rho), \quad \mathbf{x}_i \sim \text{N}(0, \mathbf{I}_p), \quad i = 1, \dots, n, \quad (31)$$

where  $\mu_i = \exp\{F(\mathbf{x}_i)\}$ .

### Setting I: when the index is known

Firstly, we study the situation that the true index parameter  $\rho$  is known when fitting models. We generate data according to the RFG model with index parameter  $\tilde{\rho} = 1.5$  and the dispersion parameter  $\tilde{\phi} = 1$  in the true model. We set the number of predictors to be  $p = 10$  and generate  $n \in \{1000, 2000, 5000\}$  observations as training sets, on which both MGCV and TDboost are fitted with  $\rho$  specified to be the true value 1.5. An additional test set of  $n' = 5000$  observations was generated for evaluating the performance of the final estimate.

Figure 3 shows simulation results for comparing the estimation performance of MGCV and TDboost, when varying the training sample size. The empirical distributions of the MADs shown as box-plots are based on 100 independent replications. We can see that in all of the cases, TDboost outperforms MGCV in terms of prediction accuracy.

We also test estimation performance on  $\mu$  when the index parameter  $\rho$  is misspecified, that is, we use a guess value  $\rho$  differing from the true value  $\tilde{\rho}$  when fitting the TDboost model. Because  $\mu$  is statistically orthogonal to  $\phi$  and  $\rho$ , meaning that the off-diagonal elements of the Fisher information matrix are zero (Jørgensen, 1997), we expect  $\hat{\mu}$  will vary very slowly as  $\rho$  changes. Indeed, using the previous simulation data with the true value  $\tilde{\rho} = 1.5$  and  $\tilde{\phi} = 1$ , we fitted TDboost models with nine guess values of  $\rho \in \{1.1, 1.2, \dots, 1.9\}$ . The resulting MADs are displayed in Figure 4, which shows the choice of the value  $\rho$  has almost no significant effect on estimation accuracy of  $\mu$ .

### Setting II: using the estimated index

Next we study the situation that the true index parameter  $\rho$  is unknown, and we use the estimated  $\rho$  obtained from the profile likelihood procedure discussed in Section 4.2 for fitting the model. The same data generation scheme is adopted as in Setting I, except now both MGCV and TDboost are fitted with  $\rho$  estimated by maximizing the profile likelihood. Figure 5 shows simulation results for comparing the estimation performance of MGCV and TDboost

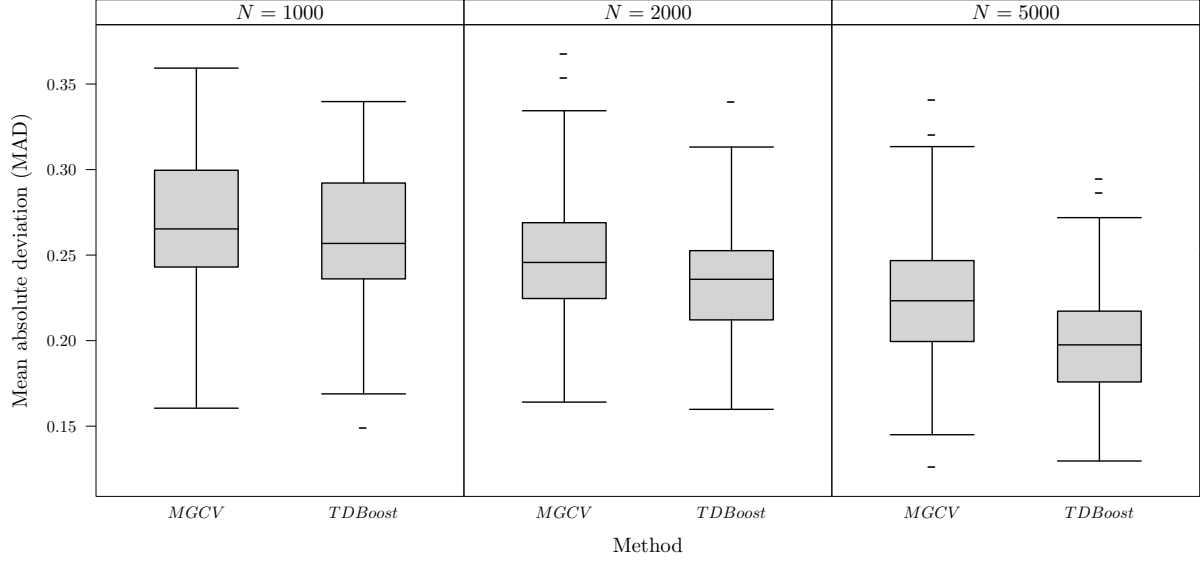


Figure 3: Simulation results for Setting I: compare the estimation performance of **MGCV** and **TDBOOST** when varying the training sample size and the dispersion parameter in the true model. Box-plots display empirical distributions of the MADs based on 100 independent replications.

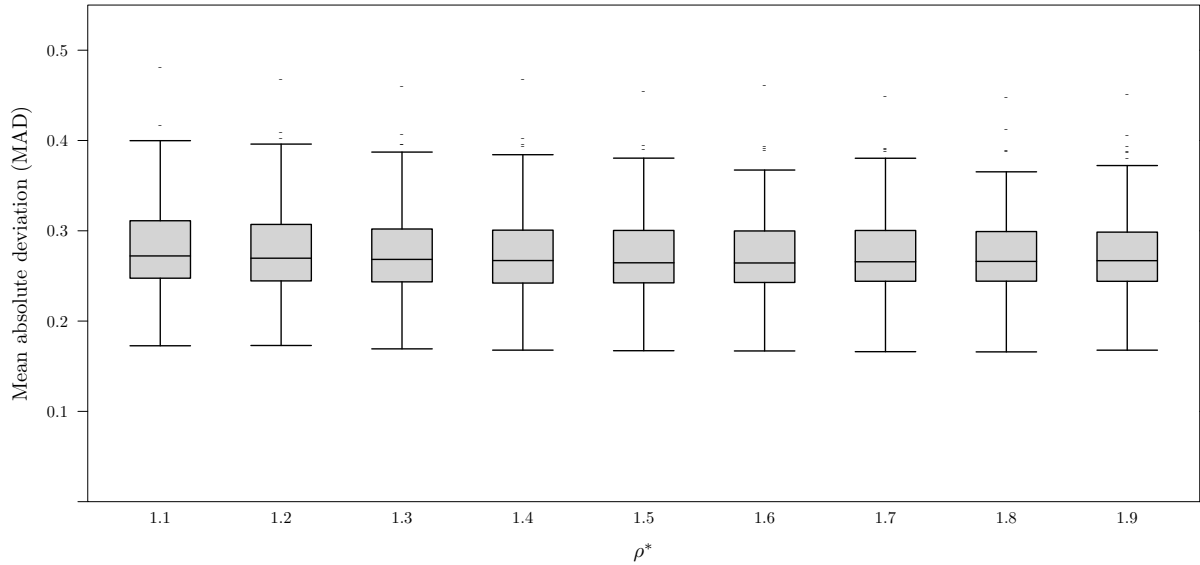


Figure 4: Simulation results for Setting I when the index is misspecified: the estimation performance of **TDBOOST** when varying the value of the index parameter  $\rho \in \{1.1, 1.2, \dots, 1.9\}$ . In the true model  $\tilde{\rho} = 1.5$  and  $\tilde{\phi} = 1$ . Box-plots show empirical distributions of the MADs based on 200 independent replications.

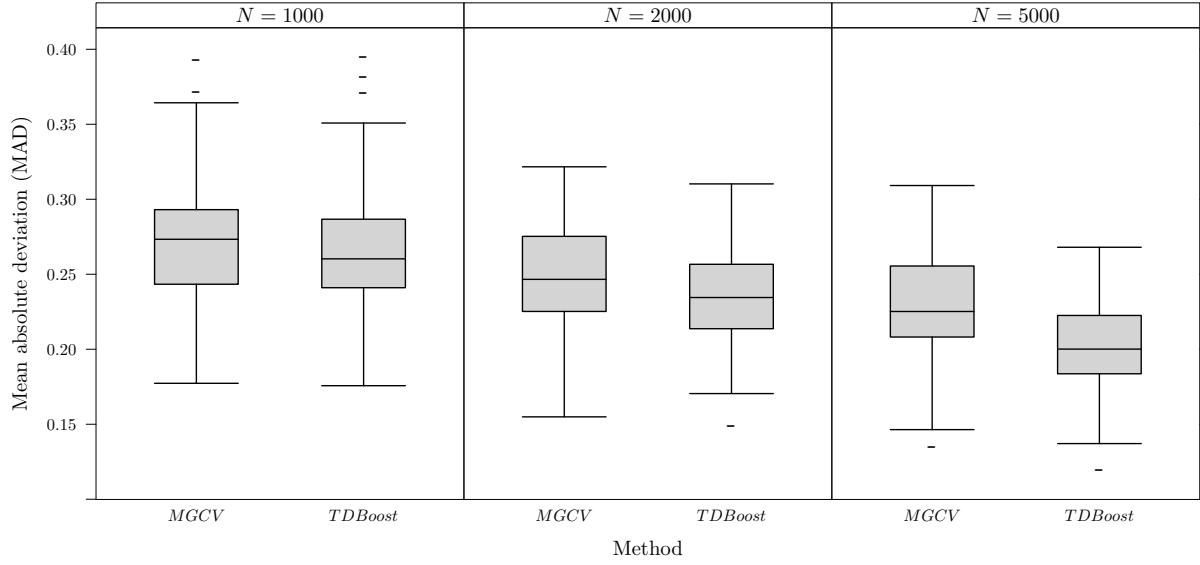


Figure 5: Simulation results for Setting II: compare the estimation performance of **MGCV** and **TDBOOST** when varying the training sample size and the dispersion parameter in the true model. Box-plots display empirical distributions of the MADs based on 100 independent replications.

in such setting. We can see that the results have no significant difference to the results of Setting I: TDBOOST still outperforms MGCV in terms of prediction accuracy when using the estimated  $\rho$  instead of the true value.

Lastly, we demonstrate our results from the estimation of the dispersion  $\phi$  and the index  $\rho$  by using the profile likelihood. A total number of 200 sets of training samples are randomly generated from a true model according to the setting (31) with  $\phi = 2$  and  $\rho = 1.7$ , each sample having 2000 observations. We fit the TDBOOST model on each sample and compute the estimates  $\phi^*$  at each of the 50 equally spaced values  $\{\rho_1, \dots, \rho_{50}\}$  on  $(1, 2)$ . The  $(\rho_j, \phi^*(\rho_j))$  corresponding to the maximal profile likelihood is the estimate of  $(\rho, \phi)$ . The estimation process is repeated 200 times. The estimated indices have mean  $\bar{\rho}^* = 1.68$  and standard error  $SE(\rho^*) = 0.026$ , so the true value  $\rho = 1.7$  is within  $\bar{\rho}^* \pm SE(\rho^*)$ . The estimated dispersions have mean  $\bar{\phi}^* = 1.82$  and standard error  $SE(\phi^*) = 0.12$ . Figure 6 shows the profile likelihood function of  $\rho$  for a single run.

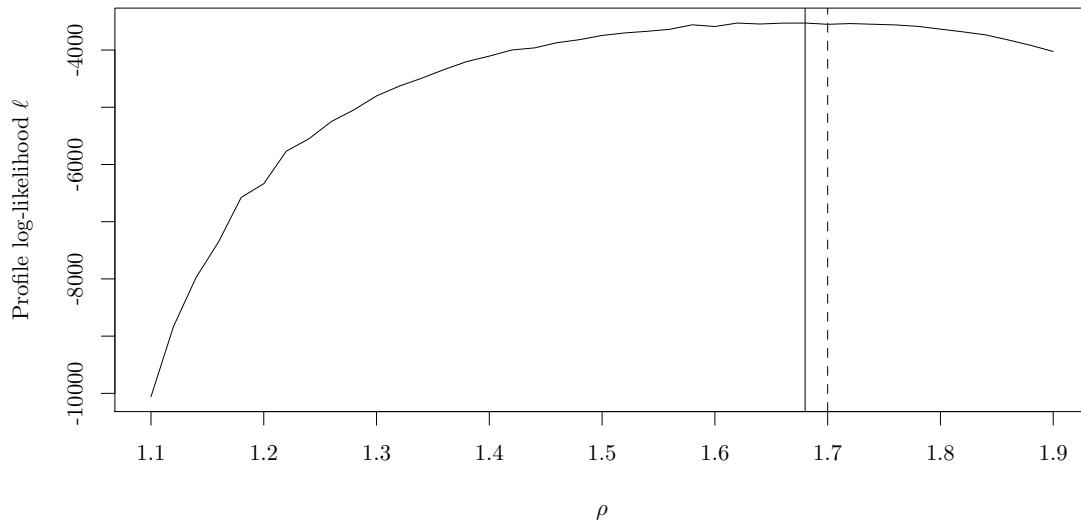


Figure 6: The curve represents the profile likelihood function of  $\rho$  from a single run. The dotted line shows the true value  $\rho = 1.7$ . The solid line shows the estimated value  $\rho^* = 1.68$  corresponding to the maximum likelihood. The associated estimated dispersion is  $\phi^* = 1.89$ .

## 6 Application: Automobile Claims

### 6.1 Dataset

We consider an auto insurance claim dataset as analyzed in Yip and Yau (2005) and Zhang and Yu (2005). The data set contains 10,296 driver vehicle records, each record including an individual driver's total claim amount ( $z_i$ ) in the last five years ( $w_i = 5$ ) and 17 characteristics  $x_i = (x_{i,1}, \dots, x_{i,17})$  for the driver and the insured vehicle. We want to predict the expected pure premium based on  $x_i$ . Table 3 summarize the data set. The descriptive statistics of the data are provided in Appendix Part D. The histogram of the total claim amounts in Figure 1 shows that the empirical distribution of these values is highly skewed. We find that approximately 61.1% of policyholders had no claims, and approximately 29.6% of the policyholders had a positive claim amount up to 10,000 dollars. Note that only 9.3% of the policyholders had a high claim amount above 10,000 dollars, but the sum of their claim amount made up to 64% of the overall sum. Another important feature of the data is that there are interactions among explanatory variables. For example, from Table 2 we can

REVOKED		AREA	
		Urban	Rural
	No Yes	3150.57 14551.62	904.70 7624.36
Difference		11401.05	6719.66

Table 2: The averaged total claim amount for different categories of the policyholders.

ID	Variable	Type	Description
1	AGE	N	Driver’s age
2	BLUEBOOK	N	Value of vehicle
3	HOMEKIDS	N	Number of children
4	KIDSDRIV	N	Number of driving children
5	MVR_PTS	N	Motor vehicle record points
6	NPOLICY	N	Number of policies
7	RETAINED	N	Number of years as a customer
8	TRAVTIME	N	Distance to work
9	AREA	C	Home/work area: Rural, Urban
10	CAR_USE	C	Vehicle use: Commercial, Private
11	CAR_TYPE	C	Type of vehicle: Panel Truck, Pickup, Sedan, Sports Car, SUV, Van
12	GENDER	C	Driver’s gender: F, M
13	JOBCLASS	C	Unknown, Blue Collar, Clerical, Doctor, Home Maker, Lawyer, Manager, Professional, Student
14	MAX_EDUC	C	Education level: High School or Below, Bachelors, High School, Masters, PhD
15	MARRIED	C	Married or not: Yes, No
16	REVOKED	C	Whether license revoked in past 7 years: Yes, No

Table 3: Explanatory variables in the claim history data set. Type N stands for numerical variable, Type C stands for categorical variable.

see that the marginal effect of the variable REVOKED on the total claim amount is much greater for the policyholders living in the urban area than those living in the rural area. The importance of the interaction effects will be confirmed later in our data analysis.

## 6.2 Models

We separate the entire dataset into a training set and a testing set with equal size. Then the TDboost model is fitted on the training set and tuned with five-fold cross validation.



For comparison, we also fit TGLM and MGCV, both of which are fitted using all the explanatory variables. In MGCV, the numerical variables AGE, BLUEBOOK, HOMEKIDS, KIDSDRIV, MVR\_PTS, NPOLICY, RETAINED and TRAVTIME are modeled by smooth terms represented using penalized regression splines. We find the appropriate smoothness for each applicable model term using Generalized Cross Validation (GCV) (Wahba, 1990). For the TDboost model, it is not necessary to carry out data transformation, since the tree-based boosting method can automatically handle different types of data. For other models, we use logarithmic transformation on BLUEBOOK, i.e.  $\log(\text{BLUEBOOK})$ , and scale all the numerical variables except for HOMEKIDS, KIDSDRIV, MVR\_PTS and NPOLICY to have mean 0 and standard deviation 1. We also create dummy variables for the categorical variables with more than two levels (CAR\_TYPE, JOBCLASS and MAX\_EDUC). For all models, we use the profile likelihood method to estimate the dispersion  $\phi$  and the index  $\rho$ , which are in turn used in fitting the final models.

### 6.3 Performance comparison

To examine the performance of TGLM, MGCV and TDboost, after fitting on the training set, we predict the pure premium  $P(\mathbf{x}) = \hat{\mu}(\mathbf{x})$  by applying each model on the independent held-out testing set. However, attention must be paid when measuring the differences between predicted premiums  $P(\mathbf{x})$  and real losses  $y$  on the testing data. The mean squared loss or mean absolute loss is not appropriate here because the losses have high proportions of zeros and are highly right skewed. Therefore an alternative statistical measure – the ordered Lorenz curve and the associated Gini index – proposed by Frees et al. (2011) are used for capturing the discrepancy between the premium and loss distributions. By calculating the Gini index, the performance of different predictive models can be compared. Here we only briefly explain the idea of the ordered Lorenz curve (Frees et al., 2011, 2013). Let  $B(\mathbf{x})$  be the “base premium”, which is calculated using the existing premium prediction model, and let  $P(\mathbf{x})$  be the “competing premium” calculated using an alternative premium prediction model. In the ordered Lorenz curve, the distribution of losses and the distribution of premiums are sorted based on the relative premium  $R(\mathbf{x}) = P(\mathbf{x})/B(\mathbf{x})$ . The ordered premium distribution is

$$\hat{D}_P(s) = \frac{\sum_{i=1}^n B(\mathbf{x}_i) I(R(\mathbf{x}_i) \leq s)}{\sum_{i=1}^n B(\mathbf{x}_i)},$$

and the ordered loss distribution is

$$\hat{D}_L(s) = \frac{\sum_{i=1}^n y_i I(R(\mathbf{x}_i) \leq s)}{\sum_{i=1}^n y_i}.$$

Two empirical distributions are based on the same sort order, which makes it possible to compare the premium and loss distributions for the same policyholder group. The ordered Lorenz curve is the graph of  $(\hat{D}_P(s), \hat{D}_L(s))$ . When the percentage of losses equals the percentage of premiums for the insurer, the curve results in a 45-degree line, known as “the line of equality”. Twice the area between the ordered Lorenz curve and the line of equality measures the discrepancy between the premium and loss distributions, and is defined as the Gini index. Curves below the line of equality indicate that, given knowledge of the relative premium, an insurer could identify the profitable contracts, whose premiums are greater than losses. Therefore, a larger Gini index (hence a larger area between the line of equality and the curve below) would imply a more favorable model.

Following Frees et al. (2013), we successively specify the prediction from each model as the base premium  $B(\mathbf{x})$  and use predictions from the remaining models as the competing premium  $P(\mathbf{x})$  to compute the Gini indices. The entire procedure of the data splitting and Gini index computation are repeated 20 times, and a matrix of the averaged Gini indices and standard errors is reported in Table 4. To pick the “best” model, we use a “minimax” strategy (Frees et al., 2013) to select the base premium model that are least vulnerable to competing premium models; that is, we select the model that provides the smallest of the maximal Gini indices, taken over competing premiums. We find that the maximal Gini index is 15.528 when using  $B(\mathbf{x}) = \hat{\mu}^{\text{TGLM}}(\mathbf{x})$  as the base premium, 12.979 when  $B(\mathbf{x}) = \hat{\mu}^{\text{MGCv}}(\mathbf{x})$ , and 4.000 when  $B(\mathbf{x}) = \hat{\mu}^{\text{TDboost}}(\mathbf{x})$ . Therefore, TDboost has the smallest maximum Gini index at 4.000, hence is the least vulnerable to alternative scores. Figure 7 also shows that when TGLM (or MGCv) is selected as the base premium, the area between the line of equality and the ordered Lorenz curve is larger when choosing TDboost as the competing premium, indicating again that the TDboost model represents the most favorable choice.

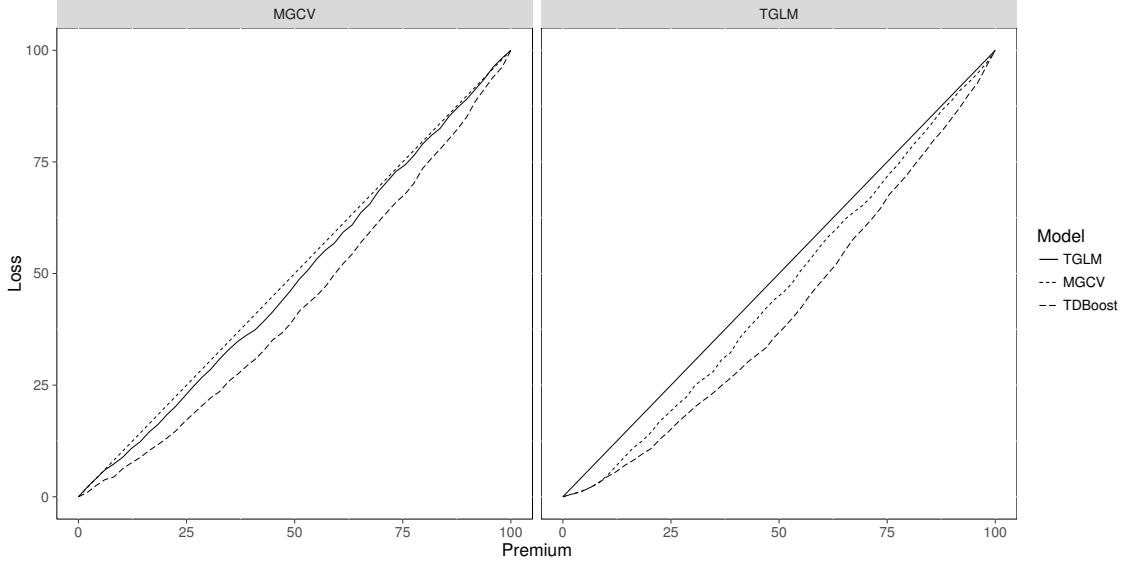


Figure 7: The ordered Lorenz curves for the auto insurance claim data.

Base Premium	Competing Premium		
	TGLM	MGCV	TDboost
TGLM	0	7.833 (0.338)	15.528 (0.509)
MGCV	3.044 (0.610)	0	12.979 (0.473)
TDboost	4.000 (0.364)	3.540 (0.415)	0

Table 4: The averaged Gini indices and standard errors in the auto insurance claim data example based on 20 random splits.

## 6.4 Interpreting the results

Next, we focus on the analysis using the TDboost model. There are several explanatory variables significantly related to the pure premium. The VI measure and the baseline value of each explanatory variable are shown in Figure 8. We find that REVOKED, MVR\_PTS, AREA and BLUEBOOK have high VI measure scores (the vertical line), and their scores all surpass the corresponding baselines (the horizontal line-length), indicating that the importance of those explanatory variables is real. We also find the variables AGE, JOBCLASS, CAR\_TYPE, NPOLICY, MAX\_EDUC, MARRIED, KIDSDRIV and CAR\_USE have larger-than-baseline VI measure scores, but the absolute scales are much less than aforementioned four variables. On the other hand, although the VI measure of, e.g., TRAVTIME is quite large, it does not significantly surpass the baseline importance.

We now use the partial dependence plots to visualize the fitted model. Figure 9 shows

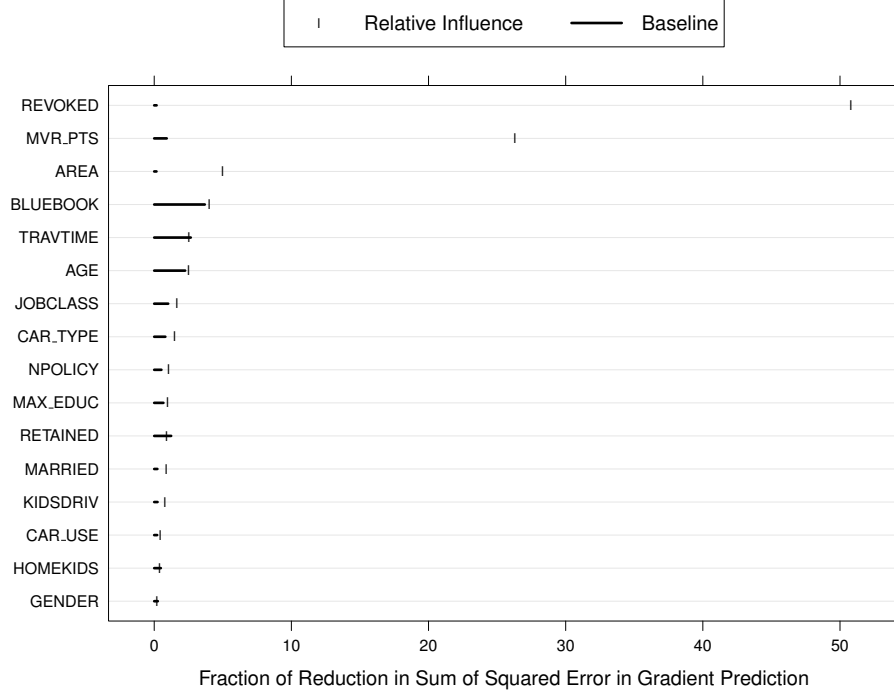


Figure 8: The variable importance measures and baselines of 17 explanatory variables for modeling the pure premium.

the main effects of four important explanatory variables on the pure premium. We clearly see that the strong nonlinear effects exist in predictors BLUEBOOK and MVR\_PTS: for the policyholders whose vehicle values are below 40K, their pure premium is negatively associated with the value of vehicle; after the value of vehicle passes 40K, the pure premium curve reaches a plateau; Additionally, the pure premium is positively associated with motor vehicle record points MVR\_PTS, but the pure premium curve reaches a plateau when MVR\_PTS exceeds six. On the other hand, the partial dependence plots suggest that a policyholder who lives in the urban area (AREA="URBAN") or with driver's license revoked (REVOKED="YES") typically has relatively high pure premium.

In our model, the data-driven choice for the tree size is  $L = 7$ , which means that our model includes higher order interactions. In Figure 10, we visualize the effects of four important second order interactions using the joint partial dependence plots. These four interactions are  $AREA \times MVR\_PTS$ ,  $AREA \times NPOLICY$ ,  $AREA \times REVOKED$  and  $AREA \times TRAVTIME$ . These four interactions all involve the variable AREA: we can see that the marginal effects of MVR\_PTS, NPOLICY, REVOKED and TRAVTIME on the pure

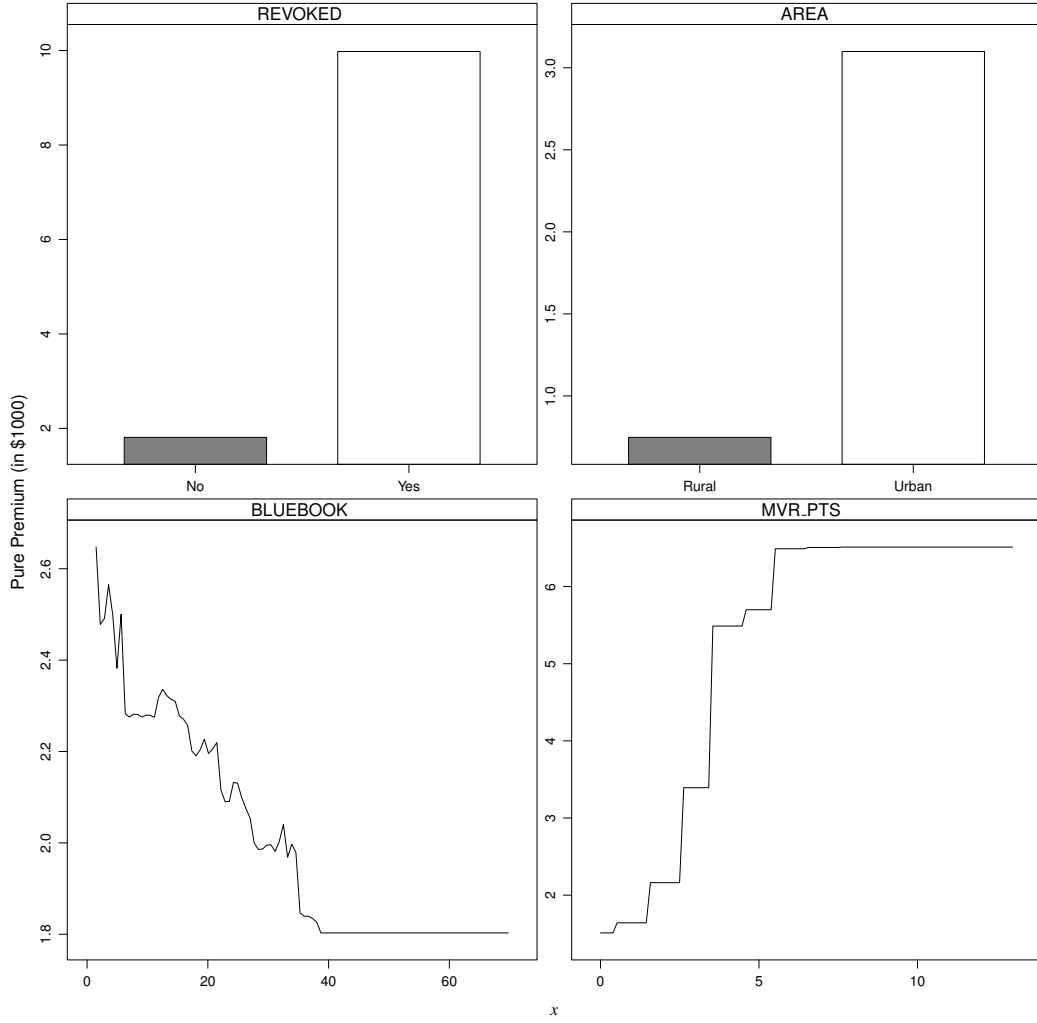


Figure 9: Marginal effects of four most significant explanatory variables on the pure premium.

premium are greater for the policyholders living in the urban area (AREA=“URBAN”) than those living in the rural area (AREA=“RURAL”). For example, a strong  $\text{AREA} \times \text{MVR\_PTS}$  interaction suggests that for the policyholders living in the rural area, motor vehicle record points of the policyholders have a weaker positive marginal effect on the expected pure premium than for the policyholders living in the urban area.

## 7 Conclusions

The need for nonlinear risk factors as well as risk factor interactions for modeling insurance claim sizes is well-recognized by actuarial practitioners, but practical tools to study them

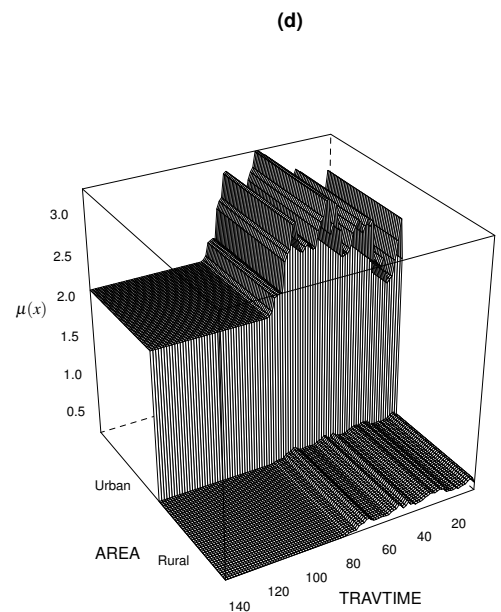
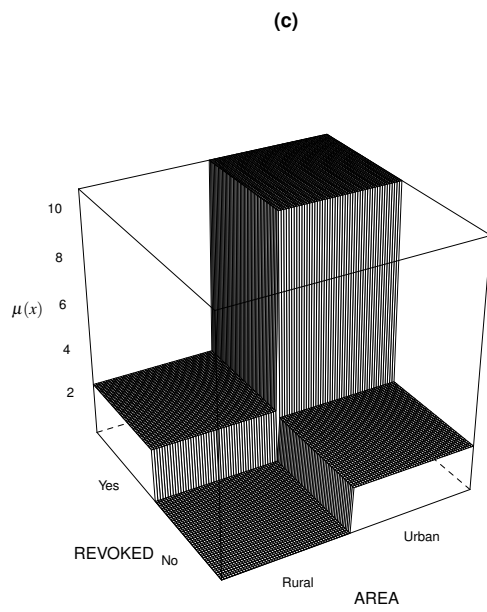
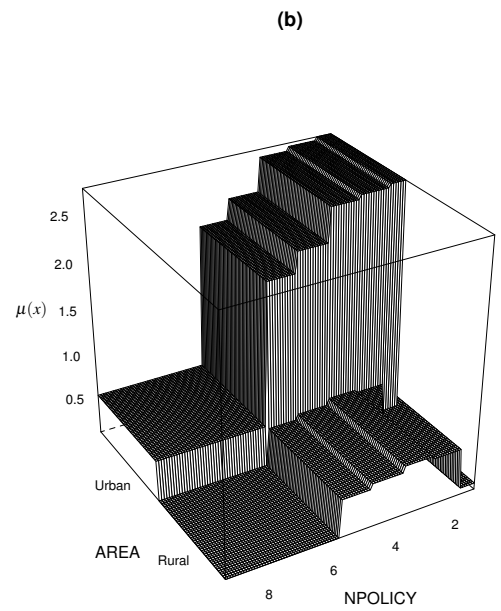
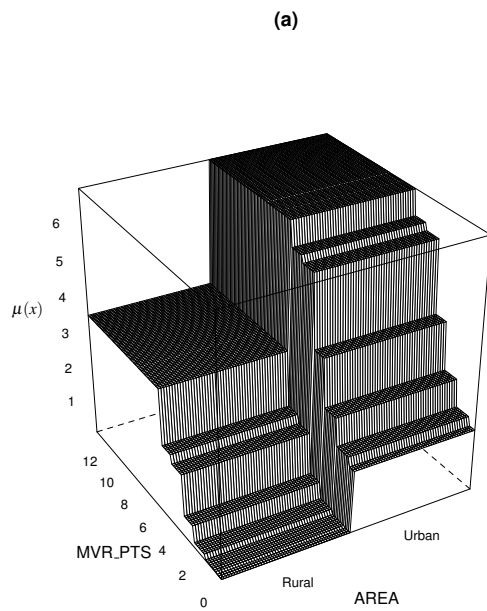


Figure 10: Four strong pairwise interactions.

are very limited. In this paper, relying on neither the linear assumption nor a pre-specified interaction structure, a flexible tree-based gradient boosting method is designed for the Tweedie model. We implement the proposed method in a user-friendly R package “TDboost” that can make accurate insurance premium predictions for complex data sets and serve as a convenient tool for actuarial practitioners to investigate the nonlinear and interaction effects. In the context of personal auto insurance, we implicitly use the policy duration as a volume measure (or exposure), and demonstrate the favorable prediction performance of TDboost for the pure premium. In cases that exposure measures other than duration are used, which is common in commercial insurance, we can extend the TDboost method to the corresponding claim size by simply replacing the duration with any chosen exposure measure.

TDboost can also be an important complement to the traditional GLM model in insurance rating. Even under the strict circumstances that the regulators demand the final model to have a GLM structure, our approach can still be quite helpful due to its ability to extract additional information such as non-monotonicity/non-linearity and important interaction. In Appendix Part E, we provide an additional real data analysis to demonstrate that our method can provide insights into the structure of interaction terms. After integrating the obtained information about the interaction terms into the original GLM model, we can much enhance the overall accuracy of the insurance premium prediction while maintaining a GLM model structure.

In addition, it is worth mentioning that the applications of the proposed method can go beyond the insurance premium prediction and be of interest to researchers in many other fields including ecology (Foster and Bravington, 2013), meteorology (Dunn, 2004) and political science (Lauderdale, 2012). See, for example, Dunn and Smyth (2005) and Qian et al. (2015) for descriptions of the broad Tweedie distribution applications. The proposed method and the implementation tool allow researchers in these related fields to venture outside the Tweedie GLM modeling framework, build new flexible models from nonparametric perspectives, and use the model interpretation tools demonstrated in our real data analysis to study their own problems of interests.

## References

- Anstey, K. J., Wood, J., Lord, S., and Walker, J. G. (2005), “Cognitive, sensory and physical factors enabling driving safety in older adults,” *Clinical psychology review*, 25, 45–65.
- Breiman, L. (1996), “Bagging predictors,” *Machine learning*, 24, 123–140.
- (1998), “Arcing classifier (with discussion and a rejoinder by the author),” *The Annals of Statistics*, 26, 801–849.
- (1999), “Prediction games and arcing algorithms,” *Neural Computation*, 11, 1493–1517.
- Breiman, L., Friedman, J., Olshen, R., Stone, C., Steinberg, D., and Colla, P. (1984), “CART: Classification and regression trees,” *Wadsworth*.
- Brent, R. P. (2013), *Algorithms for minimization without derivatives*, Courier Dover Publications.
- Bühlmann, P. and Hothorn, T. (2007), “Boosting algorithms: Regularization, prediction and model fitting,” *Statistical Science*, 22, 477–505.
- Dionne, G., Gouriéroux, C., and Vanasse, C. (2001), “Testing for evidence of adverse selection in the automobile insurance market: A comment,” *Journal of Political Economy*, 109, 444–453.
- Dunn, P. K. (2004), “Occurrence and quantity of precipitation can be modelled simultaneously,” *International Journal of Climatology*, 24, 1231–1239.
- Dunn, P. K. and Smyth, G. K. (2005), “Series evaluation of Tweedie exponential dispersion model densities,” *Statistics and Computing*, 15, 267–280.
- Elith, J., Leathwick, J. R., and Hastie, T. (2008), “A working guide to boosted regression trees,” *Journal of Animal Ecology*, 77, 802–813.
- Foster, S. D. and Bravington, M. V. (2013), “A Poisson–Gamma model for analysis of ecological non-negative continuous data,” *Environmental and ecological statistics*, 20, 533–552.



- Frees, E. W., Meyers, G., and Cummings, A. D. (2011), “Summarizing insurance scores using a Gini index,” *Journal of the American Statistical Association*, 106.
- Frees, E. W. J., Meyers, G., and Cummings, A. D. (2013), “Insurance ratemaking and a Gini index,” *Journal of Risk and Insurance*.
- Freund, Y. and Schapire, R. (1997), “A decision-theoretic generalization of on-line learning and an application to boosting,” *Journal of Computer and System Sciences*, 55, 119–139.
- Friedman, J. (2001), “Greedy function approximation: A gradient boosting machine,” *The Annals of Statistics*, 29, 1189–1232.
- Friedman, J., Hastie, T., and Tibshirani, R. (2000), “Additive logistic regression: A statistical view of boosting (With discussion and a rejoinder by the authors),” *The Annals of Statistics*, 28, 337–407.
- Friedman, J. H. (2002), “Stochastic gradient boosting,” *Computational Statistics & Data Analysis*, 38, 367–378.
- Haberman, S. and Renshaw, A. E. (1996), “Generalized linear models and actuarial science,” *Statistician*, 45, 407–436.
- Hastie, T., Tibshirani, R., and Friedman, J. (2009), *The elements of statistical learning: Data mining, inference, and prediction. Second Edition.*, Springer Series in Statistics, Springer.
- Hastie, T. J. and Tibshirani, R. J. (1990), *Generalized additive models*, vol. 43, CRC Press.
- Jørgensen, B. (1987), “Exponential dispersion models,” *Journal of the Royal Statistical Society. Series B (Methodological)*, 127–162.
- (1997), *The theory of dispersion models*, vol. 76, CRC Press.
- Jørgensen, B. and de Souza, M. C. (1994), “Fitting Tweedie’s compound Poisson model to insurance claims data,” *Scandinavian Actuarial Journal*, 1994, 69–93.
- Lauderdale, B. E. (2012), “Compound Poisson–Gamma regression models for dollar outcomes that are sometimes zero,” *Political Analysis*, 20, 387–399.

- Mildenhall, S. J. (1999), “A systematic relationship between minimum bias and generalized linear models,” in *Proceedings of the Casualty Actuarial Society*, vol. 86, pp. 393–487.
- Murphy, K. P., Brockman, M. J., and Lee, P. K. (2000), “Using generalized linear models to build dynamic pricing systems,” in *Casualty Actuarial Society Forum, Winter*, pp. 107–139.
- Nelder, J. and Wedderburn, R. (1972), “Generalized Linear Models,” *Journal of the Royal Statistical Society. Series A (General)*, 135, 370–384.
- Ohlsson, E. and Johansson, B. (2010), *Non-life insurance pricing with generalized linear models*, Springer.
- Peters, G. W., Shevchenko, P. V., and Wüthrich, M. V. (2008), “Model risk in claims reserving within Tweedie’s compound Poisson models,” *ASTIN Bulletin*, to appear.
- Qian, W., Yang, Y., and Zou, H. (2015), “Tweedie’s compound Poisson model with grouped elastic net,” *Journal of Computational and Graphical Statistics*, preprint.
- Renshaw, A. E. (1994), “Modelling the claims process in the presence of covariates,” *ASTIN Bulletin*, 24, 265–285.
- Ridgeway, G. (2007), “Generalized Boosted Regression Models,” *R package manual*.
- Sandri, M. and Zuccolotto, P. (2008), “A bias correction algorithm for the Gini variable importance measure in classification trees,” *Journal of Computational and Graphical Statistics*, 17.
- (2010), “Analysis and correction of bias in Total Decrease in Node Impurity measures for tree-based algorithms,” *Statistics and Computing*, 20, 393–407.
- Showers, V. E. and Shotick, J. A. (1994), “The effects of household characteristics on demand for insurance: A tobit analysis,” *Journal of Risk and Insurance*, 492–502.
- Smyth, G. and Jørgensen, B. (2002), “Fitting Tweedie’s compound Poisson model to insurance claims data: Dispersion modelling,” *ASTIN Bulletin*, 32, 143–157.

- Smyth, G. K. (1996), “Regression analysis of quantity data with exact zeros,” in *Proceedings of the second Australia–Japan workshop on stochastic models in engineering, technology and management*, Citeseer, pp. 572–580.
- Tweedie, M. (1984), “An index which distinguishes between some important exponential families,” in *Statistics: Applications and New Directions: Proc. Indian Statistical Institute Golden Jubilee International Conference*, pp. 579–604.
- Van de Ven, W. and van Praag, B. M. (1981), “Risk aversion and deductibles in private health insurance: Application of an adjusted tobit model to family health care expenditures,” *Health, economics, and health economics*, 125–48.
- Wahba, G. (1990), *Spline models for observational data*, vol. 59, SIAM.
- White, A. P. and Liu, W. Z. (1994), “Technical note: Bias in information-based measures in decision tree induction,” *Machine Learning*, 15, 321–329.
- Wood, S. (2001), “mgcv: GAMs and generalized ridge regression for R,” *R News*, 1, 20–25.
- (2006), *Generalized additive models: An introduction with R*, CRC press.
- Yip, K. C. and Yau, K. K. (2005), “On modeling claim frequency data in general insurance with extra zeros,” *Insurance: Mathematics and Economics*, 36, 153–163.
- Zhang, T. and Yu, B. (2005), “Boosting with early stopping: Convergence and consistency,” *The Annals of Statistics*, 1538–1579.
- Zhang, W. (2011), “cplm: Monte Carlo EM algorithms and Bayesian methods for fitting Tweedie compound Poisson linear models,” *R package*, <http://cran.r-project.org/web/packages/cplm/index.html>.