

**Technometrics** 



ISSN: 0040-1706 (Print) 1537-2723 (Online) Journal homepage: https://www.tandfonline.com/loi/utch20

## A Unified Approach to Sparse Tweedie Modeling of Multisource Insurance Claim Data

Simon Fontaine, Yi Yang, Wei Qian, Yuwen Gu & Bo Fan

To cite this article: Simon Fontaine, Yi Yang, Wei Qian, Yuwen Gu & Bo Fan (2020) A Unified Approach to Sparse Tweedie Modeling of Multisource Insurance Claim Data, Technometrics, 62:3, 339-356, DOI: 10.1080/00401706.2019.1647881

To link to this article: https://doi.org/10.1080/00401706.2019.1647881

View supplementary material 🗹



Accepted author version posted online: 26 Jul 2019. Published online: 05 Sep 2019.

Submit your article to this journal 🗹



View related articles

View Crossmark data 🗹

### A Unified Approach to Sparse Tweedie Modeling of Multisource Insurance Claim Data

Simon Fontaine<sup>a</sup>, Yi Yang<sup>b</sup>, Wei Qian<sup>c</sup>, Yuwen Gu<sup>d</sup>, and Bo Fan<sup>e</sup>

<sup>a</sup>Department of Statistics, University of Michigan, Ann Arbor, MI; <sup>b</sup>Department of Mathematics and Statistics, McGill University, Montreal, Quebec, Canada; <sup>c</sup>Department of Applied Economics and Statistics, University of Delaware, Newark, DE; <sup>d</sup>Department of Statistics, University of Connecticut, Storrs, CT; <sup>e</sup>Department of Statistics, University of Oxford, Oxford, UK

#### ABSTRACT

Actuarial practitioners now have access to multiple sources of insurance data corresponding to various situations: multiple business lines, umbrella coverage, multiple hazards, and so on. Despite the wide use and simple nature of single-target approaches, modeling these types of data may benefit from an approach performing variable selection jointly across the sources. We propose a unified algorithm to perform sparse learning of such fused insurance data under the Tweedie (compound Poisson) model. By integrating ideas from multitask sparse learning and sparse Tweedie modeling, our algorithm produces flexible regularization that balances predictor sparsity and between-sources sparsity. When applied to simulated and real data, our approach clearly outperforms single-target modeling in both prediction and selection accuracy, notably when the sources do not have exactly the same set of predictors. An efficient implementation of the proposed algorithm is provided in our R package MStweedie, which is available at https://github.com/fontaine618/MStweedie.

#### **ARTICLE HISTORY**

Received June 2018 Accepted June 2019

#### **KEYWORDS**

Backtracking line search; Groupwise proximal gradient descent; Multisource insurance data; Multitask learning; Regularization; Tweedie model

Taylor & Francis

Check for updates

Taylor & Francis Group

#### 1. Introduction

Insurance claim data are characterized by excess zeros, corresponding to insurance policies without any claims, and highly right-skewed positive values associated with nonzero claim amounts, typically in monetary value. The modeling of insurance claim data helps to predict the expected loss associated with a portfolio of policies and is widely used for premium pricing. As claim data reflect a unique mixed nature of distributions with both discrete and continuous components, there are generally two popular modeling approaches. The first type considers a frequency-severity approach where claim frequency (i.e., whether a claim exists or not) and claim amount are modeled separately (Yip and Yau 2005; Frees, Gao, and Rosenberg 2011; Shi, Feng, and Ivantsova 2015), so that the two models need to be used together for claim loss prediction. The second type uses Tweedie's compound Poisson model (or Tweedie model for short; Tweedie 1984) that considers an inherent Poisson process and models both components simultaneously. Our study will focus on the second approach that draws upon Tweedie distribution's natural structure for claim data modeling (Smyth and Jørgensen 2002; Frees, Meyers, and Cummings 2011; Zhang 2013; Shi, Feng, and Boucher 2016). It is also common practice that insurers collect and maintain external information associated with insurance policies either directly from policy holders or from third-party databases. Covariates generated from the external information can be associated with the claim loss and help improve the modeling process. Depending on the type of insurance, this information may consist of policyholder's characteristics (demographics, employment status, experience, etc.) of the insured object characteristics (car type and usage, property value, etc.) or of any other characteristic deemed relevant.

Traditionally, actuarial practitioners adopt a single-target approach that, for a given insurance product, assumes one population to be homogeneously characterized by some covariates and aims to build a single Tweedie model solely from the product's sample data. Despite the wide use and simple nature of this approach, practitioners now have access to multiple sources of insurance data. For instance, many insurers have multiple business lines such as the auto insurance and the property insurance; in umbrella coverage, claim amounts are available for multiple types of coverage and for different hazard causes of the same coverage; multiple datasets can be accumulated for a long period of time, during which business environment may have changed significantly so that earlier-year and later-year data sources may not be treated as one homogeneous population. As a result, the modern multisource insurance data may not be characterized well by a homogeneous model. With these emerging multisource insurance data problems, much attention has been drawn to addressing their modeling issues in statistics and actuarial science. Both the frequency-severity and Tweedie model approaches have been investigated in the context of multivariate regressions to model the multiple responses simultaneously (see Frees, Lee, and Yang 2016; Shi 2016 and references therein).

Variable selection is one of the most important tasks in building transparent and interpretable models for claim loss prediction. Large-scale high-dimensional sparse modeling is

CONTACT Yi Yang *yi.yang6@mcgill.ca* Department of Mathematics and Statistics, McGill University, Montreal, Quebec H3A 0G4, Canada. Color versions of one or more of the figures in the article can be found online at *www.tandfonline.com/r/TECH*. Supplementary materials for this article are available online. Please go to *www.tandfonline.com/r/TECH*.

© 2019 American Statistical Association and the American Society for Quality

commonly encountered as hundreds of covariates are often considered as candidate variables while only a few of them are believed to be associated with the claim loss or can be used in the final model production. Under the single population setting, efficient variable selection approaches designed for the Tweedie model have been developed via a shrinkagetype approach (see Qian, Yang, and Zou 2016 and references therein). The increasingly prevalent multisource data scenarios coupled with high dimensionality and large data scale pose new challenges to actuarial practitioners. To our knowledge, the corresponding variable selection issues for multisource Tweedie models have not been studied in the literature. On the one hand, simply treating all different data sources as if they were from one population is problematic due to severe model misspecification. On the other hand, it may not be ideal either to perform variable selection separately on each individual data source because it often results in a loss of estimation efficiency. In the aforementioned multiline, multitype, or multiyear scenarios, the different data sources often contain similar types of covariates and some (or all) of them can be relevant across some (or all) data sources, even if different data come from totally different sets of customers. For example, both auto and property insurance contain geographical, credit, and experience variables that may be important in both lines of business. Therefore, a proper variable selection process should ideally take advantage of the potential connections among data sources as opposed to simply treating each data source separately.

In this article, we augment the multisource claim data analysis through an integrated shrinkage-based Tweedie modeling approach that fuses different data sources to find commonly shared relevant covariates. To insure our method is plausible, we will assume that the different sources have some (if not all) covariates in common. At the same time, our method retains the ability to recover model structures and covariates unique to individual data sources. In particular, we impose a composite adaptive lasso-type penalty (Tibshirani 1996; Zou 2006; Simon et al. 2013) in the composite Tweedie model to obtain both common and source-specific variables simultaneously. We study several different candidate penalty terms for our multisource data setting and devise a new algorithm (named MStweedie) to efficiently solve the corresponding optimization problems in a unified fashion. Our proposal is closely related to the celebrated multitask lasso in Lounici et al. (2011) that intends to uncover shared information across different tasks while achieving improved estimation, selection and prediction efficiency compared to independent variable selection (Obozinski, Wainwright, and Jordan 2008; Lounici et al. 2009; Huang and Zhang 2010; Lounici et al. 2011) under least square settings. Different from the existing multitask learning and related studies that mainly focused on the least squares (see, e.g., Jenatton et al. 2010; Morales, Micchelli, and Pontil 2010; Kim and Xing 2012) or classification (see, e.g., Zhang et al. 2008; Friedman, Hastie, and Tibshirani 2010; Obozinski, Taskar, and Jordan 2010; Vincent and Hansen 2014; Qian et al. 2019) setting, our proposal solves the important challenges posed by the semicontinuous, highly right-skewed claim data with excess zeros which cannot be efficiently modeled by a Gaussian or logit distribution. In particular, we show that the MStweedie algorithm is theoretically guaranteed to converge to the optimization target with at least linear rate, and is practically flexible to handle source-specific missing covariates. In addition, we implement our proposal in an efficient and user-friendly R package called MStweedie (standing for multisource Tweedie modeling), which is available at *https://github.com/fontaine618/ MStweedie*.

The article is organized as follows. In Section 2, we introduce the sparse Tweedie model for multisource claim data and derive a general objective function. Section 3 develops a unified algorithm to efficiently optimize that objective. Section 4 provides the details of implementation and tuning parameter selection for the proposed algorithm. In Section 5, we compare the performance of our proposal to other existing methods in a series of numerical experiments on both simulated and real data. Section 6 concludes the article. The technical proofs are relegated to the appendix (supplementary materials).

#### 2. Methodology

#### 2.1. Tweedie's Compound Poisson Model

The Tweedie model is closely related to the exponential dispersion models (EDM; Jørgensen 1987)

$$f_Y(y|\theta,\phi) = a(y,\phi) \exp\left\{\frac{y\theta - \kappa(\theta)}{\phi}\right\},$$

parameterized by the natural parameter  $\theta$  and dispersion parameter  $\phi$ , where  $\kappa(\cdot)$  is the cumulant function and  $a(\cdot)$  is the normalizing function. Both  $a(\cdot)$  and  $\kappa(\cdot)$  are known functions. It can be shown that *Y* has mean  $\mu \equiv E(Y) = \dot{\kappa}(\theta)$  and variance  $\operatorname{var}(Y) = \phi \ddot{\kappa}(\theta)$ , where  $\dot{\kappa}(\theta)$  and  $\ddot{\kappa}(\theta)$  denote the first and second derivatives of  $\kappa(\theta)$ , respectively. In this article, we are primarily interested in the Tweedie EDMs, a class of EDMs that have the mean-variance relationship  $\operatorname{var}(Y) = \phi \mu^{\rho}$ , where  $\rho$  is the power parameter. Such mean-variance relation gives

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

In particular, when  $\rho \in (1, 2)$ , the Tweedie EDMs correspond to a family of distributions called the compound Poisson distributions. In the sequel, we briefly discuss the compound Poisson distributions and their connection to the Tweedie EDMs. A compound Poisson random variable can be written as the sum of a (random) Poisson number of Gamma random variables. Specifically, let  $Z_1, Z_2, \ldots, Z_N$  be N iid random variables from Gamma( $\alpha, \gamma$ ), where N follows Poisson( $\lambda$ ). We assume that the  $Z_i$ 's are independent of N. Then the sum of the  $Z_i$ 's

$$Y = \begin{cases} 0 & \text{if } N = 0, \\ Z_1 + Z_2 + \dots + Z_N & \text{if } N = 1, 2, \dots \end{cases}$$
(2)

follows the compound Poisson distribution

$$f_Y(y|\lambda, \alpha, \gamma) = P(N=0)\delta_0(y) + \sum_{j=1}^{\infty} P(N=j)f_{Y|N=j}(y)$$
$$= e^{-\lambda}\delta_0(y) + \sum_{j=1}^{\infty} \frac{\lambda^j y^{j\alpha-1} e^{-\lambda-y/\gamma}}{j!\gamma^{j\alpha} \Gamma(j\alpha)},$$

where  $\delta_0$  is the Dirac delta mass at zero,  $f_{Y|N=j}(\cdot)$  is the conditional density of *Y* given N = j, and  $\Gamma(\cdot)$  is the gamma function. The compound Poisson distributions fit into a special class of Tweedie EDMs with  $\rho \in (1, 2)$ . To see this, we reparameterize  $(\lambda, \gamma, \alpha)$  by

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

The compound Poisson model will then have the form

$$\log f_Y(y|\mu,\phi,\rho) = \frac{1}{\phi} \left( y \frac{\mu^{1-\rho}}{1-\rho} - \frac{\mu^{2-\rho}}{2-\rho} \right) + \log a_\rho(y,\phi), \quad (3)$$

where

$$a_{\rho}(y,\phi) = \begin{cases} \frac{1}{y} \sum_{j=1}^{\infty} \frac{y^{j\alpha}}{j!(\rho-1)^{j\alpha} \phi^{j(\alpha+1)} \Gamma(j\alpha)} & \text{if } y > 0, \\ 1 & \text{if } y = 0. \end{cases}$$

It can be directly seen that (3) belongs to the Tweedie EDMs. As a result, for the rest of this article, we simply refer to (3) as Tweedie's compound Poisson model (or the Tweedie model), and denote it by  $\text{Tw}(\mu, \phi, \rho)$ , where  $1 < \rho < 2$ .

This equivalence provides a very intuitive justification for the use of the Tweedie distribution in modeling insurance claim data: the random variable *N* corresponds to the number of claims during the exposure period,  $Z_1, \ldots, Z_N$  correspond to the claim amounts, and  $Y = \sum_{j=1}^N Z_j$  then corresponds to the aggregate claim amount. The case Y = 0 represents the absence of claims during the exposure period and is a frequent situation for this type of data.

#### 2.2. A Sparse Tweedie Modeling Framework for Multisource Claim Data

Suppose the claim data consist of *K* data sources (possibly from different policy products), and each data source k ( $1 \le k \le K$ ) has  $n_k$  policies. Given any policy *i* in data source *k*, assume exposure is  $w_i^{(k)}$ . For data source *k*, denote by  $\tilde{Y}_i^{(k)} = \sum_{j=1}^{N_i^{(k)}} Z_{ij}^{(k)}$  the corresponding claim loss, where  $N_i^{(k)}$  is the claim frequency and the  $Z_{ij}^{(k)}$ 's are the claim severity. The goal is to model the pure premium  $Y_i^{(k)} = \tilde{Y}_i^{(k)}/w_i^{(k)}$ . Here, the exposure is a known measure of certain risk in force (e.g., the exposure of a personal auto insurance can be the policy duration) so that in the Tweedie model, we assume  $N_i^{(k)} \sim \text{Poisson}(\lambda_i^{(k)}w_i^{(k)})$  and  $\tilde{Y}_{ij}^{(k)} | N_i^{(k)} \sim \text{Gamma}(\alpha, \gamma_i^{(k)})$ , where  $\lambda_i^{(k)}$  represents a policy-specific parameter for the expected claim frequency under unit exposure,  $\gamma_i^{(k)}$  is a policy-specific parameter for claim severity, and  $\alpha$  is a known scalar (Dunn and Smyth 2005). Further assume a mean-variance relation  $\operatorname{var}(Y_i^{*(k)}) = \phi^{(k)} \{E(Y_i^{*(k)})\}^{\rho}$ , where  $Y_i^{*(k)} = 1$ ) and  $\phi^{(k)}$  is the source-specific dispersion. Then we have  $Y_i^{(k)} \sim \operatorname{Tw}(\mu_i^{(k)}, \phi^{(k)}/w_i^{(k)}, \rho)$  with  $\mu_i^{(k)} = E(Y_i^{(k)})$  (Smyth and Jørgensen 2002; Yang, Qian, and Zou 2018).

Suppose that each policy *i* in data source *k* has *p* covariates  $\mathbf{x}_{i}^{(k)} = (x_{i1}^{(k)}, \ldots, x_{ip}^{(k)})^{\top}$ . For brevity, we assume these covariates are of the same type with equal dimension across different data

sources, but as will be discussed in our numerical studies, we can generalize this setting to handle possibly unequal dimension scenarios. We adopt the commonly used multiplicative logarithmic link

$$\log \mu_i^{(k)} = \eta_i^{(k)} = \beta_0^{(k)} + \mathbf{x}_i^{(k)\top} \boldsymbol{\beta}^{(k)}$$

where  $\boldsymbol{\beta}^{(k)} = (\beta_1^{(k)}, \dots, \beta_p^{(k)})^{\top}$  with  $\beta_j^{(k)}$  being the *j*th element of  $\boldsymbol{\beta}^{(k)}$ ,  $j = 1, \dots, p$ . Let  $\boldsymbol{\beta}_0 = (\beta_0^{(1)}, \dots, \beta_0^{(K)})^{\top}$ ,  $\boldsymbol{\beta}_j = (\beta_j^{(1)}, \dots, \beta_j^{(K)})^{\top}$ , and  $\boldsymbol{\beta} = (\boldsymbol{\beta}_1^{\top}, \dots, \boldsymbol{\beta}_p^{\top})^{\top} \in \mathbb{R}^{pK}$  be the target coefficient parameters. Assume that only a small fraction of the covariates in  $\mathbf{x}_i^{(k)}$  are relevant to  $y_i^{(k)}$  so that many elements in  $\boldsymbol{\beta}^{(k)}$  are zero. Given *K* sources of polices  $\mathbf{D}^{(k)} = \{(y_i^{(k)}, \mathbf{x}_i^{(k)}, w_i^{(k)})\}_{i=1}^{n_k}$  for  $k = 1, \dots, K$ , the multisource data setting naturally leads to a composite objective function

$$L(\boldsymbol{\beta}_{0}, \boldsymbol{\beta}) = \prod_{k=1}^{K} \prod_{i=1}^{n_{k}} f_{Y}(y_{i}^{(k)} \mid \mu_{i}^{(k)}, \phi^{(k)} / w_{i}^{(k)}, \rho), \qquad (4)$$

which, assuming independence across different data sources, becomes the likelihood function. When the independence assumption is violated, (4) can still be viewed as a composite marginal likelihood (Varin, Reid, and Firth 2011), the study of which plays an important role in allowing feasible estimation of marginal parameters (see, e.g., Chandler and Bate 2007; Shi 2016). Without loss of generality, we assume same dispersion  $\phi = \phi^{(1)} = \cdots = \phi^{(K)}$  across all data sources (otherwise, we can simply adjust  $w_i^{(K)}$ 's in (5) accordingly). Hence, we assume common dispersion parameters between subjects; it does not much affect the modeling of the mean since the two parameters,  $\mu$  and  $\phi$ , are orthogonal to each other from the compound Poisson-Gamma distribution (Shi 2016). Taking negative logarithm and omitting constant terms, we obtain the following objective function (up to a dispersion scalar)

$$\ell(\boldsymbol{\beta}_0, \boldsymbol{\beta}) = \sum_{k=1}^{K} \sum_{i=1}^{n_k} w_i^{(k)} \left\{ -\frac{y_i^{(k)} e^{(1-\rho)\eta_i^{(k)}}}{1-\rho} + \frac{e^{(2-\rho)\eta_i^{(k)}}}{2-\rho} \right\}, \quad (5)$$

which is the negative log-likelihood under the independence assumption and is a convex objective.

To take advantage of the commonly shared relevant covariates while recovering source-specific model structures, we consider the composite penalty (Zhao, Rocha, and Yu 2009)

$$P_{\alpha}(\boldsymbol{\beta}) = \sum_{j=1}^{p} v_{j} \left[ (1-\alpha) ||\boldsymbol{\beta}_{j}||_{q} + \alpha ||\boldsymbol{\beta}_{j}||_{1} \right]$$

for some  $0 \le \alpha \le 1$  and  $q \in \{2, \infty\}$ , where the  $v_j$ 's are the penalty weights. The first component in  $P_{\alpha}(\beta)$  is aimed to find common relevant covariates across data sources and the second component is intended to deal with potential betweensource differences in sparsity and to find source-specific relevant covariates. When  $\alpha = 0$ ,  $P_{\alpha}(\beta)$  simplifies to the group lasso if q = 2 (Yuan and Lin 2006), while it gives a different "group discount" if  $q = \infty$  as only the largest coefficient is penalized (Obozinski, Taskar, and Jordan 2006). In both cases, when the *j*th covariate is selected by the model, then it is selected for all sources, which means that the coefficient  $\beta_i^{(k)}$  will be nonzero for all sources k = 1, ..., K. When  $0 < \alpha < 1$  and q = 2,  $P_{\alpha}(\beta)$  becomes the sparse group lasso (Simon et al. 2013). In that case, when the *j*th feature is included in the model, the *q*-norm of  $\beta_j$  is still freed from zero, but each individual coefficient are individually penalized so that some components of  $\beta_j$  may still remain zero even though other components will be nonzero; selecting the *j*th feature only ensure that the coefficient for at least one (but not necessarily all) source is nonzero. The use of the penalty weights is motivated from the adaptive Lasso (Zou 2006) for improved variable selection performance. Our integral approach to sparse Tweedie modeling for multisource data aims to solve the regularized objective

$$f^* = \min_{\boldsymbol{\beta}_0, \boldsymbol{\beta}} f(\boldsymbol{\beta}_0, \boldsymbol{\beta}), \quad f(\boldsymbol{\beta}_0, \boldsymbol{\beta}) = \ell(\boldsymbol{\beta}_0, \boldsymbol{\beta}) + \lambda P_{\alpha}(\boldsymbol{\beta}), \quad (6)$$

where  $\lambda > 0$  is the tuning parameter. We call (6) the  $L_1/L_q(\alpha)$  regularization, and when  $\alpha = 0$ , we simply call it  $L_1/L_q$  regularization (q = 2 or  $\infty$ ).

As in most multitask learning algorithms, our proposed approach has a single-task reinterpretation. Indeed, similar to Turlach, Venables, and Wright (2005), we can model the response  $y_i^{(k)}$  through the equivalent log-link for the mean

$$\log \mu_i^{(k)} = \sum_{l=1}^K \mathbb{I}(l=k) \,\beta_0^{(k)} + \left[\mathbb{I}(l=k) \,\mathbf{x}_i^{(k)}\right]^\top \boldsymbol{\beta}^{(k)} = \tilde{\mathbf{x}}_i^{(k)\top} \boldsymbol{\beta},$$

where

$$\tilde{\mathbf{x}}_{i}^{(k)\top} = \left(0, \mathbf{0}_{K}^{\top}, \dots, 1, \underbrace{\mathbf{x}_{i}^{(k)\top}}_{\text{position } k+1}, \dots, 0, \mathbf{0}_{K}^{\top}\right),$$
$$\boldsymbol{\beta}^{\top} = \left(\boldsymbol{\beta}_{0}^{(1)}, \boldsymbol{\beta}^{(1)\top}, \dots, \boldsymbol{\beta}_{0}^{(k)}, \boldsymbol{\beta}^{(k)\top}, \dots, \boldsymbol{\beta}_{0}^{(K)}, \boldsymbol{\beta}^{(K)\top}\right).$$

Hence, our model can be seen has a single-source Tweedie model with  $K \cdot p$  features (plus the *K* intercepts) where a response  $y_i^{(k)}$  has nonzero features only for the group of feature related to the *k*th source. The penalty  $P_{\alpha}$  ( $\beta$ ) then correspond to a sparse group Lasso penalty where the groups are defined by grouping a feature across sources. The novelty of our algorithm is actually to profit from the specific structure of the groups and dummy variables to improve the efficiency of the algorithm. Fitting a single Tweedie with such a description of the groups would be inefficient computationally since a large proportion of the entries of the design matrix would artificially be zeroes.

#### 3. Algorithm

In this section, we propose an efficient algorithm to solve the penalized composite Tweedie model (6). We decompose the description of our algorithm into four parts: Section 3.1 gives a general idea on how to solve our optimization problem via the cyclic groupwise proximal gradient (GPG) descent; Section 3.2 discusses an acceleration scheme for the proposed algorithm; and Section 3.3 provides detailed solutions to the  $L_1/L_q$  regularization, which gives necessary information to introduce our complete algorithm to solve the more general  $L_1/L_q(\alpha)$  regularization in Section 3.4.

For data source k, denote the response vector by  $Y^{(k)} = (y_1^{(k)}, \ldots, y_{n_k}^{(k)})^\top$  and the  $n_k \times p$  design matrix by  $\mathbf{X}^{(k)} = (\mathbf{x}_1^{(k)}, \ldots, \mathbf{x}_{n_k}^{(k)})^\top = (X_1^{(k)}, \ldots, X_p^{(k)})$ . For consistency of notation, we also let  $X_0^{(k)} = \mathbf{1}_{n_k}$ .

#### 3.1. A Groupwise Proximal Gradient Algorithm for MStweedie

Note that the penalty term  $P_{\alpha}(\boldsymbol{\beta})$  in (6) is separable with respect to the indices of the feature sets  $j = 1, \ldots, p$ . We exploit this property and propose to iteratively update and cycle through the  $\boldsymbol{\beta}_j$ 's ( $j = 0, 1, \ldots, p$ ) via the proximal gradient (Beck and Teboulle 2009) scheme which gives rise to a cyclic GPG algorithm designed for MStweedie. Specifically, let  $\tilde{\mathbf{b}}$  be the current iterate

$$\tilde{\mathbf{b}} \equiv (\tilde{\boldsymbol{\beta}}_0, \dots, \tilde{\boldsymbol{\beta}}_{j-1}, \tilde{\boldsymbol{\beta}}_j, \tilde{\boldsymbol{\beta}}_{j+1}, \dots, \tilde{\boldsymbol{\beta}}_p)^\top,$$

and  $\tilde{\mathbf{b}}_{-i}$  be the current iterate with the *j*th group excluded

$$\tilde{\mathbf{b}}_{-j} \equiv (\tilde{\boldsymbol{\beta}}_0, \dots, \tilde{\boldsymbol{\beta}}_{j-1}, \tilde{\boldsymbol{\beta}}_{j+1}, \dots, \tilde{\boldsymbol{\beta}}_p)^\top, \quad j = 0, \dots, p.$$

Suppose we are about to update the group  $\boldsymbol{\beta}_j = (\boldsymbol{\beta}_j^{(1)}, \dots, \boldsymbol{\beta}_j^{(K)})^\top$  for some  $j \in \{0, 1, \dots, p\}$ . View the negative log-likelihood function  $\ell(\boldsymbol{\beta}_0, \boldsymbol{\beta})$  in (5) as a function of the *j*th group  $\boldsymbol{\beta}_j$ , while keeping all the other groups fixed at  $\tilde{\mathbf{b}}_{-j}$ , that is,  $\ell(\boldsymbol{\beta}_j; \tilde{\mathbf{b}}_{-j}) = \ell(\boldsymbol{\beta}_0, \boldsymbol{\beta})|_{\boldsymbol{\beta}_m = \tilde{\boldsymbol{\beta}}_m, 0 \le m \le p, m \ne j}$ . For group *j*, note that a quadratic approximation to  $\ell(\boldsymbol{\beta}_j; \tilde{\mathbf{b}}_{-j})$  around  $\tilde{\boldsymbol{\beta}}_j$  is given by

$$\ell(\boldsymbol{\beta}_{j}; \tilde{\mathbf{b}}_{-j}) \approx \ell_{Q_{j}}(\boldsymbol{\beta}_{j}; \tilde{\mathbf{b}}, t_{j})$$
(7)

$$\equiv \ell(\tilde{\mathbf{b}}) + \nabla_j \ell(\tilde{\boldsymbol{\beta}}_j; \tilde{\mathbf{b}}_{-j})^\top (\boldsymbol{\beta}_j - \tilde{\boldsymbol{\beta}}_j) + \frac{1}{2t_j} \|\boldsymbol{\beta}_j - \tilde{\boldsymbol{\beta}}_j\|_2^2, \ t_j > 0.$$

It can be seen that  $\ell_{Q_j}(\boldsymbol{\beta}_j; \tilde{\mathbf{b}}, t_j) = \ell(\boldsymbol{\beta}_j; \tilde{\mathbf{b}}_{-j})$  when  $\boldsymbol{\beta}_j = \tilde{\boldsymbol{\beta}}_j$  for any  $t_j > 0$ . To ensure the convergence of the algorithm, the value of  $t_j$  can be determined using the backtracking line search (details given later in this section). In (7), the gradient  $\nabla_j \ell(\tilde{\boldsymbol{\beta}}_j; \tilde{\mathbf{b}}_{-j})$  can be written explicitly as

$$\nabla_{j}\ell(\tilde{\boldsymbol{\beta}}_{j};\tilde{\mathbf{b}}_{-j}) \tag{8}$$
$$= \frac{\partial}{\partial\boldsymbol{\beta}_{j}}\ell(\boldsymbol{\beta}_{j};\tilde{\mathbf{b}}_{-j})\Big|_{\boldsymbol{\beta}_{j}=\tilde{\boldsymbol{\beta}}_{j}} = \left((\tilde{\boldsymbol{\eta}}^{(k)}-\tilde{\mathbf{z}}^{(k)})^{\top}\widetilde{\mathbf{W}}^{(k)}X_{j}^{(k)}\right)_{k=1}^{K},$$

where  $\tilde{\eta}^{(k)} = (\tilde{\eta}_1^{(k)}, \dots, \tilde{\eta}_{n_k}^{(k)})^\top$  with  $\tilde{\eta}_i^{(k)} = \sum_{j=0}^p x_{ij}^{(k)} \tilde{\beta}_j^{(k)}$ ,  $\tilde{\mathbf{z}}^{(k)} = (\tilde{z}_1^{(k)}, \dots, \tilde{z}_{n_k}^{(k)})^\top$  with

$$\tilde{z}_{i}^{(k)} = \tilde{\eta}_{i}^{(k)} + \frac{w_{i}^{(k)}}{\tilde{w}_{i}^{(k)}} (y_{i}^{(k)} e^{(1-\rho)\tilde{\eta}_{i}^{(k)}} - e^{(2-\rho)\tilde{\eta}_{i}^{(k)}}), \qquad (9)$$

and  $\widetilde{\mathbf{W}}^{(k)} = \operatorname{diag}(\widetilde{w}_1^{(k)}, \dots, \widetilde{w}_{n_k}^{(k)})$  with

$$\tilde{w}_i^{(k)} = w_i^{(k)} \big( (\rho - 1) y_i^{(k)} e^{(1-\rho)\tilde{\eta}_i^{(k)}} + (2-\rho) e^{(2-\rho)\tilde{\eta}_i^{(k)}} \big).$$
(10)

Now we apply the proximal gradient algorithm on  $\ell_{Q_j}(\boldsymbol{\beta}_j; \mathbf{\hat{b}}, t_j)$  to update  $\boldsymbol{\beta}_j$  as follows. For  $\tau > 0$ , define the proximal mapping

of  $h(\cdot) = (1-\alpha)||\cdot||_q + \alpha||\cdot||_1$  as the minimizer of the following problem

$$\operatorname{prox}_{\tau h}(\mathbf{u}) = \arg\min_{\mathbf{v}} \left( \tau h(\mathbf{v}) + \frac{1}{2} \|\mathbf{v} - \mathbf{u}\|_{2}^{2} \right).$$
(11)

For now, suppose that the solution to (11) is given (methods for computing the minimizer is deferred to Sections 3.3 and 3.4). We update  $\beta_i$  by minimizing the following penalized problem

$$\boldsymbol{\beta}_{j}^{+}(\tilde{\mathbf{b}}, t_{j}) = \underset{\boldsymbol{\beta}_{j}}{\operatorname{arg\,min}} \underset{\boldsymbol{\beta}_{j}}{\operatorname{arg\,min}} \ell_{Q_{j}}(\boldsymbol{\beta}_{j}; \tilde{\mathbf{b}}, t_{j}) + \lambda P_{\alpha, j}(\boldsymbol{\beta}_{j})$$

$$= \underset{\boldsymbol{\beta}_{j}}{\operatorname{arg\,min}} \frac{1}{2} \|\boldsymbol{\beta}_{j} - (\tilde{\boldsymbol{\beta}}_{j} - t_{j} \nabla_{j} \ell(\tilde{\boldsymbol{\beta}}_{j}; \tilde{\mathbf{b}}_{-j}))\|_{2}^{2} + \lambda v_{j} t_{j} h(\boldsymbol{\beta}_{j})$$

$$= \operatorname{prox}_{\lambda v_{j} t_{j} h}(\tilde{\boldsymbol{\beta}}_{j} - t_{j} \nabla_{j} \ell(\tilde{\boldsymbol{\beta}}_{j}; \tilde{\mathbf{b}}_{-j})). \qquad (12)$$

Note that in (12), when j = 0, we have  $\boldsymbol{\beta}_0^+(\tilde{\mathbf{b}}, t_0) = \tilde{\boldsymbol{\beta}}_0 - \boldsymbol{\beta}_0$  $t_0 \nabla_i \ell(\boldsymbol{\beta}_0; \mathbf{b}_{-0})$ , since the intercept term is not penalized, that is,  $P_{\alpha,0}(\boldsymbol{\beta}_0) \equiv 0$ .

To guarantee convergence, we determine the step size  $t_i$  in (12) using backtracking line search. Define

$$G_{t_j}(\tilde{\boldsymbol{\beta}}_j) = \frac{1}{t_j} \{ \tilde{\boldsymbol{\beta}}_j - \boldsymbol{\beta}_j^+(\tilde{\mathbf{b}}, t_j) \}$$
  
=  $\frac{1}{t_j} \{ \tilde{\boldsymbol{\beta}}_j - \operatorname{prox}_{\lambda \nu_j t_j h}(\tilde{\boldsymbol{\beta}}_j - t_j \nabla_j \ell(\tilde{\boldsymbol{\beta}}_j; \tilde{\mathbf{b}}_{-j})) \}.$ 

We initialize  $t_j$  with some  $t_{max} > 0$  and repeatedly shrink  $t_j$  with  $t_i \leftarrow \delta t_i$  for some prechosen  $0 < \delta < 1$  until

$$\ell(\boldsymbol{\beta}_{j}^{+}(\tilde{\mathbf{b}}, t_{j})) = \ell(\tilde{\boldsymbol{\beta}}_{j} - t_{j}G_{t_{j}}(\tilde{\boldsymbol{\beta}}_{j}))$$

$$\leq \ell(\tilde{\mathbf{b}}) - t_{j}\nabla_{j}\ell(\tilde{\boldsymbol{\beta}}_{j}; \tilde{\mathbf{b}}_{-j})^{\top}G_{t_{j}}(\tilde{\boldsymbol{\beta}}_{j}) + \frac{t_{j}}{2}\|G_{t_{j}}(\tilde{\boldsymbol{\beta}}_{j})\|_{2}^{2}.$$
(13)

Once (13) is satisfied by  $\boldsymbol{\beta}_{i}^{+}(\tilde{\mathbf{b}}, t_{j})$  for some  $t_{j}$ , we set  $\tilde{\boldsymbol{\beta}}_{j} \leftarrow$  $\boldsymbol{\beta}_{i}^{+}(\tilde{\mathbf{b}},t_{j})$  and move on to the next group j+1 and compute the update  $\boldsymbol{\beta}_{i+1}^+(\tilde{\mathbf{b}}, t_{j+1})$ . The algorithm cyclically updates groups  $j = 0, 1, \dots, p, 0, 1, \dots, p, \dots$  until convergence of  $(\boldsymbol{\beta}_0, \boldsymbol{\beta}).$ 

We summarize our proposal above with backtracking line search in Algorithm 1, and call it MStweedie-GPG for short. Moreover, we show that the proposed iterative approach is guaranteed to converge with at least linear rate in the following theorem, whose proof can be found in Appendix D (supplementary materials).

*Theorem 1.* In the MStweedie-GPG algorithm, let  $(\boldsymbol{\beta}_0^{(r)}, \boldsymbol{\beta}^{(r)})$ be the update of  $(\boldsymbol{\beta}_0, \boldsymbol{\beta})$  after the *r*th cycle,  $r \geq 0$ . The algorithm with backtracking line search converges to the global minimum  $f^*$  of (6) with at least a linear rate of convergence, that is,

$$f(\boldsymbol{\beta}_{0}^{(r+1)}, \boldsymbol{\beta}^{(r+1)}) - f^{*} \le c(f(\boldsymbol{\beta}_{0}^{(r)}, \boldsymbol{\beta}^{(r)}) - f^{*})$$

for large enough *r*, where  $c \in (0, 1)$  is a constant.

Algorithm 1: MStweedie-GPG with backtracking line search.

1 Initialize the coefficients with  $(\tilde{\boldsymbol{\beta}}_0, \tilde{\boldsymbol{\beta}})$  and choose some  $0 < \delta < 1;$ 2 repeat 3 for j = 0, 1, 2, ..., p do

Initialize  $t_i$  with  $t_{\text{max}} > 0$ ;

4 repeat 5 6 Compute  $\boldsymbol{\beta}_{i}^{+}(\tilde{\mathbf{b}}, t_{j}) = \operatorname{prox}_{\lambda v_{i} t_{i} h}(\tilde{\boldsymbol{\beta}}_{j} - t_{j} \nabla_{j} \ell(\tilde{\boldsymbol{\beta}}_{j}; \tilde{\mathbf{b}}_{-j}))$  using the proximal operator in (21), where  $\nabla_i \ell(\tilde{\boldsymbol{\beta}}_i; \tilde{\mathbf{b}}_{-i})$ is calculated from (8); Compute  $G_{t_j}(\tilde{\boldsymbol{\beta}}_j) = \frac{1}{t_i} \{ \tilde{\boldsymbol{\beta}}_j - \boldsymbol{\beta}_j^+(\tilde{\mathbf{b}}, t_j) \};$ 7 Set  $t_i \leftarrow \delta t_i$ ; 8 until  $\ell(\tilde{\boldsymbol{\beta}}_j - t_j G_{t_j}(\tilde{\boldsymbol{\beta}}_j)) <$ 9  $\ell(\tilde{\mathbf{b}}) - t_j \nabla_j \ell(\tilde{\boldsymbol{\beta}}_j; \tilde{\mathbf{b}}_{-j})^\top G_{t_j}(\tilde{\boldsymbol{\beta}}_j) + \frac{t_j}{2} \|G_{t_j}(\tilde{\boldsymbol{\beta}}_j)\|_2^2;$ Set  $\tilde{\boldsymbol{\beta}}_i \leftarrow \boldsymbol{\beta}_i^+(\tilde{\mathbf{b}}, t_i);$ 10 11 end 12 until Convergence of  $(\tilde{\boldsymbol{\beta}}_0, \tilde{\boldsymbol{\beta}})$ ; 13 Return( $\tilde{\boldsymbol{\beta}}_0, \tilde{\boldsymbol{\beta}}$ );

14 Note: when j = 0,  $\boldsymbol{\beta}_0^+(\tilde{\mathbf{b}}, t_0) = \tilde{\boldsymbol{\beta}}_0 - t_0 \nabla_0 \ell(\tilde{\boldsymbol{\beta}}_0; \tilde{\mathbf{b}}_{-0})$  and  $G_{t_0}(\tilde{\boldsymbol{\beta}}_0) = \nabla_0 \ell(\tilde{\boldsymbol{\beta}}_0; \tilde{\mathbf{b}}_{-0})$ 

#### 3.2. Accelerated MStweedie-GPG

In the vanilla MStweedie-GPG algorithm, operation (13) for backtracking is repeatedly evaluated during each groupwise update, and is thus computationally expensive. We can accelerate our algorithm by fixing the step sizes and only update them after  $(\boldsymbol{\beta}_0, \boldsymbol{\beta})$  converges in each loop. Specifically, instead of searching for a new step size to update  $\beta_i$  during each iteration within a loop, we use a fixed step size  $t_i^*$  as follows: given  $(\tilde{\beta}_0, \tilde{\beta})$ at the beginning of each loop, we set the step sizes to  $t_i^* = \sigma_i^{-1}$ for j = 0, 1, ..., p, where  $\sigma_j$  is the largest element of  $\nabla_i^2 \ell(\hat{\boldsymbol{\beta}}_j; \mathbf{b}_{-j})$ with

$$\nabla_{j}^{2}\ell(\tilde{\boldsymbol{\beta}}_{j};\tilde{\mathbf{b}}_{-j}) = \frac{\partial^{2}}{\partial\boldsymbol{\beta}_{j}\partial\boldsymbol{\beta}_{j}^{\top}}\ell(\tilde{\boldsymbol{\beta}}_{j};\tilde{\mathbf{b}}_{-j})$$
$$= \operatorname{diag}(X_{j}^{(k)\top}\widetilde{\mathbf{W}}^{(k)}X_{j}^{(k)}, k = 1,\ldots,K). \quad (14)$$

Next, we make the cyclic updates  $\tilde{\boldsymbol{\beta}}_{i} \leftarrow \boldsymbol{\beta}_{i}^{+}(\tilde{\mathbf{b}}, t_{i}^{*})$  with

$$\boldsymbol{\beta}_{j}^{+}(\tilde{\mathbf{b}}, t_{j}^{*}) = \operatorname{prox}_{\lambda \nu_{j} \sigma_{j}^{-1} h}(\tilde{\boldsymbol{\beta}}_{j} - \sigma_{j}^{-1} \nabla_{j} \ell(\tilde{\boldsymbol{\beta}}_{j}; \tilde{\mathbf{b}}_{-j})), \quad (15)$$

for j = 0, 1, ..., p, 0, 1, ..., p... until  $(\tilde{\beta}_0, \tilde{\beta})$  converges during this loop. Then we recompute step sizes  $t_i^*$  using (14) and repeat the above process. We refer to this scheme as the Accelerated MStweedie-GPG (or MStweedie-AGPG for short). We summarize this practically important acceleration strategy in Algorithm 2. It can be seen that the algorithm only updates step sizes after  $(\hat{\boldsymbol{\beta}}_0, \hat{\boldsymbol{\beta}})$  converges in the sub-iteration 2(b) of Algorithm 2. A similar technique for accelerating coordinate descent algorithms can be found in Friedman, Hastie, and Tibshirani (2010). Our empirical evidence shows that MStweedie-AGPG

344 👄 S. FONTAINE ET AL.

converges very fast and follows an overall descending trend; see Figure A1 in the appendix (supplementary materials) for an illustration. This is the algorithm we use for all our numerical studies.

Algorithm 2: Accelerated MStweedie-GPG.

1 Initialize the coefficients with  $(\tilde{\boldsymbol{\beta}}_0, \tilde{\boldsymbol{\beta}})$ ; 2 repeat Compute step sizes  $t_j^* = \sigma_j^{-1}$  for j = 0, 1, ..., p, where 3  $\sigma_i$  is defined in (14); 4 repeat for j = 0, 1, 2, ..., p do 5 Update  $\tilde{\boldsymbol{\beta}}_j \leftarrow \boldsymbol{\beta}_j^+(\tilde{\mathbf{b}}, t_j^*) =$ 6  $\operatorname{prox}_{\lambda\nu_i\sigma_i^{-1}h}(\tilde{\boldsymbol{\beta}}_j - \sigma_j^{-1}\nabla_j \ell(\tilde{\boldsymbol{\beta}}_j; \tilde{\mathbf{b}}_{-j})) \text{ with the fixed}$ step sizes  $t_i^* = \sigma_i^{-1}$ , where the proximal operator is given in (21) and  $\nabla_{j}\ell(\tilde{\boldsymbol{\beta}}_{j};\tilde{\mathbf{b}}_{-j})$  is calculated from (8); 7 end **until** Convergence of  $(\tilde{\boldsymbol{\beta}}_0, \tilde{\boldsymbol{\beta}})$ ; 8 9 until Convergence of  $(\hat{\boldsymbol{\beta}}_0, \hat{\boldsymbol{\beta}})$ ; 10 Return( $\tilde{\boldsymbol{\beta}}_0, \tilde{\boldsymbol{\beta}}$ );

#### 3.3. $L_1/L_q$ Regularization

In the unified algorithm of Section 3.1, it remains to show how to solve (12). We first discuss the  $L_1/L_q$  regularization case ( $\alpha = 0$ ), which will be used in the next subsection to derive solutions to the more general  $L_1/L_q(\alpha)$  regularization with  $\alpha \in [0, 1]$ .

The following lemma translates the proximal operator of the  $L_{\infty}$  regularization ( $q = \infty$ ) into a projection. Its proof is given in Appendix A (supplementary materials).

Lemma 1. The minimization problem

$$\boldsymbol{\beta}_{j}^{+}(\tilde{\mathbf{b}}, t_{j}) = \operatorname*{arg\,min}_{\boldsymbol{\beta}_{j}} \frac{1}{2} \|\boldsymbol{\beta}_{j} - (\tilde{\boldsymbol{\beta}}_{j} - t_{j} \nabla_{j} \ell(\tilde{\boldsymbol{\beta}}_{j}; \tilde{\mathbf{b}}_{-j}))\|_{2}^{2} + \lambda v_{j} t_{j} \|\boldsymbol{\beta}_{j}\|_{\infty}$$
(16)

is equivalent to

$$\boldsymbol{\beta}_{j}^{+}(\tilde{\mathbf{b}}, t_{j}) = \tilde{\boldsymbol{\beta}}_{j} - t_{j} \nabla_{j} \ell(\tilde{\boldsymbol{\beta}}_{j}; \tilde{\mathbf{b}}_{-j})$$

$$- \operatorname{Proj}_{B_{1}(\lambda v_{i}t_{j})}(\tilde{\boldsymbol{\beta}}_{j} - t_{j} \nabla_{j} \ell(\tilde{\boldsymbol{\beta}}_{j}; \tilde{\mathbf{b}}_{-j})),$$
(17)

where  $\operatorname{Proj}_{B_1(\tau)}(\cdot)$  is the  $L_2$ -projection onto  $B_1(\tau) = \{\mathbf{v} \mid ||\mathbf{v}||_1 \le \tau\}$ , the  $L_1$ -ball with radius  $\tau$ .

We use an extension of the algorithm suggested by Duchi et al. (2008) to perform fast projections onto the  $L_1$ -ball (see Appendix A for details in the supplementary materials). The KKT conditions of (16) can be shown (see Appendix B for details in the supplementary materials) as follows

$$\begin{cases} \|\boldsymbol{\beta}_{j} - t_{j} \nabla_{j} \ell(\boldsymbol{\beta}_{j}; \mathbf{b}_{-j})\|_{1} \leq \lambda v_{j} t_{j}, & \boldsymbol{\beta}_{j} = \mathbf{0}, \\ \|\boldsymbol{\tilde{\beta}}_{j} - t_{j} \nabla_{j} \ell(\boldsymbol{\tilde{\beta}}_{j}; \mathbf{\tilde{b}}_{-j}) - \boldsymbol{\beta}_{j}\|_{1} = \lambda v_{j} t_{j}, & \boldsymbol{\beta}_{j} \neq \mathbf{0}, \\ \boldsymbol{\tilde{\beta}}_{j}^{(k)} - t_{j} \nabla_{j} \ell(\boldsymbol{\tilde{\beta}}_{j}; \mathbf{\tilde{b}}_{-j})^{(k)} - \boldsymbol{\beta}_{j}^{(k)} = \mathbf{0}, & \boldsymbol{\beta}_{j} \neq \mathbf{0}, k \notin M(\boldsymbol{\beta}_{j}), \end{cases}$$

$$(18)$$

where  $M(\boldsymbol{\beta}_j) = \{k \in \{1, \dots, K\} : \|\boldsymbol{\beta}_j\|_{\infty} = |\boldsymbol{\beta}_j^{(k)}|\}$  is the maximizing index set.

Next, we still assume  $\alpha = 0$  and briefly discuss the  $L_2$  regularization case (q = 2) in (12). We will omit most of the details and focus only on its differences from the  $L_1/L_{\infty}$  case. The minimizer of the penalized objective

$$\boldsymbol{\beta}_{j}^{+}(\tilde{\mathbf{b}},t_{j}) = \arg\min_{\boldsymbol{\beta}_{j}} \frac{1}{2} \|\boldsymbol{\beta}_{j} - (\tilde{\boldsymbol{\beta}}_{j} - t_{j} \nabla_{j} \ell(\tilde{\boldsymbol{\beta}}_{j};\tilde{\mathbf{b}}_{-j}))\|_{2}^{2} + \lambda \nu_{j} t_{j} \|\boldsymbol{\beta}_{j}\|_{2}$$

has closed form

$$\boldsymbol{\beta}_{j}^{+}(\tilde{\mathbf{b}},t_{j}) = (\|\tilde{\boldsymbol{\beta}}_{j} - t_{j}\nabla_{j}\ell(\tilde{\boldsymbol{\beta}}_{j};\tilde{\mathbf{b}}_{-j})\|_{2} - \lambda v_{j}t_{j})_{+} \qquad (19)$$
$$\frac{\tilde{\boldsymbol{\beta}}_{j} - t_{j}\nabla_{j}\ell(\tilde{\boldsymbol{\beta}}_{j};\tilde{\mathbf{b}}_{-j})}{\|\tilde{\boldsymbol{\beta}}_{j} - t_{j}\nabla_{j}\ell(\tilde{\boldsymbol{\beta}}_{j};\tilde{\mathbf{b}}_{-j})\|_{2}},$$

and the corresponding KKT conditions are

$$\begin{cases} \|\tilde{\boldsymbol{\beta}}_{j} - t_{j} \nabla_{j} \ell(\tilde{\boldsymbol{\beta}}_{j}; \tilde{\mathbf{b}}_{-j})\|_{2} \leq \lambda v_{j} t_{j}, & \boldsymbol{\beta}_{j} = \mathbf{0}, \\ \lambda v_{j} t_{j} \frac{\boldsymbol{\beta}_{j}}{\|\boldsymbol{\beta}_{j}\|_{2}} + t_{j} \nabla_{j} \ell(\tilde{\boldsymbol{\beta}}_{j}; \tilde{\mathbf{b}}_{-j}) + (\boldsymbol{\beta}_{j} - \tilde{\boldsymbol{\beta}}_{j}) = \mathbf{0}_{K}, & \boldsymbol{\beta}_{j} \neq \mathbf{0}. \end{cases}$$

$$(20)$$

#### **3.4.** $L_1/L_q(\alpha)$ Regularization

With the  $L_1/L_q$  regularization discussed in the previous subsection to take advantage of possibly common covariates across data sources, we are now ready to discuss the more general  $L_1/L_q(\alpha)$  regularization ( $0 \le \alpha \le 1$ ) to achieve the goal of uncovering relevant covariates unique to some data source.

It could seem complicated to derive a closed form expression of the above proximal operator (the Fenchel conjugate of fcannot be derived explicitly), but it is possible to solve it with a proximal technique originally developed for the hierarchical group lasso (Jenatton et al. 2010). Specifically, we rewrite our composite penalty as a sum of  $L_q$ -norms (q = 2 or  $\infty$ ) on a set of groups  $\mathcal{G}$  that is tree-structured by noting that  $\|\boldsymbol{\beta}_j\|_1$  is separable across  $k = 1, \ldots, K$ 

$$(1 - \alpha) \|\boldsymbol{\beta}_{j}\|_{q} + \alpha \|\boldsymbol{\beta}_{j}\|_{1} = (1 - \alpha) \|\boldsymbol{\beta}_{j}\|_{q} + \alpha \sum_{k=1}^{K} |\beta_{j}^{(k)}|$$
$$= (1 - \alpha) \|\boldsymbol{\beta}_{j}\|_{q} + \sum_{k=1}^{K} \alpha \|\beta_{j}^{(k)}\|_{q},$$

where we can identify  $\mathcal{G} = \{\{1\}, \ldots, \{K\}, \{1, \ldots, K\}\}$ , which is tree-structured. Consequently, we only require the proximal operator of each norm and compose them according to the tree ordering. Let  $\mathbf{u} = \tilde{\boldsymbol{\beta}}_j - t_j \nabla_j \ell(\tilde{\boldsymbol{\beta}}_j; \tilde{\mathbf{b}}_{-j})$  and  $\tau = \lambda v_j t_j$ . It is known from Section 3.3 that the proximal operator of  $(1 - \alpha)\tau \|.\|_q$  is

$$\operatorname{prox}_{(1-\alpha)\tau\|\cdot\|_{q}}(\mathbf{u}) = \begin{cases} \mathbf{u} - \operatorname{Proj}_{B_{1}((1-\alpha)\tau)}(\mathbf{u}), & q = \infty, \\ (\|\mathbf{u}\|_{2} - (1-\alpha)\tau)_{+} \frac{\mathbf{u}}{\|\mathbf{u}\|_{2}}, & q = 2, \end{cases}$$

and the proximal operator of  $\alpha \tau | \cdot |$  is given by the soft-thresholding operator

$$\operatorname{prox}_{\alpha\tau|\cdot|}(u_k) = \operatorname{sgn}(u_k) \left(|u_k| - \alpha\tau\right)_+ =: S(u_k, \alpha\tau).$$

Defining  $S(\mathbf{u}, \alpha \tau)$  as the component-wise soft-thresholding operator, that is,  $[S(\mathbf{u}, \alpha \tau)]_k = S(u_k, \alpha \tau)$ , we get

$$\operatorname{prox}_{\tau h}(\mathbf{u}) = \operatorname{prox}_{(1-\alpha)\tau \|\cdot\|_{q}} (S(\mathbf{u}, \alpha \tau))$$
$$= \begin{cases} S(\mathbf{u}, \alpha \tau) - \operatorname{Proj}_{B_{1}((1-\alpha)\tau)}(S(\mathbf{u}, \alpha \tau)), & q = \infty; \\ (\|S(\mathbf{u}, \alpha \tau)\|_{2} - (1-\alpha)\tau)_{+} \frac{S(\mathbf{u}, \alpha \tau)}{\|S(\mathbf{u}, \alpha \tau)\|_{2}}, & q = 2, \end{cases}$$
(21)

the computation of which has been already studied in Section 3.3.

*Remark.* Although we could wish for a general algorithm for all  $q \ge 1$ , our construction is only valid for  $q \in \{2, \infty\}$ . As shown in Jenatton et al. (2010), the property used to derive the proximal operator of the composite penalty is only true when  $q \in \{2, \infty\}$ . Note also that the case q = 1 is simply the Lasso.

#### 3.5. Missing Features Properties

One of the assumptions behind our algorithm is that all sources share exactly the same set of features. In practice, distinct sets of features may be encountered from different sources. For example, if a dataset contains policies from different years where some additional information is available in the later years, we may split the data into two sources where the first source contains fewer predictors than the second one. Another example is the case where data come from—literally—different sources that do not keep track of exactly the same information on the policy.

Suppose that the *j*th feature is missing from the *k*th source. We can set  $X_j^{(k)} = \mathbf{0}$  for the corresponding *j* and *k*. It can be shown that this treatment, together with the initialization  $\beta_j^{(k)} = 0$ , keeps  $\beta_j^{(k)}$  at zero throughout the entire algorithm for all choices of  $q \in \{2, \infty\}$  and  $0 \le \alpha \le 1$ . This way, predictor *j* of source *k* is systematically excluded from the model.

Indeed, at any point of the algorithm, we have

$$\nabla_{j}\ell(\widetilde{\boldsymbol{\beta}}_{j};\widetilde{\mathbf{b}}_{-j})^{(k)} = (\widetilde{\boldsymbol{\eta}}^{(k)} - \widetilde{\mathbf{z}}^{(k)})^{\top}\widetilde{\mathbf{W}}^{(k)}X_{j}^{(k)} = 0.$$

Hence, in the proximal operator, we have  $u_k = \tilde{\beta}_j^{(k)} - t_j \nabla_j \ell(\tilde{\beta}_j; \tilde{\mathbf{b}}_{-j})^{(k)} = 0 - 0 = 0$ . Then, the soft-thresholding operator produces

$$S(u_k, \alpha \tau) = \operatorname{sgn}(u_k)(|u_k| - \alpha \tau)_+ = 0$$

for any  $0 \le \alpha \le 1$ . Thus, for q = 2, we get

$$[\operatorname{prox}_{\tau h}(\mathbf{u})]_k = (||S(\mathbf{u},\alpha\tau)||_2 - (1-\alpha)\tau)_+ \frac{S(u_k,\alpha\tau)}{||S(\mathbf{u},\alpha\tau)||_2} = 0$$

and, for  $q = \infty$ ,

$$[\operatorname{prox}_{\tau h}(\mathbf{u})]_k = -\left[\operatorname{Proj}_{B_1((1-\alpha)\tau)}(S(\mathbf{u},\alpha\tau))\right]_k$$
$$= -\operatorname{sgn}(u_k)(|u_k| - \xi)_+ = 0.$$

In any case, we obtain  $\beta_j^{(k)+} = [\operatorname{prox}_{\tau h}(\mathbf{u})]_k = 0$ . It should be pointed out that this property does not prevent the same feature from being included in the model for other sources, though.

#### 4. Implementation

#### 4.1. Regularization Path

To select the tuning parameter, we apply the MStweedie-GPG algorithm on a decreasing sequence  $(\lambda_l)_{l=1}^L$ . The sequence of the corresponding solutions produces the solution path when a fine grid of  $\lambda$  is used. We present the solution path algorithm for solving MStweedie in Algorithm 3, where we wrap the MStweedie-GPG algorithm in an outer loop over the  $\lambda$  sequence. The sequence starts at  $\lambda_1 = \lambda_{\text{max}}$ , chosen so that all coefficients except the intercepts are shrunken to zero, and iterates successively to smaller values of  $\lambda$  until the last value,  $\lambda_L$ , is reached.

The full sequence of  $\lambda$  is chosen as follows. We first compute  $\lambda_{\max}$  via the KKT conditions (see below for details) and set  $\lambda_{\min} = \varepsilon \lambda_{\max}$  for some small  $\varepsilon$  (e.g.,  $\varepsilon = 10^{-3}$ ). Then, we construct a logarithmically decreasing sequence from  $\lambda_{\max}$  to  $\lambda_{\min}$ , that is,  $\lambda_l = \lambda_{\max} (\lambda_{\min}/\lambda_{\max})^{\frac{l-1}{L-1}}$ , where  $l = 1, \ldots, L$ . Note that we want  $\tilde{\boldsymbol{\beta}}_j = \boldsymbol{0}$  for all  $j \neq 0$  when  $\lambda = \lambda_{\max}$ . From the KKT conditions, that requires  $\lambda \geq v_j^{-1} \|\nabla_j \ell(\tilde{\boldsymbol{\beta}}_j; \tilde{\mathbf{b}}_{-j})\|_1$  for all  $j \neq 0$ . Therefore, we can choose  $\lambda_{\max} = \max_{1 \leq j \leq p} v_j^{-1} \|\nabla_j \ell(\tilde{\boldsymbol{\beta}}_j; \tilde{\mathbf{b}}_{-j})\|_1$ . Now at  $\lambda = \lambda_{\max}$ , we have  $\tilde{\boldsymbol{\beta}}(\operatorname{init}) = \boldsymbol{0}$  and

$$\tilde{\beta}_{0}^{(k)}(\text{init}) = \arg\min_{\beta_{0}^{(k)}} \sum_{i=1}^{n_{k}} w_{i}^{(k)} \left\{ -y_{i}^{(k)} \frac{e^{(1-\rho)\beta_{0}^{(k)}}}{1-\rho} + \frac{e^{(2-\rho)\beta_{0}^{(k)}}}{2-\rho} \right\},$$
$$= \log \frac{\sum_{i=1}^{n_{k}} w_{i}^{(k)} y_{i}^{(k)}}{\sum_{i=1}^{n_{k}} w_{i}^{(k)}}, \qquad k = 1, \dots, K.$$
(22)

Consequently, we obtain  $\tilde{\eta}_i^{(k)} = \tilde{\beta}_0^{(k)}$  (init) and

$$\nabla_{j}\ell(\tilde{\boldsymbol{\beta}}_{j};\tilde{\mathbf{b}}_{-j})^{(k)} = \sum_{i=1}^{n_{k}} \tilde{w}_{i}^{(k)} (\tilde{\eta}_{i}^{(k)} - \tilde{z}_{i}^{(k)}) x_{ij}^{(k)}, \qquad (23)$$

which now can be used to determine  $\lambda_{max}$ .

Algorithm 3: Solution path algorithm for solving MStweedie

1. Initialize  $\tilde{\boldsymbol{\beta}}_i = \boldsymbol{0}$  and  $\tilde{\boldsymbol{\beta}}_0 = \tilde{\boldsymbol{\beta}}_0$  (init) according to (22).

2. Compute  $\nabla_j \ell(\tilde{\boldsymbol{\beta}}_j; \tilde{\mathbf{b}}_{-j})$  using (23) and set  $\lambda_{\max} = \max_{1 \le j \le p} v_j^{-1} \|\nabla_j \ell(\tilde{\boldsymbol{\beta}}_j; \tilde{\mathbf{b}}_{-j})\|_1$  and  $\lambda = \lambda_{\max}$ . 3. For l = 2, ..., L, do

(a)Increment 
$$\lambda \leftarrow \lambda \left(\frac{\lambda_{\min}}{\lambda_{\max}}\right)^{\frac{1}{L-1}}$$
,  
(b)Update  $\tilde{\boldsymbol{\beta}}$  using Algorithm 1.

Once a solution path  $\{(\tilde{\boldsymbol{\beta}}_{0}^{[l]}, \tilde{\boldsymbol{\beta}}^{[l]})\}_{l=1}^{L}$  is obtained, we could use cross-validation (CV) to perform the model selection, where the out-of-sample prediction deviance may be used as the guided criterion. The scaled deviance from a single observation is

$$\begin{split} d_i^{(k)} &= -2\phi \big\{ \log f_Y(y_i^{(k)} | \mu_i^{(k)}, \phi, \rho) - \log f_Y(y_i^{(k)} | y_i^{(k)}, \phi, \rho) \big\} \\ &= 2 \bigg\{ \frac{y_i^{(k)(2-\rho)} - y_i^{(k)} \mu_i^{(k)(1-\rho)}}{1-\rho} - \frac{y_i^{(k)(2-\rho)} - \mu_i^{(k)(2-\rho)}}{2-\rho} \bigg\}, \end{split}$$

where  $\mu_i^{(k)} = \exp(\tilde{\beta}_0^{(k)} + \mathbf{x}_i^{(k)\top} \tilde{\boldsymbol{\beta}}^{(k)})$ , and the full deviance is then the weighted sum across all observations from all sources. Often, we choose the optimal  $\lambda$  as the one that minimizes the CV deviance (call it  $\lambda_m$ ). If model simplicity and interpretability are more of a concern, one may prefer the one-standarderror rule (Hastie, Tibshirani, and Friedman 2009), that is, choose optimal  $\lambda$  as the largest  $\lambda_l$  within one standard error of  $\lambda_m$ .

#### 4.2. Further Acceleration and Stabilization Strategies

Two tricks suggested by Friedman, Hastie, and Tibshirani (2010) are added to our algorithm. First, the solution path is computed using warm starts at each iteration to increase the stability of the algorithm. This means that the initialization at  $\lambda = \lambda_l$  is chosen to be the solution  $\tilde{\mathbf{b}}^{[l-1]} = (\tilde{\boldsymbol{\beta}}_0^{[l-1]}, \tilde{\boldsymbol{\beta}}^{[l-1]})$  from previously  $\lambda = \lambda_{l-1}$ . Second, the MStweedie-GPG algorithm is augmented with the active set updates: we first run a full cycle of the updates and identify the set of active predictors  $A = \{j \in \{1, \dots, p\} | \tilde{\boldsymbol{\beta}}_j \neq 0\}$ , and then repeat the cycles only over  $j \in A$  until convergence.

Another method to speed up the calculations, similar to the active set updates, is the sequential strong rule (Tibshirani et al. 2012). Specifically, it is designed to identify an active set on which to perform the full MStweedie-GPG algorithm at each  $\lambda$ . Before entering the algorithm at  $\lambda_l$ , we check the following conditions for each j = 1, ..., p

$$\left\|\nabla_{j}\ell(\tilde{\boldsymbol{\beta}}_{j}^{[l-1]};\tilde{\mathbf{b}}_{-j}^{[l-1]})\right\|_{1} < \nu_{j}(2\lambda_{l}-\lambda_{l-1}).$$

We exclude every predictor with index *j* that meets the above condition and run the MStweedie-GPG algorithm on the remaining predictors. Once the algorithm reaches convergence with these remaining variables, we perform a final check to verify that we do not accidentally exclude a predictor that should have been included. The check is based on the KKT conditions: for each predictor *j* initially excluded, we verify the KKT condition with  $\boldsymbol{\beta}_j = 0$ , which requires  $\|\nabla_j \ell(\tilde{\boldsymbol{\beta}}_j; \tilde{\mathbf{b}}_{-j})\|_{q^*} \leq \lambda v_j$ , where  $q^* = 1$  if  $q = \infty$  and  $q^* = 2$  if q = 2. If at least one condition is violated, then the corresponding predictor is added back to the active set. This process is repeated until the KKT condition is satisfied for all excluded predictors.

The algorithm with the sequential strong rule is presented in Algorithm 4.

Algorithm 4: MStweedie sequential strong rule.

1. Do while  $V \neq \emptyset$ :

- (a) Identify  $S = \{j : \|\nabla_j \ell(\tilde{\boldsymbol{\beta}}_j; \tilde{\mathbf{b}}_{-j})\|_{q^*} \ge \nu_j (2\lambda_l \lambda_{l-1})\}$ and  $S^c = \{1, \dots, p\} \setminus S$ .
- (b) Update β̃ as in Algorithm 1 while keeping β̃<sub>j</sub> = 0 for all j ∈ S<sup>C</sup>.
- (c) Identify the violations  $V = \{j : \|\nabla_j \ell(\tilde{\boldsymbol{\beta}}_j; \tilde{\mathbf{b}}_{-j})\|_{q^*} \le \lambda \nu_j, \ j \in S^c\}.$

Note:  $q^* = 1$  if  $q = \infty$  and  $q^* = 2$  if q = 2.

#### 4.3. Adaptive MStweedie

We also consider an adaptive version of MStweedie (a-MStweedie). The a-MStweedie is motivated from Zou (2006), where the adaptive lasso is used to improve model selection performance over the regular lasso. In a-MStweedie, we first obtain  $\hat{\beta}^*$ , the cross-validated parameter estimate under equal penalty factors (i.e.,  $v_j = 1$  for all  $1 \le j \le p$ ). Then, we update the penalty factors  $v_j = ||\hat{\beta}_j^*||_q^{-\varphi}$  for some  $\varphi > 0$  (default is  $\varphi = 1$ ) and refit the model with these new penalty factors. When the initial CV yields  $\hat{\beta}_j^* = 0$  for some j, we set  $v_j$  to a large machine number to ensure that this variable is not included in the adaptive modeling.

#### 5. Numerical Studies

#### 5.1. Performance Assessment

We use deviance as the criterion to access model fit. First of all, we split the data into two parts: a training set on which a model is fit to yield the coefficient estimates, and a testing set on which these estimates are used for prediction. The train and test deviances are then obtained, respectively, from these two sets.

Three measures are considered for assessing selection performance: the percentage of variables correctly identified (*accuracy*), the percentage of identified variables that are indeed true variables (*precision*), and the percentage of true variables identified (*recall*). These three measures describe different aspects of a variable selection result and are widely used in classification and pattern recognition (see, e.g., Fawcett 2006). In terms of overall performance, accuracy is perhaps a more interesting measure as our goal is not only to find the true predictors but also to exclude those spurious ones.

#### 5.2. Synthetic Data

We consider a variety of settings under which our algorithm is tested and compared to existing ones.

#### 5.2.1. Setting 1—Unequal Coefficients, $p < n_k$

This simulation setting is inspired by Gong, Ye, and Zhang (2012), in which we set the number of sources to K = 10, the number of observations to  $n_k = 400$ , k = 1, ..., 10, and the number of covariates to p = 100. The covariates are generated from independent normal distributions. Moreover, we set the coefficient matrix  $\boldsymbol{\beta}$  to zero everywhere except the last 10 columns, which are generated from independent normal distributions of mean 0 and variance  $4^2\sigma$  with  $\sigma = 0.1$ . Finally, we generate the responses  $y_i^{(k)}$  from Tw $(\mu_i^{(k)}, \phi, \rho)$  with  $\phi = 1$  and  $\rho = 1.5$ , where  $\mu_i^{(k)} = \exp(\mathbf{x}_i^{(k)\top}\boldsymbol{\beta}^{(k)})$  for all *i* and *k*.

We randomly split the above data into two equal parts ( $n_k = 200$  for each source): the first part is used to tune the model via 10-fold CV, while the second is used for testing the model. The results are averaged over 100 replications. The following models are compared: Full Lasso ( $L_1$ -regularized Tweedie model on the full dataset, using the HDtweedie package by Qian, Yang, and Zou (2016)), Individual Lasso (Individual  $L_1$ -regularized Tweedie model for each source, also using HDtweedie), and

Table 1. Results from Setting 1 with 100 replications.

| l: Mean (stan | dard error)   |   |   |  |  |   |   |   |   |   |   |
|---------------|---|---|---|--|--|---|---|---|---|---|---|
| Full          | Lasso   | Ind.  | Lasso   | L <sub>1</sub> /   | $L_{\infty}$   | a-L <sub>1</sub>  | $/L_{\infty}$   | $L_{1}$   | /L <sub>2</sub>   | a-L <sub>1</sub> ,  | /L <sub>2</sub>   |
| 2,64<br>(553  | 2,427<br>,415)  | 17,<br>(32  | 710<br>234)   | 17,<br>(20   | 330<br>)66)  | 59<br>(5)   | 963<br>23)  | 77<br>(7  | '63<br>17)  | 559<br>(72  | 90<br>3)  |
| 1.02          | (0.36)  | 89.37   | (0.43)  | 89.46  | (0.89)   | 10.00   | (0.00)  | 40.19   | (1.34)  | 10.00   | (0.00)  |
| 89.9<br>94.0  | (0.2)<br>(1.9)  | 20.6<br>11.2  | (0.4)<br>(0.1)  | 20.5<br>11.3   | (0.9)<br>(0.1)   | 100.0<br>100.0  | (0.0)<br>(0.0)  | 69.8<br>28.3  | (1.3)<br>(1.1)  | 100.0<br>100.0  | (0.0)<br>(0.0)  |
| 4.6<br>12.76  | (1.2)<br>(0.01)   | 100.0<br>1.56   | (0.0)<br>(0.04)   | 100.0<br>1.87  | (0.0)<br>(0.06)  | 100.0<br>1.20   | (0.0)<br>(0.06)   | 100.0<br>1.30   | (0.0)<br>(0.06)   | 100.0<br>1.03   | (0.0)<br>(0.07)   |
| 1: Mean rank  | (# of times be  | st)   |   |  |  |   |   |   |   |   |   |
| Full          | Lasso   | Ind.  | Lasso   | $L_1/L_\infty$   |  | a- $L_1/L_\infty$   |   | $L_{1}/L_{2}$   |   | a-L <sub>1</sub> /L <sub>2</sub>  |   |
| 6.00          | (0)   | 4.00  | (5)   | 4.74   | (0)  | 2.08  | (16)  | 2.89  | (2)   | 1.29  | (77)  |
| 1.06          | (97)  | 5.36  | (0)   | 5.58   | (0)  | 1.97  | (3)   | 4.00  | (0)   | 1.97  | (3)   |
| 3.04          | (0)<br>(80)   | 5.36  | (0)<br>(0)  | 5.58   | (0)<br>(0)   | 1.00  | (100)   | 3.93  | (0)<br>(0)  | 1.00  | (100)   |
| 6.00<br>6.00  | (0)<br>(0)  | 1.00<br>3.74  | (100)<br>(9)  | 1.00<br>4.86   | (100)<br>(0)   | 1.00<br>1.00<br>2.29  | (100)<br>(100)<br>(9)   | 1.00<br>2.77  | (1)<br>(1)  | 1.00<br>1.00<br>1.34  | (100)<br>(100)<br>(81)  |
|               | I: Mean (stan<br>Full<br>2,64<br>(553<br>1.02<br>89.9<br>94.0<br>4.6<br>12.76<br>1: Mean rank<br>Full<br>6.00<br>1.06<br>3.04<br>1.27<br>6.00<br>6.00 | I: Mean (standard error)<br>Full Lasso<br>2,642,427<br>(553,415)<br>1.02 (0.36)<br>89.9 (0.2)<br>94.0 (1.9)<br>4.6 (1.2)<br>12.76 (0.01)<br>1: Mean rank (# of times be<br>Full Lasso<br>6.00 (0)<br>1.06 (97)<br>3.04 (0)<br>1.27 (89)<br>6.00 (0)<br>6.00 (0) | Full Lasso         Ind.           2,642,427         17,           (553,415)         (32           1.02         (0.36)         89.37           89.9         (0.2)         20.6           94.0         (1.9)         11.2           4.6         (1.2)         100.0           12.76         (0.01)         1.56           1: Mean rank (# of times best)         Ind.           6.00         (0)         4.00           1.06         (97)         5.36           3.04         (0)         5.36           6.00         (0)         1.00           6.00         (0)         1.00           1.27         (89)         5.35           6.00         (0)         1.00           6.00         (0)         1.00           6.00         (0)         3.74 | Full Lasso       Ind. Lasso         2,642,427       17,710 $(553,415)$ $(3234)$ $1.02$ $(0.36)$ $89.37$ $(0.43)$ $89.9$ $(0.2)$ $20.6$ $(0.4)$ $94.0$ $(1.9)$ $11.2$ $(0.1)$ $4.6$ $(1.2)$ $100.0$ $(0.0)$ $12.76$ $(0.01)$ $1.56$ $(0.4)$ Ind. Lasso         Ind. Lasso         6.00 $(0)$ $4.00$ $(5)$ $1.06$ $(97)$ $5.36$ $(0)$ $3.04$ $(0)$ $5.35$ $(0)$ $1.27$ $(89)$ $5.35$ $(0)$ $6.00$ $(0)$ $1.00$ $(100)$ $6.00$ $(0)$ $3.74$ $(9)$ | Full Lasso         Ind. Lasso $L_{1/}$ 2,642,427         17,710         17,           (553,415)         (3234)         (20           1.02         (0.36)         89.37         (0.43)         89.46           89.9         (0.2)         20.6         (0.4)         20.5           94.0         (1.9)         11.2         (0.1)         11.3           4.6         (1.2)         100.0         (0.0)         100.0           12.76         (0.01)         1.56         (0.04)         1.87           1: Mean rank (# of times best)         Full Lasso         Ind. Lasso $L_{1/}$ 6.00         (0)         4.00         (5)         4.74           1.06         (97)         5.36         (0)         5.58           3.04         (0)         5.36         (0)         5.57           6.00         (0)         1.00         1.00         5.57           6.00         (0)         3.74         (9)         4.86 | Full Lasso $L_1/L_{\infty}$ Full Lasso       Ind. Lasso $L_1/L_{\infty}$ 2,642,427       17,710       17,330         (553,415)       (3234)       (2066)         1.02       (0.36)       89.37       (0.43)       89.46       (0.89)         89.9       (0.2)       20.6       (0.4)       20.5       (0.9)         94.0       (1.9)       11.2       (0.1)       11.3       (0.1)         4.6       (1.2)       100.0       (0.0)       100.0       (0.0)         12.76       (0.01)       1.56       (0.04)       1.87       (0.06)         1: Mean rank (# of times best)       Full Lasso       Ind. Lasso $L_1/L_{\infty}$ 6.00       (0)       4.00       (5)       4.74       (0)         1.06       (97)       5.36       (0)       5.58       (0)         3.04       (0)       5.36       (0)       5.57       (0)         6.00       (0)       1.00       (100)       1.00       (100)         6.00       (0)       3.74       (9)       4.86       (0) | Full Lasso $L_1/L_{\infty}$ a-L_1 $2,642,427$ 17,710       17,330       59         (553,415)       (3234)       (2066)       (5         1.02       (0.36)       89.37       (0.43)       89.46       (0.89)       10.00         89.9       (0.2)       20.6       (0.4)       20.5       (0.9)       100.0         94.0       (1.9)       11.2       (0.1)       11.3       (0.1)       100.0         4.6       (1.2)       100.0       (0.0)       100.0       (0.0)       100.0         12.76       (0.01)       1.56       (0.04)       1.87       (0.06)       1.20         1: Mean rank (# of times best)       Full Lasso       Ind. Lasso $L_1/L_{\infty}$ a-L_1         6.00       (0)       4.00       (5)       4.74       (0)       2.08         1.06       (97)       5.36       (0)       5.58       (0)       1.97         3.04       (0)       5.36       (0)       5.57       (0)       1.00         1.27       (89)       5.35       (0)       5.57       (0)       1.00         6.00       (0)       1.00       < | Full Lasso $L_1/L_{\infty}$ $a-L_1/L_{\infty}$ $2,642,427$ $17,710$ $17,330$ $5963$ $(553,415)$ $(3234)$ $(2066)$ $(523)$ $1.02$ $(0.36)$ $89.37$ $(0.43)$ $89.46$ $(0.89)$ $10.00$ $(0.00)$ $89.9$ $(0.2)$ $20.6$ $(0.4)$ $20.5$ $(0.9)$ $100.0$ $(0.0)$ $94.0$ $(1.9)$ $11.2$ $(0.1)$ $11.3$ $(0.1)$ $100.0$ $(0.0)$ $4.6$ $(1.2)$ $100.0$ $(0.0)$ $100.0$ $(0.0)$ $100.0$ $(0.0)$ $12.76$ $(0.01)$ $1.56$ $(0.04)$ $1.87$ $(0.06)$ $1.20$ $(0.06)$ 1: Mean rank (# of times best)       Ind. Lasso $L_1/L_{\infty}$ $a-L_1/L_{\infty}$ Full Lasso       Ind. Lasso $L_1/L_{\infty}$ $a-L_1/L_{\infty}$ $6.00$ $(0)$ $5.36$ $(0)$ $5.58$ $(0)$ $1.97$ $(3)$ $3.04$ $(0)$ $5.35$ $(0)$ $5.57$ $(0)$ $1.00$ $(100)$ | Full Lasso $L_1/L_{\infty}$ $a \cdot L_1/L_{\infty}$ $L_1$ 2,642,427       17,710       17,330       5963       77         (553,415)       (3234)       (2066)       (523)       (7         1.02       (0.36)       89.37       (0.43)       89.46       (0.89)       10.00       (0.00)       40.19         89.9       (0.2)       20.6       (0.4)       20.5       (0.9)       100.0       (0.0)       28.3         4.6       (1.2)       100.0       (0.0)       100.0       (0.0)       100.0       100.0       100.0       100.0       100.0       100.0       100.0       100.0       1.20       100.0       100.0       1.30         1.2.76       (0.01)       1.56       (0.04)       1.87       (0.06)       1.20       (0.06)       1.30         1: Mean rank (# of times best)       Ind. Lasso $L_1/L_{\infty}$ $a \cdot L_1/L_{\infty}$ $L_1$ 6.00       (0)       4.00       (5)       4.74       (0)       2.08       (16)       2.89         1.06       (97)       5.36       (0)       5.58       (0)       1.00       1000       3.93         1.27 | Full Lasso $L_1/L_{\infty}$ $a-L_1/L_{\infty}$ $L_1/L_2$ 2,642,427       17,710       17,330       5963       7763         (553,415)       (3234)       (2066)       (523)       (717)         1.02       (0.36)       89.37       (0.43)       89.46       (0.89)       10.00       (0.00)       40.19       (1.34)         89.9       (0.2)       20.6       (0.4)       20.5       (0.9)       100.0       (0.0)       28.3       (1.1)         4.6       (1.2)       100.0       (0.0)       100.0       (0.0)       100.0       (0.0)       100.0       (0.0)       100.0       (0.0)       100.0       (0.0)       100.0       (0.0)       100.0       (0.0)       100.0       (0.0)       100.0       (0.0)       100.0       (0.0)       100.0       (0.0)       12.6       (0.01)       1.56       (0.04)       1.87       (0.06)       1.20       (0.06)       1.30       (0.06)         1: Mean rank (# of times best)       Ind. Lasso $L_1/L_{\infty}$ $a-L_1/L_{\infty}$ $L_1/L_2$ 6.00       (0)       4.00       (0)       3.93       (0)       1.97       (3)       4.00       (0)       1.2 | I: Mean (standard error)         Full Lasso       Ind. Lasso $L_1/L_{\infty}$ $a^2L_1/L_{\infty}$ $L_1/L_2$ $a^2L_1$ 2,642,427       17,710       17,30       5963       7763       559         (553,415)       (3234)       (2066)       (523)       (717)       (72         1.02       (0.36)       89.37       (0.43)       89.46       (0.89)       10.00       (0.00)       40.19       (1.34)       10.00         89.9       (0.2)       20.6       (0.4)       20.5       (0.9)       100.0       (0.0)       28.3       (1.1)       100.0         94.0       (1.9)       11.2       (0.1)       11.3       (0.1)       100.0       (0.0)       28.3       (1.1)       100.0         4.6       (1.2)       100.0       (0.0)       100.0       (0.0)       100.0       (1.30)       (0.06)       1.03         1: Mean rank (# of times best)       Ind. Lasso $L_1/L_{\infty}$ $a^2L_1/L_{\infty}$ $L_1/L_2$ $a^2L_1$ 6.00       (0)       4.00       (5)       4.74       (0)       2.08       (16)       2.89       (2)       1.29         1.06       (97)       5.36       (0) </td |

NOTE: Part (a) shows the mean values of the statistics (with their standard errors listed in the parentheses). Part (b) shows the mean rank across the six models and, in parentheses, the number of times the model is best.

MStweedie with  $L_1/L_{\infty}$ ,  $L_1/L_2$ ,  $a-L_1/L_{\infty}$  (adaptive  $L_1/L_{\infty}$ ), and  $a-L_1/L_2$  (adaptive  $L_1/L_2$ ) regularizations.

Part (a) of Table 1 lists the averages and standard errors of different statistics. The test deviance, measuring the goodness of fit of the corresponding model, shows that MStweedie with the adaptive  $L_1/L_2$  regularization is the best while the Full Lasso performs very poorly on this matter. The poor performance of Full Lasso is due to the fact that it has identical estimates across sources, which is apparently not true according to our data-generating mechanism.

If we disregard the Full Lasso (which selected no features 81 out of 100 times), the two adaptive procedures performed the best in terms of variable selection performance, where each picks exactly 10 predictors in every replication. This selection matches exactly the true active variables so that both  $a-L_1/L_2$  and  $a-L_1/L_{\infty}$  achieve perfect accuracy, precision and recall. For the other models, the number of selected variables is much larger, yielding low precision and accuracy even with perfect recall. Finally,  $a-L_1/L_2$  produces estimates that are closest (in  $L_2$  norm) to the true coefficients. For both  $L_1/L_{\infty}$  and  $L_1/L_2$ , their adaptive versions greatly increase the selection accuracy and precision by picking much fewer variables while achieving lower deviance and  $L_2$ -loss. Overall,  $L_1/L_{\infty}$  exhibits similar performance to Individual Lasso, but its adaptive version increases the performance significantly.

The results about the ranking of these methods, reported in part (b) of Table 1, gives similar conclusions. However, we can see that, although being the best on average,  $a -L_1/L_2$  is occasionally outperformed by either Individual Lasso or  $a -L_1/L_\infty$  in terms of test deviance.

#### 5.2.2. Setting 2—Equal Coefficients, $p > n_k$

In this setting, we consider the high-dimensional scenario ( $p > n_k$ ) with local correlation structure. The data are generated similarly as in Gu and Zou (2016) with  $n_k = 300$ , p = 600, and K = 5. We generate the covariates  $\mathbf{x}_i^{(k)}$  from the multivariate normal distribution with mean **0** and covariance matrix  $\boldsymbol{\Sigma} =$ 

 $(0.5^{|i-j|})_{i,j=1}^p$  and set  $\boldsymbol{\beta}_j = 0$  for all *j* except  $j \in \{2, 4, 8, 16, 32\}$  where it is set to 2 in all sources. We then simulate the responses as in Setting 1 using  $\mu_i^{(k)} = \exp(\mathbf{x}_i^{(k)\top}\boldsymbol{\beta}^{(k)}), \phi = 1$ , and  $\rho = 1.5$ . Results below are summarized from 100 replications.

Part (a) of Table 2 contains the average values and standard errors of the different statistics. The lowest average test deviance is achieved by  $a-L_1/L_{\infty}$  followed closely by Full Lasso while Individual Lasso is significantly worse. The models selected by  $L_1/L_{\infty}$  are much more complex than any other method. As in Setting 1, the two adaptive methods performed perfectly in terms of accuracy, precision and recall, since they select the five true predictors exactly. Also,  $a-L_1/L_{\infty}$  produces the best estimates in term of  $L_2$ -loss. The study of the rankings, in part (b) of Table 2, leads to the same observations except that  $a-L_1/L_2$  and Full Lasso outperform  $a-L_1/L_{\infty}$  on some occasions in terms of test deviance or  $L_2$ -loss.

#### 5.2.3. Setting 3—Within-Feature Sparsity

In multisource insurance claim data, some predictors may not be relevant to all sources. For example, property age may only help predict the property claim amount. Some information of the same policyholders, such as credit history, however, may be relevant for both sources. The model thus exhibits both *within-feature* and *between-sources* sparsity. We consider a scenario designed to generate such a model to specifically test our  $L_1/L_q(\alpha)$  regularization.

The setting is similar to Setting 2, except that we voluntarily set the coefficients of some true generating variables to zero in certain sources

$$(\boldsymbol{\beta}_2, \boldsymbol{\beta}_4, \boldsymbol{\beta}_8, \boldsymbol{\beta}_{16}, \boldsymbol{\beta}_{32}) = \begin{bmatrix} 2 & 0 & 2 & 2 & 0 \\ 0 & 2 & 2 & 0 & 0 \\ 2 & 0 & 0 & 0 & 0 \\ 2 & 2 & 0 & 0 & 2 \\ 0 & 2 & 0 & 2 & 2 \end{bmatrix},$$
$$\boldsymbol{\beta}_j = \mathbf{0}, j \notin \{2, 4, 8, 16, 32\}.$$

#### Table 2. Results from Setting 2 with 100 replications.

| (a) Setting 2  | 2: Mean (stan                 | dard error)                       |                              |                                   |                              |                                   |                                 |                                   |                               |                                   |                                  |                                   |
|--|-------------------------------|-----------------------------------|------------------------------|-----------------------------------|------------------------------|-----------------------------------|---------------------------------|-----------------------------------|-------------------------------|-----------------------------------|----------------------------------|-----------------------------------|
|  | Full                          | Lasso                             | Ind.                         | Lasso                             | L <sub>1</sub> /             | $L_{\infty}$                      | a-L <sub>1</sub>                | $/L_{\infty}$                     | L <sub>1</sub> ,              | /L <sub>2</sub>                   | a-L <sub>1</sub>                 | /L <sub>2</sub>                   |
| Test dev.  | 14<br>(1                      | 130<br>42)                        | 1,13<br>(265                 | 6,266<br>,655)                    | 30<br>(6                     | )96<br>04)                        | 11<br>(1)                       | l61<br>29)                        | 69<br>(12                     | 968<br>277)                       | 257<br>(64                       | 72<br>2)                          |
| Size   | 22.07                         | (0.76)                            | 29.73                        | (1.44)                            | 71.10                        | (2.33)                            | 5.00                            | (0.00)                            | 37.86                         | (0.94)                            | 5.00                             | (0.00)                            |
| Accuracy<br>Precision<br>Recall<br>L <sub>2</sub> loss | 97.2<br>24.5<br>100.0<br>0.43 | (0.1)<br>(0.6)<br>(0.0)<br>(0.03) | 95.7<br>22.6<br>91.2<br>8.86 | (0.2)<br>(2.0)<br>(2.5)<br>(0.07) | 89.0<br>8.1<br>100.0<br>0.70 | (0.4)<br>(0.4)<br>(0.0)<br>(0.02) | 100.0<br>100.0<br>100.0<br>0.34 | (0.0)<br>(0.0)<br>(0.0)<br>(0.01) | 94.5<br>14.2<br>100.0<br>1.03 | (0.2)<br>(0.4)<br>(0.0)<br>(0.06) | 100.0<br>100.0<br>100.0<br>0.53  | (0.0)<br>(0.0)<br>(0.0)<br>(0.05) |
| (b) Setting  | 2: Mean rank                  | (# of times be                    | est)                         |                                   |                              |                                   |                                 |                                   |                               |                                   |                                  |                                   |
|  | Full                          | Lasso                             | Ind. Lasso                   |                                   | $L_1/L_\infty$               |                                   | a- $L_1/L_\infty$               |                                   | $L_{1}/L_{2}$                 |                                   | a-L <sub>1</sub> /L <sub>2</sub> |                                   |
| Test dev.  | 2.09                          | (30)                              | 6.00                         | (0)                               | 3.89                         | (0)                               | 1.57                            | (55)                              | 4.90                          | (0)                               | 2.55                             | (15)                              |
| Size   | 3.34                          | (0)                               | 3.87                         | (7)                               | 5.97                         | (0)                               | 1.06                            | (94)                              | 4.67                          | (0)                               | 1.06                             | (94)                              |
| Accuracy<br>Precision<br>Recall<br>L <sub>2</sub> loss | 3.34<br>3.32<br>1.00<br>2.16  | (0)<br>(0)<br>(100)<br>(25)       | 4.01<br>3.94<br>1.70<br>6.00 | (0)<br>(5)<br>(86)<br>(0)         | 5.97<br>5.96<br>1.00<br>3.84 | (0)<br>(0)<br>(100)<br>(0)        | 1.00<br>1.00<br>1.00<br>1.44    | (100)<br>(100)<br>(100)<br>(61)   | 4.67<br>4.66<br>1.00<br>4.97  | (0)<br>(0)<br>(100)<br>(0)        | 1.00<br>1.00<br>1.00<br>2.59     | (100)<br>(100)<br>(100)<br>(14)   |

NOTE: Part (a) shows the mean values of the statistics and their standard errors in parentheses. Part (b) shows the mean rank across the six models and, in parentheses, the number of times the model is best.

Thus, the true model is sparse in terms of features (only five generating variables), but it is also sparse within features since some of the true features do not generate the responses in certain sources.

Under Setting 2, we see that  $a-L_1/L_{\infty}$  produces the best fit. In this setting, we compare Full Lasso (which should not perform well) and Individual Lasso to  $a-L_1/L_{\infty}(\alpha)$  for different choices of the mixing parameter  $\alpha$ . The case  $\alpha = 0$  is exactly  $a-L_1/L_{\infty}$  and we also consider  $\alpha = 0.5$ ,  $\alpha = 0.8$  and  $\alpha = 1$ . We note that there is a small difference between Individual Lasso and the case  $\alpha = 1$ : both models consider the same regularization, but the former selects a model through CV in each source, while the latter selects a model through CV for all sources simultaneously.

Under this setting, the statistics of size, accuracy, precision and recall are calculated for each component  $\beta_j^{(k)}$  instead of the vectors  $\beta_j$ . This means that the true model has size 3 + 3 + 2 + 2 + 2 = 12. Using this definition, we will see more clearly the effect of  $\alpha$  on the sparsity of the selected model. The results from 100 replications are summarized in Table 3.

The lowest test deviance is achieved by  $a - L_1/L_{\infty}(\alpha)$  with  $\alpha = 0.5$ . It is not significantly better than other values of the parameter, but clearly has improvement over the Full and Individual Lasso for out-of-sample adjustment. We also observe a decrease in the size of the model as  $\alpha$  increases: starting from 25 selected features with  $\alpha = 0$  (i.e., five features selected in five sources since there is no selection performed across sources) to less than 15 for  $\alpha = 1$ , closing in to the 12 generating features. With perfect recall for all MStweedie algorithms, this means that with  $\alpha = 1$  we achieve the best accuracy and precision. Finally, the  $L_2$  loss being the smallest under  $\alpha = 0.5$  means that its extra selected features have coefficient estimates very close to 0 and that its coefficient estimates for the true features are closer to the true values.

#### 5.2.4. Setting 4—Different Datasets

To test how our algorithm behave under circumstances where some features are missing from certain sources, we consider three simulation setups: (4A) some true generating variables are missing from certain sources, (4B) some spurious variables are missing from certain sources, and (4C) both true and spurious variables are missing from certain sources. For all cases, we generate data as in Setting 2 with K = 5,  $n_k = 300$ , p = 600and the true variable indices are {2, 4, 8, 16, 32}. In Setting 4A, we set to 0 column 32 for sources 1 and 2 and columns 16 and 32 of source 3. In Setting 4B, we set to 0 the last 100 columns of sources 1 and 2 and the last 200 columns of source 3. In Setting 4C, we consider the zero columns of Settings 4A and 4B simultaneously. For demonstration purposes, we compare Full Lasso, Individual Lasso and  $a - L_1/L_{\infty}$ . The results over 100 replications are reported in Table 4.

Under Setting 4A, where true variables are omitted in some sources, we find that  $a-L_1/L_{\infty}$  clearly outperforms both Full Lasso and Individual Lasso under all criteria. As we would expect, it does not achieve the same performance as when using the complete dataset (Setting 2) due to removal of important features.

Under Setting 4B, where only spurious variables are removed from some sources, we do not observe significant difference in any statistic compared to the models trained on the complete data.

Under Setting 4C, where both true and spurious variables are removed from some sources, we observe similar behavior as in Setting 4A, with  $a-L_1/L_{\infty}$  having slightly better performance. It seems that both Full Lasso and  $a-L_1/L_{\infty}$  are less inclined to overfit the spurious information when it is missing from some sources.

#### 5.2.5. Setting 5—Scalability Study

Under the same construction of Setting 1, we conduct a short scalability study of the influence of the number of covariates

#### Table 3. Results from Setting 3 with 100 replications.

(a) Setting 3: Mean (standard error)

|                     |                                |        |                                |        |            | a- $L_1/L_\infty(lpha)$ |            |                |             |                |             |              |  |  |  |
|---------------------|--------------------------------|--------|--------------------------------|--------|------------|-------------------------|------------|----------------|-------------|----------------|-------------|--------------|--|--|--|
|                     | Full Lasso<br>52,398<br>(5830) |        | Ind. Lasso<br>11,590<br>(1448) |        | α :        | $\alpha = 0$            |            | $\alpha = 0.5$ |             | $\alpha = 0.8$ |             | $\alpha = 1$ |  |  |  |
| Test dev.           |                                |        |                                |        | 861<br>(9) |                         | 852<br>(8) |                | 863<br>(10) |                | 866<br>(10) |              |  |  |  |
| Size                | 61.20                          | (5.52) | 30.63                          | (0.55) | 25.00      | (0.00)                  | 16.97      | (0.20)         | 15.34       | (0.18)         | 14.76       | (0.17)       |  |  |  |
| Accuracy            | 98.1                           | (0.2)  | 99.3                           | (0.0)  | 99.6       | (0.0)                   | 99.8       | (0.0)          | 99.9        | (0.0)          | 99.9        | (0.0)        |  |  |  |
| Precision           | 27.3                           | (2.4)  | 38.0                           | (0.5)  | 48.0       | (0.0)                   | 71.7       | (0.9)          | 79.4        | (1.0)          | 82.4        | (1.0)        |  |  |  |
| Recall              | 72.2                           | (3.0)  | 95.1                           | (1.0)  | 100.0      | (0.0)                   | 100.0      | (0.0)          | 100.0       | (0.0)          | 100.0       | (0.0)        |  |  |  |
| L <sub>2</sub> loss | 5.99                           | (0.05) | 2.73                           | (0.09) | 0.38       | (0.01)                  | 0.34       | (0.01)         | 0.36        | (0.01)         | 0.36        | (0.01)       |  |  |  |

(b) Setting 3: Mean rank (# of times best)

|                     |            |      |            |      | $a-L_1/L_\infty(\alpha)$ |       |                |       |                |       |              |       |  |  |
|---------------------|------------|------|------------|------|--------------------------|-------|----------------|-------|----------------|-------|--------------|-------|--|--|
|                     | Full Lasso |      | Ind. Lasso |      | $\alpha = 0$             |       | $\alpha = 0.5$ |       | $\alpha = 0.8$ |       | $\alpha = 1$ |       |  |  |
| Test dev.           | 5.96       | (0)  | 5.04       | (0)  | 2.68                     | (31)  | 2.19           | (31)  | 2.47           | (15)  | 2.66         | (23)  |  |  |
| Size                | 4.82       | (19) | 5.10       | (0)  | 4.37                     | (0)   | 2.89           | (8)   | 1.68           | (43)  | 1.23         | (78)  |  |  |
| Accuracy            | 5.56       | (0)  | 5.15       | (0)  | 4.19                     | (0)   | 2.68           | (8)   | 1.52           | (49)  | 1.08         | (92)  |  |  |
| Precision           | 5.42       | (7)  | 5.14       | (0)  | 4.19                     | (0)   | 2.75           | (8)   | 1.59           | (48)  | 1.15         | (86)  |  |  |
| Recall              | 4.27       | (34) | 2.14       | (74) | 1.00                     | (100) | 1.00           | (100) | 1.00           | (100) | 1.00         | (100) |  |  |
| L <sub>2</sub> loss | 6.00       | (0)  | 5.00       | (0)  | 2.87                     | (18)  | 2.00           | (40)  | 2.41           | (21)  | 2.72         | (21)  |  |  |

NOTE: Part (a) shows the mean values of the statistics and their standard errors in parentheses. Part (b) shows the mean rank across the five models and, in parentheses, the number of times the model is best.

Table 4. Results from Setting 4 with 100 replications: values represent mean values of the statistics and their standard errors in parentheses.

| Setting 4: M        | ean (standar                | d error)   |                        |            |               |                            |                     |            |                        |            |                     |              |
|---------------------|-----------------------------|------------|------------------------|------------|---------------|----------------------------|---------------------|------------|------------------------|------------|---------------------|--------------|
|                     | (2) Comj                    | plete data |                        |            |               | (4A) Missing true features |                     |            |                        |            |                     |              |
|                     | Full Lasso<br>1430<br>(142) |            | Ind. I                 | Ind. Lasso |               | a- $L_1/L_\infty$          |                     | Full Lasso |                        | Ind. Lasso |                     | $L_{\infty}$ |
| Test dev.           |                             |            | 1,136,266<br>(265,655) |            | 1161<br>(129) |                            | 272,840<br>(49,327) |            | 2,122,348<br>(839,704) |            | 152,315<br>(31,194) |              |
| Size                | 22.07                       | (0.76)     | 29.73                  | (1.44)     | 5.00          | (0.00)                     | 34.74               | (2.83)     | 21.38                  | (1.33)     | 11.95               | (0.44)       |
| Accuracy            | 97.2                        | (0.1)      | 95.7                   | (0.2)      | 100.0         | (0.0)                      | 94.9                | (0.5)      | 96.9                   | (0.2)      | 98.7                | (0.1)        |
| Precision           | 24.5                        | (0.6)      | 22.6                   | (2.0)      | 100.0         | (0.0)                      | 19.2                | (1.3)      | 31.5                   | (2.7)      | 44.8                | (1.8)        |
| Recall              | 100.0                       | (0.0)      | 91.2                   | (2.5)      | 100.0         | (0.0)                      | 89.4                | (1.4)      | 78.8                   | (3.5)      | 92.8                | (1.5)        |
| L <sub>2</sub> loss | 0.43                        | (0.03)     | 8.86                   | (0.07)     | 0.34          | (0.01)                     | 6.01                | (0.13)     | 9.34                   | (0.04)     | 5.15                | (0.08)       |

|  | (4B) Missing spurious features |                                   |                                      |                                   |  |                                   | (4C) Missing true and spurious features |                                   |                              |                                   |                              |                                   |  |
|--|--------------------------------|-----------------------------------|--------------------------------------|-----------------------------------|--|-----------------------------------|---|-----------------------------------|------------------------------|-----------------------------------|------------------------------|-----------------------------------|--|
|  | Full Lasso<br>1376<br>(133)    |                                   | Ind. Lasso<br>1,119,800<br>(265,508) |                                   | a- <i>L</i> <sub>1</sub> / <i>L</i> ∞<br>1164<br>(130) |                                   | Full I                                  | asso                              | Ind. Lasso                   |                                   | a- $L_1/L_\infty$            |                                   |  |
| Test dev.  |                                |                                   |                                      |                                   |  |                                   | 245,761<br>(46,198)                     |                                   | 2,116,184<br>(839,728)       |                                   | 119,081<br>(18,616)          |                                   |  |
| Size   | 19.53                          | (0.53)                            | 29.77                                | (1.33)                            | 5.00   | (0.00)                            | 31.25                                   | (2.53)                            | 21.49                        | (1.25)                            | 11.60                        | (0.42)                            |  |
| Accuracy<br>Precision<br>Recall<br>L <sub>2</sub> loss | 97.6<br>27.4<br>100.0<br>0.41  | (0.1)<br>(0.7)<br>(0.0)<br>(0.03) | 95.8<br>22.1<br>94.4<br>8.79         | (0.2)<br>(1.8)<br>(2.0)<br>(0.07) | 100.0<br>100.0<br>100.0<br>0.34                        | (0.0)<br>(0.0)<br>(0.0)<br>(0.01) | 95.5<br>21.1<br>90.4<br>5.99            | (0.4)<br>(1.5)<br>(1.3)<br>(0.12) | 96.9<br>29.9<br>80.8<br>9.33 | (0.2)<br>(2.5)<br>(3.2)<br>(0.04) | 98.8<br>46.4<br>94.6<br>5.09 | (0.1)<br>(1.7)<br>(1.5)<br>(0.08) |  |

NOTE: The four parts, respectively, show the results from Settings 2, 4A, 4B, and 4C for comparison.

#### Table 5. Description of the different parameters used in Setting 5.

| Setting | 5: Description of the scenarios   |                                    |                          |                  |                                    |
|---------|-----------------------------------|------------------------------------|--------------------------|------------------|------------------------------------|
|         | К                                 | p                                  | # of true variables      | % true variables | n <sub>k</sub>                     |
| (a)     | 20                                | $10 \times 3^{i}, i = 0, \dots, 5$ | 10                       | _                | 300                                |
| (b)     | 20                                | $10 \times 3^{i}, i = 0, \dots, 5$ | -                        | 10%              | 300                                |
| (c)     | $5 \times 2^{i}, i = 0, \dots, 5$ | 100                                | 10                       | -                | 300                                |
| (d)     | 20                                | 100                                | $2^{i}, i = 0, \dots, 5$ | _                | 300                                |
| (e)     | 5                                 | 50                                 | 10                       | _                | $5 \times 4^{i}, i = 0, \dots, 5$  |
| (f)     | 5                                 | 1000                               | 10                       | -                | $30 \times 2^{i}, i = 0, \dots, 5$ |

*p*, the number of sources *K* and sample sizes  $n_k$  on the CPU time. We consider different scenarios as shown in Table 5. The running times are averaged over 10 independent runs and are used to compare the  $L_1/L_{\infty}$  and  $L_1/L_2$  regularizations to the Individual Lasso.

Figure 1 contains the plot of the average CPU time versus the variable of interest under the four schemes considered. In parts (a) and (b), the running times of all three algorithms increase at a similar linear rate. In part (*c*), we clearly see, as we would expect, that the running time of individual regularization increases linearly with the number of sources. In contrast, the CPU times of the two MStweedie algorithms increase faster than the linear rate and seem to diminish with *K*. Note that the iteration complexity of the MStweedie algorithm is influenced by *K* mainly in the step that requires Euclidean projections. For  $L_1/L_{\infty}$  regularization, Condat (2016) pointed out that the



(a) K=20, 10 true variables

(b) K=20, 10% true variables

Figure 1. Results from the scalability study under various conditions for synthetic data. The dashed line represents what a linear relation between the CPU time and the variable of interest would follow. All axes are in logarithmic scales.

algorithm by Duchi et al. (2008) has expected and observed complexity  $\mathcal{O}(K)$ , but can be slower (up to  $\mathcal{O}(K^2)$ ) in sparse problems.

In part (d), we study the effect of sparsity by varying the proportion of true variables in the model. For all three algorithms, we note a slight increase of the computing time when the proportion increases. In parts (e) and (f), we look at the effect of the sample size  $n_k$  in the cases  $n_k > p$  and  $n_k < p$ , respectively. A linear rate can be observed for both cases with the MStweedie algorithms. In contrast, the Individual Lasso has CPU time increasing only sublinearly when  $n_k < p$ .

Overall,  $L_1/L_{\infty}$  regularization is systematically slower than  $L_1/L_2$  regularization by a multiplicative constant. Both MStweedie algorithms are slower than individual regularization only by a multiplicative constant.

#### 5.2.6. Setting 6—Correlated Responses in

Upon the request of a referee, we also study the impact of having correlated responses on the performance of our proposed algorithm. The simulation results show that all versions of our algorithm significantly produce better test deviance and the two adaptive versions clearly beats all other models. Due to limited space, we provide the simulation results in Appendix E (supplementary materials).

#### 5.3. Real Data—Automobile Insurance Claims

We apply our algorithm to the analysis of a real dataset studied in Yip and Yau (2005) and Qian, Yang, and Zou (2016). The dataset consists of many automobile insurance policy records and is available as AutoClaim in the R package cplm (Zhang 2011, 2013). A preprocessed version of the data is also available in our R package. It contains the records of 10,296 policies of which 6290 (61.1%) have no claims. We are interested

Table 6. Description of the variables in the auto insurance claim dataset.

in predicting the aggregate claim loss of the policy using the 15 predictors (along with their necessary transformations) described in Table 6. We split the dataset into two sources corresponding to potentially different types of driving license (according to whether or not the policyholder had his or her license revoked). Source 1 contains 9036 policies of which 5643 (62.5%) have no insurance claims and Source 2 contains 1260 policies of which 646 (51.3%) have no insurance claims. Figure 2 plots the histogram of the aggregate claims for both sources.

The following models are considered: the Full and Individual Lasso, and the MStweedie with both  $L_1/L_2(\alpha)$  and  $L_1/L_{\infty}(\alpha)$  regularizations as well as their adaptive counterparts under different values of the mixing parameter  $\alpha \in \{0, 0.5, 0.8, 1\}$ . We split the dataset into a training and a testing set consisting, respectively, of two-thirds and one-third of the policies of each source. The 10-fold CV is then performed to select the best model. Finally, we summarize the results by averaging them over 100 replications of training/testing random partition.

The results of the study are reported in Table 7. In terms of model fit, we note that all adaptive MStweedie methods perform very similarly while the nonadaptive procedures and the Individual Lasso are slightly worse and the Full Lasso is the worst. In terms of model sparsity, the Individual Lasso produces the simplest models on average followed by the adaptive MStweedie algorithms and then the Full Lasso. The nonadaptive MStweedie algorithms yield models that have significantly more variables.

Now, by looking at the exact variables selected within each source, we first see that MVR\_PTS and AREA are systematically included in every model except the Individual Lasso which does not include AREA for source 2. When  $\alpha$  is nonzero, there is no major difference between the models under different values of

| AutoClaim Dataset Va | riable Description |                   |  |
|----------------------|--------------------|-------------------|--|
| Variable             | Туре               | Transformation    | Description  |
| Response             |                    |                   |  |
| CLM_AMT5             | Numerical          | ×10 <sup>-3</sup> | Aggregate claim loss of policy   |
| Source identifier    |                    |                   |  |
| REVOKED              | Categorical(2)     | 1/2               | Whether the policyholder's license was (2) revoked in the past or (1) not  |
| Predictors           |                    |                   |  |
| KIDSDRIV             | Numerical          | -                 | Number of child passengers   |
| TRAVTIME             | Numerical          | -                 | Commute time   |
| CAR USE              | Categorical(2)     | 1/2               | (1) Private or (2) Commercial use  |
| BLUEBOOK             | Numerical          | log               | Car value  |
| NPOLICY              | Numerical          | -                 | Number of policies   |
| RED CAR              | Categorical(2)     | 1/2               | Whether the color of the car is (2) red or (1) not   |
| MVR PTS              | Numerical          | -                 | Number of motor vehicle record points  |
| AGE                  | Numerical          | -                 | Age of policyholder  |
| HOMEKIDS             | Numerical          | -                 | Number of children at home   |
| GENDER               | Categorical(2)     | 1/2               | Gender of policyholder: (2) male or (1) female   |
| PARENT1              | Categorical(2)     | 1/2               | Whether (2) the policyholder grew up in a single-parent family or (1) not  |
| AREA                 | Categorical(2)     | 1/2               | (1) Rural or (2) urban area  |
| CAR_TYPE             | Categorical(6)     | Dummy(5)          | Type of car: (base) Panel Truck, (2) Pickup, (3) Sedan, (4) Sports Car, (5)<br>SUV, (6) Van  |
| JOBCLASS             | Categorical(9)     | Dummy(8)          | Job class of policyholder: (base) Unknown, (2) Blue Collar, (3) Clerical, (4)<br>Doctor, (5) Home Maker, (6) Lawyer, (7) Manager, (8) Professional, (9)<br>Student |
| MAX_EDUC             | Categorical(5)     | Dummy(4)          | Maximal level of education of policyholder: (base) less than High School,<br>(2) Bachelors, (3) High School, (4) Masters, (5) PhD                                  |



# Figure 2. Frequency of the aggregate claim amounts in the AutoClaim dataset according the whether or not the policyholder's license was revoked (defining the two sources).

| Table 7. Test deviance, size of the selected model and selected variables under different regularization schemes on the AutoClaim | ı dataset. |
|---|------------|
| Auto Claiman Maran (standard amar)  |            |

| Algorithm                      | Test De | eviance | Si    | ze     | Selected variables (# of times in source 1, in source 2)   |
|--------------------------------|---------|---------|-------|--------|--|
| Full Lasso                     | 22,203  | (35)    | 5.32  | (0.30) | MVR_PTS(100,100), AREA(100,100), CAR_TYPE_4(25,25),<br>MARRIED(14,14), PARENT1(6,6), KIDSDRIV(4,4), CAR_TYPE_3(4,4),<br>JOBCLASS_6(3,3), MAX_EDUC_3(3,3), BLUEBOOK(2,2),<br>JOBCLASS_3(2,2), JOBCLASS_4(2,2), MAX_EDUC_5(1,1)  |
| Ind. Lasso                     | 19,493  | (33)    | 3.77  | (0.11) | MVR_PTS(100,100), AREA(100,36), CAR_TYPE_4(0,27),<br>CAR_USE(0,1), MARRIED(2,1), JOBCLASS_3(0,1),<br>JOBCLASS_6(1,1), MAX_EDUC_4(0,1), AGE_CAT_5(0,1)  |
| $L_1/L_\infty$                 | 19,475  | (32)    | 13.08 | (0.63) | MVR_PTS(100,100), AREA(100,100), CAR_TYPE_4(98,98),<br>JOBCLASS_3(59,59), JOBCLASS_6(59,59), CAR_TYPE_5(33,33),<br>MARRIED(25,25), JOBCLASS_7(22,22), KIDSDRIV(21,21),<br>AGE_CAT_5(20,20), AGE_CAT_2(16,16), JOBCLASS_5(15,15),<br>CAR_USE(12,12), BLUEBOOK(10,10), CAR_TYPE_6(10,10),<br>MAX_EDUC_4(10,10), JOBCLASS_4(8,8), JOBCLASS_8(7,7),<br>RED_CAR(4,4), TRAVTIME(3,3), CAR_TYPE_2(3,3),<br>CAR_TYPE_3(3,3), MAX_EDUC_2(3,3), AGE_CAT_4(3,3),<br>PARENT1(2,2), MAX_EDUC_3(2,2), AGE_CAT_3(2,2), NPOLICY(1,1),<br>GENDER(1,1), JOBCLASS_2(1,1), MAX_EDUC_5(1,1) |
| $a-L_1/L_\infty(0)$            | 19,438  | (32)    | 5.00  | (0.12) | MVR_PTS(100,100), AREA(100,100), CAR_TYPE_4(45,45),<br>JOBCLASS_6(4,4), MARRIED(1,1)   |
| $a\text{-}L_1/L_\infty(0.5)$   | 19,437  | (31)    | 4.31  | (0.05) | MVR_PTS(100,100), AREA(100,100), CAR_TYPE_4(1,28),<br>MARRIED(0,1), JOBCLASS_6(0,1)  |
| $a\text{-}L_1/L_\infty(0.8)$   | 19,431  | (32)    | 4.29  | (0.05) | MVR_PTS(100,100), AREA(100,100), CAR_TYPE_4(1,26),<br>JOBCLASS_6(0,2)  |
| $a-L_1/L_\infty(1)$            | 19,431  | (32)    | 4.29  | (0.05) | MVR_PTS(100,100), AREA(100,100), CAR_TYPE_4(0,28),<br>JOBCLASS_6(0,1)  |
| L <sub>1</sub> /L <sub>2</sub> | 19,456  | (30)    | 9.86  | (0.29) | MVR_PTS(100,100), AREA(100,100), CAR_TYPE_4(99,99),<br>JOBCLASS_3(50,50), JOBCLASS_6(44,44), CAR_TYPE_5(16,16),<br>JOBCLASS_7(16,16), JOBCLASS_5(12,12), AGE_CAT_5(11,11),<br>AGE_CAT_2(9,9), MAX_EDUC_4(8,8), CAR_USE(6,6), MARRIED(6,6),<br>BLUEBOOK(5,5), KIDSDRIV(2,2), RED_CAR(2,2), GENDER(1,1),<br>CAR_TYPE_3(1,1), CAR_TYPE_6(1,1), JOBCLASS_4(1,1),<br>JOBCLASS_8(1,1), MAX_EDUC_2(1,1), AGE_CAT_4(1,1)   |
| $a-L_1/L_2(0)$                 | 19,434  | (31)    | 5.00  | (0.11) | MVR_PTS(100,100), AREA(100,100), CAR_TYPE_4(48,48),<br>MARRIED(1,1), JOBCLASS_6(1,1)   |

#### Auto Insurance Claims

#### Table 7. Continued

| Auto Claims: Mean (standard error) |         |         |      |        |   |  |  |  |
|------------------------------------|---------|---------|------|--------|---|--|--|--|
| Algorithm                          | Test De | eviance | S    | ize    | Selected variables (# of times in source 1, in source 2)  |  |  |  |
| $a-L_1/L_2(0.5)$                   | 19,442  | (32)    | 4.61 | (0.05) | MVR_PTS(100,100), AREA(100,100), CAR_TYPE_4(1,56),<br>JOBCLASS_6(0,2), MAX_EDUC_4(0,1), AGE_CAT_5(0,1)                  |  |  |  |
| d-L1/L2(0.0)                       | 19,432  | (31)    | 4.00 | (0.03) | JOBCLASS_6(0,4), MAX_EDUC_4(0,1), AGE_CAT_5(0,1)  |  |  |  |
| $a-L_1/L_2(1)$                     | 19,428  | (31)    | 4.72 | (0.05) | MVR_PTS(100,100), AREA(100,100), CAR_TYPE_4(0,66),<br>JOBCLASS_6(0,3), MARRIED(0,1), MAX_EDUC_4(0,1),<br>AGE_CAT_5(0,1) |  |  |  |

NOTE: The results are averaged over 100 replications of the training/testing splitting.



(a) 10 folds CV train deviance

**Figure 3.** MStweedie with adaptive  $L_1/L_{\infty}$  regularization on AutoClaim data. Panel (a) shows the plot of the 10-fold CV mean deviance (and its standard error) along the  $\lambda$  sequence. Panel (b) plots the norm of the estimates,  $||\beta_j||_{\infty}$ , along the  $\lambda$  sequence. In each pane, the grey vertical line indicates the  $\lambda$  for which the CV deviance is minimal and the black vertical line indicates the  $\lambda$  value selected according to the one-standard-error rule.

 $\alpha$ , but they all behave as expected: for example, they select the variable CAR\_TYPE\_4 (corresponding to "Sports Car") only for source 2, corresponding to between-sources sparsity that the  $\alpha = 0$  method cannot uncover.

CV is used to select the optimal value of  $\lambda$ . We plot the CV deviance as well as its standard error along the sequence of  $\lambda$  values and display the minimal value as well as the selected  $\lambda$  according to the one-standard-error rule in Figure 3. The figure also contains the plot of the norm of the estimated coefficients for  $a - L_1/L_{\infty}$ . It provides an excellent example of why the one-standard-error rule is often favored in practice: its selected model does not have a significantly different model fit than the one minimizing the CV error, but it is considerably sparser.

Furthermore, to have a real data example that approaches more real-world situations, we artificially increase the proportion of zeros of the dataset by sub-sampling the nonzero responses. We consider target proportions between 65% and 95% and remove enough nonzero claim amounts observation from the dataset to reach the given proportion. The new datasets will be smaller in size and the proportion of zeroes may differ between the two sources: the simple random sampling ensures that the disproportion remains the same on average. Table 8 contains details on the new datasets.

The same experimental methodology as with the original data is performed; Figure 4 contains the normalized test deviance and model size, both averaged over 10 replications of the sampling, of the two base algorithms and of  $a-L_1/L_{\infty}(\alpha)$  for  $\alpha \in \{0.8, 1\}$ , which performed the best on the original dataset. Uniformly over the range of proportion of zeros, our algorithm exhibit performance similar or better than individual Lasso and significantly better than 
 Table 8.
 Number of observations and proportion of zeros in the whole dataset and within each sources for the datasets sampled from the AutoClaim dataset to yield a target global proportion of zeros.

| AutoClaim: resampled datasets |        |          |          |  |  |  |
|-------------------------------|--------|----------|----------|--|--|--|
|                               |        |          |          |  |  |  |
| Ν                             | Global | Source 1 | Source 2 |  |  |  |
| 10,296                        | 61.6   | 62.5     | 51.3     |  |  |  |
| 9677                          | 65.0   | 66.3     | 55.5     |  |  |  |
| 8986                          | 70.0   | 71.2     | 60.8     |  |  |  |
| 8387                          | 75.0   | 76.1     | 66.5     |  |  |  |
| 7863                          | 80.0   | 81.0     | 72.4     |  |  |  |
| 7400                          | 85.0   | 85.8     | 78.9     |  |  |  |
| 6989                          | 90.0   | 90.6     | 85.4     |  |  |  |
| 6622                          | 95.0   | 95.2     | 92.8     |  |  |  |

the Lasso on the complete dataset. Except for a proportion of zeros of 95%, the adjustment to test data and the sparsity are essentially the same between the multisource algorithms and the independent Lasso. For a proportion of 95%, the multisource algorithm produce significantly sparser models with very similar adjustment. Hence, sharing information between sources for variable selection seem to allow the algorithm to discard more efficiently faint signals from particular features.

#### 5.4. Discussion—Choosing the Regularization

Before fitting the model, it is not obvious which regularization should be used. All simulations as well as the real data example indicate that adaptive regularization should always be preferred to increase the prediction, estimation and selection efficiency of our model. In our experiments, it seems that  $L_1/L_{\infty}$  performs better than  $L_1/L_2$  when the coefficients are the same across tasks for a same feature while the converse is true when the coefficients differ between sources. Additionally, the sparse penalty only helps reduce the size of the model for variable selection (it does not improve the fit) and it seems that most of the gain comes from relatively small  $\alpha$  values.

Hence, we can outline some general guideline in selecting the regularization. If the user suspect that the coefficient will vary wildly between sources, then  $L_1/L_2$  should be preferred while  $L_1/L_\infty$  should be preferred when the coefficient are thought to be roughly the same. A sparse penalty never seem to hurt, so using  $\alpha = 0.5$  may uncover some additional sparsity. Nonetheless, it is difficult to predict what the optimal solution would look like and it might be useful to study the results of a trial run to tune the regularization of the real fit.

#### 6. Conclusion

In this article, we develop a unified algorithm for sparse learning of multisource insurance data using the MStweedie method. The Mstweedie-GPG algorithm we proposed cyclically updates each group of coefficients via the proximal gradient descent scheme and enjoys fast convergence guarantee. This procedure is embedded in a solution path algorithm to achieve the best balance between goodness of fit and model sparsity.

Experiments on simulated data show that our approach clearly outperforms simpler methods in prediction and selection accuracy. It is particularly effective for datasets having distinct structures across the sources. The various regularization schemes behave as expected and thus provide additional flexibility for our algorithm to allow user specification of the desired type of sparsity. While our implementation scales well



Figure 4. For the two base algorithms and the two best multisource algorithms: (a) normalized test deviance and (b) number of variables in the model, both averaged over 10 replications of the sampling of the AutoClaim dataset to yield target global proportions of zeros. The error bars represent one standard error around the mean.

with the number of observations and variables in a dataset, we caution that an increasing number of sources may slow down the calculation because of the increased number of Euclidean projections required. When applied to real data constituted of aggregate claim amount of the automobile insurance, our procedure convey similar messages to those from the simulated experiments. We also note that although our approach is specifically designed for the Tweedie model with actuarial applications, it is possible to develop similar algorithms for alternative model choices.

In addition, though beyond the scope of our work, a promising approach is to use the multivariate copula to account for the conditional correlation between data sources. For example, Shi (2016) and Frees et al. (2018) proposed multivariate Tweedie copula models, Czado et al. (2012) used a copula on the frequency-severity pair of a single claim with Gamma severity (see also Shi and Zhang 2013 and Shi, Feng, and Ivantsova 2015), and Frees, Shi, and Valdez (2009) use a copula to jointly model a single frequency with hierarchical Generalized Beta claim amounts (see also Frees and Valdez 2008). There is also work on joint modeling of multivariate claim counts (e.g., Bermúdez and Karlis 2011; Nikoloulopoulos 2013; Shi and Valdez 2014). See Frees, Lee, and Yang (2016) and many references therein for a comprehensive review of multivariate insurance claim data modeling. Accordingly, variable selection of multisource data within a multivariate copula model framework can be a promising topic and we leave it for further investigation.

Last but not least, we note that the Tweedie model has a wide range of applications well beyond the scope of our presentation in this article. Examples of nonnegative valued data with excess zeros can also be found in other actuarial settings (Tong, Mues, and Thomas 2013; Frees, Jin, and Lin 2013; Frees, Gao, and Rosenberg 2011; Lauderdale 2012), and in ecology (Blakey et al. 2016; Foster and Bravington 2013; Zhang 2011), fishery (Ancelet et al. 2010; Shono 2008), meteorology (Dunn 2004; Smyth 1996; Swan 2006), and health (Buu et al. 2011; Moger and Aalen 2005; Smyth 1996), to name a few. We hope that this work builds new and useful research tool for many of these promising applications.

#### **Supplementary Materials**

- **MSTweedie** The R package implementing our proposed methods available at the address *https://github.com/fontaine618/MSTweedie*.
- AutoClaim The dataset used in the real data experiment is available within MSTweedie package.
- **Appendices** This appendix file contains additional numerical examples and results not shown in the main article. (appendix.pdf)

#### Acknowledgments

We sincerely thank the editor, the associate editor, and two anonymous reviewers for their valuable comments.

#### Funding

Yang's research is partially supported by NSERC RGPIN-2016-05174 and FRQ-NT NC-205972. Qian's research is partially supported by NSF DMS-1916376 and JMPC Faculty Fellowship.

#### References

- Ancelet, S., Etienne, M.-P., Benoît, H., and Parent, E. (2010), "Modelling Spatial Zero-Inflated Continuous Data With an Exponentially Compound Poisson Process," *Environmental and Ecological Statistics*, 17, 347– 376. [355]
- Beck, A., and Teboulle, M. (2009), "A Fast Iterative Shrinkage-Thresholding Algorithm for Linear Inverse Problems," SIAM Journal on Imaging Sciences, 2, 183–202. [342]
- Bermúdez, L., and Karlis, D. (2011), "Bayesian Multivariate Poisson Models for Insurance Ratemaking," *Insurance: Mathematics and Economics*, 48, 226–236. [355]
- Blakey, R. V., Law, B. S., Kingsford, R. T., Stoklosa, J., Tap, P., and Williamson, K. (2016), "Bat Communities Respond Positively to Large-Scale Thinning of Forest Regrowth," *Journal of Applied Ecology*, 53, 1694–1703. [355]
- Buu, A., Johnson, N. J., Li, R., and Tan, X. (2011), "New Variable Selection Methods for Zero-Inflated Count Data With Applications to the Substance Abuse Field," *Statistics in Medicine*, 30, 2326–2340. [355]
- Chandler, R. E., and Bate, S. (2007), "Inference for Clustered Data Using the Independence Loglikelihood," *Biometrika*, 94, 167–183. [341]
- Condat, L. (2016), "Fast Projection Onto the Simplex and the  $\ell_1$  Ball," Mathematical Programming, 158, 575–585. [350]
- Czado, C., Kastenmeier, R., Brechmann, E. C., and Min, A. (2012), "A Mixed Copula Model for Insurance Claims and Claim Sizes," *Scandinavian Actuarial Journal*, 2012, 278–305. [355]
- Duchi, J., Shalev-Shwartz, S., Singer, Y., and Chandra, T. (2008), "Efficient Projections Onto the  $\ell_1$ -Ball for Learning in High Dimensions," in *Proceedings of the 25th International Conference on Machine Learning*, ACM, pp. 272–279. [344,351]
- Dunn, P. K. (2004), "Occurrence and Quantity of Precipitation can be Modelled Simultaneously," *International Journal of Climatology*, 24, 1231– 1239. [355]
- Dunn, P. K., and Smyth, G. K. (2005), "Series Evaluation of Tweedie Exponential Dispersion Model Densities," *Statistics and Computing*, 15, 267–280. [341]
- Fawcett, T. (2006), "An Introduction to ROC Analysis," *Pattern Recognition Letters*, 27, 861–874. [346]
- 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. [355]
- Frees, E. W., Bolancé, C., Guillen, M., and Valdez, E. (2018), "Joint Models of Insurance Lapsation and Claims," arXiv no. 1810.04567. [355]
- Frees, E. W., Gao, J., and Rosenberg, M. A. (2011), "Predicting the Frequency and Amount of Health Care Expenditures," North American Actuarial Journal, 15, 377–392. [339,355]
- Frees, E. W., Jin, X., and Lin, X. (2013), "Actuarial Applications of Multivariate Two-Part Regression Models," *Annals of Actuarial Science*, 7, 258– 287. [355]
- Frees, E. W., Lee, G., and Yang, L. (2016), "Multivariate Frequency-Severity Regression Models in Insurance," *Risks*, 4, 4. [339,355]
- Frees, E. W., Meyers, G., and Cummings, A. D. (2011), "Summarizing Insurance Scores Using a Gini Index," *Journal of the American Statistical Association*, 106, 1085–1098. [339]
- Frees, E. W., Shi, P., and Valdez, E. A. (2009), "Actuarial Applications of a Hierarchical Insurance Claims Model," ASTIN Bulletin: The Journal of the IAA, 39, 165–197. [355]
- Frees, E. W., and Valdez, E. A. (2008), "Hierarchical Insurance Claims Modeling," *Journal of the American Statistical Association*, 103, 1457– 1469. [355]
- Friedman, J., Hastie, T., and Tibshirani, R. (2010), "Regularization Paths for Generalized Linear Models via Coordinate Descent," *Journal of Statisti*cal Software, 33, 1. [340,343,346]
- Gong, P., Ye, J., and Zhang, C. (2012), "Robust Multi-Task Feature Learning," in Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, ACM, pp. 895–903. [346]
- Gu, Y., and Zou, H. (2016), "High-Dimensional Generalizations of Asymmetric Least Squares Regression and Their Applications," *The Annals of Statistics*, 44, 2661–2694. [347]

- Hastie, T., Tibshirani, R., and Friedman, J. (2009), *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, Springer Series in Statistics (2nd ed.), New York: Springer. [346]
- Huang, J., and Zhang, T. (2010), "The Benefit of Group Sparsity," *The Annals of Statistics*, 38, 1978–2004. [340]
- Jenatton, R., Mairal, J., Bach, F. R., and Obozinski, G. R. (2010), "Proximal Methods for Sparse Hierarchical Dictionary Learning," in *Proceedings of* the 27th International Conference on Machine Learning (ICML-10), pp. 487–494. [340,344,345]
- Jørgensen, B. (1987), "Exponential Dispersion Models," Journal of the Royal Statistical Society, Series B, 49, 127–162. [340]
- Kim, S., and Xing, E. P. (2012), "Tree-Guided Group Lasso for Multi-Response Regression With Structured Sparsity, With an Application to EQTL Mapping," *The Annals of Applied Statistics*, 6, 1095–1117. [340]
- Lauderdale, B. E. (2012), "Compound Poisson–Gamma Regression Models for Dollar Outcomes That Are Sometimes Zero," *Political Analysis*, 20, 387–399. [355]
- Lounici, K., Pontil, M., Tsybakov, A. B., and Van De Geer, S. (2009), "Taking Advantage of Sparsity in Multi-Task Learning," arXiv no. 0903.1468. [340]
- Lounici, K., Pontil, M., Van De Geer, S., and Tsybakov, A. B. (2011), "Oracle Inequalities and Optimal Inference Under Group Sparsity," *The Annals* of Statistics, 39, 2164–2204. [340]
- Moger, T. A., and Aalen, O. O. (2005), "A Distribution for Multivariate Frailty Based on the Compound Poisson Distribution With Random Scale," *Lifetime Data Analysis*, 11, 41–59. [355]
- Morales, J., Micchelli, C. A., and Pontil, M. (2010), "A Family of Penalty Functions for Structured Sparsity," in Advances in Neural Information Processing Systems, pp. 1612–1623. [340]
- Nikoloulopoulos, A. K. (2013), "Copula-Based Models for Multivariate Discrete Response Data," in *Copulae in Mathematical and Quantitative Finance*, Berlin, Heidelberg: Springer, pp. 231–249. [355]
- Obozinski, G., Taskar, B., and Jordan, M. (2006), "Multi-Task Feature Selection," Statistics Department, UC Berkeley, Tech. Rep. 2. [341]
- (2010), "Joint Covariate Selection and Joint Subspace Selection for Multiple Classification Problems," *Statistics and Computing*, 20, 231–252.
   [340]
- Obozinski, G., Wainwright, M. J., and Jordan, M. I. (2008), "Union Support Recovery in High-Dimensional Multivariate Regression," in 2008 46th Annual Allerton Conference on Communication, Control, and Computing, IEEE, pp. 21–26. [340]
- Qian, W., Li, W., Sogawa, Y., Fujimaki, R., Yang, X., and Liu, J. (2019), "An Interactive Greedy Approach to Group Sparsity in High Dimensions," *Technometrics*, 61, 409–421. [340]
- Qian, W., Yang, Y., and Zou, H. (2016), "Tweedie's Compound Poisson Model With Grouped Elastic Net," *Journal of Computational and Graphical Statistics*, 25, 606–625. [340,346,351]
- Shi, P. (2016), "Insurance Ratemaking Using a Copula-Based Multivariate Tweedie Model," *Scandinavian Actuarial Journal*, 2016, 198–215. [339,341,355]
- Shi, P., Feng, X., and Boucher, J.-P. (2016), "Multilevel Modeling of Insurance Claims Using Copulas," *The Annals of Applied Statistics*, 10, 834– 863. [339]
- Shi, P., Feng, X., and Ivantsova, A. (2015), "Dependent Frequency–Severity Modeling of Insurance Claims," *Insurance: Mathematics and Economics*, 64, 417–428. [339,355]
- Shi, P., and Valdez, E. A. (2014), "Multivariate Negative Binomial Models for Insurance Claim Counts," *Insurance: Mathematics and Economics*, 55, 18–29. [355]
- Shi, P., and Zhang, W. (2013), "Managed Care and Health Care Utilization: Specification of Bivariate Models Using Copulas," North American Actuarial Journal, 17, 306–324. [355]

- Shono, H. (2008), "Application of the Tweedie Distribution to Zero-Catch Data in CPUE Analysis," *Fisheries Research*, 93, 154–162. [355]
- Simon, N., Friedman, J., Hastie, T., and Tibshirani, R. (2013), "A Sparse-Group Lasso," *Journal of Computational and Graphical Statistics*, 22, 231– 245. [340,342]
- 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. [355]
- Smyth, G. K., and Jørgensen, B. (2002), "Fitting Tweedie's Compound Poisson Model to Insurance Claims Data: Dispersion Modelling," ASTIN Bulletin: The Journal of the IAA, 32, 143–157. [339,341]
- Swan, T. (2006), "Generalized Estimating Equations When the Response Variable has a Tweedie Distribution: An Application for Multi-Site Rainfall Modelling," Ph.D. dissertation, University of Southern Queensland. [355]
- Tibshirani, R. (1996), "Regression Shrinkage and Selection via the Lasso," *Journal of the Royal Statistical Society*, Series B, 58, 267–288. [340]
- Tibshirani, R., Bien, J., Friedman, J., Hastie, T., Simon, N., Taylor, J., and Tibshirani, R. J. (2012), "Strong Rules for Discarding Predictors in Lasso-Type Problems," *Journal of the Royal Statistical Society*, Series B, 74, 245– 266. [346]
- Tong, E. N., Mues, C., and Thomas, L. (2013), "A Zero-Adjusted Gamma Model for Mortgage Loan Loss Given Default," *International Journal of Forecasting*, 29, 548–562. [355]
- Turlach, B. A., Venables, W. N., and Wright, S. J. (2005), "Simultaneous Variable Selection," *Technometrics*, 47, 349–363. [342]
- 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. [339]
- Varin, C., Reid, N., and Firth, D. (2011), "An Overview of Composite Likelihood Methods," *Statistica Sinica*, 21, 5–42. [341]
- Vincent, M., and Hansen, N. R. (2014), "Sparse Group Lasso and High Dimensional Multinomial Classification," *Computational Statistics & Data Analysis*, 71, 771–786. [340]
- Yang, Y., Qian, W., and Zou, H. (2018), "Insurance Premium Prediction via Gradient Tree-Boosted Tweedie Compound Poisson Models," *Journal of Business & Economic Statistics*, 36, 456–470. [341]
- 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. [339,351]
- Yuan, M., and Lin, Y. (2006), "Model Selection and Estimation in Regression With Grouped Variables," *Journal of the Royal Statistical Society*, Series B, 68, 49–67. [341]
- Zhang, H. H., Liu, Y., Wu, Y., and Zhu, J. (2008), "Variable Selection for the Multicategory SVM via Adaptive Sup-Norm Regularization," *Electronic Journal of Statistics*, 2, 149–167. [340]
- Zhang, W. (2011), "cplm: Monte Carlo EM Algorithms and Bayesian Methods for Fitting Tweedie Compound Poisson Linear Models," R Package, available at http://cran.r-project.org/web/packages/cplm/index. html. [351,355]
- Zhang, Y. (2013), "Likelihood-Based and Bayesian Methods for Tweedie Compound Poisson Linear Mixed Models," *Statistics and Computing*, 23, 743–757. [339,351]
- Zhao, P., Rocha, G., and Yu, B. (2009), "The Composite Absolute Penalties Family for Grouped and Hierarchical Variable Selection," *The Annals of Statistics*, 37, 3468–3497. [341]
- Zou, H. (2006), "The Adaptive Lasso and Its Oracle Properties," *Journal of the American Statistical Association*, 101, 1418–1429. [340,342,346]

# Appendices for "A Unified Approach to Sparse Tweedie Modeling of Multi-Source Insurance Claim Data"

Simon Fontaine<sup>\*</sup>, Yi Yang<sup>†</sup>, Wei Qian<sup>‡</sup>, Yuwen Gu<sup>§</sup>, Bo Fan<sup>¶</sup>

July 19, 2019

#### **Appendix A.** Projection onto the $L_1$ -Ball

Proof of Lemma 1. Note that (16) can be written as

$$\operatorname{prox}_{\tau h}(\mathbf{u}) = \operatorname*{arg\,min}_{\boldsymbol{\beta}_j} \frac{1}{2} \|\boldsymbol{\beta}_j - (\tilde{\boldsymbol{\beta}}_j - t_j \nabla_j \ell(\tilde{\boldsymbol{\beta}}_j; \tilde{\mathbf{b}}_{-j}))\|_2^2 + \lambda v_j t_j \|\boldsymbol{\beta}_j\|_{\infty},$$

where  $\tau = \lambda v_j t_j$ ,  $h = || \cdot ||_{\infty}$  and  $\mathbf{u} = \tilde{\boldsymbol{\beta}}_j - t_j \nabla_j \ell(\tilde{\boldsymbol{\beta}}_j; \tilde{\mathbf{b}}_{-j})$ . By the Moreau decomposition (Parikh et al., 2014), we have

$$\operatorname{prox}_{h}(\mathbf{u}) = \mathbf{u} - \operatorname{prox}_{h^{*}}(\mathbf{u}),$$

where  $h^*$  denotes the convex conjugate of h. We want to derive a similar identity for  $\tau h$ ,  $\tau > 0$ . The convex conjugate of  $\tau h$  is

$$(\tau h)^*(\mathbf{u}) = \sup_{\mathbf{v}} \left( \mathbf{u}^\top \mathbf{v} - \tau h(\mathbf{v}) \right) = \tau \sup_{\mathbf{v}} \left( \frac{1}{\tau} \mathbf{u}^\top \mathbf{v} - h(\mathbf{v}) \right) = \tau h^* \left( \frac{\mathbf{u}}{\tau} \right).$$

<sup>\*</sup>Department of Mathematics and Statistics, University of Montreal (fontaines@dms.umontreal.ca)

<sup>&</sup>lt;sup>†</sup>Corresponding author, Department of Mathematics and Statistics, McGill University (yi.yang6@mcgill.ca)

<sup>&</sup>lt;sup>‡</sup>Department of Applied Economics and Statistics, University of Delaware (weiqian@udel.edu)

<sup>&</sup>lt;sup>§</sup>Department of Statistics, University of Connecticut (yuwen.gu@uconn.edu)

<sup>&</sup>lt;sup>¶</sup>Department of Statistics, University of Oxford (bo.fan@lmh.ox.ac.uk)

Then, we get

$$prox_{(\tau h)^*}(\mathbf{u}) = \arg\min_{\mathbf{v}} (\tau h)^*(\mathbf{v}) + \frac{1}{2} \|\mathbf{v} - \mathbf{u}\|_2^2$$
  

$$= \arg\min_{\mathbf{v}} \tau h^* \left(\frac{\mathbf{v}}{\tau}\right) + \frac{1}{2} \|\mathbf{v} - \mathbf{u}\|_2^2$$
  

$$= \arg\min_{\mathbf{v}} h^* \left(\frac{\mathbf{v}}{\tau}\right) + \frac{1}{2\tau} \|\mathbf{v} - \mathbf{u}\|_2^2 \qquad (\mathbf{v} = \tau \mathbf{z})$$
  

$$= \arg\min_{\tau \mathbf{z}} h^*(\mathbf{z}) + \frac{1}{2\tau} \|\tau \mathbf{z} - \mathbf{u}\|_2^2$$
  

$$= \tau \arg\min_{\mathbf{z}} h^*(\mathbf{z}) + \frac{1}{2\frac{1}{\tau}} \|\mathbf{z} - \frac{\mathbf{u}}{\tau}\|_2^2$$
  

$$= \tau \operatorname{prox}_{\frac{1}{\tau}h^*} \left(\frac{\mathbf{u}}{\tau}\right),$$

so we have the identity

$$\operatorname{prox}_{\tau h}(\mathbf{u}) = \mathbf{u} - \operatorname{prox}_{(\tau h)^*}(\mathbf{u}) = \mathbf{u} - \tau \operatorname{prox}_{\frac{1}{\tau}h^*}\left(\frac{\mathbf{u}}{\tau}\right).$$

For  $h = || \cdot ||_{\infty}$ , it can be shown that  $\tau \operatorname{prox}_{\frac{1}{\tau}h^*}\left(\frac{\mathbf{u}}{\tau}\right)$  is equivalent to the  $L_2$ -projection of  $\mathbf{u}$  onto an  $L_1$ -ball  $B_1(\tau)$ ,

$$au \operatorname{prox}_{\frac{1}{\tau}h^*}\left(\frac{\mathbf{u}}{\tau}\right) = \operatorname{Proj}_{B_1(\tau)}(\mathbf{u}).$$

To see this, note that the convex conjugate  $h^*$  of  $h = ||\cdot||_\infty$  is

$$h^{*}(\mathbf{u}) = I_{\{\mathbf{u}: ||\mathbf{u}||_{1} \le 1\}} = \begin{cases} 0, & ||\mathbf{u}||_{1} \le 1, \\ +\infty, & ||\mathbf{u}||_{1} > 1, \end{cases}$$

and

$$2\tau h^*\left(\frac{\mathbf{z}}{\tau}\right) = \begin{cases} 0, & ||\mathbf{z}||_1 \le \tau, \\ +\infty, & ||\mathbf{z}||_1 > \tau. \end{cases}$$

Then

$$\begin{aligned} \tau \operatorname{prox}_{\frac{1}{\tau}h^*} \left( \frac{\mathbf{u}}{\tau} \right) &= \tau \operatorname{arg\,min}_{\mathbf{v}} h^*(\mathbf{v}) + \frac{\tau}{2} \left\| \mathbf{v} - \frac{\mathbf{u}}{\tau} \right\|_2^2 \\ &= \operatorname{arg\,min}_{\mathbf{z}} h^* \left( \frac{\mathbf{z}}{\tau} \right) + \frac{\tau}{2} \left\| \frac{\mathbf{z}}{\tau} - \frac{\mathbf{u}}{\tau} \right\|_2^2 \\ &= \operatorname{arg\,min}_{\mathbf{z}} h^* \left( \frac{\mathbf{z}}{\tau} \right) + \frac{1}{2\tau} \left\| \mathbf{z} - \mathbf{u} \right\|_2^2 \\ &= \operatorname{arg\,min}_{\mathbf{z}} 2\tau h^* \left( \frac{\mathbf{z}}{\tau} \right) + \left\| \mathbf{z} - \mathbf{u} \right\|_2^2. \end{aligned}$$

The objective function is minimized at where  $2\tau h^*\left(\frac{\mathbf{z}}{\tau}\right)$  is finite, i.e.,  $||\mathbf{z}||_1 \leq \tau$ . Hence, we get

$$\tau \operatorname{prox}_{\frac{1}{\tau}h^*} \left( \frac{\mathbf{u}}{\tau} \right) = \operatorname*{arg\,min}_{\mathbf{z}:||\mathbf{z}||_1 \le \tau} \|\mathbf{z} - \mathbf{u}\|_2^2 = \operatorname{Proj}_{B_1(\tau)}(\mathbf{u}).$$

If  $||\mathbf{u}||_1 \leq \tau$ , we obviously have  $\operatorname{Proj}_{B_1(\tau)}(\mathbf{u}) = \mathbf{u}$ . Otherwise, we have to solve

$$\sum_{k=1}^{K} (|u_k| - \xi)_+ = \tau$$

for  $\xi$  and compute

$$\left[\operatorname{Proj}_{B_1(\tau)}(\mathbf{u})\right]_k = \operatorname{sgn}(u_k)\left(|u_k| - \xi\right)_+.$$

Duchi et al. (2008) suggest a linear time algorithm to perform projection onto the simplex that can be easily extended to projection onto the  $L_1$ -ball. Algorithm 5 summarizes the procedure.

#### Appendix B. KKT Conditions

Denote  $\mathbf{u} = \tilde{\boldsymbol{\beta}}_j - t_j \nabla_j \ell(\tilde{\boldsymbol{\beta}}_j | \tilde{\mathbf{b}}_{-j})$ . Note that

$$||\mathbf{u}||_{\infty} = \max_{k} |u_{k}| = \max_{k} |\mathbf{e}_{k}^{\top}\mathbf{u}|,$$

where  $\mathbf{e}_k = (I(j = k), 1 \le j \le K)^{\top}$ . For each individual  $|\mathbf{e}_k^{\top} \mathbf{u}|$ , we have

$$\partial |\mathbf{e}_k^\top \mathbf{u}| = \mathbf{e}_k \partial |\mathbf{e}_k^\top \mathbf{u}| = \mathbf{e}_k \cdot s_k,$$

Algorithm 5: Linear time projection of  $\mathbf{y} \in \mathbb{R}^n$  onto the  $L_1$ -ball of radius z > 0 (Duchi et al., 2008)

- 1. Consider  $\mathbf{v} = (|y_1|, ..., |y_n|)^{\top};$
- 2. Project  $\mathbf{v}$  onto the simplex:
  - (a) Initialize  $U = \{1, ..., n\}, s = 0, \rho = 0;$
  - (b) While  $U \neq \emptyset$ , do:
    - i. Pick  $k \in U$  at random;
    - ii. Partition  $U = G \cup L$ , where  $G = \{j \in U | v_j \ge v_k\}$  and  $L = U \setminus G$ ;
    - iii. Compute  $\Delta \rho = |G|$  and  $\Delta s = \sum_{j \in G} v_j$ ;
    - iv. If  $(s + \Delta s) (\rho + \Delta \rho)v_k < z$ , then set  $s \leftarrow s + \Delta s$ ,  $\rho \leftarrow \rho + \Delta \rho$  and  $U \leftarrow L$ . Otherwise, set  $U \leftarrow G \setminus \{k\}$ ;
  - (c) Set  $\theta = (s z)/\rho$ ;
  - (d) Compute the projection onto the simplex  $\mathbf{w} = (w_1, \dots, w_n)^\top$ , where  $w_i = \max(v_i \theta, 0);$
- 3. Output  $\mathbf{x} = (x_1, \dots, x_n)^{\top}$ , the projection onto the  $L_1$ -Ball, where  $x_i = w_i \cdot \operatorname{sgn}(y_i)$ .

where

$$s_k = \begin{cases} \{1\} & \mathbf{e}_k^\top \mathbf{u} > 0, \\ \{-1\} & \mathbf{e}_k^\top \mathbf{u} < 0, \\ [-1,1] & \mathbf{e}_k^\top \mathbf{u} = 0. \end{cases}$$

Thus we can obtain the sub-differential for  $||\mathbf{u}||_{\infty}$ 

$$\partial ||\mathbf{u}||_{\infty} = \operatorname{conv} \bigcup_{k \in M(\mathbf{u})} \{\mathbf{e}_k \cdot s_k\},$$

where  $M(\mathbf{u}) = \{k : |\mathbf{e}_k^\top \mathbf{u}| = ||\mathbf{u}||_{\infty}\}$  is the maximizing indices set and conv denotes the convex hull. This implies that an optimal solution needs to satisfy the condition:  $\mathbf{0} \in \nabla_j \ell(\tilde{\boldsymbol{\beta}}_j; \tilde{\mathbf{b}}_{-j}) + t_j^{-1}(\boldsymbol{\beta}_j - \tilde{\boldsymbol{\beta}}_j) + \lambda v_j \partial ||\boldsymbol{\beta}_j||_{\infty}$ , i.e.,

$$\frac{1}{\lambda v_j t_j} \left( \tilde{\boldsymbol{\beta}}_j - t_j \nabla_j \ell(\tilde{\boldsymbol{\beta}}_j; \tilde{\mathbf{b}}_{-j}) \right) - \frac{1}{\lambda v_j t_j} \boldsymbol{\beta}_j \in \operatorname{conv} \bigcup_{k \in M(\boldsymbol{\beta}_j)} \{ \mathbf{e}_k \cdot s_k \}.$$
(24)

If  $\beta_j = 0$ , then  $M(\beta_j) = \{1, ..., K\}$  resulting in a convex hull equal to the  $L_1$  unit ball formed by  $\{\mathbf{e}_k \cdot s\}_{k=1}^K$ . Thus, from (24), we require  $\|\tilde{\beta}_j - t_j \nabla_j \ell(\tilde{\beta}_j; \tilde{\mathbf{b}}_{-j})\|_1 \leq \lambda v_j t_j$ . In practice, our algorithm builds the model upwards: it will never exclude a feature from the model (i.e., by setting  $\beta_j = 0$ ) once it is already included (i.e.,  $\tilde{\beta}_j \neq 0$  for some previous iteration) so that these two inequalities will be equivalent.

For  $\beta_i \neq 0$ , we need to verify the above inclusion directly. If (24) holds, then we must have

$$\frac{1}{\lambda v_j t_j} \left( \tilde{\beta}_j^{(k)} - t_j \nabla_j \ell(\tilde{\boldsymbol{\beta}}_j; \tilde{\mathbf{b}}_{-j})^{(k)} \right) - \frac{1}{\lambda v_j t_j} \beta_j^{(k)} = 0$$

for all  $k \notin M(\beta_j)$ , i.e.,  $|\beta_j^{(k)}| \neq ||\beta_j||_{\infty}$ , while  $||t_j^{-1}\tilde{\beta}_j - \nabla_j \ell(\tilde{\beta}_j; \tilde{\mathbf{b}}_{-j}) - t_j^{-1}\beta_j||_1 = \lambda v_j$  since the convex hull must be a subset of the boundary of the  $L_1$  ball of radius  $\lambda v_j$ . These two conditions are also sufficient for (24) to hold.

#### Appendix C. Algorithm Verification

To check the validity of our algorithm, we consider the modeling under  $L_1/L_{\infty}$  regularization of simulated data with K = 5, p = 20,  $n_k = 200$  and 4 true variables in setting 1.

In Section 3.1, we have seen that the inner loop of the algorithm (the MStweedie-GPG algorithm) should feature the strict descent property. We can plot the difference in the objective function  $\ell_Q(\tilde{\beta}_0, \tilde{\beta}) - \ell_Q(\beta_0, \beta)$  and check whether this value is positive for every cycle of the MStweedie-GPG algorithm. The theoretical solution should always exhibit the descent property where a numerical solution will possibly violate that check. Figure A1 displays this verification for the current example. Except minor violations, we can see that this property is satisfied by our implementation.

The KKT conditions are at the heart of minimizing the penalized likelihood  $\ell(\beta_0, \beta) + \lambda P_{\alpha}(\beta)$ . Along the solution path, the KKT conditions in (18) should always be verified by the theoretical solution. However, a numerical solution could only approach this analytical value within certain precision and therefore may fail the KKT check. Thus, we can plot the values of these conditions for both zero and non-zero estimates and check how far they deviate from their theoretical values. Figure A2 shows these conditions for every  $j = 1, \ldots, p$  along the sequence of  $\lambda$  values. There are exactly no violations of the condition on excluded variables and the condition on included variables is never violated by a large value.

#### Descent property check



**Figure A1:** Verification of the descent property in the MStweedie-GPG algorithm with synthetic data: the difference in objective function is plotted versus the iteration number (representing one MStweedie-GPG cycle). The vertical dotted lines represent new  $\lambda$  values in the solution path.



**Figure A2:** Verification of the KKT conditions with synthetic data. The curves in each panel trace the path of the value  $||\beta_j||_1/v_j - \lambda$  for one j. In part (a), we verify the condition on non-zero estimates, i.e. variables included in the model for a given  $\lambda$ , where we expect the value to be 0. In part (b), we verify the condition on zero estimates, i.e. variables excluded from the model, where we expect the value to be below 0.

#### Appendix D. Convergence of MStweedie-GPG with Line Search

**Lemma 2.** For each  $j \in \{0, 1, ..., p\}$ ,  $\nabla_j \ell(\beta_j; \tilde{\mathbf{b}}_{-j})$  is uniformly Lipschitz continuous in the sublevel set  $\mathcal{L}_0 = \{(\beta_0, \beta) : f(\beta_0, \beta) \leq f(0, 0)\}$ , where  $f(\beta_0, \beta) = \ell(\beta_0, \beta) + \lambda P_{\alpha}(\beta)$ . In other words, there exists  $M_j \in (0, \infty)$  such that the inequality

$$\|\nabla_j \ell(\boldsymbol{\beta}_j; \tilde{\mathbf{b}}_{-j}) - \nabla_j \ell(\boldsymbol{\beta}'_j; \tilde{\mathbf{b}}_{-j})\|_2 \le M_j \|\boldsymbol{\beta}_j - \boldsymbol{\beta}'_j\|_2$$

holds for any  $\beta_j, \beta'_j$  and  $\tilde{\mathbf{b}}_{-j}$  such that  $(\beta_j, \tilde{\mathbf{b}}_{-j}) \in \mathcal{L}_0$  and  $(\beta'_j, \tilde{\mathbf{b}}_{-j}) \in \mathcal{L}_0$ . Moreover,  $\nabla \ell(\beta_0, \beta)$  is uniformly Lipschitz continuous with constant  $M \in (0, \infty)$ , i.e., for all  $(\beta_0, \beta), \ (\beta'_0, \beta') \in \mathcal{L}_0$ ,

$$\|\nabla \ell(\boldsymbol{\beta}_0, \boldsymbol{\beta}) - \nabla \ell(\boldsymbol{\beta}_0', \boldsymbol{\beta}')\|_2 \le M \|(\boldsymbol{\beta}_0, \boldsymbol{\beta}) - (\boldsymbol{\beta}_0', \boldsymbol{\beta}')\|_2.$$

#### **Proof of Lemma 2**

*Proof.* As will be shown in the proof of Theorem 1, the MStweedie-GPG algorithm is descending along its iterations and we can thus restrict the domain of  $(\boldsymbol{\beta}_0, \boldsymbol{\beta})$  to the sublevel set  $\mathcal{L}_0$ . Without loss of generality, assume not all  $y_i^{(k)}$ 's are zero. Define  $\eta_i^{(k)} = \beta_0^{(k)} + \mathbf{x}_i^{(k)\top} \boldsymbol{\beta}^{(k)}, i = 1, \dots, n_k, k = 1, \dots, K$ . It follows that the set

$$\mathcal{C}_0 = \{ \boldsymbol{\eta} = (\eta_i^{(k)}, 1 \le i \le n_k, 1 \le k \le K) : (\boldsymbol{\beta}_0, \boldsymbol{\beta}) \in \mathcal{L}_0 \}$$

is convex compact. Therefore, for all  $(\beta_0, \beta) \in \mathcal{L}_0, \eta_i^{(k)}$  is bounded by  $\eta_{\max}$ , where

$$\eta_{\max} = \max_{1 \le i \le n_k, 1 \le k \le K} \sup_{(\beta_0, \beta) \in \mathcal{L}_0} |\eta_i^{(k)}| < \infty.$$

Also,  $\boldsymbol{w}_i^{(k)}$  and  $\boldsymbol{y}_i^{(k)}$  are bounded, respectively, by

$$w_{\max} = \max_{1 \le i \le n_k, 1 \le k \le K} w_i^{(k)}$$
 and  $y_{\max} = \max_{1 \le i \le n_k, 1 \le k \le K} y_i^{(k)}$ .

Let

$$\overline{w}_i^{(k)} = w_i^{(k)} \big( (\rho - 1) y_i^{(k)} e^{(1-\rho)\eta_i^{(k)}} + (2-\rho) e^{(2-\rho)\eta_i^{(k)}} \big).$$

Note that  $\overline{w}_i^{(k)}$  is bounded by

$$\max_{1 \le i \le n_k, 1 \le k \le K} \sup_{(\beta_0, \beta) \in \mathcal{L}_0} |\overline{w}_i^{(k)}| \le w_{\max} (y_{\max}(\rho - 1)e^{(\rho - 1)\eta_{\max}} + (2 - \rho)e^{(2 - \rho)\eta_{\max}}) \equiv C.$$

Let  $M_j = C \max_{1 \le k \le K} ||X_j^{(k)}||_2^2$ . We can see that

$$\nabla_{j}^{2}\ell(\boldsymbol{\beta}_{j};\tilde{\mathbf{b}}_{-j}) = \frac{\partial^{2}}{\partial\boldsymbol{\beta}_{j}\partial\boldsymbol{\beta}_{j}^{\top}}\ell(\boldsymbol{\beta}_{j};\tilde{\mathbf{b}}_{-j})$$
  
= diag $\left(X_{j}^{(k)\top}[\operatorname{diag}(\overline{w}_{1}^{(k)},\ldots,\overline{w}_{n_{k}}^{(k)})]X_{j}^{(k)},k=1,\ldots,K\right)$   
 $\leq M_{j}\mathbf{I}_{K}, \quad \forall(\boldsymbol{\beta}_{j};\tilde{\mathbf{b}}_{-j}) \in \mathcal{L}_{0}.$ 

It follows from the mean-value theorem that  $\nabla_j \ell(\beta_j; \tilde{\mathbf{b}}_{-j})$  is uniformly Lipschitz continuous on the sublevel set  $\mathcal{L}_0$ . Indeed, the inequality

$$\|\nabla_j \ell(\boldsymbol{\beta}_j; \tilde{\mathbf{b}}_{-j}) - \nabla_j \ell(\boldsymbol{\beta}'_j; \tilde{\mathbf{b}}_{-j})\|_2 \le M_j \|\boldsymbol{\beta}_j - \boldsymbol{\beta}'_j\|_2$$

holds for any  $\beta_j, \beta'_j$  and  $\tilde{\mathbf{b}}_{-j}$  satisfying  $(\beta_j, \tilde{\mathbf{b}}_{-j}) \in \mathcal{L}_0$  and  $(\beta'_j, \tilde{\mathbf{b}}_{-j}) \in \mathcal{L}_0$ . Now let

$$M = \max_{1 \le k \le K} C\Lambda_{\max}(\hat{\mathbf{X}}^{(k)\top} \hat{\mathbf{X}}^{(k)}),$$

where  $\hat{\mathbf{X}}^{(k)} = (\mathbf{1}_{n_k}, \mathbf{X}^{(k)})$  and  $\Lambda_{\max}(\cdot)$  denotes the largest eigenvalue of the enclosed matrix. We can similarly show that  $\nabla \ell(\boldsymbol{\beta}_0, \boldsymbol{\beta})$  is uniformly Lipschitz continuous with constant M for all  $(\boldsymbol{\beta}_0, \boldsymbol{\beta}) \in \mathcal{L}_0$ .

#### **Proof of Theorem 1**

*Proof.* To simplify notation, let  $\mathbf{b} = (\boldsymbol{\beta}_0, \boldsymbol{\beta})$  such that  $\mathbf{b}_j = \boldsymbol{\beta}_j, 0 \leq j \leq p$ . Also, let  $\ell(\mathbf{b}) = \ell(\boldsymbol{\beta}_0, \boldsymbol{\beta}), h(\mathbf{b}) = \lambda P_{\alpha}(\boldsymbol{\beta})$  and  $f(\mathbf{b}) = \ell(\mathbf{b}) + h(\mathbf{b})$ . Since h is separable in  $\mathbf{b}$ , we let  $h_j(\mathbf{b}_j) = \lambda P_{\alpha,j}(\mathbf{b}_j), 0 \leq j \leq p$ . Denote by  $\nabla \ell(\mathbf{b}) = \partial \ell(\mathbf{b})/\partial \mathbf{b}$  the gradient of  $\ell$  and by  $\nabla_j \ell(\mathbf{b}) = \partial \ell(\mathbf{b})/\partial \mathbf{b}_j$  the groupwise gradient of  $\ell$ . Let  $\nabla_j^2 \ell(\mathbf{b}) = \partial^2 \ell(\mathbf{b})/(\partial \mathbf{b}_j \partial \mathbf{b}_j^{\mathsf{T}})$  be the Hessian matrix of  $\ell(\cdot)$  for group j. In Lemma 2, we have shown that  $\nabla \ell(\cdot)$  is uniformly Lipschitz continuous on the sublevel set  $\mathcal{L}_0$  with constant M and  $\nabla_j \ell(\cdot)$  is uniformly Lipschitz continuous on the sublevel set  $\mathcal{L}_0$  with constant  $M_j, 0 \leq j \leq p$ . Moreover, from (10), it can be shown that  $\overline{w}_i^{(k)}$  is lower-bounded

in the sublevel set  $\mathcal{L}_0$ . First, we have

$$\overline{w}_i^{(k)} \ge \left(\frac{\rho - 1}{2 - \rho}\right)^{3 - 2\rho} w_i^{(k)} (y_i^{(k)})^{2 - \rho} I(y_i^{(k)} > 0) + (2 - \rho) e^{-(2 - \rho)\eta_{\max}} I(y_i^{(k)} = 0) > 0$$

for all  $\mathbf{b} \in \mathcal{L}_0$  and  $1 \le i \le n_k, 1 \le k \le K$ . Let

$$w_{\min} = \min\left\{ \left(\frac{\rho - 1}{2 - \rho}\right)^{3 - 2\rho} \min_{i,k:y_i^{(k)} > 0} w_i^{(k)} (y_i^{(k)})^{2 - \rho}, \ (2 - \rho)e^{-(2 - \rho)\eta_{\max}} \right\}.$$

Then we can see that  $\overline{w}_i^{(k)} \ge w_{\min} > 0$ . Therefore

$$\nabla_j^2 \ell(\mathbf{b}) \succeq \operatorname{diag} \left( X_j^{(k)\top} [\operatorname{diag}(\overline{w}_1^{(k)}, \dots, \overline{w}_{n_k}^{(k)})] X_j^{(k)}, k = 1, \dots, K \right)$$
$$\succeq w_{\min} \operatorname{diag} \left( \|X_j^{(k)}\|_2^2, k = 1, \dots, K \right).$$

As long as none of  $\hat{\mathbf{X}}^{(k)}$ 's columns are zero (otherwise we simply remove that column and the corresponding group variable), this implies that  $\ell(\cdot)$  is groupwise strongly convex in  $\mathcal{L}_0$ .

Let  $t_j^{r+1}$  be the first step size that satisfies (13) when updating group  $\mathbf{b}_j$  in the (r+1)-st cycle of MStweedie-GPG. We claim that

$$\frac{\delta}{M_j} \le t_j^{r+1} \le t_{\max}, \ 0 \le j \le p.$$
(25)

Indeed, recall that in the line search,  $t_j$  starts with  $t_{max}$ . The search then continues by scaling  $t_j$  down with the factor  $\delta \in (0, 1)$ . Therefore, the last inequality holds in (25). Denote

$$G_{t_j}(\tilde{\mathbf{b}}) = G_{t_j}(\tilde{\boldsymbol{\beta}}_j; \tilde{\mathbf{b}}_{-j}) = \frac{\tilde{\boldsymbol{\beta}}_j - \operatorname{prox}_{\lambda v_j t_j h}(\tilde{\boldsymbol{\beta}}_j - t_j \nabla_j \ell(\tilde{\boldsymbol{\beta}}_j; \tilde{\mathbf{b}}_{-j}))}{t_j} = \frac{\tilde{\boldsymbol{\beta}}_j - \boldsymbol{\beta}_j^+}{t_j}.$$

By the definition of  $M_j$ , we can see that

$$\ell(\boldsymbol{\beta}_{j}^{+}; \tilde{\mathbf{b}}_{-j}) \leq \ell(\tilde{\mathbf{b}}) + \nabla_{j}\ell(\tilde{\boldsymbol{\beta}}_{j}; \tilde{\mathbf{b}}_{-j})^{\top}(\boldsymbol{\beta}_{j}^{+} - \tilde{\boldsymbol{\beta}}_{j}) + \frac{M_{j}}{2} \|\boldsymbol{\beta}_{j}^{+} - \tilde{\boldsymbol{\beta}}_{j}\|_{2}^{2}$$
$$= \ell(\tilde{\mathbf{b}}) - t_{j}\nabla_{j}\ell(\tilde{\boldsymbol{\beta}}_{j}; \tilde{\mathbf{b}}_{-j})^{\top}G_{t_{j}}(\tilde{\mathbf{b}}) + \frac{M_{j}t_{j}^{2}}{2} \|G_{t_{j}}(\tilde{\mathbf{b}})\|_{2}^{2}$$

holds for any  $t_j$ . Compared to (13), the above inequality implies that (13) can be satisfied by all  $t_j \in [0, M_j^{-1}]$ . Consequently, the first inequality holds in (25). Now let  $t_{\min} = \delta/(\max_{0 \le j \le p} M_j)$ , we conclude that  $t_j^{r+1} \in [t_{\min}, t_{\max}]$  for all j and r.

In the cyclic MStweedie-GPG algorithm, let  $b^r$  be the update of b after the *r*-th cycle. For notational convenience, define the following auxiliary variables

$$\mathbf{B}_{j}^{r+1} \equiv (\mathbf{b}_{0}^{r+1}, \dots, \mathbf{b}_{j-1}^{r+1}, \mathbf{b}_{j}^{r}, \mathbf{b}_{j+1}^{r}, \dots, \mathbf{b}_{p}^{r})^{\top}, j = 0, \dots, p, \mathbf{B}_{-j}^{r+1} \equiv (\mathbf{b}_{0}^{r+1}, \dots, \mathbf{b}_{j-1}^{r+1}, \mathbf{b}_{j+1}^{r}, \dots, \mathbf{b}_{p}^{r})^{\top}, j = 0, \dots, p,$$

For  $\mathbf{z} \in \mathbb{R}^{K}$ , let

$$(\mathbf{z}; \mathbf{B}_{-j}^{r+1}) \equiv (\mathbf{b}_0^{r+1}, \dots, \mathbf{b}_{j-1}^{r+1}, \mathbf{z}, \mathbf{b}_{j+1}^r, \dots, \mathbf{b}_p^r)^\top$$

Clearly we have  $\mathbf{B}_0^{r+1} = \mathbf{b}^r$  and  $\mathbf{B}_{p+1}^{r+1} = \mathbf{b}^{r+1}$ , and we have

$$\mathbf{B}_{j}^{r+1} = (\mathbf{b}_{j}^{r}; \mathbf{B}_{-j}^{r+1}), \qquad \mathbf{B}_{j+1}^{r+1} = (\mathbf{b}_{j}^{r+1}; \mathbf{B}_{-j}^{r+1}).$$

Under the new notation, (13) can be rewritten as

$$\ell(\mathbf{B}_{j+1}^{r+1}) = \ell(\mathbf{b}_{j}^{r+1}; \mathbf{B}_{-j}^{r+1}) \le \ell(\mathbf{B}_{j}^{r+1}) - t_{j}^{r+1} \nabla_{j} \ell(\mathbf{B}_{j}^{r+1})^{\top} G_{t_{j}^{r+1}}(\mathbf{B}_{j}^{r+1}) + \frac{t_{j}^{r+1}}{2} \|G_{t_{j}^{r+1}}(\mathbf{B}_{j}^{r+1})\|_{2}^{2},$$
(26)

where

$$G_{t_j^{r+1}}(\mathbf{B}_j^{r+1}) \equiv G_{t_j^{r+1}}(\mathbf{b}_j^r; \mathbf{B}_{-j}^{r+1}) = -\frac{\mathbf{b}_j^{r+1} - \mathbf{b}_j^r}{t_j^{r+1}}.$$
(27)

Next, we show that for any  $\mathbf{z} \in \mathbb{R}^{K}$ ,

$$f(\mathbf{B}_{j+1}^{r+1}) \le f(\mathbf{z}; \mathbf{B}_{-j}^{r+1}) + G_{t_j^{r+1}}(\mathbf{B}_j^{r+1})^\top (\mathbf{b}_j^r - \mathbf{z}) - \frac{t_j^{r+1}}{2} \|G_{t_j^{r+1}}(\mathbf{B}_j^{r+1})\|_2^2.$$
(28)

Let

$$\ell_{Q_j}(\mathbf{B}_{j+1}^{r+1}) = \ell_{Q_j}(\mathbf{b}_j^{r+1}; \mathbf{B}_{-j}^{r+1}) = \ell(\mathbf{B}_j^{r+1}) + \nabla_j \ell(\mathbf{B}_j^{r+1})^\top (\mathbf{b}_j^{r+1} - \mathbf{b}_j^r) + \frac{1}{2t_j^{r+1}} \|\mathbf{b}_j^{r+1} - \mathbf{b}_j^r\|_2^2.$$

The gradient of  $\ell_{Q_j}$  is

$$\nabla_{j}\ell_{Q_{j}}(\mathbf{B}_{j+1}^{r+1}) = \nabla_{j}\ell(\mathbf{B}_{j}^{r+1}) + \frac{\mathbf{b}_{j}^{r+1} - \mathbf{b}_{j}^{r}}{t_{j}} = \nabla_{j}\ell(\mathbf{B}_{j}^{r+1}) - G_{t_{j}^{r+1}}(\mathbf{B}_{j}^{r+1}).$$
(29)

By subgradient optimality condition, we have

$$\mathbf{0} \in \nabla_j \ell_{Q_j}(\mathbf{B}_{j+1}^{r+1}) + \partial h_j(\mathbf{b}_j^{r+1}),$$

thus

$$G_{t_j^{r+1}}(\mathbf{B}_j^{r+1}) - \nabla_j \ell(\mathbf{B}_j^{r+1}) \in \partial h_j(\mathbf{b}_j^{r+1}).$$
(30)

Now by convexity of  $\ell$ 

$$\ell(\mathbf{z}; \mathbf{B}_{-j}^{r+1}) \ge \ell(\mathbf{b}_j^r; \mathbf{B}_{-j}^{r+1}) + \nabla_j \ell(\mathbf{B}_j^{r+1})^\top (\mathbf{z} - \mathbf{b}_j^r),$$
(31)

and the convexity of h

$$h(\mathbf{z}; \mathbf{B}_{-j}^{r+1}) = h_j(\mathbf{z}) + \sum_{0 \le m \le p, m \ne j} h_m(\mathbf{b}_m^{r+I(m \le j)}) \ge h(\mathbf{b}_j^{r+1}; \mathbf{B}_{-j}^{r+1}) + \partial h_j(\mathbf{b}_j^{r+1})^\top (\mathbf{z} - \mathbf{b}_j^{r+1})$$
(32)

and (13), we have that for any  $\mathbf{z} \in \mathbb{R}^{K}$ ,

$$\begin{split} f(\mathbf{B}_{j+1}^{r+1}) &= f(\mathbf{b}_{j}^{r+1}; \mathbf{B}_{-j}^{r+1}) = \ell(\mathbf{b}_{j}^{r+1}; \mathbf{B}_{-j}^{r+1}) + h(\mathbf{b}_{j}^{r+1}; \mathbf{B}_{-j}^{r+1}) \\ \stackrel{(26)}{\leq} \ell(\mathbf{B}_{j}^{r+1}) - t_{j}^{r+1} \nabla_{j} \ell(\mathbf{B}_{j}^{r+1})^{\top} G_{t_{j}^{r+1}}(\mathbf{B}_{j}^{r+1}) + \frac{t_{j}^{r+1}}{2} \|G_{t_{j}^{r+1}}(\mathbf{B}_{j}^{r+1})\|_{2}^{2} + h(\mathbf{b}_{j}^{r+1}; \mathbf{B}_{-j}^{r+1}) \\ \stackrel{(31)(32)}{\leq} \ell(\mathbf{z}; \mathbf{B}_{-j}^{r+1}) + \nabla_{j} \ell(\mathbf{B}_{j}^{r+1})^{\top} (\mathbf{b}_{j}^{r} - \mathbf{z}) - t_{j}^{r+1} \nabla_{j} \ell(\mathbf{B}_{j}^{r+1})^{\top} G_{t_{j}^{r+1}}(\mathbf{B}_{j}^{r+1}) \\ &+ \frac{t_{j}^{r+1}}{2} \|G_{t_{j}^{r+1}}(\mathbf{B}_{j}^{r+1})\|_{2}^{2} + h_{j}(\mathbf{z}) + \partial h_{j}(\mathbf{b}_{j}^{r+1})^{\top} (\mathbf{b}_{j}^{r+1} - \mathbf{z}) + \sum_{0 \leq m \leq p, m \neq j} h_{m}(\mathbf{b}_{m}^{r+l(m < j)}) \\ \stackrel{(30)}{=} \ell(\mathbf{z}; \mathbf{B}_{-j}^{r+1}) + \nabla_{j} \ell(\mathbf{B}_{j}^{r+1})^{\top} (\mathbf{b}_{j}^{r} - \mathbf{z}) - t_{j}^{r+1} \nabla_{j} \ell(\mathbf{B}_{j}^{r+1})^{\top} G_{t_{j}^{r+1}}(\mathbf{B}_{j}^{r+1}) \\ &+ \frac{t_{j}^{r+1}}{2} \|G_{t_{j}^{r+1}}(\mathbf{B}_{j}^{r+1})\|_{2}^{2} + h_{j}(\mathbf{z}) + (G_{t_{j}^{r+1}}(\mathbf{B}_{j}^{r+1}) - \nabla_{j} \ell(\mathbf{B}_{j}^{r+1}))^{\top} (\mathbf{b}_{j}^{r+1} - \mathbf{z}) \\ &+ \sum_{0 \leq m \leq p, m \neq j} h_{m}(\mathbf{b}_{m}^{r+l(m < j)}) \\ &= \ell(\mathbf{z}; \mathbf{B}_{-j}^{r+1}) + h(\mathbf{z}; \mathbf{B}_{-j}^{r+1}) + \nabla_{j} \ell(\mathbf{B}_{j}^{r+1})^{\top} (\mathbf{b}_{j}^{r} - \mathbf{b}_{j}^{r+1}) - t_{j}^{r+1} \nabla_{j} \ell(\mathbf{B}_{j}^{r+1})^{\top} G_{t_{j}^{r+1}}(\mathbf{B}_{j}^{r+1}) \\ &+ \frac{t_{j}^{r+1}}{2} \|G_{t_{j}^{r+1}}(\mathbf{B}_{j}^{r+1})\|_{2}^{2} + G_{t_{j}^{r+1}}(\mathbf{B}_{j}^{r+1})^{\top} (\mathbf{b}_{j}^{r} - \mathbf{b}_{j}^{r+1}) - t_{j}^{r+1} \nabla_{j} \ell(\mathbf{B}_{j}^{r+1})^{\top} G_{t_{j}^{r+1}}(\mathbf{B}_{j}^{r+1}) \\ &+ \frac{t_{j}^{r+1}}{2} \|G_{t_{j}^{r+1}}(\mathbf{B}_{j}^{r+1})\|_{2}^{2} + G_{t_{j}^{r+1}}(\mathbf{B}_{j}^{r+1})^{\top} (\mathbf{b}_{j}^{r+1} - \mathbf{b}_{j}^{r} + \mathbf{b}_{j}^{r} - \mathbf{z}) \\ &= \ell(\mathbf{z}; \mathbf{B}_{-j}^{r+1}) + h(\mathbf{z}; \mathbf{B}_{j}^{r+1}) \|_{2}^{2} + G_{t_{j}^{r+1}}(\mathbf{B}_{j}^{r+1})^{\top} (\mathbf{b}_{j}^{r} - \mathbf{z}) - \frac{t_{j}^{r+1}}{2} \|G_{t_{j}^{r+1}}(\mathbf{B}_{j}^{r+1})\|_{2}^{2}, \end{split}$$

which proves (28).

Now taking  $\mathbf{z} = \mathbf{b}_j^r$  in (28), we have

$$f(\mathbf{B}_{j}^{r+1}) - f(\mathbf{B}_{j+1}^{r+1}) \ge \frac{t_{j}^{r+1}}{2} \|G_{t_{j}^{r+1}}(\mathbf{B}_{j}^{r+1})\|_{2}^{2} = \frac{1}{2t_{j}^{r+1}} \|\mathbf{b}_{j}^{r} - \mathbf{b}_{j}^{r+1}\|_{2}^{2} \ge \frac{1}{2t_{\max}} \|\mathbf{b}_{j}^{r} - \mathbf{b}_{j}^{r+1}\|_{2}^{2},$$

which implies that the MStweedie-GPG algorithm is descending. Moreover, we have the descent property of MStweedie-GPG over the cycles

$$f(\mathbf{b}^{r}) - f(\mathbf{b}^{r+1}) = \sum_{j=0}^{p} [f(\mathbf{B}_{j}^{r+1}) - f(\mathbf{B}_{j+1}^{r+1})] \ge (2t_{\max})^{-1} \|\mathbf{b}^{r} - \mathbf{b}^{r+1}\|_{2}^{2}.$$
 (33)

Now let  $\mathcal{X}^* := {\mathbf{b}^* \in \mathcal{L}_0 : f(\mathbf{b}^*) = \min_{\mathbf{b} \in \mathcal{L}_0} f(\mathbf{b})}$  be the optimal solution set of problem (6) and define  $d_{\mathcal{X}^*}(\mathbf{b}) := \min_{\mathbf{b}^* \in \mathcal{X}^*} ||\mathbf{b} - \mathbf{b}^*||_2$  to be the minimum distance from  $\mathbf{b}$  to  $\mathcal{X}^*$ . Let  $\mathbf{b}^{r*}$  be the point in  $\mathcal{X}^*$  such that  $||\mathbf{b}^r - \mathbf{b}^{r*}||_2 = d_{\mathcal{X}^*}(\mathbf{b}^r)$ . We also have  $f(\mathbf{b}^{r*}) = f^* := \min_{\mathbf{b} \in \mathcal{L}_0} f(\mathbf{b})$ . By the mean value theorem, there exists  $\mu \in [0, 1]$  and  $\boldsymbol{\zeta}^r = \mu \mathbf{b}^{r+1} + (1 - \mu)\mathbf{b}^{r*}$  such that

$$\ell(\mathbf{b}^{r+1}) - \ell(\mathbf{b}^{r*}) = (\nabla \ell(\boldsymbol{\zeta}^r))^{\top} (\mathbf{b}^{r+1} - \mathbf{b}^{r*}).$$

It follows that

$$\begin{split} f(\mathbf{b}^{r+1}) - f^* &= f(\mathbf{b}^{r+1}) - f(\mathbf{b}^{r*}) \\ &= \ell(\mathbf{b}^{r+1}) - \ell(\mathbf{b}^{r*}) + \sum_{j=0}^p [h_j(\mathbf{b}^{r+1}_j) - h_j(\mathbf{b}^{r*}_j)] \\ &= \sum_{j=0}^p [\nabla_j \ell(\boldsymbol{\zeta}^r)^\top (\mathbf{b}^{r+1}_j - \mathbf{b}^{r*}_j) + h_j(\mathbf{b}^{r+1}_j) - h_j(\mathbf{b}^{r*}_j)] \\ &= \sum_{j=0}^p [\nabla_j \ell(\mathbf{B}^{r+1}_j)^\top (\mathbf{b}^{r+1}_j - \mathbf{b}^{r*}_j) + h_j(\mathbf{b}^{r+1}_j) - h_j(\mathbf{b}^{r*}_j) \\ &+ (\nabla_j \ell(\boldsymbol{\zeta}^r) - \nabla_j \ell(\mathbf{B}^{r+1}_j))^\top (\mathbf{b}^{r+1}_j - \mathbf{b}^{r*}_j)]. \end{split}$$

By convexity of h, we have

$$\begin{split} \nabla_{j}\ell(\mathbf{B}_{j}^{r+1})^{\top}(\mathbf{b}_{j}^{r+1} - \mathbf{b}_{j}^{r*}) + h_{j}(\mathbf{b}_{j}^{r+1}) - h_{j}(\mathbf{b}_{j}^{r*}) \\ &\leq \nabla_{j}\ell(\mathbf{B}_{j}^{r+1})^{\top}(\mathbf{b}_{j}^{r+1} - \mathbf{b}_{j}^{r*}) - \partial h_{j}(\mathbf{b}_{j}^{r+1})^{\top}(\mathbf{b}_{j}^{r*} - \mathbf{b}_{j}^{r+1}) \\ \stackrel{(30)}{=} \nabla_{j}\ell(\mathbf{B}_{j}^{r+1})^{\top}(\mathbf{b}_{j}^{r+1} - \mathbf{b}_{j}^{r*}) - (G_{t_{j}^{r+1}}(\mathbf{B}_{j}^{r+1}) - \nabla_{j}\ell(\mathbf{B}_{j}^{r+1}))^{\top}(\mathbf{b}_{j}^{r*} - \mathbf{b}_{j}^{r+1}) \\ &= -G_{t_{j}^{r+1}}(\mathbf{B}_{j}^{r+1})(\mathbf{b}_{j}^{r*} - \mathbf{b}_{j}^{r+1}) \\ &= \frac{1}{t_{j}^{r+1}}(\mathbf{b}_{j}^{r+1} - \mathbf{b}_{j}^{r})(\mathbf{b}_{j}^{r*} - \mathbf{b}_{j}^{r} + \mathbf{b}_{j}^{r} - \mathbf{b}_{j}^{r+1}) \\ &\leq \frac{1}{t_{j}^{r+1}}[(\mathbf{b}_{j}^{r+1} - \mathbf{b}_{j}^{r})^{\top}(\mathbf{b}_{j}^{r*} - \mathbf{b}_{j}^{r}) - ||\mathbf{b}_{j}^{r+1} - \mathbf{b}_{j}^{r}||_{2}^{2}] \\ &\leq \frac{1}{2t_{j}^{r+1}}[||\mathbf{b}_{j}^{r*} - \mathbf{b}_{j}^{r}||_{2}^{2} + ||\mathbf{b}_{j}^{r} - \mathbf{b}_{j}^{r+1}||_{2}^{2}] \\ &\leq \frac{1}{2t_{\min}}[||\mathbf{b}_{j}^{r*} - \mathbf{b}_{j}^{r}||_{2}^{2} + ||\mathbf{b}_{j}^{r} - \mathbf{b}_{j}^{r+1}||_{2}^{2}]. \end{split}$$

Moreover, by the Lipschitz continuity of  $\nabla\ell(\cdot)$  and the Cauchy–Schwarz inequality, we have

$$\begin{split} & \left(\sum_{j=0}^{p} (\nabla_{j}\ell(\boldsymbol{\zeta}^{r}) - \nabla_{j}\ell(\mathbf{B}_{j}^{r+1}))^{\top}(\mathbf{b}_{j}^{r+1} - \mathbf{b}_{j}^{r*})\right)^{2} \\ & \leq \left(\sum_{j=0}^{p} \|\nabla\ell(\boldsymbol{\zeta}^{r}) - \nabla\ell(\mathbf{B}_{j}^{r+1})\|_{2}^{2}\right) \left(\sum_{j=0}^{p} \|\mathbf{b}_{j}^{r+1} - \mathbf{b}_{j}^{r*}\|_{2}^{2}\right) \\ & \leq \left(\sum_{j=0}^{p} M^{2} \|\boldsymbol{\zeta}^{r} - \mathbf{B}_{j}^{r+1}\|_{2}^{2}\right) \|\mathbf{b}^{r+1} - \mathbf{b}^{r*}\|_{2}^{2} \\ & = \left(\sum_{j=0}^{p} M^{2} \sum_{j'=0}^{p} \|\mu(\mathbf{b}_{j'}^{r+1} - \mathbf{b}_{j'}^{r}) + (1 - \mu)(\mathbf{b}_{j'}^{r*} - \mathbf{b}_{j'}^{r}) + \mathbf{b}_{j'}^{r} - \mathbf{b}_{j'}^{r+I(j'\leq j)}\|_{2}^{2}\right) \\ & \cdot 2(\|\mathbf{b}^{r+1} - \mathbf{b}^{r}\|_{2}^{2} + \|\mathbf{b}^{r*} - \mathbf{b}^{r}\|_{2}^{2}) \\ & \leq \left(2\sum_{j=0}^{p} M^{2}\|\mathbf{b}^{r+1} - \mathbf{b}^{r}\|_{2}^{2} + \|\mathbf{b}^{r*} - \mathbf{b}^{r}\|_{2}^{2}\right) \cdot 2(\|\mathbf{b}^{r+1} - \mathbf{b}^{r}\|_{2}^{2} + \|\mathbf{b}^{r*} - \mathbf{b}^{r}\|_{2}^{2}) \\ & \leq 4(p+1)M^{2}(\|\mathbf{b}^{r+1} - \mathbf{b}^{r}\|_{2}^{2} + \|\mathbf{b}^{r*} - \mathbf{b}^{r}\|_{2}^{2})^{2}. \end{split}$$

Altogether these imply

$$f(\mathbf{b}^{r+1}) - f^* \leq \sum_{j=0}^{p} \frac{1}{2t_{\min}} [\|\mathbf{b}_j^{r*} - \mathbf{b}_j^{r}\|_2^2 + \|\mathbf{b}_j^{r} - \mathbf{b}_j^{r+1}\|_2^2] + 2M\sqrt{p+1} (\|\mathbf{b}^{r+1} - \mathbf{b}^{r}\|_2^2 + \mathbf{d}_{\mathcal{X}^*}^2(\mathbf{b}^{r})) \leq \left(\frac{1}{2t_{\min}} + 2M\sqrt{p+1}\right) (\|\mathbf{b}^{r+1} - \mathbf{b}^{r}\|_2^2 + \mathbf{d}_{\mathcal{X}^*}^2(\mathbf{b}^{r})).$$
(34)

According to our algorithm,

$$\mathbf{b}_{j}^{r+1} = \underset{\mathbf{z} \in \mathbb{R}^{K}}{\arg\min} \, \ell_{Q_{j}}(\mathbf{z}; \mathbf{B}_{j}^{r+1}) + h_{j}(\mathbf{z})$$
$$= \underset{\mathbf{z} \in \mathbb{R}^{K}}{\arg\min} \, \ell(\mathbf{B}_{j}^{r+1}) + \nabla_{j} \ell(\mathbf{B}_{j}^{r+1})^{\top}(\mathbf{z} - \mathbf{b}_{j}^{r}) + \frac{1}{2t_{j}^{r+1}} \|\mathbf{z} - \mathbf{b}_{j}^{r}\|_{2}^{2} + h_{j}(\mathbf{z}).$$
(35)

By the optimality condition of  $\mathbf{b}_{j}^{r+1}$  in (35), we have

$$\mathbf{b}_{j}^{r+1} = \operatorname{prox}_{t_{j}^{r+1}h_{j}}(\mathbf{b}_{j}^{r+1} - t_{j}^{r+1}\nabla_{j}\ell_{Q_{j}}(\mathbf{b}_{j}^{r+1}; \mathbf{B}_{j}^{r+1})).$$

Now let  $c_0 = \min(1, t_{\max})$ . It follows from Lemma 4.3 of Kadkhodaie et al. (2014) that

$$\begin{split} \|\mathbf{b}_{j}^{r} - \operatorname{prox}_{h_{j}}(\mathbf{b}_{j}^{r} - \nabla_{j}\ell(\mathbf{b}^{r}))\|_{2} \\ &\leq \frac{1}{\max(1, 1/t_{j}^{r+1})} \|\mathbf{b}_{j}^{r} - \operatorname{prox}_{t_{j}^{r+1}h_{j}}(\mathbf{b}_{j}^{r} - t_{j}^{r+1}\nabla_{j}\ell(\mathbf{b}^{r}))\|_{2} \\ &= \min(1, t_{j}^{r+1}) \|\mathbf{b}_{j}^{r} - \operatorname{prox}_{t_{j}^{r+1}h_{j}}(\mathbf{b}_{j}^{r} - t_{j}^{r+1}\nabla_{j}\ell(\mathbf{b}^{r}))\|_{2} \\ &\leq c_{0} \|\mathbf{b}_{j}^{r} - \operatorname{prox}_{t_{j}^{r+1}h_{j}}(\mathbf{b}_{j}^{r} - t_{j}^{r+1}\nabla_{j}\ell(\mathbf{b}^{r})) + \mathbf{b}_{j}^{r+1} - \mathbf{b}_{j}^{r+1}\|_{2} \\ &\leq c_{0} \|\|\mathbf{b}_{j}^{r+1} - \operatorname{prox}_{t_{j}^{r+1}h_{j}}(\mathbf{b}_{j}^{r} - t_{j}^{r+1}\nabla_{j}\ell(\mathbf{b}^{r}))\|_{2} + \|\mathbf{b}_{j}^{r+1} - \mathbf{b}_{j}^{r}\|_{2} ] \\ &\leq c_{0} \|\|\mathbf{prox}_{t_{j}^{r+1}h_{j}}(\mathbf{b}_{j}^{r+1} - t_{j}^{r+1}\nabla_{j}\ell(\mathbf{b}^{r}))\|_{2} + \|\mathbf{b}_{j}^{r+1} - \mathbf{b}_{j}^{r}\|_{2} ] \\ &\leq c_{0} \|\|\mathbf{prox}_{t_{j}^{r+1}h_{j}}(\mathbf{b}_{j}^{r+1} - t_{j}^{r+1}\nabla_{j}\ell(\mathbf{b}^{r+1}))\|_{2} + \|\mathbf{b}_{j}^{r+1} - \mathbf{b}_{j}^{r}\|_{2} ] \\ &\leq c_{0} \|\|\mathbf{prox}_{t_{j}^{r+1}h_{j}}(\mathbf{b}_{j}^{r+1} - t_{j}^{r+1}\nabla_{j}\ell(\mathbf{b}^{r+1}))\|_{2} + \|\mathbf{b}_{j}^{r+1} - \mathbf{b}_{j}^{r}\|_{2} ] \\ &\leq c_{0} \|\|\mathbf{b}_{j}^{r+1} - \mathbf{b}_{j}^{r}\|_{2} + c_{0}t_{j}^{r+1}\|\nabla_{j}\ell(\mathbf{B}_{j}^{r+1}) + \frac{1}{t_{j}^{r+1}}(\mathbf{b}_{j}^{r+1} - \mathbf{b}_{j}^{r})\|_{2} \\ &\leq 3c_{0} \|\mathbf{b}_{j}^{r+1} - \mathbf{b}_{j}^{r}\|_{2} + c_{0}t_{\max}\|\nabla_{j}\ell(\mathbf{B}_{j}^{r+1}) - \nabla_{j}\ell(\mathbf{b}^{r})\|_{2} \\ &\leq 3c_{0} \|\mathbf{b}_{j}^{r+1} - \mathbf{b}_{j}^{r}\|_{2} + c_{0}t_{\max}\|\nabla_{j}\ell(\mathbf{B}_{j}^{r+1}) - \nabla_{j}\ell(\mathbf{b}^{r})\|_{2} \\ &\leq 3c_{0} \|\mathbf{b}_{j}^{r+1} - \mathbf{b}_{j}^{r}\|_{2} + c_{0}t_{\max}\|\nabla_{\ell}\ell(\mathbf{B}_{j}^{r+1}) - \nabla_{\ell}\ell(\mathbf{b}^{r})\|_{2} \\ &\leq 3c_{0} \|\mathbf{b}_{j}^{r+1} - \mathbf{b}_{j}^{r}\|_{2} + c_{0}t_{\max}\|\nabla_{\ell}\ell(\mathbf{B}_{j}^{r+1}) - \nabla_{\ell}\ell(\mathbf{b}^{r})\|_{2} \end{aligned}$$

It follows that

$$\|\mathbf{b}^{r} - \operatorname{prox}_{h}(\mathbf{b}^{r} - \nabla \ell(\mathbf{b}^{r}))\|_{2} \le (3c_{0} + c_{0}t_{\max}M\sqrt{p+1})\|\mathbf{b}^{r+1} - \mathbf{b}^{r}\|_{2}.$$
 (36)

Note that

$$\ell(\boldsymbol{\eta}) = \sum_{k=1}^{K} \sum_{i=1}^{n_k} w_i^{(k)} \left\{ -\frac{y_i^{(k)} e^{(1-\rho)\eta_i^{(k)}}}{1-\rho} + \frac{e^{(2-\rho)\eta_i^{(k)}}}{2-\rho} \right\}$$

is strongly convex in  $\eta \in C_0$  and  $\eta$  is an affine transformation of  $(\beta_0, \beta)$ , i.e.,  $\eta_i^{(k)} = \beta_0^{(k)} + \mathbf{x}_i^{(k)\top} \boldsymbol{\beta}^{(k)}$ . It follows from Zhang et al. (2013) that for any given  $\xi \ge f^* = \min_{\mathbf{b} \in \mathcal{L}_0} f(\mathbf{b})$ , there exists  $\kappa, \epsilon > 0$  such that, for all  $\mathbf{b} \in \mathcal{L}_0$  satisfying  $f(\mathbf{b}) \le \xi$  and  $\|\mathbf{b} - \operatorname{prox}_h(\mathbf{b} - \nabla \ell(\mathbf{b}))\|_2 \le \epsilon$ , we have

$$d_{\mathcal{X}^*}(\mathbf{b}) \le \kappa \|\mathbf{b} - \operatorname{prox}_h(\mathbf{b} - \nabla \ell(\mathbf{b}))\|_2.$$
(37)

From (33), we can see that

$$\sum_{i=0}^{r} \|\mathbf{b}^{i} - \mathbf{b}^{i+1}\|_{2}^{2} \le 2t_{\max} \sum_{i=0}^{r} \left[ f(\mathbf{b}^{i}) - f(\mathbf{b}^{i+1}) \right] = 2t_{\max} \left[ f(\mathbf{b}^{0}) - f(\mathbf{b}^{r+1}) \right] \le 2t_{\max} f(\mathbf{b}^{0}) < \infty,$$

then we must have  $\|\mathbf{b}^{r+1} - \mathbf{b}^r\|_2 \to 0$  as  $r \to \infty$ . Thus, it follows from (36) that as  $r \to \infty$ ,  $\|\mathbf{b}^r - \operatorname{prox}_h(\mathbf{b}^r - \nabla \ell(\mathbf{b}^r))\|_2 \to 0$ , and further by (37), this implies that  $d_{\mathcal{X}^*}(\mathbf{b}^r) \to 0$  as  $r \to \infty$ . Consequently, from (34) it follows that  $f(\mathbf{b}^r) \to f^*$ , which proves that the MStweedie-GPG algorithm converges to the global minimum. Let  $\Delta^r = f(\mathbf{b}^r) - f^*$ ,  $c_1 = \frac{1}{2t_{\min}} + 2M\sqrt{p+1}$ . By (37) and (34) again, we have for large enough r,

$$\begin{aligned} \Delta^{r+1} &= f(\mathbf{b}^{r+1}) - f^* \leq c_1 [\mathrm{d}^2_{\mathcal{X}^*}(\mathbf{b}^r) + \|\mathbf{b}^{r+1} - \mathbf{b}^r\|_2^2] \\ &\leq c_1 \kappa^2 \|\mathbf{b}^r - \mathrm{prox}_h(\mathbf{b}^r - \nabla \ell(\mathbf{b}^r))\|_2^2 + c_1 \|\mathbf{b}^{r+1} - \mathbf{b}^r\|_2^2 \\ &\leq (c_1 \kappa^2 (3c_0 + c_0 t_{\max} M \sqrt{p+1})^2 + c_1) \|\mathbf{b}^{r+1} - \mathbf{b}^r\|_2^2 \\ &\leq (c_1 \kappa^2 (3c_0 + c_0 t_{\max} M \sqrt{p+1})^2 + c_1) \cdot 2t_{\max} [f(\mathbf{b}^r) - f(\mathbf{b}^{r+1})] \\ &= 2t_{\max} (c_1 \kappa^2 (3c_0 + c_0 t_{\max} M \sqrt{p+1})^2 + c_1) (\Delta^r - \Delta^{r+1}). \end{aligned}$$

This implies that

$$\Delta^{r+1} \le \frac{c_2}{1+c_2} \Delta^r,\tag{38}$$

where  $c_2 = 2t_{\max}(c_1\kappa^2(3c_0 + c_0t_{\max}M\sqrt{p+1})^2 + c_1)$ . Let  $c_3 = c_2/(1+c_2)$ . From (38), we can see that  $f(\mathbf{b}^r)$  approaches  $f^*$  with linear rate  $O(c_3^r)$ . By (33) this further implies that  $\{\mathbf{b}^r, r \ge 0\}$  converges at least linearly.

#### Appendix E. Numerical Studies on Correlated Responses

#### **Setting 6 – Correlated responses**

In this simulation setting, we study the impact of having correlated responses on the performance of our proposed algorithm. Correlation is introduced using two compounded techniques. We consider a simultaneous setting, i.e. an observation consists of a vector of features **x** which is used to predict all K responses. When the coefficients are similar across tasks, then there will be correlation induced from the fact that the means  $\mu^{(k)} = \exp\left(\mathbf{x}\boldsymbol{\beta}^{(k)}\right)$ ,  $k = 1, \ldots, K$ , will be related. If we simply generate K Tweedie variables from these means, then the random variables will be independent, conditionally on the vector of means. To introduce additional correlation, we consider the following

setup inspired from a claim count decomposition suggested by Bermúdez and Karlis (2011). We generate K' > K independent Tweedie variables with means  $\mu^{(k)} = \exp\left(\mathbf{x}\boldsymbol{\beta}^{(k)}\right), k = 1, \dots, K'$ , for some choice of coefficients  $\boldsymbol{\beta}^{(k)}$  and produce the responses by taking a linear combination of these independent Tweedie random variables. In particular, we consider  $\tilde{Y}^{(k)}, k = 1, \dots, 6$ , the independent Tweedie random variables generating the K = 3 observed responses as follows:

$$\begin{bmatrix} Y^{(1)} \\ Y^{(2)} \\ Y^{(3)} \end{bmatrix} = \underbrace{\begin{bmatrix} 1 & 0 & 0 & 1 & 0 & 1 \\ 0 & 1 & 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 0 & 1 & 1 \end{bmatrix}}_{=:A} \begin{bmatrix} \widetilde{Y}^{(1)} & \widetilde{Y}^{(2)} & \widetilde{Y}^{(3)} & \widetilde{Y}^{(4)} & \widetilde{Y}^{(5)} & \widetilde{Y}^{(6)} \end{bmatrix}^{\top}.$$

The correlation depends on the mean of each independent Tweedie, but it is clear there will be correlation introduced in that way. Indeed, if  $\tilde{Y}^{(4)} > 0$ , then both  $Y^{(1)}$  and  $Y^{(2)}$  will be non-zero.

This construction actually has a real interpretation. Suppose the  $Y^{(k)}$  represent different aspect of a car insurance policy (1: personal injury, 2: property damage, 3: third party). Then, the random variable  $\tilde{Y}^{(4)}$  can be seen as the total claim amount that is common to personal injuries and property damages but without third party damages, while the difference between those aspects is captured by  $\tilde{Y}^{(1)}$  and  $\tilde{Y}^{(2)}$ , which are independent.

We consider three experiments under this setting. In the first two cases, we set  $\rho = 1.5$ ,  $\phi = 40$  and  $n^{(k)} = 1000$  and consider p = 50 features of which only the first 5 are truly generating the data. Each  $x_{ij}^{(k)}$  is produced from a standard normal distribution. In the experiment 6A, we consider equal contribution of the features across the sources so that the Lasso on the full dataset should be sufficient:

$$\begin{bmatrix} \boldsymbol{\beta}^{(1)} \quad \boldsymbol{\beta}^{(2)} \quad \boldsymbol{\beta}^{(3)} \quad \boldsymbol{\beta}^{(4)} \quad \boldsymbol{\beta}^{(5)} \quad \boldsymbol{\beta}^{(6)} \end{bmatrix} = \begin{bmatrix} 0.2 & 0.2 & 0.2 & 0.8 & 0.8 & 0.8 \\ 0.2 & 0.2 & 0.2 & 0.2 & 0.8 & 0.8 & 0.8 \\ 0.2 & 0.2 & 0.2 & 0.2 & 0.8 & 0.8 & 0.8 \\ -0.2 & -0.2 & -0.2 & -0.8 & -0.8 & -0.8 \\ -0.2 & -0.2 & -0.2 & -0.8 & -0.8 & -0.8 \\ -0.2 & -0.2 & -0.2 & -0.8 & -0.8 & -0.8 \\ 0_{45} \quad \boldsymbol{0}_{45} \quad \boldsymbol{0}_{45} \quad \boldsymbol{0}_{45} & \boldsymbol{0}_{45} & \boldsymbol{0}_{45} \end{bmatrix}$$

Upon generating 100 replications of the experiment, we obtain the following empirical correlation

| (a) Setting | (a) Setting 6A: Mean (standard error) |                 |                |                   |               |              |  |  |
|-------------|---------------------------------------|-----------------|----------------|-------------------|---------------|--------------|--|--|
|             | Full Lasso                            | Ind. Lasso      | $L_1/L_\infty$ | a- $L_1/L_\infty$ | $L_{1}/L_{2}$ | $a-L_1/L_2$  |  |  |
| Test dev.   | 48.13 (0.28)                          | 52.11 (0.43)    | 49.82 (0.33)   | 48.63 (0.29)      | 49.85 (0.35)  | 48.32 (0.29) |  |  |
| Size        | 11.47 (0.48)                          | 10.02 (0.35)    | 6.74 (0.20)    | 5.54 (0.13)       | 7.40 (0.21)   | 5.44 (0.10)  |  |  |
| Accuracy    | 87.06 (0.96)                          | 89.92 (0.70)    | 96.52 (0.39)   | 98.92 (0.25)      | 95.20 (0.43)  | 99.08 (0.21) |  |  |
| Precision   | 50.31 (1.82)                          | 55.07 (1.67)    | 78.96 (1.77)   | 93.26 (1.37)      | 72.76 (1.89)  | 94.02 (1.23) |  |  |
| Recall      | 100.00 (0.00)                         | 99.80 (0.20)    | 100.00 (0.00)  | 100.00 (0.00)     | 100.00 (0.00) | 99.80 (0.20) |  |  |
| (b) Setting | 6A: Mean rank                         | (nb. times best | )              |                   |               |              |  |  |
|             | Full Lasso                            | Ind. Lasso      | $L_1/L_\infty$ | a- $L_1/L_\infty$ | $L_{1}/L_{2}$ | $a-L_1/L_2$  |  |  |
| Test dev.   | 1.75 (50)                             | 5.89 (0)        | 4.39 (1)       | 2.58 (20)         | 4.28 (0)      | 2.11 (29)    |  |  |
| Size        | 5.21 (1)                              | 4.83 (3)        | 2.50 (29)      | 1.26 (83)         | 3.20 (21)     | 1.18 (87)    |  |  |
| Accuracy    | 5.20 (1)                              | 4.84 (3)        | 2.49 (30)      | 1.25 (84)         | 3.19 (22)     | 1.21 (86)    |  |  |
| Precision   | 5.20(1)                               | 4.84 (3)        | 2.49 (30)      | 1.25 (84)         | 3.19 (22)     | 1.18 (87)    |  |  |
| Recall      | 1.00 (100)                            | 1.04 (99)       | 1.00 (100)     | 1.00 (100)        | 1.00 (100)    | 1.04 (99)    |  |  |

**Table A1:** Results from Setting 6A with 100 replications. Part (a) shows the mean values of the statistics and their standard errors in parentheses. Part (b) shows the mean rank across the six models and, in parentheses, the number of times the model is best.

matrix,

$$\begin{bmatrix} 1.00 & 0.76 & 0.73 \\ 0.76 & 1.00 & 0.74 \\ 0.73 & 0.74 & 1.00 \end{bmatrix},$$

and the three sources respectively have 73.6%, 74.5% and 74.9% of zeroes.

Since the mean is exponential in the coefficients, we cannot compute the true coefficients generating the real responses so that it is impossible to produce the  $L_2$ -loss measure of performance. However, the selection accuracy measures (accuracy, precision and recall) are still relevant. The results of training and testing the usual six models are contained in Table A1. The two adaptive versions of our algorithm (especially  $a - L_1/L_2$ ) achieve test deviance values similar to that of Full Lasso, but using far less features. The accuracy and precision are therefore much better with a similar fit to the test data. This suggests that our proposal is quite robust to correlated data and can actually benefit from it.

In experiment 6B, we consider unequal contribution of the coefficients to the means of the

independent random variables:

$$\begin{bmatrix} \boldsymbol{\beta}^{(1)} \quad \boldsymbol{\beta}^{(2)} \quad \boldsymbol{\beta}^{(3)} \quad \boldsymbol{\beta}^{(4)} \quad \boldsymbol{\beta}^{(5)} \quad \boldsymbol{\beta}^{(6)} \end{bmatrix} = \begin{bmatrix} 0.5 & 0.5 & 0 & 1.5 & 3.0 & 1.5 \\ 0.5 & 0 & 0.5 & 3.0 & 1.5 & 1.5 \\ 0 & 0.5 & 0.5 & 1.5 & 1.5 & 3.0 \\ 0 & 0 & 0 & 0 & -3.0 & 0 \\ 0 & 0 & 0 & -3.0 & 0 & 0 \\ 0 & 45 & \mathbf{0}_{45} & \mathbf{0}_{45} & \mathbf{0}_{45} & \mathbf{0}_{45} \end{bmatrix}.$$

Hence, the Full Lasso should not perform well in this case, but the individual Lasso should be able to capture the differences between sources. Upon generating 100 replications of the experiment, we obtain the following empirical correlation matrix,

$$\begin{bmatrix} 1.00 & 0.83 & 0.21 \\ 0.83 & 1.00 & 0.48 \\ 0.21 & 0.48 & 1.00 \end{bmatrix},$$

and the three sources respectively have 73.0%, 69.3% and 73.3% of zeroes. The results of training and testing the usual six models are contained in Table A2. We get that all our models systematically out-performs both the Full Lasso and independent Lasso in term of test deviance. While the independent Lasso can assign different parameter values in each sources, it does not benefit from the sharing of information between sources and is thus more prone to over-fit. This is what we observe through the poor model fit to test data and larger number of variables in the model.

In experiment 6C, rather than considering a linear combination of independent random variables, we consider a product:

$$Y^{(1)} = \widetilde{Y}^{(1)}\widetilde{Y}^{(4)}\widetilde{Y}^{(6)},$$
  

$$Y^{(2)} = \widetilde{Y}^{(2)}\widetilde{Y}^{(4)}\widetilde{Y}^{(5)},$$
  

$$Y^{(3)} = \widetilde{Y}^{(3)}\widetilde{Y}^{(5)}\widetilde{Y}^{(6)}.$$

This new construction allows us to compute the true generating coefficients in each task and to produce the  $L_2$ -loss measure of performance. Indeed, we find that they are given by the sub of the

| (a) Setting 6B: Mean (standard error) |               |                   |                |                   |               |              |  |
|---------------------------------------|---------------|-------------------|----------------|-------------------|---------------|--------------|--|
|                                       | Full Lasso    | Ind. Lasso        | $L_1/L_\infty$ | a- $L_1/L_\infty$ | $L_{1}/L_{2}$ | $a-L_1/L_2$  |  |
| Test dev.                             | 39.58 (0.80)  | 43.13 (1.81)      | 36.22 (0.78)   | 31.69 (0.53)      | 36.55 (0.86)  | 31.51 (0.60) |  |
| Size                                  | 7.48 (0.23)   | 16.55 (0.49)      | 10.63 (0.37)   | 5.62 (0.13)       | 10.29 (0.38)  | 5.57 (0.13)  |  |
| Accuracy                              | 94.80 (0.44)  | 76.90 (0.98)      | 88.70 (0.73)   | 98.72 (0.27)      | 89.34 (0.76)  | 98.74 (0.25) |  |
| Precision                             | 71.28 (1.83)  | 32.72 (0.96)      | 51.53 (1.47)   | 91.96 (1.42)      | 53.85 (1.67)  | 92.23 (1.38) |  |
| Recall                                | 98.80 (0.48)  | 100.00 (0.00)     | 99.80 (0.20)   | 99.80 (0.20)      | 99.60 (0.28)  | 99.40 (0.34) |  |
| (b) Setting                           | 6B: Mean ranl | x (nb. times best | )              |                   |               |              |  |
|                                       | Full Lasso    | Ind. Lasso        | $L_1/L_\infty$ | a- $L_1/L_\infty$ | $L_{1}/L_{2}$ | $a-L_1/L_2$  |  |
| Test dev.                             | 5.10 (0)      | 5.33 (0)          | 3.67 (2)       | 1.67 (39)         | 3.76 (0)      | 1.47 (59)    |  |
| Size                                  | 2.79 (21)     | 5.86 (0)          | 4.21 (0)       | 1.32 (78)         | 3.99 (2)      | 1.25 (82)    |  |
| Accuracy                              | 2.82 (20)     | 5.86 (0)          | 4.20 (0)       | 1.28 (82)         | 4.00 (2)      | 1.25 (82)    |  |
| Precision                             | 2.82 (20)     | 5.86 (0)          | 4.21 (0)       | 1.29 (81)         | 4.01 (2)      | 1.24 (83)    |  |
| Recall                                | 1.25 (94)     | 1.00 (100)        | 1.03 (99)      | 1.04 (99)         | 1.06 (98)     | 1.11 (97)    |  |

**Table A2:** Results from Setting 6B with 100 replications. Part (a) shows the mean values of the statistics and their standard errors in parentheses. Part (b) shows the mean rank across the six models and, in parentheses, the number of times the model is best.

coefficients of the independent random variables of the product. For example,

$$\mathbb{E}\left\{Y^{(1)}\right\} = \mu^{(1)}\mu^{(4)}\mu^{(6)} = \exp\left\{\mathbf{x}\left(\boldsymbol{\beta}^{(1)} + \boldsymbol{\beta}^{(4)} + \boldsymbol{\beta}^{(6)}\right)\right\}.$$

Hence the true coefficients is the product between the matrix of independent coefficients and the structure matrix *A*:

$$\boldsymbol{\beta}^{\text{true}} = \boldsymbol{\beta} A^{\top} = \begin{bmatrix} 1 & 1 & 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 & 1 & 1 \\ 0 & -1 & -1 & -1 & -1 & -1 \\ 0 & 0 & 0 & 0 & -1 & 0 \\ 0_{45} & \mathbf{0}_{45} & \mathbf{0}_{45} & \mathbf{0}_{45} & \mathbf{0}_{45} & \mathbf{0}_{45} \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 1 & 0 \\ 0 & 1 & 1 \\ 1 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 3 & 3 & 2 \\ 3 & 2 & 3 \\ -2 & -3 & -3 \\ 0 & -1 & -1 \\ -1 & -1 & 0 \\ \mathbf{0}_{45} & \mathbf{0}_{45} & \mathbf{0}_{45} \end{bmatrix}$$

•

In this setting, pairs of responses are often zero at the same: e.g.  $\tilde{Y}^{(4)} = 0$  implies  $Y^{(1)} = 0$  and  $Y^{(2)} = 0$ . Upon generating 100 replications of the experiment with  $\phi = 10$ , we obtain the following

| (a) Setting 6C: Mean (standard error) |               |                  |                |                   |               |              |  |
|---------------------------------------|---------------|------------------|----------------|-------------------|---------------|--------------|--|
|                                       | Full Lasso    | Ind. Lasso       | $L_1/L_\infty$ | a- $L_1/L_\infty$ | $L_{1}/L_{2}$ | $a-L_1/L_2$  |  |
| Test dev.                             | 7.44 (3.86)   | 8.10 (1.29)      | 6.13 (1.84)    | 2.99 (0.88)       | 7.11 (2.00)   | 3.45 (1.24)  |  |
| Size                                  | 7.91 (0.31)   | 17.50 (0.64)     | 9.97 (0.54)    | 5.60 (0.32)       | 9.99 (0.46)   | 5.08 (0.24)  |  |
| Accuracy                              | 89.34 (0.51)  | 72.80 (1.16)     | 86.38 (0.92)   | 94.88 (0.50)      | 86.50 (0.73)  | 95.80 (0.36) |  |
| Precision                             | 54.26 (1.93)  | 28.55 (1.01)     | 49.90 (2.01)   | 83.26 (2.20)      | 48.59 (1.90)  | 86.32 (1.87) |  |
| Recall                                | 75.80 (1.49)  | 89.00 (1.25)     | 81.60 (1.63)   | 80.40 (1.63)      | 82.40 (1.74)  | 79.80 (1.65) |  |
| L2 loss                               | 4.32 (0.08)   | 5.09 (0.10)      | 4.32 (0.08)    | 2.86 (0.08)       | 4.38 (0.10)   | 2.86 (0.09)  |  |
| (b) Setting                           | 6C: Mean ranl | k (nb. times bes | t)             |                   |               |              |  |
|                                       | Full Lasso    | Ind. Lasso       | $L_1/L_\infty$ | a- $L_1/L_\infty$ | $L_{1}/L_{2}$ | $a-L_1/L_2$  |  |
| Test dev.                             | 3.68 (6)      | 5.39 (1)         | 4.21 (2)       | 1.77 (39)         | 4.27 (3)      | 1.68 (49)    |  |
| Size                                  | 3.19 (17)     | 5.66 (0)         | 3.84 (9)       | 1.69 (60)         | 3.88 (6)      | 1.42 (68)    |  |
| Accuracy                              | 3.43 (9)      | 5.68 (0)         | 3.89 (4)       | 1.52 (69)         | 3.84 (5)      | 1.33 (75)    |  |
| Precision                             | 3.42 (13)     | 5.71 (0)         | 3.93 (6)       | 1.54 (70)         | 3.82 (7)      | 1.29 (77)    |  |
| Recall                                | 2.95 (36)     | 1.32 (84)        | 1.99 (57)      | 2.16 (54)         | 1.88 (60)     | 2.13 (50)    |  |
| L2 loss                               | 4.05 (0)      | 5.42 (0)         | 4.12 (1)       | 1.60 (54)         | 4.19 (1)      | 1.62 (44)    |  |

**Table A3:** Results from Setting 6C with 100 replications. Part (a) shows the mean values of the statistics and their standard errors in parentheses. Part (b) shows the mean rank across the six models and, in parentheses, the number of times the model is best.

empirical correlation matrix,

| 1.00 | 0.35 | 0.38 |   |
|------|------|------|---|
| 0.35 | 1.00 | 0.41 | , |
| 0.38 | 0.41 | 1.00 |   |

and the three sources all have 92.9% of zeroes. Table A3 contain the results for the six models. All versions of our algorithm significantly produce better test deviance and the two adaptive versions clearly beats all other models. Also, the adaptive versions have smaller models and therefore much improved selection accuracy. Finally, the estimated coefficients by the adaptive algorithms are much closer to the truth.