

Efficient and proper Generalised Linear Models with power link functions

Vali Asimit

Bayes Business School, City, University of London, 106 Bunhill Row, EC1Y 8TZ, UK

Alexandru Badescu

Department of Mathematics and Statistics, University of Calgary, Calgary, AB T2N 1N4, Canada

Feng Zhou

University of Cambridge, Cambridge, CB2 0SR, UK and Bayes Business School, City, University of London, 106 Bunhill Row, EC1Y 8TZ, UK

Abstract

The generalised linear model is a flexible predictive model for observational data that is widely used in practice as it extends linear regression models to non-Gaussian data. In this paper we introduce the concept of a properly defined generalised linear model by requiring the conditional mean of the response variable to be properly mapped through the chosen link function and the log-likelihood function to be concave. We provide a comprehensive classification of proper generalised linear models for the Tweedie family and its popular subclasses under different link function specifications. Our main theoretical findings show that most Tweedie generalised linear models are not proper for canonical and log link functions, and identify a rich class of proper Tweedie generalised linear models with power link functions. Using self-concordant log-likelihoods and linearisation techniques, we provide novel algorithms for estimating several special cases of proper and not proper Tweedie generalised linear models with power link functions. The effectiveness of our methods is determined through an extensive numerical comparison of our estimates and those obtained using three built-in packages, MATLAB *fitglm*, R *glm2* and Python *sm.GLM* libraries, which are all implemented based on the standard Iteratively Reweighted Least Squares method. Overall, we find that our algorithms consistently outperform these benchmarks in terms of both accuracy and efficiency, the largest improvements being documented for high-dimensional settings.

Keywords: Exponential dispersion family, proper generalised linear model, Tweedie regression, power link function, self-concordance.

JEL classification: C13, C35, C44

1. Introduction

1.1. Literature Review and Main Goals

Generalised linear modelling (GLM) is a predictive model for observational data which creates a bridge between statistics and machine/statistical learning. That is, GLM provides not only statistical goodness of fit evidence (Nelder and Wedderburn, 1972; McCullagh et al., 1989; Bickel and Doksum, 2015) but also machine/statistical learning evidence such as feature/variable selection (Kuo and Mallick, 1998; Hastie et al., 2001; Fouskakis, 2012).¹

The basic GLM requires assumptions about two key quantities, the underlying parametric distribution and the choice of *link function (LF)*. The estimation procedure is based on an optimisation algorithm if the most common estimation method is chosen, i.e. *maximum likelihood estimation (MLE)*. The asymptotic theory of M-estimators requires a concave log-likelihood function, which is the ideal setting so that efficient and stable estimates are obtained; the existence and uniqueness of the MLE estimator is an essential assumption that requires some regularity conditions (Wedderburn, 1976; Mäkeläinen et al., 1981). Consequently, we introduce the concept of a *proper* GLM which requires the conditional mean of the response variable to be properly mapped through the chosen LF and for the log-likelihood function to be well-defined and concave. Since the GLM literature typically relies on exponential dispersion models (Jørgensen, 1987), *our first main goal* is to provide a classification of proper GLMs under this modelling assumption for different LF specifications. This allows the modeller to reduce the numerical issues and understand which combination of parametric family and LF would provide the best possible setting for implementation purposes. The most common LFs belong to the class of *log* or *power* functions, see e.g. McCullagh et al. (1989) and Bickel and Doksum (2015), and thus the main focus will be on these choices.

The most popular algorithms for fitting exponential dispersion GLMs are *Iteratively Reweighted Least Squares (IRLS)*, *Broyden-Fletcher-Goldfarb-Shanno (BFGS)* and *Limited-memory BFGS (L-BFGS)*. IRLS is the standard algorithm which is reasonably scalable when the number of covariates/features is smaller than the sample size. However, IRLS requires inverting the Hessian matrix at every step, which is computationally challenging in non-sparse problems when either the number of features/covariates or the sample size is small. A remedy for this is given by either BFGS or L-BFGS, where the inverse of the Hessian is approximated so that it is feasible to solve higher-dimensional GLM Regressions.² *The second main goal* of the paper is to identify viable alternative estimation algorithms to IRLS. Given that the underlying distribution of the response variable is parametrised according to an exponential dispersion family, the MLE could also be obtained via the vanilla Newton’s method, which by design is the same as IRLS if the *canonical* LF is in place; the application of Newton’s method is also known as the Fisher Scoring method in the GLM literature. Our aim is to improve this estimation method for both

¹GLMs have been successfully implemented in different research fields. For example, in actuarial science GLM applications include mortality modelling (Debón et al., 2008), default prediction (Breedon, 2016), cyber risk modelling (Eling and Wirfs, 2019), insurance pricing (DeLong et al., 2021), failure prediction modelling (van Staden et al., 2022), etc.

²The standard benchmark for high-dimensional problems is to have the number of features/covariates greater than 500, and these cases are typically implemented using penalised GLM. However, feature/variable selection is beyond the scope of this paper.

proper and not proper GLM settings, by making use of the mathematical properties of *power* LFs. For convex problems, Newton’s algorithm can be further refined if the objective function is in addition *self concordant (SC)*, i.e. a convex function whose third derivative is bounded relative to the second derivative in the interior of its domain.³ This property allows defining an augmented Newton’s method which requires a fewer number of iterations for convergence to the optimal solution, see e.g. [Boyd and Vandenberghe \(2004\)](#) or [Nesterov \(2004\)](#) for further details on SC and their fast convergence iterative methods. Since the log-likelihood associated to special cases of Tweedie GLMs (e.g. Poisson and Gamma) equipped with some particular *power* LF specifications is an SC function, we rely on this method for implementing them.⁴ For non-convex problems, which is typically the case for many exponential dispersion GLMs (e.g. Inverse Gaussian with *power* LFs), the use of standard IRLS-type algorithms leads to significant computational problems, as illustrated in the next subsection. In such cases, one could either construct bespoke optimisation algorithms designed to tackle a specific problem or rely on mainstream optimisation tools (e.g. generic interior-point methods) if the former are not available. In this paper we also aim to identify tractable solutions for non-convex GLM instances by exploring linearisation techniques, see e.g. [Boyd et al. \(2011\)](#).

1.2. Motivation and Contributions

The impact of using standard IRLS-based built-in packages on fitting not proper exponential dispersion GLMs is illustrated in the following motivational example. Specifically, using synthetic data, we compare the estimates of an Inverse Gaussian GLM based on the *log* LF, which is an example of a not proper GLM due to the non-concavity of its log-likelihood function, obtained with either MATLAB’s *fitglm* library or the non-linear optimisation solver provided by MATLAB’s *fmincon* function. Figure 1 displays box plots of the ratio between the L_1 distance (from the true value) of the estimates obtained with the latter method and those computed using MATLAB’s *fitglm* values. The results suggest that the *fmincon*-based estimation significantly outperforms the *fitglm* counterpart, especially for large size problems, which indicates that IRLS is not designed to perform well for not proper GLM settings.

To summarize, for any GLM implementation, one should not only consider a proper framework, but also construct bespoke algorithms to deal with the optimisation problem when possible. *Our contributions* address both these fundamental issues. *First*, we provide a comprehensive characterisation of proper MLE-based GLMs for a variety of exponential dispersion models, including the Tweedie family and its well-known special cases, under various LF specifications. Our main theoretical findings indicate that most of Tweedie generalised linear models are not proper for *canonical* and *log link* functions, and identify a rich class of proper Tweedie gen-

³In a GLM context, a modified version of the SC property with a different control of the third derivative has been used by [Bach \(2010\)](#) for analyzing the statistical properties of Logistic Regressions.

⁴We should note that the augmented Newton’s method for SC objective functions still requires the inverse of the Hessian matrix, but in a much lower number of iterations, which reduces the computational time. If the size of the GLM is large, then one may need compromises like those given by BFGS and L-BFGS algorithms where the inverse of the Hessian is efficiently computed, although we do not recommend this choice unless the augmented Newton’s method is overwhelmed by the size of the problem. In conclusion, the SC objective functions are expected to bring an improvement to IRLS, and large sized problems could be combined with the Hessian inverse approximations brought by L-BFGS or BFGS.

eralised linear models with power link functions. Consequently, using the Tweedie family for GLM implementation needs a careful approach, since, despite its very flexible parametrisation, the non-standard (Tweedie) models may lead to serious computational issues. *Second*, for a few standard Tweedie GLMs equipped with special cases of *power* LFs, we provide efficient and accurate bespoke algorithms for solving high-dimensional problems which cannot be properly tackled with standard IRLS-type methods. Specifically, we propose the *Newton’s method for Self-Concordant problems (NSC)* for solving Poisson and Gamma Regressions and the *Alternating Linearisation Methods (ALM)* algorithm for Inverse Gaussian Regressions. We provide a comprehensive comparison between these algorithms and those available in the standard built-in GLM libraries from various software, such as, MATLAB *fitglm*, R *glm2* and Python statsmodels *sm.GLM*. We find that our methods outperform these benchmarks in terms of both accuracy and efficiency, the largest improvements being documented for high-dimensional problems.

The remainder of the paper is organized as follows. Section 2 introduces the notion of *proper* GLMs for exponential dispersion models and reviews the LF candidates. Section 3 provides a comprehensive classification of proper Tweedie GLMs and its subclasses. Section 4 introduces the NSC and ALM algorithms for solving Poisson and Gamma, and Inverse Gaussian Regressions, respectively. The numerical comparison between these algorithms and the standard built-in libraries from MATLAB, R and Python is illustrated in Section 5. Section 6 concludes the paper.

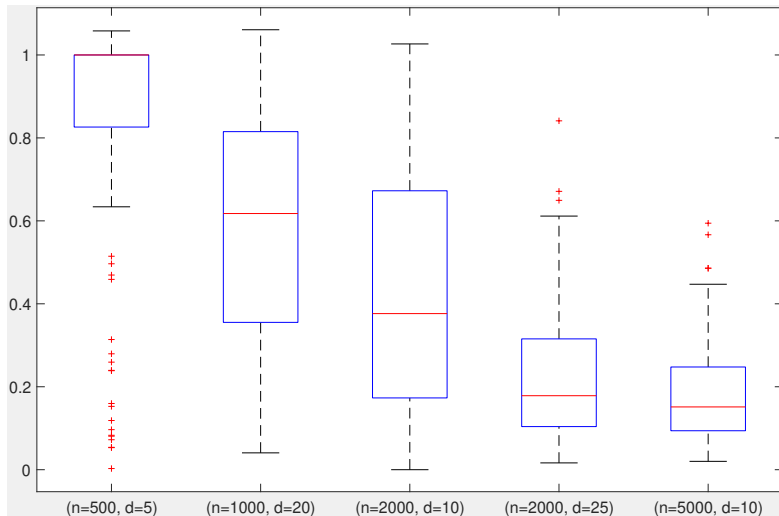


Figure 1: Box plots of MATLAB *fmincon* vs *fitglm* for Inverse Gaussian GLM

Notes: This figure shows the box plots of the ratio between the L_1 distance (from the true value) of the MLE-based GLM solutions obtained with MATLAB’s *fmincon* function and the IRLS-based GLM solution obtained with MATLAB’s *fitglm* library. Each box plot is constructed based on $N = 500$ simulations according to the DGP scheme outlined in Appendix C, for different specifications for the number of observations and the number of covariates. All GLMs are fitted with *log* LFs, i.e. a non-proper GLM.

2. Proper GLMs and LF candidates for exponential dispersion models

A univariate GLM setting assumes that the response variable Y , defined on $\mathcal{Y} \subseteq \mathbb{R}$, is explained by covariates/features \mathbf{X} defined on $\mathcal{X} \subseteq \mathbb{R}^d$. Let $\{P_{\theta, \phi} : \theta \in \Theta \subseteq \mathbb{R}, \phi \in \Phi \subseteq \mathbb{R}\}$ be the parametric set of distributions for Y , which is assumed to be an *exponential dispersion model*

characterised by the following probability density/mass function:⁵

$$\log(f_Y(y; \theta, \phi)) = \frac{\theta y - b(\theta)}{a(\phi)} + c(y, \phi). \quad (2.1)$$

Here, $a(\cdot)$, $b(\cdot)$ and $c(\cdot, \cdot)$ are real-valued functions defined on Φ , Θ and $\mathcal{Y} \times \Phi$, respectively, and ϕ is the dispersion parameter. When ϕ is fixed, (2.1) resembles an exponential family with *canonical* parameter θ . Under standard regularity conditions, the mean and variance of Y are

$$\mathbb{E}[Y] = b'(\theta) \quad \text{and} \quad \text{Var}[Y] = a(\phi)b''(\theta). \quad (2.2)$$

The GLM consists of n independent r.v.'s (observations) Y_1, \dots, Y_n with Y_i distributed according to (2.1) with parameters θ_i and ϕ , and functions $a_i(\phi)$, $b(\theta_i)$ and $c(y_i, \phi)$, and conditional mean linked through a linear predictor $\eta_i = \mathbf{x}_i^\top \boldsymbol{\beta}$ via a real-valued function h , so that

$$\mathbb{E}[Y_i | \mathbf{X}_i = \mathbf{x}_i] = h(\mathbf{x}_i^\top \boldsymbol{\beta}). \quad (2.3)$$

Here, \mathbf{x}_i is a d -dimensional vector of realized features/covariates for any $i = 1, \dots, n$.⁶

The inverse function of h , provided that it exists, is known as the *link function* (LF) and it is denoted by $g = h^{-1}$. The standard GLM literature differentiates the GLMs by the parametric choice made in (2.1) and the preferred LF g . However, from the maximum likelihood estimation (MLE) perspective, the function h is more relevant than g , and thus, the remaining results are described in terms of the former. If the dispersion parameter ϕ is known (otherwise it is estimated through the variance function from (2.2)), the MLE associated with the GLM defined in Equations (2.1) and (2.3) is obtained by solving the following non-linear optimisation problem

$$\hat{\boldsymbol{\beta}} = \arg \max_{\boldsymbol{\beta} \in \mathbb{R}^d} \ell(\boldsymbol{\beta}) = \sum_{i=1}^n \frac{\theta_i y_i - b(\theta_i)}{a_i(\phi)} \quad \text{with} \quad \theta_i = (b'^{-1} \circ h)(\mathbf{x}_i^\top \boldsymbol{\beta}). \quad (2.4)$$

Without loss of generality, we let $a_i(\phi) = a(\phi)$.⁷ The above optimisation problem is well-defined and admits a (unique) solution if the functions a , b and h satisfy certain regularity conditions. These constraints formalise the concept of a *proper GLM* and are summarised below.

Definition 2.1. *The GLM defined in Equations (2.1) and (2.3) is said to be proper if the following two conditions are satisfied:*

C1. *The conditional mean relationship from (2.3) is properly mapped, i.e. $h : \mathfrak{R} \rightarrow b'(\Theta) \subseteq$*

⁵Although the univariate assumption for the response variable Y is not essential, it simplifies the exposition.

⁶Note that although the linear predictor suggests observing d covariates/features (since $\mathcal{X} \subseteq \mathbb{R}^d$), in fact we only assume $d - 1$ covariates as we impose $x_{i,1} = 1$ for any $i = 1, \dots, n$ almost surely. This convention simplifies the notation, so that the linear predictor becomes $\eta_i = \mathbf{x}_i^\top \boldsymbol{\beta} = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_{d-1} x_{i,d-1}$.

⁷A popular choice in the GLM literature is to consider $a_i(\phi) = a(\phi)/w_i$ with $a(\phi) = \phi$ and w_i non-negative fixed weights for all $i = 1, \dots, n$. Under this assumption, the non-linear optimisation from Equation (2.4) is equivalent to solving a weighted MLE for a GLM where the response variable follows a canonical one-parameter exponential family distribution. While this could simplify the estimation of $\boldsymbol{\beta}$, and some bespoke model adequacy is typically available to check whether the predefined weights w_i are acceptable, in reality, this is more like a trial error approach which is often resolved by relying on domain knowledge.

$Conv(\mathcal{Y})$ with $b' : \Theta \rightarrow b'(\Theta)$ an injective function.⁸

C2. Assume that the likelihood function is well-defined in (2.4). The individual likelihood contribution is a (strictly) concave function, i.e.

$$\begin{cases} \text{sgn}(a(\phi)) \cdot (y \cdot (b'^{-1} \circ h)(\eta) - (b \circ b'^{-1} \circ h)(\eta)) \text{ is (strictly) concave} \\ \text{in } \eta \text{ on } \mathfrak{R} \text{ for any given } y \in \mathcal{Y}, \end{cases}$$

where sgn is the signum function.

Condition **C1** ensures that the GLM estimation is well-defined. More specifically, we require the function b' to be injective, so that it admits an inverse.⁹ Condition **C2** implies that the likelihood function ℓ defined in (2.4) is a concave function in $\eta \in \mathfrak{R}$, since the composition of a concave function with an affine mapping is concave and the sum of concave functions is also concave; in other words, (2.4) is a concave programming instance. Consequently, under the constraints from Definition 2.1, the optimisation problem in (2.4) leads to solutions which are global maximum (see e.g. [Boyd and Vandenberghe \(2004\)](#)). Note that the asymptotic distribution of $\hat{\beta}$ – like any M-estimator – requires Equation (2.4) to have a unique solution, which is not always guaranteed. However, this condition is always satisfied if the function from Condition **C2** is strictly concave.¹⁰ Note that we exclude from our analysis the cases in which $n \leq d$ or $d/n \rightarrow \delta \in (0, 1)$ with n being large. The latter case leads to potentially biased M-estimators and the asymptotic normality of such estimators may fail; see e.g. the discussions in [El Karoui et al. \(2013\)](#) and [Sur and Candès \(2019\)](#) focused on Linear and Logistic Regressions. The standard choice for solving (2.4) is to assume the function h satisfies

$$h(\eta) = b'(\eta), \quad \eta \in \mathfrak{R}. \quad (2.5)$$

Under the specification from (2.5), its equivalent LF g is known as the *canonical LF*. The sufficient conditions for a proper *canonical LF*-based GLM are summarised in the lemma below.

Lemma 2.2. *Let a GLM be equipped with its canonical LF. The MLE-based GLM is proper if $\Theta = \mathfrak{R}$ and b is strictly convex (concave) on Θ provided that $a(\phi) > 0$ ($a(\phi) < 0$) for all $\phi \in \Phi$.*

Although the *canonical LF* has useful mathematical/statistical properties, it does not always satisfy the conditions from Lemma 2.2, and therefore leads to not proper GLMs. Below, we

⁸Note that $Conv$ is the convex-hull of a set. In addition, $Conv(\mathcal{Y})$ should be read as \mathcal{Y} when Y is continuously distributed, while the convex hull operator makes a difference when Y is a discrete random variable (see e.g. Bernoulli and Poisson families).

⁹The function b' is automatically surjective since the codomain coincides with its image $b'(\Theta)$.

¹⁰The technical conditions for existence and uniqueness of the MLE estimate are well-known (see e.g. [Wedderburn \(1976\)](#) and [Mäkeläinen et al. \(1981\)](#)), and are standard in the literature, i.e. the log-likelihood function is strictly concave and some boundary conditions are satisfied. The MLE solutions could be on the boundary of the parameter space, which makes the estimation quite problematic, but we exclude such extreme cases from our analysis. For example, the latter is observed in the Logistic Regression when there exists a hyperplane that perfectly separates the ‘0’/‘1’ classes, which is also known as *complete separation*; this means that there is a continuum of points on the boundary where the absolute maximum is attained (see e.g. [Albert and Anderson \(1984\)](#)).

briefly introduce two of the most popular alternative choices in the literature, namely the *log* and *power* classes of LFs.¹¹ The *log* LF is defined by taking

$$h(\eta) = e^\eta, \quad \eta \in \mathfrak{R}. \quad (2.6)$$

Similar to the previous case, this choice may fail to produce a proper GLM in certain situations, but a general classification as in Lemma 2.2 for such models is not available. Moreover, *log* LFs has been further associated to computationally unstable MLE procedures, which leads us to considering the following family of LFs which could address some of these issues due to their appealing mathematical properties.¹² The *power* LF is defined via the following expression

$$h(\eta) = \eta^\gamma, \quad \eta \in \mathfrak{R} \quad \text{and} \quad \gamma \in \mathfrak{R}^*. \quad (2.7)$$

Popular cases of *power* LFs used in numerical applications are the *identity*, *square* and *square-root* functions which are obtained by taking $\gamma = 1, 1/2$ and 2 in (2.7), respectively. Furthermore, the *reciprocal* versions of these cases (i.e. *reciprocal identity*, *reciprocal square* and *reciprocal square-root*) are obtained by letting $\gamma = -1, -1/2$ and -2 , respectively.

Lemma 2.3 provides the sufficient conditions for **C1** to be satisfied under the choice from (2.7).¹³

Lemma 2.3. *Let a GLM with a power LF be chosen. Condition C1 in Definition 2.1 is satisfied if either of the following conditions are satisfied:*

- (i) γ is a non-zero even integer and $b(\Theta) = \mathfrak{R}_+ \subseteq \text{Conv}(\mathcal{Y})$ such that $b' : \Theta \rightarrow \mathfrak{R}_+$ is an injective mapping.
- (ii) γ is an odd integer and $b'(\Theta) = \text{Conv}(\mathcal{Y}) = \mathfrak{R}$ such that $b' : \Theta \rightarrow \mathfrak{R}$ is an injective mapping.

The above result helps us identify when a GLM is not proper due to Condition **C1** violation. For example, a direct consequence of Lemma 2.3 is that *power* LFs are not appropriate choices for GLMs where the function b' has a bounded image; this is the case of Logistic Regression (see Appendix B.2 for more details).

One way to tackle the not proper GLM issue for *power* LFs is to consider restrictions and/or modifications to these functions. For this purpose, we first introduce the class of *half-power* LFs which corresponds to taking

$$h(\eta) = \begin{cases} \eta^\gamma, & \eta > 0, \\ +\infty, & \eta \leq 0, \end{cases} \quad (2.8)$$

¹¹Note that both these functions are also *canonical* LFs for certain GLM cases. A detailed characterisation of these LFs within the context of a proper GLM is provided in Section 3 for several well-known cases of exponential dispersion models. Other classes of LFs such as *probit* and *complementary log-log* are introduced and discussed in Appendix B.3 for Logistic Regressions.

¹²Generally speaking, *power* LFs are useful for constructing convex optimisation algorithms for estimating GLMs in an accurate and efficient way. Examples of such algorithms are provided in Section 4.

¹³Note that a general characterisation for Condition **C2** cannot be provided for the *power* LF. The proofs of Lemmas 2.2 and 2.3 follow immediately from Definition 2.1.

with $\gamma \in \mathfrak{R}^*$.¹⁴ Finally, one can consider the *negative* versions of the *power/half-power* functions, called *negative power/negative half-power*, respectively, which are obtained by multiplying h from (2.7)/(2.8) by -1 .

3. Special examples of GLMs and main results

This section provides a classification of proper MLE-based GLM for a variety of exponential dispersion models and discusses the potential issues associated with the use of the different LFs introduced in Section 2. Specifically, we focus on the more general Tweedie family, together with three of its most popular special cases, namely the Poisson, Gamma and Inverse Gaussian distributions.¹⁵ A summary of proper GLMs is provided at the end of the section.

3.1. Poisson Regression – Poisson family

We assume $Y \sim \text{Poisson}(\theta)$ with probability mass function given by

$$\log(f_Y(y; \theta, \phi)) = \theta y - e^\theta - \log(y!), \quad (y, \theta, \phi) \in \mathbb{N} \times \mathfrak{R} \times \{1\}.$$

The above expression is obtained as a special case of (2.1) by taking

$$a(\phi) = \phi = 1, \quad b(\theta) = e^\theta, \quad c(y, \phi) = -\log(y!).$$

In addition, $b'(\Theta) = \mathfrak{R}_+^*$ and $b'^{-1}(\mu) = \log(\mu)$. Proposition 3.1 provides a characterisation of a proper Poisson Regression model according to our Definition 2.1.

Proposition 3.1. *Assume that $Y \sim \text{Poisson}(\theta)$. The Poisson GLM is proper if and only if $h : \mathfrak{R} \rightarrow \mathfrak{R}_+^*$, and*

$$-y \log(h(\eta)) + h(\eta) \quad \text{is convex in } \eta \text{ on } \mathfrak{R} \text{ for any given } y \in \mathbb{N}. \quad (3.1)$$

The Poisson *canonical* LF is the *log* function and this choice leads to a proper GLM due to either Lemma 2.2 or Proposition 3.1. The *power* LF does not satisfy the conditions from Proposition 3.1 unless $\gamma = 2k$ with $k \in \mathbb{N}^*$; specifically, Condition **C1** does not hold unless γ is a non-zero even integer, while Condition **C2** requires $\gamma \geq 1$. The *half-power* LF satisfies the conditions stated in Proposition 3.1 for any $\gamma \in [1, \infty)$. Thus, the simplified Poisson regression (i.e. $\phi = 1$) with a proper *half-power* LF, obtained by taking any $\gamma \geq 1$ in (2.8), leads to solving

$$\hat{\boldsymbol{\beta}} = \arg \max_{\boldsymbol{\beta} \in \mathfrak{R}^d} \ell(\boldsymbol{\beta}) = \sum_{i=1}^n \left(\gamma y_i \log(\mathbf{x}_i^\top \boldsymbol{\beta}) - (\mathbf{x}_i^\top \boldsymbol{\beta})^\gamma \right). \quad (3.2)$$

While these *half-power* LFs lead to proper GLMs that could be solved via a general convex programming algorithm, the *half-identity* and *half-square-root* cases can be solved via a compu-

¹⁴Note that the special cases for γ that we considered for the standard *power* LFs are defined in the same way for the *half-power* scenarios, e.g. we use the term *reciprocal half-square-root* for h following Equation (2.8) with $\gamma = -2$.

¹⁵In addition, the Linear and Logistic Regression models are also illustrated in Section Appendix B though we mention that only the Linear Regression is a special case of the Tweedie Regression.

tationally efficient algorithm, as outlined in Section 4.1. Finally, note that the *half-square-root* and the standard *square-root* LFs are closely related, but the latter does not satisfy (3.1) because Condition **C2** does not hold in this case. Essentially, the *half-square-root* case optimises the strictly concave instance in (3.2) on the \mathfrak{R}^d cone such that $\mathbf{x}_i^\top \boldsymbol{\beta} > 0$ for all $i = 1, \dots, n$, while the *square-root* solves a similar problem to (3.2) (where $\log(\mathbf{x}_i^\top \boldsymbol{\beta})$ is replaced by $\log|\mathbf{x}_i^\top \boldsymbol{\beta}|$) on \mathfrak{R}^d , but its objective function is not concave on the entire feasibility set, namely \mathfrak{R}^d . An analogous differentiation between the *half-identity* and *identity* LFs can be formulated as well. Finally, Condition **C1** is not satisfied for any *negative power* LF or *negative half-power* LF, which are not proper for Poisson GLM.

3.2. Gamma Regression – Gamma family

We assume $Y \sim \text{Gamma}(\theta, \phi)$ with probability distribution function given by

$$\log(f_Y(y; \theta, \phi)) = \frac{\theta y + \log(-\theta)}{\phi} + \frac{1-\phi}{\phi} \log(y) - \log\left(\phi^{\frac{1}{\phi}} \Gamma\left(\frac{1}{\phi}\right)\right), \quad (y, \theta, \phi) \in \mathfrak{R}_+^* \times \mathfrak{R}_-^* \times \mathfrak{R}_+^*.$$

The above expression is obtained as a special case of (2.1) by taking

$$a(\phi) = \phi, \quad b(\theta) = -\log(-\theta), \quad c(y, \phi) = \frac{1-\phi}{\phi} \log(y) - \log\left(\phi^{\frac{1}{\phi}} \Gamma\left(\frac{1}{\phi}\right)\right).$$

In addition, $b'(\Theta) = \mathfrak{R}_+^*$ and $b'^{-1}(\mu) = -\mu^{-1}$. Proposition 3.2 provides a characterisation of a proper Gamma Regression model according to our Definition 2.1.

Proposition 3.2. *Assume that $Y \sim \text{Gamma}(\theta, \phi)$. The Gamma GLM is proper if and only if $h: \mathfrak{R} \rightarrow \mathfrak{R}_+^*$, and*

$$\frac{y}{h(\eta)} + \log(h(\eta)) \quad \text{is convex in } \eta \text{ on } \mathfrak{R} \text{ for any given } y \in \mathfrak{R}_+^*. \quad (3.3)$$

The *canonical* LF associated to the Gamma GLM is the *reciprocal identity* function. This function does not satisfy the conditions stated in Lemma 2.2 or Proposition 3.2, since Condition **C1** does not hold, and therefore, unlike in the Poisson case, the *canonical* GLM is not proper. A popular alternative for Gamma GLM is represented by the *log* LF; this choice satisfies the conditions stated in Proposition 3.2 and is thus appropriate for Gamma GLM. As in Section 3.1, we now discuss the impact of using *power/half-power* LFs in Gamma GLM. First, a *power* LF does not satisfy the conditions from Proposition 3.2 unless $\gamma = -2k$, with $k \in \mathbb{N}^*$; specifically, Condition **C1** does not hold unless γ is a non-zero even integer, while Condition **C2** requires $\gamma \leq -1$. Second, one could find that *half-power* LFs always satisfy Condition **C1**, but Condition **C2** holds if and only if $\gamma \leq -1$, leading to proper Gamma GLM in this case. Note that the simplified Gamma GLM (i.e. $\phi = 1$) with such proper *half-power* LF is equivalent to solving

$$\hat{\boldsymbol{\beta}} = \arg \max_{\boldsymbol{\beta} \in \mathfrak{R}^d} \ell(\boldsymbol{\beta}) = \sum_{i=1}^n \left(-\gamma \log(\mathbf{x}_i^\top \boldsymbol{\beta}) - y_i (\mathbf{x}_i^\top \boldsymbol{\beta})^{-\gamma} \right), \quad (3.4)$$

where $\gamma \leq -1$. While *half-power* LFs with $\gamma \leq -1$ lead to proper GLMs that could be solved via a general convex programming algorithm, the *half-reciprocal identity* and *half-reciprocal-square-*

root cases could be solved via a computationally efficient algorithm, as outlined in Section 4.1. Finally, Condition **C1** is not satisfied for any *negative power* LF or *negative half-power* LF, which are not proper for Poisson GLM.

3.3. Inverse Gaussian Regression – Inverse Gaussian (IG) family

We assume $Y \sim IG(\theta, \phi)$ with probability distribution function given by

$$\log(f_Y(y; \theta, \phi)) = \frac{\theta y - \sqrt{2\theta}}{-1/\phi} + \frac{1}{2} \left(\log \left(\frac{\phi}{2\pi y^3} \right) - \frac{\phi}{y} \right), \quad (y, \theta, \phi) \in \mathfrak{R}_+^* \times \mathfrak{R}_+^* \times \mathfrak{R}_+^*.$$

The above function is also a special case of (2.1) where

$$a(\phi) = -\frac{1}{\phi}, \quad b(\theta) = \sqrt{2\theta}, \quad c(y, \phi) = \frac{1}{2} \left(\log \left(\frac{\phi}{2\pi y^3} \right) - \frac{\phi}{y} \right).$$

In addition, $b'(\Theta) = \mathfrak{R}_+^*$ and $b'^{-1}(\mu) = \frac{1}{2}\mu^{-2}$. Proposition 3.3 provides the characterisation of a proper Inverse Gaussian Regression model according to our Definition 2.1.

Proposition 3.3. *Assume that $Y \sim IG(\theta, \phi)$. The Inverse Gaussian GLM is proper if and only if $h : \mathfrak{R} \rightarrow \mathfrak{R}_+^*$, and*

$$\frac{y}{2h^2(\eta)} - \frac{1}{h(\eta)} \quad \text{is convex in } \eta \text{ on } \mathfrak{R} \text{ for any given } y \in \mathfrak{R}_+^* \quad (3.5)$$

The *canonical* LF for the GLM based on the *IG* distribution is the *reciprocal square* function. Similar to the Gamma scenario, this function does not satisfy the conditions stated in Lemma 2.2 or Proposition 3.3, namely Condition **C1**, and therefore, it is not a proper GLM. Under the *log* LF assumption Condition **C1** is satisfied, but Condition **C2** is violated since (3.5) does not hold. The effect of non-convexity is depicted in our motivational example from Figure 1.

As before, we also investigate the *power* and *half-power* LFs in the context of an *IG* GLM. *First*, we notice that there is no *power* LF that satisfies the conditions in Proposition 3.3; specifically, Condition **C1** does not hold unless γ is a non-zero even integer, while Condition **C2** is satisfied if and only if $\gamma \in [-1, -1/2]$. *Second*, one could find that *half-power* LFs always satisfy Condition **C1**, but Condition **C2** holds if and only if $\gamma \in [-1, -1/2]$, concluding that *half-power* LF leads to a proper GLM only in this case. Given the previous findings, running *IG* Regressions with *power* or *half-power* LFs would require a compromise. That is, the *power* LF with $\gamma = 2k$, $k \in \mathbb{Z}^*$ is the best possible choice so that constrained programming is avoided (for proper *IG* GLM with *half-power* LFs such that $\gamma \in [-1, -1/2]$ for which n linear inequality constraints are needed), which is computationally undesirable for large samples. Such choice require an efficient algorithm to solve the non-concave log-likelihood function optimisation. We show how to achieve this in Section 4.2 for the *reciprocal-square-root* LF.

3.4. Main results on Tweedie Regression – Tweedie family

In this section we focus our analysis on a more general class of GLMs based on the Tweedie family, which includes the previous distributions as special/limiting cases. As before, our main

goal is to investigate if the Tweedie distribution leads to proper GLMs. Assume that $Y \sim \text{Tweedie}(\theta, \phi)$ with probability distribution function defined below

$$\log(f_Y(y; \theta, \phi)) = \frac{\theta y - K_p(\theta)}{\phi} + \log(\mu'_\phi((-\infty, y])), \quad (y, \theta, \phi) \in \mathcal{Y} \times \Theta \times \mathfrak{R}_+^*, \quad (3.6)$$

where $\Theta \subseteq \mathfrak{R}$, μ_ϕ is a Radon measure on $\mathcal{Y} \subseteq \mathfrak{R}$ and the function K_p is given by

$$K_p(\theta) := \begin{cases} \frac{\alpha - 1}{\alpha} \left(\frac{\theta}{\alpha - 1} \right)^\alpha, & p \in (-\infty, 0] \cup (1, \infty) \setminus \{2\}, \\ e^\theta, & p = 1, \\ -\log(-\theta), & p = 2, \end{cases}$$

with $\alpha = \frac{p-2}{p-1}$. The expression from (3.6) is obtained as a special case of (2.1) by taking

$$a(\phi) = \phi, \quad b(\theta) = K_p(\theta), \quad c(y, \phi) = \log(\mu'_\phi((-\infty, y])).$$

Moreover, the Poisson, Gamma and Inverse Gaussian families are obtained as special cases by taking $p = 1$ with $\mathcal{Y} = \mathbb{N}$ and $\Theta = \mathfrak{R}$, $p = 2$ with $\mathcal{Y} = \mathfrak{R}_+^*$ and $\Theta = \mathfrak{R}_-^*$, and $p = 3$ with $\mathcal{Y} = \mathfrak{R}_+^*$ and $\Theta = \mathfrak{R}_-^*$, respectively.¹⁶

Without loss of generality we henceforth assume that $p \neq \{1, 2\}$, since these two cases have already been investigated in Sections 3.1 and 3.2. Note that one should carefully choose Θ, \mathcal{Y} and p so that $K_p(\cdot)$ is well-defined on Θ . In this section, we assume that $\Theta \in \{\mathfrak{R}, \mathfrak{R}^*, \mathfrak{R}_+^*, \mathfrak{R}_-^*\}$, and thus, the function b' is well-defined and bijective on Θ only under the three settings considered in the theorem below. Extensions to subsets of these sets are obtainable at the expense of the exposition, and for this reason, we proceed with this simplification.

We now provide a characterisation of proper Tweedie GLMs, where we exclude the previous cases investigated in Sections 3.1 and 3.2 and Appendix B.1. First, we identify in Theorem 3.4 all possible settings under which Condition C1 from Definition 2.1 is satisfied.

Theorem 3.4. *Let $Y \sim \text{Tweedie}(\theta, \phi)$ parameterised as in (3.6) with $p \in (-\infty, 0) \cup (1, 2) \cup (2, \infty)$ (or equivalently, $\alpha \in (-\infty, 2) \setminus \{0, 1\}$) such that $\mathcal{Y}, \Theta \in \{\mathfrak{R}, \mathfrak{R}^*, \mathfrak{R}_+^*, \mathfrak{R}_-^*\}$. Then, Condition C1 is only satisfied for the following settings:*

- a) $\Theta = b'(\Theta) = \mathfrak{R}_+^*$ (or \mathfrak{R}_+), $\mathcal{Y} \in \{\mathfrak{R}_+^*, \mathfrak{R}\}$ (or $\mathcal{Y} \in \{\mathfrak{R}_+, \mathfrak{R}\}$) and $1 < \alpha < 2$ (which is equivalent to $p < 0$), with $h : \mathfrak{R} \rightarrow \mathfrak{R}_+^*$ (or $h : \mathfrak{R} \rightarrow \mathfrak{R}_+$);
- b) $\Theta = \mathfrak{R}_-^*$, $b'(\Theta) = \mathfrak{R}_+^*$, $\mathcal{Y} \in \{\mathfrak{R}_+^*, \mathfrak{R}_+, \mathfrak{R}\}$ and $\alpha \in (-\infty, 1) \setminus \{0\}$ (which is equivalent to $p \in (1, \infty) \setminus \{2\}$), with $h : \mathfrak{R} \rightarrow \mathfrak{R}_+^*$;
- c) $\Theta = \mathfrak{R}$, $b'(\Theta) = \mathfrak{R}_+^*$, $\mathcal{Y} \in \{\mathfrak{R}_+^*, \mathfrak{R}_+, \mathfrak{R}^*\}$, $\alpha \in \{-2l + 1 : l \in \mathbb{N}^*\}$, with $h : \mathfrak{R} \rightarrow \mathfrak{R}_+^*$.
- d) $\Theta = \mathfrak{R}$, $b'(\Theta) = \mathfrak{R}^*$, $\mathcal{Y} \in \{\mathfrak{R}^*, \mathfrak{R}\}$, $\alpha \in \{-2l : l \in \mathbb{N}^*\}$, with $h : \mathfrak{R} \rightarrow \mathfrak{R}^*$.

Setting a) includes a pedantic reference on whether the response variable could or could not include $y = 0$, and thus, we made a difference between the cases $\Theta = \mathfrak{R}_+^*$ and $\Theta = \mathfrak{R}_+$. Note

¹⁶Other notable examples are Gaussian ($p = 0$ with $\mathcal{Y} = \Theta = \mathfrak{R}$), Compound Poisson-Gamma ($1 < p < 2$ with $\mathcal{Y} = \Theta = \mathfrak{R}_+$) and Positive stable ($p > 2$ with $\mathcal{Y} = \Theta = \mathfrak{R}_+$).

that the generic Condition **C1** in Definition 2.1 requires the range of $E[Y]$, namely $b'(\Theta)$, to be a subset of $Conv(\mathcal{Y})$, though a more practical condition would be $b'(\Theta) = Conv(\mathcal{Y})$, which we assume henceforth. Setting c) is a subcase of setting b) from the implementation point of view, since the modeller chooses the Tweedie models so that \mathcal{Y} matches the data range of values. However, our classification in Theorem 3.4 has to differentiate between models with different parameter sets Θ . The next results focus on the validity of Condition **C2** from Definition 2.1 for the above Tweedie GLM settings under the LF specifications introduced in Section 2. The *power LF* class, together with its restrictions/modifications, is investigated in Theorem 3.5 below.

Theorem 3.5. *Let $Y \sim Tweedie(\theta, \phi)$ parameterised as in (3.6) with $b'(\Theta) = \mathcal{Y}$, for which condition **C1** is satisfied. Then, Condition **C2** is not satisfied by settings a)–d), for any*

(i) *power LF, except for the following cases:*

- *setting b) with $0 < \alpha < 1$ and $\gamma = -2k$, for any $k \in \mathbb{N}^*$, with $(1 - \gamma)\alpha \leq 1$,*
- *setting b) with $\alpha < 0$ and $\gamma = 2k$, for any $k \in \mathbb{Z}^*$,*
- *setting c) and $\gamma = 2k$, for any $k \in \mathbb{Z}^*$.*

(ii) *half-power LF, except for the following cases:*

- *setting a) with $1 < \alpha < 2$ and $\frac{\alpha-1}{\alpha} \leq \gamma \leq \alpha - 1$,*
- *setting b) with $0 < \alpha < 1$ and $\frac{\alpha-1}{\alpha} \leq \gamma \leq \alpha - 1$,*
- *setting b) with $\alpha < 0$ and $\gamma \leq \alpha - 1$ or $\frac{\alpha-1}{\alpha} \leq \gamma$,*
- *setting c) with $\alpha \in \{-2l + 1 : l \in \mathbb{N}^*\}$ and $\gamma \leq \alpha - 1$ or $\frac{\alpha-1}{\alpha} \leq \gamma$.*

(iii) *negative power or negative half-power LF.*

We notice that the above results are in agreement with our previous findings. For example, one could recover our discussion from Section 3.3 on proper IG GLMs, which is a special case of Theorem 3.5 if we take $p = 3$ (or equivalently $\alpha = 1/2$), where we found that proper IG GLMs with *half-power LF* are achieved if and only if $\gamma \in [-1, -1/2]$. In addition, Theorem 3.5 provides necessary and sufficient conditions for proper GLMs under other distributional assumptions. For example, Tweedie GLMs based on Positive stable distributions (i.e. $p > 2$ or equivalently $0 < \alpha < 1$) are proper only for *power LFs* with $\gamma = -2k$, $k \in \mathbb{N}^*$, with $(1 - \gamma)\alpha \leq 1$ and *half-power LFs* with $\frac{\alpha-1}{\alpha} \leq \gamma \leq \alpha - 1$. Similarly, the Compound Poisson-Gamma GLM (i.e. $1 < p < 2$ or equivalently $\alpha < 0$) is proper only for *power LFs* with $\gamma = 2k$, $k \in \mathbb{Z}^*$ or *half-power LFs* with $\gamma \leq \alpha - 1$ or $\frac{\alpha-1}{\alpha} \leq \gamma$. A complete summary of proper Tweedie GLMs is illustrated in Table 1 of Section 3.5.

Note that if $p \in (-\infty, 0] \cup (1, \infty) \setminus \{2\}$, which is equivalent to $\alpha \in (-\infty, 2] \setminus \{0, 1\}$, then the simplified Tweedie regression (i.e. $\phi = 1$) with LF h is equivalent to solving

$$\hat{\boldsymbol{\beta}} = \arg \max_{\boldsymbol{\beta} \in \mathbb{R}^d} \ell(\boldsymbol{\beta}) = \sum_{i=1}^n \left(y_i(\alpha - 1) \left(h(\mathbf{x}_i^\top \boldsymbol{\beta}) \right)^{\frac{1}{\alpha-1}} - \frac{\alpha-1}{\alpha} \left(h(\mathbf{x}_i^\top \boldsymbol{\beta}) \right)^{\frac{\alpha}{\alpha-1}} \right). \quad (3.7)$$

A few comments on (3.7) would help understanding the issues with deploying Tweedie GLMs. *First*, one may discard Condition **C2** at the expense of losing all useful properties of the M-estimators (MLE is only a special case), such as the asymptotic distribution, which questions the asymptotic bias and variance of these estimators. If that is the case, one can only hope for the numerical optimisation to behave well, but this is possible from case to case, and one would need to perform extensive numerical implementations to check whether the optimisation algorithm shows a reasonable performance for specific choices of (α, \mathcal{Y}, h) . Such compromise is done in Algorithm 2 for solving (4.5), where $\alpha = \frac{1}{2}$ as $p = 3$, $\mathcal{Y} = \mathfrak{R}_+$, and *reciprocal square-root* LF; one could recover (4.5) from (3.7) for this particular choice of (α, \mathcal{Y}, h) . *Second*, there are other parametrisations other than the one in Algorithm 2 for which Condition **C2** is not satisfied while all other regularity conditions in Definition 2.1 hold. In these instances, one has to rely on non-convex optimisation, but more importantly, has to accept that some (possibly all) statistical properties of the MLE estimator may not hold. The modeller needs to identify stable computational methods (as in Algorithm 2) instead of assuming that the general purpose GLM solvers are indeed computationally stable. Finally, we notice that the proper GLMs identified in Theorem 3.5 (ii) require solving a constrained optimisation problem on the convex cone $\{\boldsymbol{\beta} \in \mathfrak{R}^d : \mathbf{x}_i^\top \boldsymbol{\beta} \geq 0, i = 1, \dots, n\}$. Unfortunately, this is computationally expensive for large values of n , which is a negative attribute. These optimisations could be solved via convex programming and not via off-the-shelf GLM packages that relies on IRLS which cannot be adapted when such constraints are needed.

The classification of proper Tweedie GLMs based on *canonical* and *log* LF is illustrated below.

Theorem 3.6. *Let $Y \sim \text{Tweedie}(\theta, \phi)$ parameterised as in (3.6) with $b'(\Theta) = \mathcal{Y}$, for which condition **C1** is satisfied. Then, Condition **C2** is not satisfied by settings a)–d), for any*

(i) *canonical LF.*

(ii) *log LF, except for setting b) with $\alpha < 0$ or setting c).*

Theorem 3.6 shows that there are no proper Tweedie GLMs if the *canonical* LF is chosen. In addition, we notice that the Compound Poisson-Gamma GLM is proper for any *log* LF.

3.5. Summary results

Table 1 summarises our findings discussed in Section 3 and Appendix B. *First*, we recall that the *canonical* LFs, which are the standard choices in all built-in GLM implementations (available in MATLAB, Python, R, etc.), lead to not proper Tweedie GLMs, except for the Gaussian and Poisson cases. *Second*, *log* LFs tend to have the similar limitations to *canonical* LFs for Tweedie modelling. *Third*, the *power* and *half-power* LFs allow more flexibility than *log* LFs to GLM modelling when proper GLM are sought.

4. Alternative algorithms for GLMs with power LFs

The goal of this section is to not only provide efficient methods for solving high-dimensional problems while addressing the potential numerical issues in the optimisation stage, but to also create tractable models for dealing with non-convex instances, which cannot be tackled with

Table 1: Summary of proper GLMs and violations of Conditions C1 and C2

Regression model	LF	Predictor $(\hat{y} = h(\mathbf{x}^\top \hat{\boldsymbol{\beta}}))$	Violations
Gaussian/Linear	<i>identity</i> (canonical)	$\mathbf{x}^\top \hat{\boldsymbol{\beta}}$	No
Logistic	<i>logit</i> (canonical)	$(1 + \exp(-\mathbf{x}^\top \hat{\boldsymbol{\beta}}))^{-1}$	No
	<i>probit</i>	$\Phi(\mathbf{x}^\top \hat{\boldsymbol{\beta}})$	No
	<i>complementary log-log</i>	$1 - \exp(-\exp(\mathbf{x}^\top \hat{\boldsymbol{\beta}}))$	No
Poisson	<i>log</i> (canonical)	$\exp(\mathbf{x}^\top \hat{\boldsymbol{\beta}})$	No
	<i>power</i>	$(\mathbf{x}^\top \hat{\boldsymbol{\beta}})^\gamma$	No, if $\gamma = 2k, k \in \mathbb{N}^*$
	<i>half-power</i>	$(\mathbf{x}^\top \hat{\boldsymbol{\beta}})^\gamma \cdot I_{\{\mathbf{x}^\top \hat{\boldsymbol{\beta}} > 0\}}$	No, if $\gamma \geq 1$
Gamma	<i>reciprocal identity</i> (canonical)	$(\mathbf{x}^\top \hat{\boldsymbol{\beta}})^{-1}$	C1
	<i>log</i>	$\exp(\mathbf{x}^\top \hat{\boldsymbol{\beta}})$	No
	<i>power</i>	$(\mathbf{x}^\top \hat{\boldsymbol{\beta}})^\gamma$	No, if $\gamma = -2k, k \in \mathbb{N}^*$
	<i>half-power</i>	$(\mathbf{x}^\top \hat{\boldsymbol{\beta}})^\gamma \cdot I_{\{\mathbf{x}^\top \hat{\boldsymbol{\beta}} > 0\}}$	No, if $\gamma \leq -1$
Inverse Gaussian	<i>reciprocal square</i> (canonical)	$(\mathbf{x}^\top \hat{\boldsymbol{\beta}})^{-1/2}$	C1
	<i>log</i>	$\exp(\mathbf{x}^\top \hat{\boldsymbol{\beta}})$	C2
	<i>power</i>	$(\mathbf{x}^\top \hat{\boldsymbol{\beta}})^\gamma$	C1 , if $\gamma \neq 2k, k \in \mathbb{Z}^*$, and C2 , if $\gamma \notin [-1, -1/2]$
	<i>half-power</i>	$(\mathbf{x}^\top \hat{\boldsymbol{\beta}})^\gamma \cdot I_{\{\mathbf{x}^\top \hat{\boldsymbol{\beta}} > 0\}}$	No, if $\gamma \in [-1, -1/2]$
Tweedie (except of some of the above)	<i>canonical</i>	$((1-p) \cdot \mathbf{x}^\top \hat{\boldsymbol{\beta}})^{1/(1-p)}$	see Theorem 3.6
	<i>log</i>	$\exp(\mathbf{x}^\top \hat{\boldsymbol{\beta}})$	see Theorem 3.6
special cases:	<i>power</i> or <i>negative power</i>	$(\mathbf{x}^\top \hat{\boldsymbol{\beta}})^\gamma$ or $-(\mathbf{x}^\top \hat{\boldsymbol{\beta}})^\gamma, \gamma \in \mathbb{R}^*$	see Theorem 3.5
Gaussian, Poisson and Gamma)	<i>half-power</i>	$(\mathbf{x}^\top \hat{\boldsymbol{\beta}})^\gamma \cdot I_{\{\mathbf{x}^\top \hat{\boldsymbol{\beta}} > 0\}}, \gamma \in \mathbb{R}^*$	see Theorem 3.5
	<i>negative half-power</i>	$-(\mathbf{x}^\top \hat{\boldsymbol{\beta}})^\gamma \cdot I_{\{\mathbf{x}^\top \hat{\boldsymbol{\beta}} > 0\}}, \gamma \in \mathbb{R}^*$	see Theorem 3.5

Notes: This table presents a summary of proper GLMs equipped with the LFs discussed in Section 3 and Appendix B, and the potential violations of Conditions C1 and C2 from Definition 2.1 associated with these regressions. Φ stands for the $\mathcal{N}(0, 1)$ cumulative distribution function and I_A represents the indicator function for set A .

standard built-in GLM algorithms. In this sense, we introduce the *Newton’s method for Self-Concordant problems (NSC)* for Poisson and Gamma regressions equipped with some bespoke *half-power* LFs, and the *Alternating Linearisation Method (ALM)* for solving Inverse Gaussian regressions based on the *reciprocal-square-root* LF.¹⁷

4.1. The NSC algorithm for Poisson and Gamma Regressions

The explicit structure of such *self-concordant* functions allows defining a refined Newton’s method which is generally more efficient due to a reduced number of iterations.¹⁸ First, we introduce the definition of a *self-concordant* function, which was first provided by Nesterov (2004), although a simplified version is provided in Boyd and Vandenberghe (2004), which we follow in this paper.

¹⁷This is also known as *inverse-square-root* LF, but we avoid referring to ‘inverse’ since the GLM uses the inverse of a function to identify the functional estimator h with the LF g .

¹⁸For further details on SC problems and their fast convergence iterative methods, see Boyd and Vandenberghe (2004); Nesterov (2004).

Definition 4.1. Let $f : \Omega \rightarrow \mathfrak{R}$ be a closed convex function¹⁹ where $\Omega = \text{dom}(f)$ is an open set in \mathfrak{R}^d and $f \in \mathcal{C}^3(\text{dom}(f))$. The function f is self-concordant on Ω if the function $g(t) := f(\mathbf{u} + t\mathbf{v})$ satisfies $|g'''(t)| \leq 2(g''(t))^{3/2}$ for any $t \in \text{dom}(g) \subseteq \mathfrak{R}$, $\mathbf{u} \in \text{dom}(f)$, and $\mathbf{v} \in \mathfrak{R}^d$ such that $\mathbf{u} + t\mathbf{v} \in \text{dom}(f)$.

Note that the constant 2 in Definition 4.1, see $|g'''(t)| \leq 2(g''(t))^{3/2}$, is chosen for convenience and helps to identify an explicit upper bound for the total number of iterations required by the Newton's method for SC functions. If constant 2 is replaced by M , i.e. $|g'''(t)| \leq M(g''(t))^{3/2}$, then we say that its equivalent function f is SC with constant M ; e.g., if f is SC with constant M , then it is not difficult to show that $\tilde{f}(\cdot) := \frac{M^2}{4}f(\cdot)$ is SC with constant 2.

We explore the Poisson and Gamma Regressions based on some special choices of *half-power* LFs by solving (3.2) and (3.4), since the associated negative log-likelihoods are not only convex (actually strictly convex in those two cases), but also *self-concordant*. This is illustrated in Theorem 4.2 below, where the *half-identity* and *half-square-root* LFs for Poisson Regression are explored in Theorem 4.2 a), while the *half-reciprocal identity* and *half-reciprocal-square-root* LFs for Gamma Regression are explored in Theorem 4.2 b).

Theorem 4.2. Let $\{(y_i, \mathbf{x}_i) : 1 \leq i \leq n\}$ be a sample of size n drawn from (Y, \mathbf{X}) , where $\mathbf{X} = (X_1, X_2, \dots, X_d)$ with $d \geq 1$ and define $\Omega := \bigcup_{i=1}^n \{\boldsymbol{\beta} \in \mathfrak{R}^d : \mathbf{x}_i^\top \boldsymbol{\beta} > 0\}$. The following statements hold:

- a) The MLE-based Poisson GLM equipped with the half-power LF from (2.8) with either $\gamma = 2$ (and $\gamma = 1$) is self-concordant, and it leads to an optimisation problem with a self-concordant objective function f_P (\check{f}_P) on Ω , where

$$\min_{\boldsymbol{\beta} \in \Omega} f_P(\boldsymbol{\beta}) := \sum_{i=1}^n \left(\frac{1}{2} (\mathbf{x}_i^\top \boldsymbol{\beta})^2 - y_i \log(\mathbf{x}_i^\top \boldsymbol{\beta}) \right), \quad (4.1)$$

$$\min_{\boldsymbol{\beta} \in \Omega} \check{f}_P(\boldsymbol{\beta}) := \sum_{i=1}^n \left(\mathbf{x}_i^\top \boldsymbol{\beta} - y_i \log(\mathbf{x}_i^\top \boldsymbol{\beta}) \right). \quad (4.2)$$

- b) The MLE-based Gamma GLM equipped with the half-power LF from (2.8) with $\gamma = -2$ (and $\gamma = -1$) is self-concordant, and it leads to an optimisation problem with a self-concordant objective function f_G (\check{f}_G) on Ω , where

$$\min_{\boldsymbol{\beta} \in \Omega} f_G(\boldsymbol{\beta}) := \sum_{i=1}^n \left(\frac{y_i}{2} (\mathbf{x}_i^\top \boldsymbol{\beta})^2 - \log(\mathbf{x}_i^\top \boldsymbol{\beta}) \right), \quad (4.3)$$

$$\min_{\boldsymbol{\beta} \in \Omega} \check{f}_G(\boldsymbol{\beta}) := \sum_{i=1}^n \left(y_i \cdot \mathbf{x}_i^\top \boldsymbol{\beta} - \log(\mathbf{x}_i^\top \boldsymbol{\beta}) \right). \quad (4.4)$$

As previously mentioned, the constant of an SC function does not have any impact in the actual iterative algorithm, and could change only the upper bound of the total number of steps (that

¹⁹A function $f : A \subseteq \mathfrak{R}^d \rightarrow B$ is closed convex if f is convex and closed on A , where f is closed if for any $\alpha \in \mathfrak{R}$, $\{\mathbf{x} \in \text{dom}(f) : f(\mathbf{x}) \leq \alpha\}$ is a closed set.

is in an explicit form for SC functions; for details, see the Newton's step in Algorithm 1). One may show that a tighter bound could be obtained for (4.1) and (4.2), i.e. the objective function is SC with constant M_P and \check{M}_P , respectively, where

$$M_P = \check{M}_P := 2 \max_{1 \leq i \leq n} \left\{ y_i^{-1/2} I_{\{y_i > 0\}} + I_{\{y_i = 0\}} \right\},$$

which satisfies $M_P \leq 2$. However, no tighter bound (tighter than 2) is possible for the Gamma GLMs in either (4.3) and (4.4).

Theorem 4.2 allows us to use the standard SC algorithm which is detailed in (Nesterov, 2004; Boyd and Vandenberghe, 2004), and is provided here as Algorithm 1.

Algorithm 1: Standard SC algorithm for solving (4.1) and (4.3)

Result: $\mathbf{z}^{(k^*)}$ which approximates \mathbf{z}^* , the global optimum of $\min_{\mathbf{z} \in \Omega} f(\mathbf{z})$ with $f(\cdot)$ being SC on Ω , where k^* is the termination step.

Choose $\mathbf{z}^{(0)} \in \text{dom}(f)$, $\epsilon > 0$, and $\lambda^* \in (0, \tilde{\lambda})$ where $\tilde{\lambda} = \frac{3-\sqrt{5}}{2}$;

Let $\nabla f(\cdot)$ and $\nabla^2 f(\cdot)$ be the gradient and Hessian, respectively, of f on Ω ;

Define the *step/search direction* function $\Delta(\cdot) := [\nabla^2 f(\cdot)]^{-1} \nabla f(\cdot)$ on Ω ;

Define $\lambda_f(\cdot) := \left(\nabla f(\cdot)^\top [\nabla^2 f(\cdot)]^{-1} \nabla f(\cdot) \right)^{1/2}$ on Ω ;

Step 1: Damped phase

(i) If $\lambda_f(\mathbf{z}^{(0)}) < \lambda^*$ go to Step 2;

(ii) While $\lambda_f(\mathbf{z}^{(k)}) \geq \lambda^*$ do $\mathbf{z}^{(k+1)} = \mathbf{z}^{(k)} - \frac{1}{1 + \lambda_f(\mathbf{z}^{(k)})} \Delta(\mathbf{z}^{(k)})$ for all $k \geq 0$;

Step 2: Newton (or quadratically convergence) phase

While $\lambda_f(\mathbf{z}^{(k)}) > \epsilon$ do $\mathbf{z}^{(k+1)} = \mathbf{z}^{(k)} - \Delta(\mathbf{z}^{(k)})$ for all $k \geq k_{DP}^*$, where k_{DP}^* is the termination step in Step 1.

This algorithm can be viewed as a modification of the Newton's method and consists of two phases that help reducing the number of iterations. More specifically, Step 1, called the *damped phase*, guarantees that $f(\mathbf{z}^{(k)}) - f(\mathbf{z}^{(k+1)}) \geq \omega(\lambda^*)$ and in turn, the number of iterations in Step 1, denoted by N_{DP} , is bounded with

$$N_{DP} \leq \frac{f(\mathbf{z}^{(0)}) - f(\mathbf{z}^*)}{\omega(\lambda^*)}, \quad \text{where } \omega(\lambda) := \lambda - \log(1 + \lambda) \text{ on } \mathfrak{R}_+.$$

This represents the advantage of Algorithm 1 as compared to relying only on the Newton's method, see Theorem 4.1.10 of Nesterov (2004) or Section 9.6.4 of Boyd and Vandenberghe (2004) for further details on this issue.²⁰ The total number of iterations in Step 2 is $\log_2 \log_2(1/\epsilon)$ if an accuracy of $f(\mathbf{z}^{(k^*)}) - f(\mathbf{z}^*) \leq \epsilon$ is sought. The latter bound is very small, e.g., 4.32 and 5.82 for $\epsilon = 10^{-6}$ and $\epsilon = 10^{-17}$, respectively. Note that $\epsilon = 10^{-17}$ is the MATLAB machine

²⁰More formal convergence measures for Step 1 that are compared to the equivalent convergence measures of the standard Newton's method are available in Theorems 4.1.11 and 4.1.12 of Nesterov (2004).

epsilon, which is the top end tolerance level benchmark in MATLAB.

Remark 4.3. *Inverting the Hessian is often challenging, and an alternative solution to computing the step/search direction, i.e. computing $\Delta(\mathbf{z}) := [\nabla^2 f(\mathbf{z})]^{-1} \nabla f(\mathbf{z})$ for a given \mathbf{z} , is to solve $\nabla^2 f(\mathbf{z}) \mathbf{t} = \nabla f(\mathbf{z})$ in \mathbf{t} , which is a linear system of equations. If we denote by $t_f^*(\mathbf{z})$ the latter solution, we have $\Delta(\mathbf{z}^{(k)}) = t_f^*(\mathbf{z}^{(k)})$ and*

$$\lambda_f(\mathbf{z}^{(k)}) = \sqrt{\nabla f(\mathbf{z}^{(k)})^\top [\nabla^2 f(\mathbf{z}^{(k)})]^{-1} \nabla f(\mathbf{z}^{(k)})} = \sqrt{\nabla f(\mathbf{z}^{(k)})^\top t_f^*(\mathbf{z}^{(k)})}.$$

4.2. The ALM algorithm for the Inverse Gaussian Regression

We showed in Section 3.3 that the Inverse Gaussian Regression model is not proper for any power LF. However, it is still possible to create a tractable model for this parametric family for a particular power LF. Indeed, we assume a *reciprocal-square-root* LF (i.e. power LF from (2.7) with $\gamma = -2$) which satisfies Condition C1 but not Condition C2 of Definition 2.1. This choice leads to solving the following (non-linear) optimisation problem:

$$\min_{\beta \in \Omega} f_{IG}(\beta) = \sum_{i=1}^n \left(\frac{y_i}{2} (\mathbf{x}_i^\top \beta)^4 - (\mathbf{x}_i^\top \beta)^2 \right). \quad (4.5)$$

The advantage of using the *reciprocal-square-root* LF is that (4.5) has a tractable solution via the *Alternating Linearisation Method (ALM)*, see e.g. Boyd et al. (2011) for further details. More specifically, the variable β can be split into two variables, so that the ALM reformulation of (4.5) is given by:

$$\min_{(\mathbf{z}, \mathbf{t}) \in \mathbb{R}^d \times \mathbb{R}^d} G(\mathbf{z}, \mathbf{t}) = \sum_{i=1}^n \left(\frac{y_i}{2} (\mathbf{x}_i^\top \mathbf{z})^2 (\mathbf{x}_i^\top \mathbf{t})^2 - (\mathbf{x}_i^\top \mathbf{z}) (\mathbf{x}_i^\top \mathbf{t}) \right) \quad \text{so that } \mathbf{z} = \mathbf{t}. \quad (4.6)$$

The iterative algorithm that efficiently solves (4.6) is given as Algorithm 2 and is an *Alternating Linearisation Method with backtracking (ALM-bktr)*, i.e. a bespoke ALM algorithm. This algorithm provides an approximation for β^* , which denotes a local optimum of (4.5), by generating two sequences $\{\mathbf{z}_s : s \geq 0\}$ and $\{\mathbf{t}_s : s \geq 0\}$ such that $\mathbf{z}_s \rightarrow \beta^*$ and/or $\mathbf{t}_s \rightarrow \beta^*$. The main idea is to solve a two-block variant of (4.6), which is a *convex quadratic programming (QP)* instance in \mathbf{z} for any given \mathbf{t} that could be efficiently solved, and the same holds if \mathbf{z} and \mathbf{t} are interchanged. The ALM algorithm relies on replacing the function G by their linearisation and an additional regularization factor in order to obtain an approximation to the initial objective function f_{IG} from (4.5). Thus, we define the following functions

$$\begin{aligned} H_1(\mathbf{z}, \mathbf{t}; \mu) &:= G(\mathbf{z}, \mathbf{t}) + \langle G_2(\mathbf{t}, \mathbf{t}), \mathbf{z} - \mathbf{t} \rangle + \frac{1}{2\mu} \|\mathbf{z} - \mathbf{t}\|_2^2, \\ H_2(\mathbf{z}, \mathbf{t}; \mu) &:= G(\mathbf{z}, \mathbf{t}) + \langle G_1(\mathbf{z}, \mathbf{z}), \mathbf{t} - \mathbf{z} \rangle + \frac{1}{2\mu} \|\mathbf{z} - \mathbf{t}\|_2^2, \end{aligned}$$

where $\|\cdot\|_2$ is the L^2 norm on \mathfrak{R}^d , μ is a positive constant, and G_1 and G_2 are the partial derivatives of G given below:

$$G_1(\mathbf{z}, \mathbf{t}) := \frac{\partial G}{\partial \mathbf{z}} = \sum_{i=1}^n \left(y_i \left(\mathbf{x}_i^\top \mathbf{z} \right) \left(\mathbf{x}_i^\top \mathbf{t} \right)^2 - \left(\mathbf{x}_i^\top \mathbf{t} \right) \right) \mathbf{x}_i,$$

$$G_2(\mathbf{z}, \mathbf{t}) := \frac{\partial G}{\partial \mathbf{t}} = \sum_{i=1}^n \left(y_i \left(\mathbf{x}_i^\top \mathbf{z} \right)^2 \left(\mathbf{x}_i^\top \mathbf{t} \right) - \left(\mathbf{x}_i^\top \mathbf{z} \right) \right) \mathbf{x}_i.$$

Algorithm 2 for solving (4.5), and therefore (4.6), is described below.²¹

Algorithm 2: Standard ALM algorithm for solving (4.5)

Result: $(\mathbf{z}_{s^*}, \mathbf{t}_{s^*})$ that approximates $\boldsymbol{\beta}^*$, a local optimum of (4.5), where s^* is the termination step.

Choose $\mu_{1,0} = \mu_{2,0} = \mu_0 > 0$, $b \in (0, 1)$, and $\mathbf{z}_0 = \mathbf{t}_0 \in \mathfrak{R}^d$;

for $s \in \{0, 1, \dots\}$ **do**

$\mathbf{z}_{s+1} := \arg \min_{\mathbf{z} \in \mathfrak{R}^d} H_1(\mathbf{z}, \mathbf{t}_s; \mu_{1,s});$

if $f_{IG}(\mathbf{z}_{s+1}) \leq H_1(\mathbf{z}_{s+1}, \mathbf{t}_s; \mu_{1,s})$ **then**

| choose $\mu_{1,s+1} \geq \mu_{1,s}$;

else

| find the lowest $n_{1,s} \geq 1$ such that $f_{IG}(\mathbf{u}_{1,s}) \leq H_1(\mathbf{u}_{1,s}, \mathbf{t}_s; \mu_{1,s}^*)$, where

$\mu_{1,s}^* = \mu_{1,s} b^{n_{1,s}}$ and $\mathbf{u}_{1,s} := \arg \min_{\mathbf{z} \in \mathfrak{R}^d} H_1(\mathbf{z}, \mathbf{t}_s; \mu_{1,s}^*);$

| $\mu_{1,s+1} := \mu_{1,s}^*/b$ and $\mathbf{z}_{s+1} := \mathbf{u}_{1,s}$;

end

$\mathbf{t}_{s+1} := \arg \min_{\mathbf{t} \in \mathfrak{R}^d} H_2(\mathbf{z}_{s+1}, \mathbf{t}; \mu_{2,s});$

if $f_{IG}(\mathbf{t}_{s+1}) \leq H_2(\mathbf{z}_{s+1}, \mathbf{t}_{s+1}; \mu_{2,s})$ **then**

| choose $\mu_{2,s+1} \geq \mu_{2,s}$;

else

| find the lowest $n_{2,s} \geq 1$ such that $f_{IG}(\mathbf{u}_{2,s}) \leq H_2(\mathbf{z}_{s+1}, \mathbf{u}_{2,s}; \mu_{2,s}^*)$, where

$\mu_{2,s}^* = \mu_{2,s} b^{n_{2,s}}$ and $\mathbf{u}_{2,s} := \arg \min_{\mathbf{t} \in \mathfrak{R}^d} H_2(\mathbf{z}_{s+1}, \mathbf{t}; \mu_{2,s}^*);$

| $\mu_{2,s+1} := \mu_{2,s}^*/b$ and $\mathbf{t}_{s+1} := \mathbf{u}_{2,s}$;

end

end

5. Numerical Examples and Analyses

This section presents several numerical experiments to determine the efficiency and accuracy of the proposed algorithms, and investigates to what extent they can improve the standard built-in GLM libraries from various software. Specifically, we implement the *NSC* Algorithm 1 introduced in Section 4.1 for the Poisson (with *half-square-root* LF) and Gamma Regressions

²¹The algorithm stops whenever $\sum \frac{|\mathbf{z}_{s+1} - \mathbf{t}_{s+1}|}{|\mathbf{z}_{s+1}|}$ reaches the user's defined value (e.g. the default value in our numerical examples is taken to be 10^{-4} to balance the speed and precision with other benchmark algorithms). Once the process is stopped, we use \mathbf{z}_{s+1} (or \mathbf{t}_{s+1}) if H_1 is smaller (or larger) than H_2 .

(with *half-reciprocal-square-root* LF), and the *ALM* Algorithm 2 introduced in Section 4.2 for solving Inverse Gaussian Regressions (with *reciprocal-square-root* LF).

For each specification of the number of observations n and number of covariates d , we synthetically construct N data generating processes (henceforth called DGP) and perform the above GLM estimations using both algorithms.²² The effectiveness of our methods is determined by comparing our estimates with the “true” regression coefficients β_k , for any $k = 1, \dots, N$, obtained by using three standard built-in packages: MATLAB *fitglm*, R *glm2* and Python statsmodels *sm.GLM* libraries.²³ To assess the accuracy of Algorithms 1 and 2 relative to these benchmarks we consider two performance indicators. First, we compute the Absolute Error Ratio (*AER*) and its mean (*MAER*), defined as:

$$MAER = \frac{1}{N} \sum_{k=1}^N AER_k \quad \text{with} \quad AER_k = \frac{AE(\hat{\beta}_k^{alg})}{AE(\hat{\beta}_k^{benchmark})}, \quad k = 1, \dots, N, \quad (5.1)$$

Here, the Absolute Error (AE) associated to each estimator $\hat{\beta}_k$ is defined by the L^1 -norm:

$$AE(\hat{\beta}_k) = \sum_{j=1}^d |\hat{\beta}_{k,j} - \beta_{k,j}^{true}|, \quad (5.2)$$

where $\beta_{k,j}^{true}$ is the j^{th} component of the k^{th} simulated “true” regression coefficient according to the DGP scheme outlined in Appendix C, and $\hat{\beta}_{k,j}^{alg}$ and $\hat{\beta}_{k,j}^{benchmark}$ are their corresponding estimated values obtained with Algorithms 1 and 2, and the three software benchmark packages, respectively. The performance of our approach is further evaluated by computing the log-likelihood ratio statistics, which compare the GLM with the saturated model. Thus, we introduce below the Deviance Ratio (*DR*) and its mean (*MDR*):

$$MDR = \frac{1}{N} \sum_{k=1}^N DR_k \quad \text{with} \quad DR_k = \frac{D(\hat{\beta}_k^{alg})}{D(\hat{\beta}_k^{benchmark})}, \quad k = 1, \dots, N. \quad (5.3)$$

Here, the Deviance (D) of each GLM is defined by:

$$D(\hat{\beta}_k) = -2\phi\left(\ell(\hat{\beta}_k) - \ell_s\right), \quad (5.4)$$

where $\ell(\hat{\beta}_k)$ is the log-likelihood function corresponding to the fitted GLM for the k^{th} simulated DGP scenario, while ℓ_s is the maximum value of the log-likelihood of the saturated model that is computed using the same function as in (2.4) with $\theta_i = b^{i-1}(y_i)$. Explicit expressions for the

²²Note that unlike in the theoretical presentation, d represents here the number of covariates excluding the trivial one corresponding to the intercept β_0 , so that the full matrix of explanatory variables is obtained by adding the n -dimensional unit vector to \mathbf{X} . Details on the DGP simulation are illustrated in Appendix C.

²³We remark that all three softwares rely on the Iteratively Reweighted Least Squares (IRLS) method to estimate the regression coefficients. Generally speaking, R *glm2* provides an improvement over the standard R *glm* package by using the step-halving approach in order to improve the convergence properties of IRLS (see e.g. Marschner (2011)).

deviance of all GLMs considered in our numerical experiments are provided in [Appendix D](#). Note that an *MAER* or *MDR* value smaller than 1 indicates that our approach is more accurate on average than the benchmark with respect to the corresponding performance measure. The efficiency of our algorithms relative to their benchmarks is also investigated by reporting the Mean Computational Time Ratio (*MCTR*) introduced as:

$$MCTR = \frac{1}{N} \sum_{k=1}^N CTR_k \quad \text{with} \quad CTR_k = \frac{CT(\hat{\beta}_k^{alg})}{CT(\hat{\beta}_k^{benchmark})}, \quad k = 1, \dots, N. \quad (5.5)$$

Here, $CT(\hat{\beta}_k^{alg})$ and $CT(\hat{\beta}_k^{benchmark})$ are the Algorithm 1 and benchmark computational times recorded for the k^{th} simulated DGP scenario, respectively. It follows that our algorithms are faster on average whenever $MCTR < 1$.²⁴ For a consistent and fair comparison of the computational time efficiency, all benchmarks have been implemented using their corresponding default starting values and the same specifications in the optimisation procedure, i.e. maximum number of iterations = 10,000 and tolerance level = 10^{-6} . Since Algorithms 1 and 2 are coded in MATLAB, we use the MATLAB *fitglm* starting values for our estimations.

The performance indicators *MAER* and *MDR* (both in **bold**), and *MCTR* are computed based on $N = 500$ replicates. Note that Algorithms 1 and 2 always converge within a very reasonable number of iterations, which is not the case for the three benchmarks. Therefore, the number of replicates (out of 500 simulations) for which the optimisation problem (associated to the benchmarks) do not converge within the allocated maximum number of iterations is illustrated as #NaN in our tables. Consequently, these cases are discarded from the computation of our performance indicators so that the benchmarks' performance are computed in the most advantageous possible to those benchmarks.

Table 2 presents the results for the Poisson GLM regression. We first notice that in terms of accuracy, Algorithm 1 consistently outperforms both MATLAB *fitglm* and Python *sm.GLM* libraries for all cases considered. The improvements relative to Python *sm.GLM* are quite significant with respect to both *MAER* and *MDR* with the largest augmentations being noticed for larger scale settings when the ratio between the sample size and the number of covariates/features decreases; for example, when $n/d = 5$, the improvements for both indicators are on average of around 15%, 37% and 53% for $n = 100, 500$ and $1,000$, respectively. The *MAER* and *MDR* for the MATLAB *fitglm* benchmark are closer to 1, but unlike in the previous case, there are many scenarios when the *fitglm* MLE does not converge. This typically happens for the bigger scale problems, as it is the case when $n = 1,000$ and $d = 200$ (our largest setting) where convergence was not achieved in half of the cases. Unlike the MATLAB and Python libraries, R *glm2* seems to perform very similarly to our Algorithm 1 for the Poisson GLM, the *MAER/MDR* values being typically slightly above/below 1. The *MCTR* values indicate that Algorithm 1 is always more efficient than both Python *sm.GLM* and MATLAB *fitglm*, with the largest improvements observed for small dimension settings. The smallest differences in runtime happen when $n = 1,000$ and $d = 50$, when our algorithm is five and seven times

²⁴Note that for streamline purposes we only report the *MCTR* values for Algorithm 1.

faster than the aforementioned benchmarks, respectively. However, while R *glm2* is also slower when $n = 100$ than our Algorithm 1, it becomes more efficient for larger values of n .

Table 2: MAER, MCTR and MDR for Poisson GLM

		n = 100			n = 500			n = 1,000		
		d = 5	d = 10	d = 20	d = 25	d = 50	d = 100	d = 50	d = 100	d = 200
MATLAB <i>fitglm</i>	MAER	0.9730	0.9620	0.9523	0.9685	0.9721	0.9713	0.9758	0.9782	0.9816
	MDR	0.9947	0.9935	0.9883	0.9977	0.9986	0.9970	0.9998	1.0002	1.0021
	MCTR	0.0134	0.0169	0.0272	0.0630	0.0625	0.0762	0.1446	0.1012	0.1069
	#NaN	16	32	58	37	67	182	46	87	256
Python <i>sm.GLM</i>	MAER	0.9393	0.8998	0.8431	0.9002	0.8463	0.6227	0.8838	0.8131	0.4723
	MDR	0.9177	0.8972	0.8518	0.9093	0.8553	0.6268	0.8915	0.8166	0.4721
	MCTR	0.0065	0.0082	0.0129	0.0551	0.0531	0.0340	0.2016	0.1022	0.0531
	#NaN	0	0	0	0	0	0	0	0	0
R <i>glm2</i>	MAER	0.9999	0.9967	1.0014	1.0082	1.0085	1.0161	1.0087	1.0168	1.0376
	MDR	0.9579	0.9708	0.9858	0.9832	0.9882	0.9950	0.9870	0.9911	1.0057
	MCTR	0.2553	0.2815	0.5043	1.5819	1.5695	1.0513	3.3093	2.0093	1.3328
	#NaN	0	0	0	0	0	0	0	0	0

Notes: This table reports the **Mean Absolute Error Ratio (MAER)**, **Mean Deviance Ratio (MDR)** and Mean Computational Time Ratio (*MCTR*) of Algorithm 1 from Section 4.1 relative to its benchmarks, **MATLAB *fitglm***, **Python *sm.GLM*** and **R *glm2***, for the Poisson GLM equipped with the *half-power* LF from (2.8) with $\gamma = 2$. These indicators are computed based on the MLE values obtained from $N = 500$ simulations according to the DGP scheme outlined in Appendix C, for different specifications for the number of observations n and number of covariates d . The number of replicates (out of 500 simulations) that the benchmarks cannot converge is shown as #NaN. All benchmarks are implemented using the same starting values with a maximum of 10,000 iterations and 10^{-6} tolerance level.

The Gamma GLM results are illustrated in Table 3. First we notice that Algorithm 1 consistently outperforms all benchmarks in terms of both accuracy and efficiency. Unlike in the Poisson case, our method performs significantly better than R *glm2* with respect to both accuracy indicators, with an average improvement ranging from 40% – 77% and 40% – 68% for *MAER* and *MDR*, respectively, when $n = 1,000$. We further notice a reverse situation regarding the MATLAB *fitglm* and Python *sm.glm* GLM libraries when comparing to the results from Table 2. Specifically, on the one hand, the MLE procedure from Python *sm.glm* does not converge in many instances, but when it converges, the estimates are very close to those obtained via Algorithm 1. On the other hand, despite always converging, the MATLAB *fitglm* optimisation produces *MAER* and *MDR* values which are significantly lower than 1, with the lowest values recorded when $n = 1,000$. The reported average computational times favours again our methodology; the only *MCTR* values greater than 1 are spotted for the larger scale settings for R *glm2*, which provided inaccurate estimates in all these cases.

In summary, based on our DGP for Poisson and Gamma GLMs, we can argue that overall, our Algorithm 1 provides the most accurate and efficient estimation approach relative to the three benchmarks, while R *glm2* is the second best, generally speaking being more stable than the

Table 3: MAER, MCTR and MDR for Gamma GLM

		n = 100			n = 500			n = 1,000		
		d = 5	d = 10	d = 20	d = 25	d = 50	d = 100	d = 50	d = 100	d = 200
MATLAB <i>fitglm</i>	MAER	0.9216	0.9449	0.9722	0.6554	0.7141	0.8469	0.5547	0.5734	0.7167
	MDR	0.9534	0.9511	0.9687	0.6713	0.6753	0.8061	0.5065	0.4424	0.6202
	MCTR	0.0579	0.0270	0.0404	0.2549	0.1142	0.0995	0.5530	0.1991	0.1954
	#NaN	0	0	0	0	0	0	0	0	0
Python <i>sm.GLM</i>	MAER	0.9831	0.9930	0.9989	0.9962	0.9999	1.0000	0.9932	1.0000	1.0000
	MDR	0.9953	0.9980	1.0000	0.9998	1.0000	1.0000	0.9997	1.0000	1.0000
	MCTR	0.0700	0.2049	0.2314	1.6705	0.9847	0.5505	3.7401	2.0635	0.8492
	#NaN	78	55	21	406	268	124	471	373	206
R <i>glm2</i>	MAER	0.9450	0.9608	0.9850	0.5843	0.7216	0.8928	0.4018	0.5434	0.7679
	MDR	0.9496	0.9621	0.9878	0.5887	0.6859	0.8585	0.3944	0.4643	0.6840
	MCTR	0.2945	0.5550	0.5101	6.5451	3.4493	1.6073	12.2574	5.1892	1.5737
	#NaN	0	0	0	0	0	0	0	0	0

Notes: This table reports the **Mean Absolute Error Ratio (MAER)**, **Mean Deviance Ratio (MDR)** and **Mean Computational Time Ratio (MCTR)** of Algorithm 1 from Section 4.1 relative to its benchmarks, **MATLAB *fitglm***, **Python *sm.GLM*** and **R *glm2***, for the Gamma GLM equipped with the *half-power* LF from (2.8) with $\gamma = -2$. These indicators are computed based on the MLE values obtained from $N = 500$ simulations according to the DGP scheme outlined in Appendix C, for different specifications for the number of observations n and number of covariates d . The number of replicates (out of 500 simulations) that the benchmarks cannot converge is shown as #NaN. All benchmarks are implemented using the same starting values with a maximum of 10,000 iterations and 10^{-6} tolerance level.

We next turn our attention to the implementation results of the *ALM* Algorithm 2 for solving Inverse Gaussian Regressions based on the *reciprocal-square-root* LF. The benchmark chosen in our analysis is the MATLAB *fitglm* package and we only focus on the accuracy of our methodology. Figure 2 illustrates the box plots of the MATLAB *fitglm*-based *AER* and *DR* for the same values of n and d as in the previous tables. First, we notice (in all nine cases) that the *AER* indicators are more or less symmetrically distributed around 1, with a median value smaller (but closer) to 1, suggesting that our Algorithm 2 slightly outperforms MATLAB *fitglm* relative to this performance measure. However, our method performs much better in terms of the deviance measure, as almost all *DR* values are below 1, with the most significant differences being documented for larger dimension problems and the smallest n/d ratio (i.e. $n/d = 5$). Furthermore, for each value of n , we notice a decreasing trend in the median of *DRs* as the number of covariates increases. These observations are consistent with the previous findings on Algorithm 1 regarding the significant improvements in accuracy for bigger datasets.

²⁵Note that these conclusions are drawn solely based on our DGP and a limited number of experiments, so further implementations may be needed to further investigate this problem.

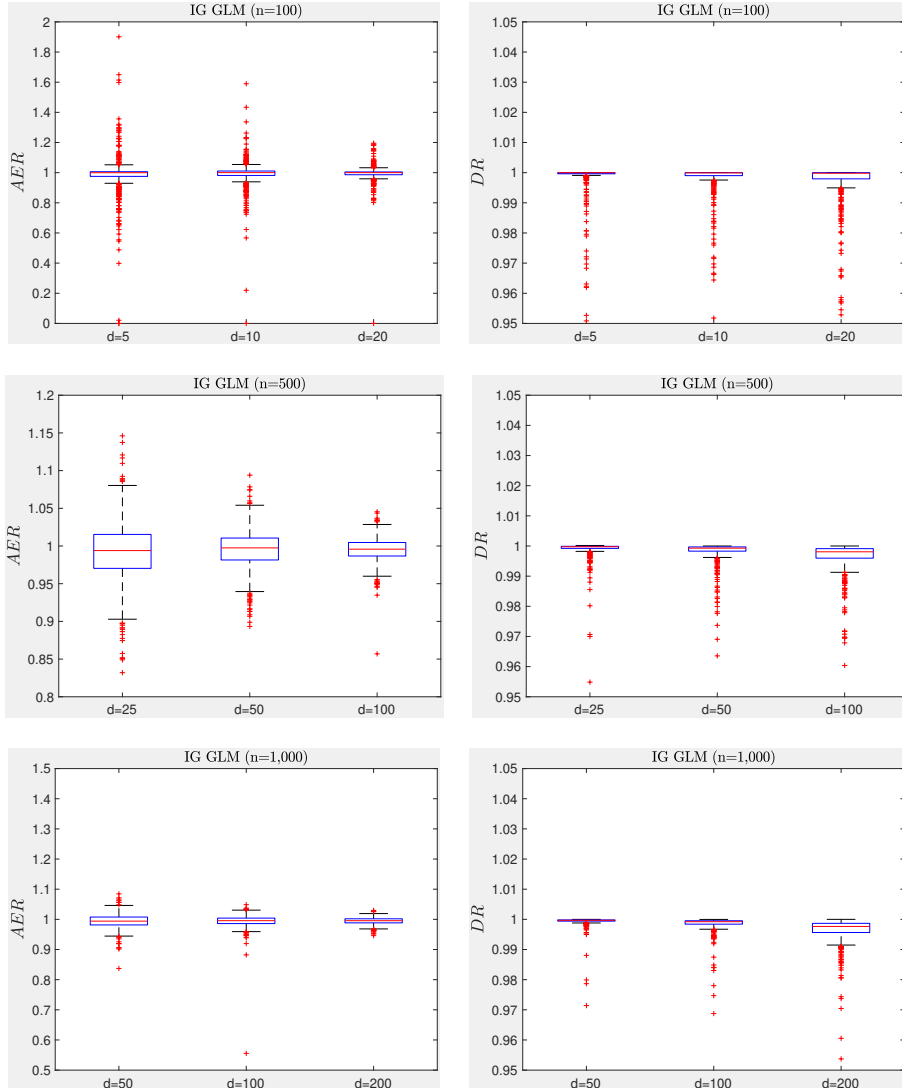


Figure 2: Absolute Error Ratio (AER) and Deviance Ratio (DR) for Inverse Gaussian GLM based on MATLAB *fitglm*.

Notes: This figure shows the box plots of Absolute Error Ratio (AER) in left panel and Deviance Ratio (DR) in right panel of Algorithm 2 from Section 4.2 relative to the MATLAB *fitglm* benchmark for the Inverse Gaussian GLM based on the *reciprocal-square-root* LF. Each box plot is constructed using AER s and DR s computed based on MLE values obtained from $N = 500$ simulations according to the DGP scheme outlined in Appendix C, for different specifications for the number of observations n and the number of covariates d . All implementations use the same starting value with a maximum of 10,000 iterations and 10^{-6} tolerance level.

6. Conclusions

This paper makes two important contributions to the GLM literature. First, we provide a general characterisation of proper GLMs for various exponential dispersion models, including the Tweedie family. The main finding is that although most Tweedie GLMs are not proper for *canonical* and *log* LFs, a rich class of proper Tweedie GLMs can be identified for *power* LFs. Second, we propose specialized optimisation algorithms for implementing several instances of Tweedie GLMs under *power* LFs. These algorithms outperform standard methods in terms of accuracy and efficiency, particularly in high-dimensional scenarios, as demonstrated via a thorough comparison with existing libraries like MATLAB *fitglm*, R *glm2*, and Python *sm.GLM*.

References

- Albert, A. and J. A. Anderson (1984). On the existence of maximum likelihood estimates in logistic regression models. *Biometrika* 71(1), 1–10.
- Bach, F. (2010). Self-concordant analysis for logistic regression. *Electronic Journal of Statistics* 4(none), 384 – 414.
- Bickel, P. J. and K. A. Doksum (2015). *Mathematical Statistics: Basic Ideas and Selected Topics, Volume I*. Second ed., Chapman and Hall/CRC.
- Boyd, S. and L. Vandenberghe (2004). *Convex optimization*. Cambridge university press.
- Boyd, S. P., N. Parikh, E. Chu, B. Peleato, and J. Eckstein (2011). Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends in Machine Learning* 3(1), 1–122.
- Breeden, J. L. (2016). Incorporating lifecycle and environment in loan-level forecasts and stress tests. *European Journal of Operational Research* 255(2), 649–658.
- Debón, A., F. Montes, and F. Puig (2008). Modelling and forecasting mortality in spain. *European Journal of Operational Research* 189(3), 624–637.
- Delong, L., M. Lindholm, and M. V. Wüthrich (2021). Making tweedie’s compound poisson model more accessible. *European Actuarial Journal* 11(1), 185–226.
- El Karoui, N., D. Bean, P. J. Bickel, C. Lim, and B. Yu (2013). On robust regression with high-dimensional predictors. *Proceedings of the National Academy of Sciences* 110(36), 14557–14562.
- Eling, M. and J. Wirfs (2019). What are the actual costs of cyber risk events? *European Journal of Operational Research* 272(3), 1109–1119.
- Fouskakis, D. (2012). Bayesian variable selection in generalized linear models using a combination of stochastic optimization methods. *European journal of operational research* 220(2), 414–422.
- Hastie, T., R. Tibshirani, and J. Friedman (2001). *The Elements of Statistical Learning*. Springer Series in Statistics. New York, NY, USA: Springer New York Inc.
- Jørgensen, B. (1987). Exponential dispersion models. *Journal of the Royal Statistical Society, Series B (Methodology)* 49(2), 127–145.
- Kuo, L. and B. Mallick (1998). Variable selection for regression models. *Sankhyā: The Indian Journal of Statistics, Series B* 60(1), 65–81.
- Mäkeläinen, T., K. Schmidt, and G. Styan (1981). On the existence and uniqueness of the maximum likelihood estimate of a vector-valued parameter in fixed-size samples. *Annals of Statistics* 9(4), 758–567.

- Marschner, I. C. (2011). Glm2: Fitting Generalized Linear Models with Convergence Problems. *The R Journal* 3(2), 12–15.
- McCullagh, P., J. Nelder, and R. Wedderburn (1989). *Generalized Linear Models*. Second ed., Chapman and Hall/CRC.
- Nelder, J. and R. Wedderburn (1972). Generalized linear models. *Journal of the Royal Statistical Society: Series A* 135(3), 370–384.
- Nesterov, Y. E. (2004). *Introductory Lectures on Convex Optimization - A Basic Course*, Volume 87 of *Applied Optimization*. Springer.
- Sur, P. and E. J. Candès (2019). A modern maximum-likelihood theory for high-dimensional logistic regression. *Proceedings of the National Academy of Sciences* 116(29), 14516–14525.
- van Staden, H. E., L. Deprez, and R. N. Boute (2022). A dynamic “predict, then optimize” preventive maintenance approach using operational intervention data. *European Journal of Operational Research* 302(3), 1079–1096.
- Wedderburn, R. (1976). On the existence and uniqueness of the maximum likelihood estimates for certain generalized linear models. *Biometrika* 63(1), 27–32.

Appendix A. Proofs

Appendix A.1. Proof of Propositions 3.1-3.3

The proofs follow easily by verifying the conditions in Definition 2.1 for the Poisson, Gamma and Inverse Gaussian families, respectively.

Appendix A.2. Proof of Theorem 3.4

The identification of the three classes of Tweedie GLM that are well-defined is not difficult, and thus, we only outline some arguments without further details that are quite obvious. Clearly, $b'(\theta) = (\theta/(\alpha - 1))^{\alpha-1}$ for all $\theta \in \mathfrak{R}$. Since $\alpha < 2$, then setting a) is readily true and we require $\alpha \in (1, 2)$, which is equivalent to $p < 0$, whenever $\Theta \in \{\mathfrak{R}_+^*, \mathfrak{R}_+\}$. Setting b) is the mirror case of setting a), and the proof is very similar. Settings c) and d) are similar to the previous ones, and the analysis depends if $\alpha - 1$ is an odd or even negative integer.

Appendix A.3. Proof of Theorem 3.5

First, we investigate parts (i) and (iii) (the *negative power* LF case) together, and therefore assume only *power* or *negative power* LFs. Condition **C2** requires

$$y(\alpha - 1)(h(\eta))^{\frac{1}{\alpha-1}} - \frac{\alpha - 1}{\alpha}(h(\eta))^{\frac{\alpha}{\alpha-1}} \quad \text{to be concave in } \eta \text{ on } \mathfrak{R} \text{ for all } y \in \mathcal{Y}. \quad (\text{A.1})$$

Setting a) is first justified, but only for *power* LFs since the image of h is \mathfrak{R}_+^* , and in turn, $\gamma = 2k, k \in \mathbb{Z}^*$. Denote $a_1 = y(\alpha - 1)$, $a_2 = \frac{1-\alpha}{\alpha}$ and $\gamma' = \frac{1}{\alpha-1}$. Equation (A.1) is equivalent to

$$\xi(\eta; y) := a_1 \eta^{\gamma\gamma'} + a_2 \eta^{\gamma(\gamma'+1)} \quad \text{is concave in } \eta \text{ on } \mathfrak{R} \text{ for all } y \in \mathcal{Y}. \quad (\text{A.2})$$

Note that $\mathcal{Y} = \mathfrak{R}_+^*$ is assumed. Since $1 < \alpha < 2$ and $y > 0$, then $a_1 > 0$ and $a_2 < 0$, and in turn, (A.2) holds if and only if $\gamma\gamma' \in [0, 1]$ and $\gamma(\gamma' + 1) \notin (0, 1)$. This is equivalent to having $\gamma \geq 0$, $\gamma\gamma' \leq 1$ and $\gamma(\gamma' + 1) \geq 1$, since $\gamma' > 1$ in this case, which is further equivalent to $\frac{\alpha-1}{\alpha} \leq \gamma \leq \alpha - 1$. The latter cannot hold since $\alpha - 1 \in (0, 1)$, $\alpha - 1 - \frac{\alpha-1}{\alpha} \in (0, 1/2)$ and $\gamma = 2k, k \in \mathbb{Z}^*$, which concludes that no proper GLM model is possible for setting a).

Setting b) is now justified, but only for *power* LFs since the image of h is \mathfrak{R}_+^* , and thus, $\gamma = 2k, k \in \mathbb{Z}^*$. We split this in two subcases, setting b1) and setting b2) for $0 < \alpha < 1$ and $\alpha < 0$, respectively.

Setting b1) holds if and only if $\gamma\gamma' \notin (0, 1)$ and $\gamma(\gamma' + 1) \in [0, 1]$, since $a_1 < 0$ and $a_2 > 0$, which is equivalent to having $\gamma \leq 0$, $\gamma\gamma' \geq 1$ and $\gamma(\gamma' + 1) \geq 1$ as $\gamma' < -1$, and in turn, $\frac{\alpha-1}{\alpha} \leq \gamma \leq \alpha - 1$. The later is true if and only if $\gamma = -2k$ for any $k \in \mathbb{N}^*$ and $(1 - \gamma)\alpha \leq 1$ since $0 < \alpha < 1$, which concludes setting b1).

Setting b2) implies that $a_1, a_2 < 0$ and $\gamma' \in (-1, 0)$. Therefore, setting b2) holds if and only if $\gamma\gamma' \notin (0, 1)$ and $\gamma(\gamma' + 1) \notin (0, 1)$, which is equivalent to having $\gamma \geq 0$ and $\gamma(\gamma' + 1) \geq 1$ or $\gamma \leq 0$ and $\gamma\gamma' \geq 1$, and in turn, $\frac{\alpha-1}{\alpha} \leq \gamma$ or $\gamma \leq \alpha - 1$ must hold, which concludes setting b2).

Setting c) is similar to setting b2), and we thus skip its proof. Setting d) requires for *power* and *negative power* LFs to having $\gamma' \in \mathbb{Z}$ so that the likelihood function is well-defined in (2.4) (and thus, in (A.1)), but also γ to be an odd integer so that the image of h is \mathfrak{R}^* . These do

not hold since $\gamma' \in (-1, 0)$, which justifies our claim for setting d). This concludes parts (i) and (iii) (the *negative power* LF case).

The proof of parts (ii) and (iii) (the *negative half-power* LF case) follows in a similar way, with one small difference. That is, *half-power* LFs require $\gamma \in \mathfrak{R}^*$ instead of $\gamma = 2k, k \in \mathbb{Z}^*$, but everything else does not significantly change. For these reasons, we do not provide additional details on this proof.

Appendix A.4. Proof of Theorem 3.6

We first show part (i), and assume *canonical* LFs. Note first $h(\eta) = b'(\eta) = (\eta/(\alpha - 1))^{\alpha-1}$, which implies that $\alpha \in \mathbb{Z} \setminus \{1\}$. This implies that amongst settings a)–c), only setting b2), which was introduced in [Appendix A.3](#), might hold while all other settings are clearly infeasible. The image of h is \mathfrak{R}_+^* and therefore, α is an odd negative integer, which is a *power* LF with an odd parameter γ . This contradicts our findings in the proof of part (i) from [Theorem 3.4](#) for setting b2), and concludes that no *canonical* LF leads to proper GLM in settings a)–c). Setting d) requires α to be an even negative integer and $\gamma' \in \mathbb{Z}$ as explained in the previous proof, which are infeasible conditions. Thus, no *canonical* LF leads to proper GLM in setting d). This concludes part (i).

We now show part (ii) and assume *log* LFs. Using the same notations as in [Appendix A.3](#), Equation (A.1) is equivalent to

$$\xi(\eta; y) := a_1 e^{\eta\gamma'} + a_2 e^{\eta(\gamma'+1)} \quad \text{is concave in } \eta \text{ on } \mathfrak{R} \text{ for all } y \in \mathcal{Y}, \quad (\text{A.3})$$

which requires $a_1, a_2 \leq 0$ due to the convexity property of $e^{\eta\gamma}$ in η on \mathfrak{R} , for any $\gamma \in \mathfrak{R}$. The latter explains that only setting b2) is feasible amongst settings a)–c). Setting d) is infeasible since the the image of h is \mathfrak{R}^* , which is impossible for a *log* LF. The proof is now complete.

Appendix A.5. Proof of Theorem 4.2

We proceed by showing part a), but only for (4.1), since (4.2) could be argued similarly. Let

$$f_{i,P}(\boldsymbol{\beta}) = \left(\frac{1}{2} \left(\mathbf{x}_i^\top \boldsymbol{\beta} \right)^2 - y_i \log \left(\mathbf{x}_i^\top \boldsymbol{\beta} \right) \right) \quad \text{for all } 1 \leq i \leq n, \quad (\text{A.4})$$

so that $f_P(\boldsymbol{\beta}) = \sum_{i=1}^n f_{i,P}(\boldsymbol{\beta})$. First, we show that f_P is a closed convex function on Ω . From (A.4), $f_{i,P}$ is convex (and therefore, continuous) on Ω , and since $\text{dom}(f_P) = \Omega$ is an open set and $\lim_{\boldsymbol{\beta} \rightarrow \boldsymbol{\beta}_0} f_{i,P}(\boldsymbol{\beta}) = \infty$ for all $\boldsymbol{\beta}_0 \in \partial \text{dom}(f_P)$, it follows that $f_{i,P}$ is closed convex on Ω . The closed convex property of f_P follows from the fact that it is a sum of closed convex functions.

We next prove that f_P is self-concordant on Ω . For any $t \in \mathfrak{R}$, $\mathbf{u} \in \Omega$ and $\mathbf{v} \in \mathfrak{R}^d$, such that $\mathbf{u} + t\mathbf{v} \in \Omega$, we define the function $g_{i,P}(t) = f_{i,P}(\mathbf{u} + t\mathbf{v})$, or any $i = 1, \dots, n$, and let $g_P(t) = \sum_{i=1}^n g_{i,P}(t)$. Next, we show that

$$|g_{i,P}'''(t)| \leq 2 (g_{i,P}''(t))^{3/2}. \quad (\text{A.5})$$

Note that

$$g''_{i,P}(t) = \left(\mathbf{x}_i^\top \mathbf{v}\right)^2 + \frac{y_i \left(\mathbf{x}_i^\top \mathbf{v}\right)^2}{\left(\mathbf{x}_i^\top \mathbf{u} + t \mathbf{x}_i^\top \mathbf{v}\right)^2} \quad \text{and} \quad g'''_{i,P}(t) = -\frac{2y_i \left(\mathbf{x}_i^\top \mathbf{v}\right)^3}{\left(\mathbf{x}_i^\top \mathbf{u} + t \mathbf{x}_i^\top \mathbf{v}\right)^3}.$$

Clearly, (A.5) holds whenever $\mathbf{x}_i^\top \mathbf{v} = 0$, and thus, we further assume that $\mathbf{x}_i^\top \mathbf{v} \neq 0$. Now,

$$\left|g'''_{i,P}(t)\right| \left(g''_{i,P}(t)\right)^{-3/2} = 2y_i \left(y_i + \left(\mathbf{x}_i^\top \mathbf{u} + t \mathbf{x}_i^\top \mathbf{v}\right)^2\right)^{-3/2} \leq 2,$$

since $y_i \leq y_i^{3/2} \leq (y_i + \epsilon_i)^{3/2}$ for any non-negative integer y_i and any $\epsilon_i \geq 0$ (recall that $y_i \in \mathbb{N}$ as the sampling distribution is Poisson). The self-concordant property of f_P follows from

$$\left|g'''_P(t)\right| = \left|\sum_{i=1}^n g'''_{i,P}(t)\right| \leq \sum_{i=1}^n \left|g'''_{i,P}(t)\right| \leq 2 \sum_{i=1}^n \left(g''_{i,P}(t)\right)^{3/2} \leq 2 \left(\sum_{i=1}^n g''_{i,P}(t)\right)^{3/2} = 2 \left(g''_P(t)\right)^{3/2}.$$

Note that the first inequality follows from the triangle inequality, the second from (A.5), and the last one from the fact that the p -norm on \mathfrak{R}^n , $\|\mathbf{x}\|_p := \left(\sum_{i=1}^n |x_i|^p\right)^{1/p}$ is a decreasing function in p on \mathfrak{R}_+^* for any $\mathbf{x} \in \mathfrak{R}^n$, and thus, $\|\mathbf{x}\|_1 \leq \|\mathbf{x}\|_{2/3}$. This completes the proof for part a). The proof of part b) follows in a similar way, and thus, we only provide the main steps. As before, we only show (4.3) since its proof is very similar to the proof of (4.4). We denote

$$f_{i,G}(\boldsymbol{\beta}) = \left(\frac{y_i}{2} \left(\mathbf{x}_i^\top \boldsymbol{\beta}\right)^2 - \log \left(\mathbf{x}_i^\top \boldsymbol{\beta}\right)\right) \quad \text{for all } 1 \leq i \leq n,$$

so that $f_G(\boldsymbol{\beta}) = \sum_{i=1}^n f_{i,G}(\boldsymbol{\beta})$. Following the same arguments as in part a), we may show that f_G is a closed convex function on Ω . The proof that f_G is self-concordant on Ω follows in a similar way by defining the function $g_{i,G}(t) = f_{i,G}(\mathbf{u} + t\mathbf{v})$ and $g_G(t) = \sum_{i=1}^n g_{i,G}(t)$ for any $t \in \mathfrak{R}$, $\mathbf{u} \in \Omega$ and $\mathbf{v} \in \mathfrak{R}^d$, such that $\mathbf{u} + t\mathbf{v} \in \Omega$, and showing that $\left|g'''_{i,G}(t)\right| \leq 2 \left(g''_{i,G}(t)\right)^{3/2}$. The second and third order derivatives of $g_{i,G}$ are given by

$$g''_{i,G}(t) = y_i \left(\mathbf{x}_i^\top \mathbf{v}\right)^2 + \frac{\left(\mathbf{x}_i^\top \mathbf{v}\right)^2}{\left(\mathbf{x}_i^\top \mathbf{u} + t \mathbf{x}_i^\top \mathbf{v}\right)^2} \quad \text{and} \quad g'''_{i,G}(t) = -\frac{2 \left(\mathbf{x}_i^\top \mathbf{v}\right)^3}{\left(\mathbf{x}_i^\top \mathbf{u} + t \mathbf{x}_i^\top \mathbf{v}\right)^3}.$$

Clearly, the required inequality holds if $\mathbf{x}_i^\top \mathbf{v} = 0$, and thus, $\mathbf{x}_i^\top \mathbf{v} \neq 0$ is further assumed. Now,

$$\left|g'''_{i,G}(t)\right| \left(g''_{i,G}(t)\right)^{-3/2} = 2 \left(y_i \left(\mathbf{x}_i^\top \mathbf{u} + t \mathbf{x}_i^\top \mathbf{v}\right)^2 + 1\right)^{-3/2} \leq 2,$$

since $(1 + y_i \epsilon_i)^{-3/2} \leq 1$ for any $y_i > 0$ and $\epsilon_i \geq 0$ (recall that $y_i \in \mathfrak{R}_+^*$ as the sampling distribution is Gamma). This completes the proof.

Appendix B. Other special cases of GLMs

Appendix B.1. Linear Regression – Gaussian family

Assume that $Y \sim N(\theta, \phi^2)$ with probability distribution function given by

$$\log(f_Y(y; \theta, \phi)) = \frac{\theta y - \frac{\theta^2}{2}}{\phi} - \frac{1}{2} \left(\frac{y^2}{\phi} + \log(2\pi\phi) \right), \quad (y, \theta, \phi) \in \mathfrak{R} \times \mathfrak{R} \times \mathfrak{R}_+^*.$$

The above pdf is obtained as a special case of (2.1) by taking

$$a(\phi) = \phi, \quad b(\theta) = \frac{\theta^2}{2}, \quad c(y, \phi) = -\frac{1}{2} \left(\frac{y^2}{\phi} + \log(2\pi\phi) \right).$$

In addition, $b'(\Theta) = \mathfrak{R}$ and $b'^{-1}(\mu) = \mu$. Proposition Appendix B.1 provides a characterisation of the LFs under which the Gaussian GLM is properly defined according to Definition 2.1.

Proposition Appendix B.1. *Assume that $Y \sim N(\theta, \phi^2)$. The Gaussian GLM is proper if and only if $h : \mathfrak{R} \rightarrow \mathfrak{R}$ and*

$$-yh(\eta) + \frac{h^2(\eta)}{2} \quad \text{is convex in } \eta \text{ on } \mathfrak{R} \text{ for any given } y \in \mathfrak{R}. \quad (\text{B.1})$$

Proof. The proof follows from verifying conditions C1 and C2 from Definition 2.1. ■

Corollary Appendix B.2 identifies the only class of LFs which satisfies Equation (B.1).

Corollary Appendix B.2. *The Gaussian GLM is proper if and only if the LF is linear.*

Proof. Since any convex real function defined on a finite open set I is continuous with non-decreasing left (and right) derivatives, then (B.1) implies that

$$h'_+(\eta_1) h(\eta_1) - yh'_+(\eta_1) \leq h'_+(\eta_2) h(\eta_2) - yh'_+(\eta_2) \quad \text{for all } y \in \mathfrak{R}, \quad (\text{B.2})$$

and any reals $\eta_1 < \eta_2$ from I , where h'_+ is the right derivative of h . Assume now that h is not linear on \mathfrak{R} , and thus, not linear on I . Then, there exists $\eta_1 < \eta_2$ from I such that $h'_+(\eta_2) - h'_+(\eta_1) \neq 0$. The latter contradicts (B.2), and in turn, we must have h linear on \mathfrak{R} , and no other possible LF leads to a MLE-based Gaussian GLM. ■

The *canonical* LF for Gaussian GLMs is the *identity* function. Corollary Appendix B.2 implies the *canonical* LF leads to a proper GLM and it is the only *power* function with this property.

Appendix B.2. Logistic Regression – Bernoulli family

Assume that $Y \sim \text{Bernoulli}(\theta)$ with probability mass function given by

$$\log(f_Y(y; \theta, \phi)) = \theta y - \log(1 + e^\theta) \quad \text{with } (y, \theta, \phi) \in \{0, 1\} \times \mathfrak{R} \times \{1\}.$$

The above function is obtained as a special case of (2.1) by taking

$$a(\phi) = 1, \quad b(\theta) = \log(1 + e^\theta), \quad c(y, \phi) = 0.$$

In addition, $b'(\Theta) = (0, 1)$ and $b^{-1}(\mu) = \log \frac{\mu}{1-\mu}$. Proposition [Appendix B.3](#) provides a brief characterisation of a proper Logistic regression model.

Proposition Appendix B.3. *Assume that $Y \sim \text{Bernoulli}(\theta)$. The Bernoulli GLM is proper if and only if $h : \mathfrak{R} \rightarrow (0, 1)$, and*

$$y \log(h(\eta)) + (1-y) \log(1-h(\eta)) \quad \text{is concave in } \eta \text{ on } \mathfrak{R} \text{ for any given } y = \{0, 1\}. \quad (\text{B.3})$$

Proof. The proof follows easily by verifying the conditions **C1** and **C2** from Definition 2.1. ■
A direct consequence of the above is that the MLE-based Bernoulli GLM is proper if and only if $h(\eta)$ and $h(1-\eta)$ are log-concave functions²⁶ on \mathfrak{R} . Three standard choices for h have been proposed for this family in the literature, and all of them lead to proper GLMs:

- (i) *logit* LF, which corresponds to having $h(\eta) = \frac{1}{1+e^{-\eta}}$, which is also the Bernoulli *canonical* LF that satisfies the conditions in Proposition 2.2 since b is strictly convex on \mathfrak{R} .
- (ii) *probit* LF, which corresponds to having $h(\eta) = \Phi(\eta)$, where Φ is the cdf of a standard Gaussian random variable. In this case, it is not difficult to show that h satisfies the characterisation from Proposition [Appendix B.3](#).
- (iii) *complementary log-log* LF, which corresponds to having $h(\eta) = 1 - \exp(-\exp(-\eta))$. It is not difficult to show that h satisfies the conditions in Proposition [Appendix B.3](#).

Finally, it is clear that no power LF satisfies the conditions in Proposition [Appendix B.3](#).

Appendix C. Data Generation Process

This section briefly outlines the DGPs for the Poisson, Gamma and Inverse Gaussian GLMs.

- **Step 1:** Generate the matrix of covariates $\mathbf{X} = \{X_{i,j}\}_{i=1,j=1}^{n,d}$, from a Gaussian distribution with mean μ and unit standard deviation, $X_{i,j} \sim \mathcal{N}(\mu, 1)$. Note that for each GLM, we let μ to be a function of d , such that the expected value of the response variable is within reasonable bounds in order to avoid exaggerating the parameter values when generating \mathbf{Y} in Step 3, which typically affect the estimation procedure for the benchmarks.²⁷
- **Step 2:** Generate the regression coefficient $\boldsymbol{\beta} = \{\beta_j\}_{j=0}^d$ by setting $\beta_j = j/d$.
- **Step 3:** For any $i = 1, \dots, n$, let $\theta_i = \beta_0 + \sum_{j=1}^d \beta_j x_{i,j}$ and generate the response variable $\mathbf{Y} = \{Y_i\}_{i=1}^n$ by simulating each Y_i from *Poisson* (θ_i^2) for the Poisson GLM, *Gamma* ($\theta_i^2, 1$) for the Gamma GLM and *IG* ($\theta_i^{-2}, 1$) for the Inverse Gaussian GLM.

Appendix D. Deviance for Poisson, Gamma and Inverse Gaussian GLMs

- Poisson GLM with *half-square-root* LF

$$D(\hat{\boldsymbol{\beta}}) = \sum_{i=1}^n \left(4y_i \log \left(\frac{\sqrt{y_i}}{\mathbf{x}_i^\top \hat{\boldsymbol{\beta}}} \right) + 2 \left(\left(\mathbf{x}_i^\top \hat{\boldsymbol{\beta}} \right)^2 - y_i \right) \right) \cdot I_{\mathbf{x}_i^\top \hat{\boldsymbol{\beta}} > 0} + 0 \cdot I_{\mathbf{x}_i^\top \hat{\boldsymbol{\beta}} = y_i = 0} + \infty \cdot I_{\text{else}}$$

²⁶A function $f : A \rightarrow B$ is log-concave on A if $\log(f(\alpha x + (1-\alpha)y)) \geq \alpha \log(f(x)) + (1-\alpha) \log(f(y))$ for all $x, y \in A$ and $0 < \alpha < 1$.

²⁷Note that in such cases all standard benchmarks fail to converge in most scenarios.

- Gamma GLM with *half-reciprocal-square-root* LF

$$D(\hat{\boldsymbol{\beta}}) = \sum_{i=1}^n \left(2y_i \left((\mathbf{x}_i^\top \hat{\boldsymbol{\beta}})^2 - y_i^{-1} \right) - 2 \log(y_i) - 4 \log(\mathbf{x}_i^\top \hat{\boldsymbol{\beta}}) \right) \cdot I_{\mathbf{x}_i^\top \hat{\boldsymbol{\beta}} > 0} + \infty \cdot I_{\mathbf{x}_i^\top \hat{\boldsymbol{\beta}} \leq 0}$$

- Inverse Gaussian GLM with *reciprocal-square root* LF

$$D(\hat{\boldsymbol{\beta}}) = \phi^2 \sum_{i=1}^n \left(y_i \left((\mathbf{x}_i^\top \hat{\boldsymbol{\beta}})^4 - y_i^{-2} \right) - 2 \left((\mathbf{x}_i^\top \hat{\boldsymbol{\beta}})^2 - y_i^{-1} \right) \right) \cdot I_{\mathbf{x}_i^\top \hat{\boldsymbol{\beta}} > 0} + \infty \cdot I_{\mathbf{x}_i^\top \hat{\boldsymbol{\beta}} \leq 0}$$