Efficient Estimation of COM-Poisson Regression and Additive Model

Suneel Babu Chatla^{a,*}, Galit Shmueli^a

^aInstitute of Service Science, National Tsing Hua University, Hsinchu 30013, Taiwan

Abstract

The Conway-Maxwell-Poisson (CMP) or COM-Poison regression is a popular model for count data due to its ability to capture both under dispersion and over dispersion. However, CMP regression is limited when dealing with complex nonlinear relationships. With today's wide availability of count data, especially due to the growing collection of data on human and social behavior, there is need for count data models that can capture complex nonlinear relationships. One useful approach is additive models; but, there has been no additive model implementation for the CMP distribution. To fill this void, we first propose a flexible estimation framework for CMP regression based on iterative reweighed least squares (IRLS) and then extend this model to allow for additive components using a penalized splines approach. Because the CMP distribution belongs to the exponential family, convergence of IRLS is guaranteed under some regularity conditions. Further, it is also known that IRLS provides smaller standard errors compared to gradient-based methods. We illustrate the usefulness of this approach through extensive simulation studies and using real data from a bike sharing system in Washington, DC.

Keywords: IRLS, Smoothing Splines, backfitting, Over and under dispersion, Time series

^{*}Corresponding author

Email addresses: suneel.chatla@iss.nthu.edu.tw (Suneel Babu Chatla), galit.shmueli@iss.nthu.edu.tw (Galit Shmueli)

1. Introduction

Count data have become popular dependent variables in studies in various areas, especially due to the growing availability of data on human and social behavior. Examples include the number of crimes in each neighborhood, number of accidents at an intersection, number of Facebook comments, ridership in bike sharing programs, etc. The wide availability of count data and the need for modeling such data as a function of other factors to establish causal relationships or to quantify correlated relationships has led to the widespread use of count data models.

The most commonly used regression models for cross-sectional count data are Poisson regression and Negative-Binomial regression. In addition, the Conway-Maxwell-Poisson (CMP) distribution (also known as the COM-Poisson distribution) has gained increasing popularity for its flexibility and ability to handle both over and under dispersed data. Revived by Shmueli et al. [29], the CMP distribution is a two-parameter generalization of the Poisson, Bernoulli, and Geometric distributions. Suppose Y is a random variable that follows a CMP distribution, then the probability mass function (p.m.f.) for $Y \in \{0, 1, 2, ...\}$ is defined as

$$P(Y = y) = \frac{\lambda^y}{(y!)^{\nu} \zeta(\lambda, \nu)}, \quad \text{where} \quad \zeta(\lambda, \nu) = \sum_{s=0}^{\infty} \frac{\lambda^s}{(s!)^{\nu}}$$

for the parameters $\lambda, \nu > 0$ and $0 < \lambda < 1, \nu = 0$.

The CMP distribution includes three well-known distributions as special cases: Poisson ($\nu = 1$), Geometric ($\nu = 0, \lambda < 1$), and Bernoulli ($\nu \to \infty$ with probability $\frac{\lambda}{1+\lambda}$). Due to the additional parameter ν , the CMP distribution is flexible enough to handle both over dispersion ($\nu < 1$) and under dispersion ($\nu > 1$) which are common in count data [27]. For more details on the distributional properties please refer to [4].

One of the major limitations of the CMP distribution is that the normalizing constant $\zeta(\lambda,\nu)$, which is an infinite series, does not have a closed form representation, and therefore there is no closed form representation available for the mean. This makes it difficult to model the mean directly as a function of covariates, as in standard models such as Poisson and Logistic regression. However, the CMP distribution belongs to the exponential family and thus has the properties and advantages of that family. Defining $\boldsymbol{\theta} = (\lambda, \nu)$, the CMP likelihood has the following form of an exponential family [17]:

$$P_Y(y|\boldsymbol{\theta}) = h(y) \exp\bigg(\sum_{i=1}^2 \eta_i(\boldsymbol{\theta}) T_i(y) - A(\boldsymbol{\theta})\bigg),$$

where the natural parameters are $\eta_1(\boldsymbol{\theta}) = \ln \lambda$ and $\eta_2(\boldsymbol{\theta}) = -\nu$ with corresponding statistics $T_1(y) = y, T_2(y) = \ln(y!)$ and $A(\boldsymbol{\theta}) = \ln \zeta(\lambda, \nu), h(y) = 1$, as mentioned in [29].

Although CMP regression is flexible in terms of handling both over and under dispersion, it is sometimes too restrictive for modeling nonlinear relationships or time series data. At the same time, additive models are widely used for modeling nonlinear relationships such as time series [5, 32]. Additive models have the advantage of being parsimonious while at the same time providing more flexibility to capture complicated relationships. Currently, there exists no additive model implementation for the CMP regression. Motivated by the need for flexible count data regression models for applications such as bike sharing, which can assist service providers in better management of their resources, we develop an additive model for CMP regression. Existing additive model implementations are heavily dependent upon the iterative reweighted least squares (IRLS) estimation framework, which currently does not exist for CMP regression. In this study, we propose and implement an IRLS estimation framework for CMP regression and then extend that to additive models.

The outline of this paper is as follows: In Section 2, we describe the CMP regression and the problems associated with IRLS implementation. In Section 3, we develop an IRLS framework for estimating a CMP regression by providing theory and the pseudo algorithm. We evaluate our proposed IRLS methodology with the existing methods using an extensive simulation study in Section 4. In Section 5, we use the IRLS framework to develop an additive model for the CMP distribution, and again evaluate its performance using a simulation in Section 6. In Section 7, we use our proposed additive model to draw valuable insights from a bike sharing application. Section 8 presents conclusions and future directions.

2. CMP Regression

Assume that we have a random sample of n observations $\{y_i, \boldsymbol{x}_i^T, \boldsymbol{z}_i^T\}_{i=1}^n$, where $\boldsymbol{x}_i^T = [1, x_{i1}, \cdots, x_{ip}]$ and $\boldsymbol{z}_i^T = [1, z_{i1}, \cdots, z_{iq}]$. In matrix notation, let

 $Y = [y_1, \ldots, y_n]^T$, $X = [\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n]^T$ and $Z = [\boldsymbol{z}_1, \ldots, \boldsymbol{z}_n]^T$ with the parameter vectors $\boldsymbol{\lambda} = (\lambda_1, \ldots, \lambda_n)^T$, $\boldsymbol{\nu} = (\nu_1, \ldots, \nu_n)^T$ and $\boldsymbol{\zeta} = (\zeta_1, \ldots, \zeta_n)^T$. We also denote mean and variance functions as $E[\cdot], V[\cdot]$ respectively.

When needed, we use the vector notation. With a slight abuse of notation, we extend the operations on scalars to operations on vectors. For example, we write $\ln(Y!) = (\ln(y_1!), \ldots, \ln(y_n!))^T$, $\ln(\boldsymbol{\lambda}) = (\ln(\lambda_1), \ldots, \ln(\lambda_n))^T$ and $\frac{\partial \ln \boldsymbol{\zeta}}{\partial \ln \boldsymbol{\lambda}} = (\frac{\partial \ln \zeta_1}{\partial \ln \lambda_1}, \ldots, \frac{\partial \ln \zeta_n}{\partial \ln \lambda_n})^T$. Unless otherwise stated, any operation on a vector simply denotes an extension of that operation to each component of that vector.

The CMP regression can be formulated as

$$\ln(\boldsymbol{\lambda}) = X\boldsymbol{\beta} \tag{1}$$

$$\ln(\boldsymbol{\nu}) = Z\boldsymbol{\gamma} \tag{2}$$

where $\boldsymbol{\beta} \in \mathbb{R}^{p+1}, \boldsymbol{\gamma} \in \mathbb{R}^{q+1}$.

The log link is used for the λ model. As mentioned in Sellers and Shmueli [27], this choice of log link is useful for two reasons. First, it coincides with the link function in two well-known cases: in Poisson regression, it reduces to $E[y_i] = \lambda_i$; in logistic regression, where $p_i = \frac{\lambda_i}{1+\lambda_i}$, it reduces to $logit(p_i) = \ln \lambda_i$. The second advantage of using a log link function is that it leads to elegant estimation, inference, and diagnostics. At the same time, we deliberately consider a log link for the $\boldsymbol{\nu}$ model, although the canonical link is identity, to restrict model predictions to the range $(0, \infty)$. This is important because while $\boldsymbol{\gamma}$ is unconstrained, ν_i $(i = 1, \ldots, n)$ can only take positive values and we cannot use the identity link between $\boldsymbol{\nu}$ and $\boldsymbol{\gamma}$.

In applications, it is common to treat $\boldsymbol{\nu}$ as nuisance parameter. For this reason, usually the data matrix Z contains only the intercept. Yet, since the $\boldsymbol{\nu}$ parameter models the dispersion, it is always better to include covariates that can potentially control for it [28]. In theory, one could use the same predictors for modeling both parameters. However, in practice, to avoid collinearity issues, it is better to have at least one different covariate in either the $\ln(\boldsymbol{\lambda})$ or the $\ln(\boldsymbol{\nu})$ model.

Using this model formulation, the log likelihood for the i^{th} observation can be written as

$$\ell_i(y_i, \boldsymbol{\beta}, \boldsymbol{\gamma}) = y_i \boldsymbol{x}_i^T \boldsymbol{\beta} - \ln(y_i!) \exp\{\boldsymbol{z}_i^T \boldsymbol{\gamma}\} - \ln \zeta_i(\exp\{\boldsymbol{x}_i^T \boldsymbol{\beta}\}, \exp\{\boldsymbol{z}_i^T \boldsymbol{\gamma}\}),$$

which yields the following score equations:

$$\frac{\partial \ell_i}{\partial \boldsymbol{\beta}^T} = \boldsymbol{x}_i \left(y_i - \frac{\partial \ln \zeta_i}{\partial \ln \lambda_i} \right) = \boldsymbol{x}_i (y_i - E[y_i]),
\frac{\partial \ell_i}{\partial \boldsymbol{\gamma}^T} = \boldsymbol{z}_i \left[\left(-\ln(y_i!) - \frac{\partial \ln \zeta_i}{\partial \nu_i} \right) \nu_i \right] = \boldsymbol{z}_i \left[\left(-\ln(y_i!) + E[\ln(y_i!)] \right) \right] \nu_i.$$
(3)

Since the derivatives of ζ_i (i = 1, ..., n) do not have closed form representations, the score equations in (3) cannot be solved as in standard generalized linear models (GLM) such as Poisson. For this reason, the existing implementations of CMP regression either use numerical gradient-based methods or Markov Chain Monte Carlo (MCMC), but do not use IRLS, which is the workhorse routine for estimation of all the standard GLMs. Although gradient-based methods have a faster convergence rate than IRLS, they are not efficient because they use the observed information matrix, and are not robust to outliers. In contrast, IRLS is more efficient and robust but is slower than gradient-based methods [9].

Another advantage of the IRLS algorithm is that modeling extensions such as additive models and lasso can be implemented easily [41]. To the best of our knowledge, there is no implementation of an IRLS algorithm for CMP regression. While Sellers and Shmueli [27] briefly outlined the IRLS algorithm, they did not implement it. Their approach is based on solving the following weighted least squares (WLS) equation at the m^{th} iteration (in matrix notation):

$$\begin{bmatrix} X^T \\ (g(Y) * Z)^T \end{bmatrix} W \begin{bmatrix} X & (g(Y) * Z) \end{bmatrix} \begin{bmatrix} \boldsymbol{\beta}^{(m)} \\ \boldsymbol{\gamma}^{(m)} \end{bmatrix} = \begin{bmatrix} X^T \\ (g(Y) * Z)^T \end{bmatrix} W \quad T$$

where $W = \text{diag}(V[y_1], \ldots, V[y_n]), g(Y) = (g(y_1), \ldots, g(y_n))^T$ with $g(y_i) = \frac{-\ln(y_i!) + E[\ln(y_i!)]}{y_i - E[y_i]} \nu_i$ and the adjusted dependent variable is $T = (t_1, \ldots, t_n)^T$ with

$$t_i = \boldsymbol{x}_i^T \boldsymbol{\beta}^{(m-1)} + g(y_i) \boldsymbol{z}_i^T \boldsymbol{\gamma}^{(m-1)} + \frac{y_i - E[y_i]}{V[y_i]}$$

Sellers and Shmueli [27] used only an intercept in the Z matrix and did not use a log link function for $\boldsymbol{\nu}$. Here we generalize their approach using a log link function for $\boldsymbol{\nu}$ and not restricting Z to have only an intercept. While the approach looks reasonable, it has the following two drawbacks: 1. The formulation by Sellers and Shmueli [27] does not use the expected information matrix. For example, based on their WLS formulation, the information for the intercept term in the $\ln(\nu)$ model can be written as:

$$\sum_{i=1}^{n} g(y_i)^2 V[y_i] = \sum_{i=1}^{n} \left[\frac{-\ln(y_i!) + E[\ln(y_i!)]}{y_i - E[y_i]} \right]^2 \nu_i^2 V[y_i]$$

$$\neq \sum_{i=1}^{n} V[\ln(y_i!)] \nu_i^2.$$
(4)

The value in the right hand side of the inequality is the expected information for the intercept term in the $\ln(\nu)$ model using the score equations (3). Clearly, there is some discrepancy as the information evaluated from the Sellers and Shmueli [27] derivation is not equal to the expected information. What we know is that the expected information is efficient and we do not know whether the Sellers and Shmueli [27] formula achieves the same efficiency (at least asymptotically). The Sellers and Shmueli [27] formulation matches the expected information only if

$$\left[\frac{-\ln(y_i!) + E[\ln(y_i!)]}{y_i - E[y_i]}\right]^2 = \frac{E\left[-\ln(y_i!) + E[\ln(y_i!)]\right]^2}{E\left[y_i - E[y_i]\right]^2}.$$

Based on the well known *Cramèr - Rao inequality* [16], the variance of any unbiased estimator is bounded by the inverse of the expected (Fisher) information. In general, the IRLS should use the expected information.

2. The idea of combining both models into a single WLS framework is computationally attractive. However, since both β and γ are dependent on each other, updating both of them in single model is problematic especially with least squares. When we implemented this approach, most of the time the algorithm remained close to the initial values and sometimes it chose very small values of ν_i (i = 1, ..., n) irrespective of the true values.

To overcome these limitations, we propose a two step IRLS algorithm with guaranteed convergence that uses the expected information matrix for updates. Our approach also makes it easier to extend the CMP regression for the estimation of additive components.

3. IRLS Framework for CMP Regression

To implement the IRLS method, we must first calculate the cumulants $E[y_i]$, $E[\ln(y_i!)]$, $V[y_i]$ and $V[\ln(y_i!)]$ for i = 1, ..., n.

3.1. Calculation of Cumulants

The standard way of calculating cumulants is to use the p.m.f. Since the p.m.f. for the CMP distribution involves an infinite series, a simple approach is to truncate the infinite series in such a way that the error is bounded $(\epsilon = 10^{-6})$ [29].

Another way of calculating the cumulants is by using the properties of the canonical parameter of the exponential family. The t^{th} cumulants for y_i and $\ln(y_i!)$ for $i = 1, \ldots, n$ can be obtained as:

$$\kappa_t[y_i] = \frac{\partial^{(t)} \ln \zeta_i(\lambda_i, \nu_i)}{\partial^{(t)} \ln \lambda_i}, \quad \kappa_t[\ln(y_i!)] = -\frac{\partial^{(t)} \ln \zeta_i(\lambda_i, \nu_i)}{\partial^{(t)} \nu_i}.$$

There has been some active research trying to approximate the ζ_i function using a closed form representation. Shmueli et al. [29] provided that for fixed positive integer ν_i the following asymptotic approximation holds:

$$\zeta_i(\lambda_i, \nu_i) = \frac{e^{\nu_i \lambda_i^{1/\nu_i}}}{\lambda_i^{\frac{\nu_i - 1}{2\nu_i}} (2\pi)^{\frac{\nu_i - 1}{2}} \sqrt{\nu_i}} (1 + \mathcal{O}(\lambda_i^{\frac{-1}{\nu_i}})).$$
(5)

Earlier, Olver [22] had derived the same leading term in the asymptotic expansion (5) and proved that it is valid for $0 < \nu_i \leq 4$. Gillispie and Green [8] built on the work of [22] to confirm that Equation (5) holds for all $\nu_i > 0$.

For higher order cumulants (t > 1) the cumulant generating function has the following form:

$$\kappa_t[y_i] = \nu_i \lambda_i^{1/\nu_i} (e^{t/\nu_i} - 1).$$
(6)

Although this approximation is appealing theoretically, it has limited practical value. To get a better approximation with the formulation in Equation (6) we should have larger λ_i^{1/ν_i} values, i.e., larger counts.

Recently, Gaunt et al. [7] further improved the asymptotic approximation in Equation (5) by providing lower order terms :

$$\zeta_{i}(\lambda_{i},\nu_{i}) = \frac{e^{\nu_{i}\lambda_{i}^{1/\nu_{i}}}}{\lambda_{i}^{\frac{\nu_{i}-1}{2\nu_{i}}}(2\pi)^{\frac{\nu_{i}-1}{2}}\sqrt{\nu_{i}}} \left(1 + c_{1}(\nu_{i}\lambda_{i}^{1/\nu_{i}})^{-1} + c_{2}(\nu_{i}\lambda_{i}^{1/\nu_{i}})^{-2} + \mathcal{O}(\lambda_{i}^{\frac{-3}{\nu_{i}}})\right),$$
(7)

where $c_1 = \frac{\nu_i^2 - 1}{24}$, $c_2 = \frac{\nu_i^2 - 1}{48} + \frac{c_1^2}{2}$. Daly and Gaunt [4] derived the leading term in the asymptotic approximation of all cumulants. However, since Gaunt et al. [7] provided the expressions for the cumulants for $E[y_i]$ and $V[y_i]$ including the first two correction terms, we use their results to approximate the mean and variance. Define $\alpha_i = \lambda_i^{1/\nu_i}$, then

$$E[y_i] = \alpha_i - \frac{\nu_i - 1}{2\nu_i} - \frac{\nu_i^2 - 1}{24\nu_i^2} \alpha_i^{-1} - \frac{\nu_i^2 - 1}{24\nu_i^3} \alpha_i^{-2} + \mathcal{O}(\alpha_i^{-3}),$$
(8)

$$V[y_i] = \frac{\alpha_i}{\nu_i} + \frac{\nu_i^2 - 1}{24\nu_i^3} \alpha_i^{-1} + \frac{\nu_i^2 - 1}{12\nu_i^4} \alpha_i^{-2} + \mathcal{O}(\alpha_i^{-3}).$$
(9)

Since we also need the first two cumulants for $\ln(y_i!)$, we used the asymptotic expression in (7) to calculate both $E[\ln(y_i!)]$ and $V[\ln(y_i!)]$:

$$E[\ln(y_i!)] = \alpha_i \left(\frac{\ln\lambda_i}{\nu_i} - 1\right) + \frac{\ln\lambda_i}{2\nu_i^2} + \frac{1}{2\nu_i} + \frac{\ln 2\pi}{2} - \frac{\alpha_i^{-1}}{24} \left(1 + \frac{1}{\nu_i^2} + \frac{\ln\lambda_i}{\nu_i} - \frac{\ln\lambda_i}{\nu_i^3}\right)$$
(10)
$$- \frac{\alpha_i^{-2}}{24} \left(\frac{1}{\nu_i^3} + \frac{\ln\lambda_i}{\nu_i^2} - \frac{\ln\lambda_i}{\nu_i^4}\right) + \mathcal{O}(\alpha_i^{-3}),$$
$$V[\ln(y_i!)] = \alpha_i \frac{(\ln\lambda_i)^2}{\nu_i^3} + \frac{\ln\lambda_i}{\nu_i^3} + \frac{1}{2\nu_i^2} + \frac{\alpha_i^{-1}}{24\nu_i^5} \left(-2\nu_i^2 + 4\nu_i\ln\lambda_i + (-1 + \nu_i^2)(\ln\lambda_i)^2\right) + \frac{\alpha_i^{-2}}{24\nu_i^6} \left(-3\nu_i^2 - 2\nu_i(-3 + \nu_i^2)\ln\lambda_i + 2(-1 + \nu_i^2)(\ln\lambda_i)^2\right) + \mathcal{O}(\alpha_i^{-3}).$$
(11)

While Gaunt et al. [7] showed that the asymptotic approximation is much better after including two more terms, it is still not close to the true value at least for the parameter range $0 < \lambda_i < 2$. Hence we consider the remaining range to check whether the asymptotic approximation is accurate or not. To illustrate this, we computed the $\ln \zeta_i$ function for the parameter range $2 \leq \lambda_i \leq 20, 0.2 \leq \nu_i \leq 10$ using both the truncated infinite series (bounding error=10⁻⁶) and the asymptotic expression in (7). We plotted the differences between the two sets of values in Figure 1. It can be observed that the new approximation is reasonably good when $\lambda_i \geq 2$ and $\nu_i \leq 1$, while for higher



Figure 1: The differences between $\ln \zeta_i$ calculated using the truncated infinite series and the asymptotic expression in (7) for the parameter range: $2 \le \lambda_i \le 20, 0.2 \le \nu_i \le 10$.

values of ν_i (> 2) the asymptotic approximation tends to over estimate the true value. For this reason we use the cumulants derived from the asymptotic expression only when $\lambda_i \geq 2$, $\nu_i \leq 1$, while for other values we use the p.m.f. to calculate cumulants recursively with some bounding error. Although the approximation works for a limited range, this is very helpful because the asymptotic series converges very slowly when $\nu_i < 1$ and this approximation eases the computational burden significantly.

Similarly, for $\ln(y_i!)$ the values are computed recursively until $y_i < 254$ and after that Stirling's approximation is used as it is reasonably close [1].

3.2. Two step method

Let us define $\boldsymbol{u}_{(p+1)\times 1} = \sum_{i=1}^{n} \frac{\partial \ell_i}{\partial \boldsymbol{\beta}^T}$ and $\boldsymbol{v}_{(q+1)\times 1} = \sum_{i=1}^{n} \frac{\partial \ell_i}{\partial \boldsymbol{\gamma}^T}$. From Equation (3), the full information matrix I can be written as

$$I_{(p+1)\times(q+1)} = E\begin{bmatrix} \begin{pmatrix} \boldsymbol{u} \\ \boldsymbol{v} \end{pmatrix} (\boldsymbol{u}^T \boldsymbol{v}^T) \end{bmatrix} = \begin{bmatrix} E[\boldsymbol{u}\boldsymbol{u}^T] & E[\boldsymbol{u}\boldsymbol{v}^T] \\ E[\boldsymbol{v}\boldsymbol{u}^T] & E[\boldsymbol{v}\boldsymbol{v}^T] \end{bmatrix} = \begin{bmatrix} I_{11} & I_{12} \\ I_{21} & I_{22} \end{bmatrix},$$

with

$$E[\boldsymbol{u}\boldsymbol{u}^{T}] = I_{11} = X^{T}\Sigma_{Y}X,$$

$$E[\boldsymbol{u}\boldsymbol{v}^{T}] = I_{12} = I_{21}^{T} = -\boldsymbol{\nu} * X^{T}\Sigma_{Y,\ln(Y!)}Z,$$

$$E[\boldsymbol{v}\boldsymbol{v}^{T}] = I_{22} = \boldsymbol{\nu}^{2} * Z^{T}\Sigma_{\ln(Y!)}Z,$$

where * denotes element-wise multiplication, $\Sigma_Y = \text{diag}(V[y_1], \ldots, V[y_n]),$ $\Sigma_{Y,\ln(Y!)} = [\text{cov}(y_i, \ln(y_j!))]_{1 \le i,j \le n}, \text{ and } \Sigma_{\ln(Y!)} = \text{diag}(V[\ln(y_1!)], \ldots, V[\ln(y_n!)]).$ Using the information matrix I, the IRLS algorithm has the following

Using the information matrix I, the IRLS algorithm has the following form for the m^{th} update:

$$\begin{bmatrix} \boldsymbol{\beta} \\ \boldsymbol{\gamma} \end{bmatrix}^{(m)} = \begin{bmatrix} \boldsymbol{\beta} \\ \boldsymbol{\gamma} \end{bmatrix}^{(m-1)} + I^{-1} \begin{bmatrix} \boldsymbol{u} \\ \boldsymbol{v} \end{bmatrix},$$

which implies the following equations:

$$X^{T} \Sigma_{Y} X \boldsymbol{\beta}^{(m)} - X^{T} \Sigma_{Y, \ln(Y!)} \boldsymbol{\nu} * Z \boldsymbol{\gamma}^{(m)} =$$

$$X^{T} \Sigma_{Y} X \boldsymbol{\beta}^{(m-1)} - X^{T} \Sigma_{Y, \ln(Y!)} \boldsymbol{\nu} * Z \boldsymbol{\gamma}^{(m-1)} + X^{T} (Y - E[Y])$$
(12)

and

$$-\boldsymbol{\nu} * Z^{T} \Sigma_{Y,\ln(Y!)} X \boldsymbol{\beta}^{(m)} + \boldsymbol{\nu}^{2} * Z^{T} \Sigma_{\ln(Y!)} Z \boldsymbol{\gamma}^{(m)} = -\boldsymbol{\nu} Z^{T} \Sigma_{Y,\ln(Y!)} X \boldsymbol{\beta}^{(m-1)} + \boldsymbol{\nu}^{2} * Z^{T} \Sigma_{\ln(Y!)} Z \boldsymbol{\gamma}^{(m-1)}$$
(13)
$$+\boldsymbol{\nu} * Z^{T} (-\ln(Y!) + E[\ln(Y!)]).$$

Each of the two equations in (12) and (13) are complicated and contain updates for both parameters β and γ . However, if we fix one parameter in each equation, a nice closed form expression appears for the updates. When we fix γ in equation (12) the equation reduces to

$$X^T \Sigma_Y X \boldsymbol{\beta}^{(m)} = X^T \Sigma_Y X \boldsymbol{\beta}^{(m-1)} + X^T (Y - E[Y]).$$
(14)

This equation is nothing but WLS of X on Y with weights Σ_Y . Similarly, if we fix β in equation (13) then the equation reduces to

$$\boldsymbol{\nu}^2 * Z^T \Sigma_{\ln(Y!)} Z \boldsymbol{\gamma}^{(m)} = \boldsymbol{\nu}^2 * Z^T \Sigma_{\ln(Y!)} Z \boldsymbol{\gamma}^{(m-1)} + \boldsymbol{\nu} * Z^T (-\ln(Y!) + E[\ln(Y!)])$$
(15)

Again this is a WLS of $\boldsymbol{\nu} * Z$ on $\ln(Y!)$ with weights $\Sigma_{\ln(Y!)}$.

The two update equations (14) and (15) are elegant and can be easily estimated with WLS methods. This approach is not only convenient for estimation but also helpful for generalizing to other modeling extensions such as additive models and the lasso.

3.3. Proof of Convergence of the Two Step Method

To prove the convergence properties of our proposed two step algorithm, we start with the following assumptions. Consider the parameter space $\Theta \in (0, \infty) \times (0, \infty)$, and the likelihood function L.

(A1). Let $\hat{\boldsymbol{\theta}}_{\mathbf{0}} = (\hat{\boldsymbol{\lambda}}_{\mathbf{0}}, \hat{\boldsymbol{\nu}}_{\mathbf{0}})^T \in \Theta$ be a starting value, such that $D_0 = \{\boldsymbol{\theta} = (\boldsymbol{\lambda}, \boldsymbol{\nu})^T \in \Theta | L(\boldsymbol{\theta}) \geq L(\hat{\boldsymbol{\theta}}_{\mathbf{0}}) \}$ is compact.

(A2). The function L is uniquely maximized over D_0 for $\boldsymbol{\theta} = \boldsymbol{\theta}$.

(A3). Suppose that we have given parameter functions $\psi_i : D_0 \to \Theta_i$ (i = 1, 2) and let $M_i(\theta), \theta \in D_0$ be the corresponding sections: $M_i(\theta) = \{ \boldsymbol{\eta} \in D_0 | \psi_i(\boldsymbol{\eta}) = \psi_i(\theta) \}$ (i = 1, 2). Then we assume that, for i = 1, 2 and $\boldsymbol{\theta} \in D_0$, L is maximized uniquely by $T_i(\boldsymbol{\theta})$ on the section $M_i(\boldsymbol{\theta})$ and that $T_i(\boldsymbol{\theta})$ is continuous on D_0 .

(A4). The point of global maximum $\hat{\boldsymbol{\theta}}$ is uniquely determined by the condition that it is the partial maximum along each section $M_i(\boldsymbol{\theta})$. In other words,

$$\sup_{\boldsymbol{\eta}\in M_i(\boldsymbol{\theta})} L(\boldsymbol{\eta}) = L(\boldsymbol{\theta}), \quad i = 1, 2$$

implies $\boldsymbol{\theta} = \hat{\boldsymbol{\theta}}$, or equivalently, $T_i(\boldsymbol{\theta}) = \boldsymbol{\theta}$ implies $\boldsymbol{\theta} = \hat{\boldsymbol{\theta}}$.

Assumptions A1 and A2 are based on the fact that the CMP distribution is unimodal and it has a log-concave p.m.f. [12, 31]. The remaining assumptions A3 and A4 follow from the properties of exponential family distributions. It is well known that the marginal distributions in a regular k-variate exponential family also belong to an exponential family [16, 17]. It means that for a distribution that belongs to an exponential family like the CMP distribution, the estimates obtained from maximizing the marginal likelihood are the same as the estimates obtained from maximizing the full likelihood.

Theorem 1. Under assumptions A1-A4, the two step IRLS algorithm

$$\hat{\boldsymbol{\theta}}_{n+1} = T_1(T_2(\hat{\boldsymbol{\theta}}_n))$$

converges to $\hat{\boldsymbol{\theta}}$ for any starting value in D_0 .

The proof is similar to Jensen et al. [15]. The authors showed that under the above assumptions any partial maximization algorithm converges to the true value for a given starting value.

Algorithm 1 IRLS Framework for CMP distribution

1: Set initial values for $\nu_i^{(0)}$ and $\lambda_i^{(0)} = (y_i + 0.1)^{\nu_i^{(0)}}$ for i = 1, ..., n. 2: Compute $\eta_{i1}^{(0)} = \ln(\lambda_i^{(0)})$ and $\eta_{i2}^{(0)} = \ln(\nu_i^{(0)})$ for i = 1, ..., n. 3: Compute $D^{(0)}(\boldsymbol{\lambda}^{(0)}, \boldsymbol{\nu}^{(0)}) = -2\sum_{i=1}^n \ell(\lambda_i^{(0)}, \nu_i^{(0)})$. 4: Compute $E[y_i]^{(0)}$ and $V[y_i]^{(0)}$ for i = 1, ..., n. 5: for k in 1:maxIter do Compute the adjusted dependent variable for each $i = 1, \ldots, n$: $t_{i1}^{(k)} = \eta_{i1}^{(k-1)} + \frac{y_i - E[y_i]^{(k-1)}}{V[y_i]^{(k-1)}}$. 6: Perform a weighted least squares regression of $T_1^{(k)} = (t_{11}, \ldots, t_{n1})^T$ on X with weights $W_1 = \text{diag}(V(y_1)^{(k-1)}, \ldots, V(y_n)^{(k-1)})$ to obtain $\boldsymbol{\beta}^{(k)}$. 7: $\begin{array}{l} \begin{array}{c} (\gamma_{11}, \dots, \gamma_{n1}) \\ (\gamma_{11}, \dots$ 8: 9: 10: Perform a weighted least squares regression of $T_2^{(k)} = (t_{12}, \ldots, t_{n2})^T$ on $\boldsymbol{\nu}^{(k-1)} * Z$ with weights $W_2 = \text{diag}(V[\ln(y_1!)]^{(k-1)}, \ldots, V[\ln(y_n!)]^{(k-1)})$ 11: to obtain $\boldsymbol{\gamma}^{(k)}$ Update $\eta_{i2}^{(k)} = \boldsymbol{z}_i^T \boldsymbol{\gamma}^{(k)}$ and $\nu_i^{(k)} = \exp(\eta_{i2}^{(k)})$ for each i = 1, ..., n. Compute $D^{(k)}(\boldsymbol{\lambda}^{(k)}, \boldsymbol{\nu}^{(k)}) = -2\sum_{i=1}^n \ell(\lambda_i^{(k)}, \nu_i^{(k)})$. 12:13:if $\frac{D^{(k)}-D^{(k-1)}}{D^{(k)}} > 10^{-6}$ then 14:Initiate step size optimization. 15:end if 16:if $\left|\frac{D^{(k)}-D^{(k-1)}}{D^{(k)}}\right| < 10^{-6}$ then 17:Convergence achieved. Break the loop. 18:19:else Compute $E[y_i]^{(k)}$ and $V[y_i]^{(k)}$ for each i = 1, ..., n. 20:end if 21: 22: end for

3.4. Practical issues

3.4.1. Initial Values

Unlike other nonlinear optimization algorithms, IRLS does not require initial values for the parameters β and γ but it does require initial values for λ and ν . We can provide suitable initial values based on the approximate method of a moments estimator for λ_i such as $(y_i + 0.1)^{\nu_i^{(o)}}$ for $i = 1, \ldots, n$. However, we do not have a closed form expression for $\nu_i^{(o)}$. In practice we observed that starting close to zero (e.g., $\nu_i = 0.2$ or 0.5) yields satisfactory results.

3.4.2. Stopping Criterion

The standard IRLS algorithm uses the deviance as a stopping criterion. If the absolute relative change in the deviance is below some tolerance threshold, the algorithm stops. In general, the deviance for the i^{th} observation is defined as:

$$D_i = -2(\ell(y_i; \hat{\lambda}_i, \hat{\nu}_i) - \ell(y_i; \hat{\lambda}_{i,sat}, \hat{\nu}_{i,sat})).$$

The estimates for both $\hat{\lambda}_{i,sat}, \hat{\nu}_{i,sat}$ depend on each other and we do not have closed forms especially for the estimate $\hat{\nu}_{i,sat}$. For this reason, we consider only the term $-2\sum \ell(y_i; \hat{\lambda}_i, \hat{\nu}_i)$ and use it as our stopping criterion. Since the likelihood for the saturated model is constant across all the iterations, ignoring the term $\ell(y_i; \hat{\lambda}_{i,sat}, \hat{\nu}_{i,sat})$ does not impact our stopping criterion. In addition, this function is monotonic and if the algorithm is converging, it will decrease with every iteration.

3.4.3. Step Size

It is common for IRLS to exhibit convergence problems [19]. To avoid non-convergence issues we used the step-halving approach suggested by Marschner [19]. The algorithm invokes step-halving either at the boundary or if the deviance is increasing. This step-halving makes sure that the algorithm remains in the interior space which is required for convergence.

3.5. Inference

Proposition 1. Under regularity conditions [16][p.158], the maximum likelihood estimators $\hat{\theta} = (\hat{\beta}, \hat{\gamma})^T$ are consistent and asymptotically normal:

$$\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \xrightarrow{d} \mathcal{N}_{p+q+2}(0, I^{-1}(\boldsymbol{\theta})).$$

The proof is an immediate consequence of the result from Keener [16][p.158]. Since the algorithm estimates each parameter vector separately while keeping the other parameter vector fixed, it only provides the marginal information for the respective parameters. The conditional information matrices can be straightforwardly obtained by using the matrix operations [18] as following:

$$\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) | \hat{\boldsymbol{\gamma}} \xrightarrow{d} \mathcal{N}_{p+1}(0, (I_{11} - I_{12}I_{22}^{-1}I_{21})^{-1})$$
$$\sqrt{n}(\hat{\boldsymbol{\gamma}} - \boldsymbol{\gamma}) | \hat{\boldsymbol{\beta}} \xrightarrow{d} \mathcal{N}_{q+1}(0, (I_{22} - I_{21}I_{11}^{-1}I_{12})^{-1}).$$

We note that the estimates for both β and γ are not independent and inferences on one parameter will be influenced by the other estimate. As it was mentioned earlier, for most practical applications, inference on the parameter β is of primary interest and usually the parameter γ will be treated as a nuisance parameter.

4. Simulation Study for the CMP Regression

We conducted an extensive simulation study to evaluate the performance of our proposed IRLS algorithm in comparison to existing gradient-based methods for estimating the CMP regression model.

At present there are two R packages (*CMPRegression* by [26]; *CompGLM* by [23]) for fitting the CMP regression model. Both use general purpose optimization functions to maximize the likelihood function. While these two R packages are technically the same, they differ in terms of their implementations. From now on, we denote these packages as Opt_1 and Opt_2 and our implementation as *IRLS*. While Opt_1 does not use a log link for the ν model, Opt_2 does use a log link and allows the user to model the ν parameter as well. Similarly, while Opt_1 only provides the log likelihood to the optimization function, Opt_2 also provides the gradients. More importantly, in terms of computational issues, the support functions for Opt_1 were implemented in R and the support functions for Opt_2 were implemented in C++ which makes it work much faster.

Both Opt_1 and Opt_2 methods have some limitations. One obvious limitation is their inability to handle larger counts. Even with a single large value in the data both methods produce errors. Since both methods supply the likelihood to an external optimization routine in R and the source code is not available for the external function, we were not able to identify the reasons as to why these methods fail with larger counts. However, we believe



Figure 2: Distribution of the simulated y corresponding to the results in Table 1 ($\nu = 0.5, 1, 2.5$ and 4, respectively).

that the problem is related to computing numerical derivatives which is often problematic with the CMP distribution because of the normalizing constant.

In order to provide a clear comparison of our method with the two aforementioned gradient-based methods, we carefully constructed simulated datasets without any large counts, as shown in Figure 2, so that none of these methods face any convergence issues. Although the distributions of the dependent variable in Figure 2, look similar they have different counts (please note the scale on the x-axis). After exploring several coefficient values, we reached these datasets for which both Opt_1 and Opt_2 did not show any problems. For the data simulated from the other sets of coefficient values only *IRLS* provided results. We considered sample size n = 500 and chose 4 covariates to be included in the model. The covariates are simulated from normal and uniform distributions and also allow mild correlation between one pair of covariates $(x_1 \sim U(0, 1), x_2 \sim N(0, 0.5), x_3 \sim N(0, 0.1)$ and $x_4 = 0.2x_3 + N(0, 0.5)$).

We considered four different values for ν in order to capture over dispersion ($\nu = 0.5$), equi dispersion ($\nu = 1$) and under dispersion ($\nu = 2.5, 4$) scenarios. The true values for the regression coefficients and their estimated values using 20 bootstrap replications are reported in Table 1. From the results, it can be observed that *IRLS* performs equally well to the existing gradient-based methods, especially Opt_1 . While the three methods are indistinguishable for the over dispersion case ($\nu = 0.5$), we observe that there are some clear discrepancies for the under dispersion case ($\nu = 2.5$ or 4). In particular, Opt_2 has some issues when there is under dispersion in the data.

For a couple of models from our simulation study ($\nu = 0.5, 4$) we compared the computation times of *IRLS* with both Opt_1 and Opt_2 with increasing sample sizes. It is well known that the convergence speed of the IRLS algorithm depends on its starting value. We therefore take a sub sample of data and then run the algorithm to get an approximated value for ν and feed it as the initial value for the estimation using the full data. We call this *IRLS*₂ and use *IRLS*₁ to denote the original algorithm which always starts at $\nu = 0.2$.

The computation times for the four methods are shown in Figures 3 and 4. The computation times for $IRLS_2$ include the time for estimating the model for the sub sample to obtain an initial value for the parameter estimates. From the results, it can be observed that Opt_2 is superior to $IRLS_1$, $IRLS_2$ and Opt_1 . Opt_1 is painfully slow and often takes many minutes where both $IRLS_2$ and Opt_2 take a few seconds. As expected, the $IRLS_2$ algorithm performs much faster than the original $IRLS_1$ and at times it even works faster than Opt_2 . It is also worth mentioning that while Opt_2 is very fast, it has some issues when there is under dispersion as we have already seen in the simulation results¹.

Ideally, we would want theoretical computation times which can provide a more rigorous comparative study. For the IRLS algorithm, it is easy to obtain the theoretical computation times given the number of iterations needed for the convergence, because for each IRLS iteration we fit two least squares models and we know the theoretical computation times for the least squares regression. However, both Opt_1 and Opt_2 algorithms use external optimization functions to optimize the likelihood. The theoretical computation times for these algorithms are not easy to obtain because of the complicated nature of the algorithms and also due to the computations involved in estimating the normalizing constant in the CMP distribution.

¹Although not reported here, we also examined the performance of the WLS formulation of Sellers and Shmueli [27] which was described in Section 2. Most of the time, their algorithm converged to the wrong value, typically with ν close to zero irrespective of the true ν value used to simulate the data.

		ν =	= 0.5	
	θ	$\hat{ heta}_{IRLS}$ (sd)	$ \hat{\theta}_{Opt_1} \\ (\text{sd}) $	$\hat{ heta}_{Opt_2} \ (ext{sd})$
β_0	0.05	0.05	0.04	0.04
, 0		(0.07)	(0.07)	(0.07)
β_1	0.5	0.52	0.52	0.52^{-1}
		(0.07)	(0.07)	(0.07)
β_2	-0.5	-0.52	-0.52	-0.52
		(0.06)	(0.06)	(0.06)
β_3	0.25	0.31	0.31	0.31
		(0.21)	(0.21)	(0.21)
β_4	-0.25	-0.25	-0.25	-0.25
		(0.06)	(0.06)	(0.06)
log(u)	-0.69	-0.65	-0.67	-0.67
		(0.09)	(0.09)	(0.09)
		$\nu =$: 2.5	
		$\frac{\nu}{\hat{\theta}_{IBLS}}$	$\frac{2.5}{\hat{\theta}_{Ont_1}}$	$\hat{\theta}_{Ont_2}$
	θ	$\nu = \frac{\hat{\theta}_{IRLS}}{(\text{sd})}$	$\begin{array}{c} \begin{array}{c} 2.5 \\ \hline \\ \hat{\theta}_{Opt_1} \\ (\text{sd}) \end{array}$	$\frac{\hat{\theta}_{Opt_2}}{(\mathrm{sd})}$
Ba	θ 1	$\nu = \frac{\hat{\theta}_{IRLS}}{(sd)}$	$\begin{array}{c} \begin{array}{c} 2.5 \\ \hline \hat{\theta}_{Opt_1} \\ (\text{sd}) \end{array}$	$\begin{array}{c} \hat{\theta}_{Opt_2} \\ \text{(sd)} \end{array}$
β_0	θ 1	$\nu = \frac{\hat{\theta}_{IRLS}}{(\text{sd})}$ 1.02 (0.13)	$\begin{array}{c} 2.5\\ \hline \hat{\theta}_{Opt_1}\\ (\text{sd})\\ \hline 1.02\\ (0.13) \end{array}$	$\begin{array}{c} \hat{\theta}_{Opt_2} \\ \text{(sd)} \end{array}$ 0.73 (0.16)
β_0	θ 1 3	$\nu = \frac{\hat{\theta}_{IRLS}}{(\text{sd})}$ 1.02 (0.13) 3.09	$ \begin{array}{c} $	
$egin{array}{c} eta_0 \ eta_1 \end{array}$	θ 1 3	$\nu = \frac{\hat{\theta}_{IRLS}}{(sd)}$ 1.02 (0.13) 3.09 (0.22)	$\begin{array}{c} \vdots 2.5 \\ \hline \hat{\theta}_{Opt_1} \\ \text{(sd)} \\ \hline 1.02 \\ (0.13) \\ 3.09 \\ (0.23) \end{array}$	$ \begin{array}{c} \hat{\theta}_{Opt_2} \\ (sd) \end{array} $ 0.73 (0.16) 2.43 (0.23)
β_0 β_1 β_2	θ 1 3 -3	$\nu = \frac{\hat{\theta}_{IRLS}}{(\text{sd})}$ 1.02 (0.13) 3.09 (0.22) -3.10	$\begin{array}{c} \vdots 2.5 \\ \hline \hat{\theta}_{Opt_1} \\ (\text{sd}) \\ \hline 1.02 \\ (0.13) \\ 3.09 \\ (0.23) \\ \mathbf{-3.10} \end{array}$	
β_0 β_1 β_2	θ 1 3 -3	$\nu = \frac{\hat{\theta}_{IRLS}}{\text{(sd)}}$ 1.02 (0.13) 3.09 (0.22) -3.10 (0.20)	$\begin{array}{c} \vdots 2.5 \\ \hline \hat{\theta}_{Opt_1} \\ (\text{sd}) \\ \hline 1.02 \\ (0.13) \\ 3.09 \\ (0.23) \\ \mathbf{-3.10} \\ (0.20) \end{array}$	$ \hat{\theta}_{Opt_2} (sd) 0.73 (0.16) 2.43 (0.23) -2.42 (0.24) $
$egin{array}{c} eta_0 \ eta_1 \ eta_2 \ eta_3 \end{array}$	$\begin{array}{c} \\ \theta \\ \hline 1 \\ 3 \\ -3 \\ 2 \end{array}$	$\nu = \frac{\hat{\theta}_{IRLS}}{(\text{sd})}$ 1.02 (0.13) 3.09 (0.22) -3.10 (0.20) 2.17	$\begin{array}{c} \vdots 2.5 \\ \hline \hat{\theta}_{Opt_1} \\ (\text{sd}) \\ \hline 1.02 \\ (0.13) \\ 3.09 \\ (0.23) \\ \mathbf{-3.10} \\ (0.20) \\ 2.17 \end{array}$	
eta_0 eta_1 eta_2 eta_3	θ 1 3 -3 2	$\nu = \frac{\hat{\theta}_{IRLS}}{(\text{sd})}$ 1.02 (0.13) 3.09 (0.22) -3.10 (0.20) 2.17 (0.39)	$\begin{array}{c} \vdots 2.5 \\ \hline \hat{\theta}_{Opt_1} \\ (\text{sd}) \\ \hline 1.02 \\ (0.13) \\ 3.09 \\ (0.23) \\ \mathbf{-3.10} \\ (0.20) \\ 2.17 \\ (0.39) \\ \end{array}$	
$egin{array}{c} eta_0 \ eta_1 \ eta_2 \ eta_3 \ eta_4 \end{array}$	θ 1 3 -3 2 -2	$\nu = \frac{\hat{\theta}_{IRLS}}{(\text{sd})}$ 1.02 (0.13) 3.09 (0.22) -3.10 (0.20) 2.17 (0.39) -2.06	$\begin{array}{c} \begin{array}{c} 2.5 \\ \hline \hat{\theta}_{Opt_1} \\ (\text{sd}) \\ \hline 1.02 \\ (0.13) \\ 3.09 \\ (0.23) \\ \mathbf{-3.10} \\ (0.20) \\ 2.17 \\ (0.39) \\ \mathbf{-2.06} \end{array}$	
β_0 β_1 β_2 β_3 β_4	θ 1 3 -3 2 -2	$\nu = \frac{\hat{\theta}_{IRLS}}{(\text{sd})}$ 1.02 (0.13) 3.09 (0.22) -3.10 (0.20) 2.17 (0.39) -2.06 (0.17)	$\begin{array}{c} \begin{array}{c} 2.5 \\ \hline \hat{\theta}_{Opt_1} \\ (\text{sd}) \\ \hline 1.02 \\ (0.13) \\ 3.09 \\ (0.23) \\ \mathbf{-3.10} \\ (0.20) \\ 2.17 \\ (0.39) \\ \mathbf{-2.06} \\ (0.17) \\ \end{array}$	$\begin{array}{c} \hat{\theta}_{Opt_2} \\ (\text{sd}) \\ \hline \textbf{0.73} \\ (0.16) \\ \textbf{2.43} \\ (0.23) \\ \textbf{-2.42} \\ (0.24) \\ \textbf{1.76} \\ (0.37) \\ \textbf{-1.61} \\ (0.16) \end{array}$
$egin{array}{c} eta_0 & & & & & & & & & & & & & & & & & & &$	θ 1 3 -3 2 -2 0.91	$\nu = \frac{\hat{\theta}_{IRLS}}{(\text{sd})}$ 1.02 (0.13) 3.09 (0.22) -3.10 (0.20) 2.17 (0.39) -2.06 (0.17) 0.95	$\begin{array}{c} \begin{array}{c} 2.5 \\ \hline \hat{\theta}_{Opt_1} \\ (\text{sd}) \\ \hline 1.02 \\ (0.13) \\ 3.09 \\ (0.23) \\ \mathbf{-3.10} \\ (0.20) \\ 2.17 \\ (0.39) \\ \mathbf{-2.06} \\ (0.17) \\ 0.94 \end{array}$	$ \hat{\theta}_{Opt_2} \\ (sd) \\ \hline \textbf{0.73} \\ (0.16) \\ \textbf{2.43} \\ (0.23) \\ \textbf{-2.42} \\ (0.24) \\ \textbf{1.76} \\ (0.37) \\ \textbf{-1.61} \\ (0.16) \\ \textbf{0.69} \\ \hline \textbf{0.69} $

Table 1: Comparison of the estimated parameters from three methods $(\hat{\theta}_{IRLS}, \hat{\theta}_{Opt_1}, \hat{\theta}_{Opt_2})$. Results are obtained using 20 bootstrap replications. θ denotes the true parameter values. Values in parentheses are standard errors.



Figure 3: Comparison of the methods in terms of their computational timings for a data with $\nu = 0.5$ with increasing sample sizes. While $IRLS_1$ is initialized at $\nu = 0.2$, $IRLS_2$ is initialized with a ν computed from a sample model. Right panel removes Opt_1 to provide clearer separation of other methods.



Figure 4: Comparison of the methods in terms of their computational timings for a data with $\nu = 4$ with increasing sample sizes. While $IRLS_1$ starts at $\nu = 0.2$, $IRLS_2$ starts with a ν computed from a sample model. Right panel removes Opt_1 to provide clearer separation of other methods.

5. A CMP Generalized Additive Model

5.1. Background

A generalized additive model (GAM) [14] is a generalized linear model (GLM; [20]) with a linear predictor involving smooth functions of covariates:

$$g\{E[y_i]\} = \boldsymbol{x}_i^* \boldsymbol{\theta}^* + \sum_{j=1}^p f_j(x_{ij}), \quad i = 1, \dots, n$$
 (16)

where $g(\cdot)$ is a smooth monotonic and twice differentiable link function, \boldsymbol{x}_i^* is the *i*th row of X^* , which is the model matrix for the parametric model components, $\boldsymbol{\theta}^*$ is the parameter vector, and f_j are the smooth functions of the covariate \boldsymbol{x}_j and they are subject to identifiability constraints, such as $\sum_{i=1}^{n} f_j(x_{ij}) = 0$ for all j. There exist multiple methods for estimating the smooth functions f_j [14, 13, 35, 25, 24, 10]. Among those the two most popular approaches that use spline bases are smoothing splines [13, 14] and penalized splines [35].

The smoothing splines approach uses the *backfitting* algorithm to estimate the smooth functions. The algorithm can be used within the IRLS framework by incorporating another inner loop to estimate smooth functions at every iteration. The *backfitting* algorithm is elegant as it has the flexibility to incorporate a wide variety of smoothing methods for component estimation. The convergence of the backfitting algorithm and its related properties can be found in [2]. However, as suggested by Wood [35], Gu and Wahba [11], it is not easy to efficiently integrate the estimation of the smoothing parameter into the model estimation framework. Traditional methods such as cross validation are often prohibitive because of the high computational cost involved in the search for multiple smoothing parameters.

The penalized splines approach has become a popular choice for fitting additive models due to the availability of a variety of methods with efficient implementations [40]. The idea is to represent each f_j with intermediate rank spline-type basis expansions, in which case the model becomes the GLM. In order to avoid overfitting, the model is estimated by penalized likelihood maximization. In practice, the penalized maximum likelihood is maximized by penalized iteratively re-weighted least squares (P-IRLS). In particular, the GAM is fitted by iteratively minimizing

$$\|\sqrt{W^{(k)}}(T^{(k)} - X\boldsymbol{\beta})\|^2 + \sum_j \eta_j \boldsymbol{\beta}^T S_j \boldsymbol{\beta} \text{ w.r.t. } \boldsymbol{\beta}.$$
 (17)

 $T^{(k)}$ denotes the adjusted response variable and $W^{(k)}$ denotes the weights at the *k*th iteration of the P-IRLS algorithm. The S_j are matrices of known coefficients such that $\boldsymbol{\beta}^T S_j \boldsymbol{\beta}$ measures the roughness of f_j . The η_j are smoothing parameters that control the trade-off between fit and smoothness and their selection can be achieved by minimizing the Generalized Cross Validation (GCV) score, AIC, or another criterion [35, 34, 37].

There are two types of computational methods available for the estimation of η_j . (i) *Performance iteration* uses the fact that at each P-IRLS step a working penalized linear model is estimated and the smoothing parameter estimation can be performed on each such working model. (ii) In *outer iteration* the P-IRLS algorithm is iterated to convergence for each trial set of smoothing parameters and the GCV or AIC scores are only evaluated on convergence [35].

5.2. Implementation of the CMP Generalized Additive Model

Similar to the CMP regression, the CMP generalized additive model can be formulated as

$$\ln(\lambda_i) = \boldsymbol{x}_i^* \boldsymbol{\theta}^* + \sum_{j=1}^p f_j(x_{ij}), \qquad (18)$$

$$\ln(\nu_i) = \mathbf{z}_i^* \boldsymbol{\delta}^* + \sum_{j=1}^k m_j(z_{ij}),$$
(19)

for i = 1, ..., n, where $\boldsymbol{\theta}^*$ and $\boldsymbol{\delta}^*$ are the parameter vectors for the parametric part of $\ln(\lambda_i)$ and $\ln(\nu_i)$ respectively. The smooth functions f_j and m_j are the smooth functions for the covariates \boldsymbol{x}_j and \boldsymbol{z}_j and are subject to the identifiability constraints. For the sake of easy presentation, we omit the strictly parametric part from now onwards.

We consider the *performance iteration* method as it is very efficient and computes faster than the *outer iteration* method. Although not common, there is some evidence that *performance iteration* faces some convergence issues because the objective function for the smoothing parameters keeps changing with every iteration of P-IRLS [35, 38]. In contrast, the *outer iteration* method is more robust to convergence related issues but usually takes longer time to compute. More importantly, on convergence, it requires some derivatives to estimate the smooth parameters. The only way to get derivatives for the CMP distribution is to use numerical derivatives and they are often prone to errors due to the normalizing constant in the likelihood which is an infinite series. For this reason, we consider only the *performance iteration* method to estimate CMP GAM and leave the *outer iteration* method for future research.

Now, each P-IRLS iteration involves minimizing the following two objective functions:

$$\|\sqrt{W_1^{(k)}}(T_1^{(k)} - X\boldsymbol{\beta})\|^2 + \boldsymbol{\beta}^T H_1\boldsymbol{\beta} + \sum_j \eta_{1j}\boldsymbol{\beta}^T S_{1j}\boldsymbol{\beta} \text{ w.r.t. } \boldsymbol{\beta}, \qquad (20)$$

$$\|\sqrt{W_2^{(k)}}(T_2^{(k)} - Z\boldsymbol{\gamma})\|^2 + \boldsymbol{\gamma}^T H_2 \boldsymbol{\gamma} + \sum_j \eta_{2j} \boldsymbol{\gamma}^T S_{2j} \boldsymbol{\gamma} \text{ w.r.t. } \boldsymbol{\gamma}.$$
(21)

 W_1, W_2 are the weight matrices, η_{1j}, η_{2j} are the smooth parameters for the $\ln(\lambda_i)$ and $\ln(\nu_i)$ models. The matrices H_1, H_2 are fixed positive semi definite penalty matrices which allow for multiple extensions to the GAMs such as ridge penalties under suitable constraints [34]. $T_1^{(k)}$ and $T_2^{(k)}$ are defined similarly as in the Algorithm 1.

Given the smoothing parameters η_{1j} , η_{2j} the objective functions in (21) are solved using any penalized least squares type of methodology. However, the smoothing parameters need to be estimated here. We consider the GCV method, in which the smoothing parameters are chosen to minimize

$$V_{1g} = \frac{n \parallel T_1 - A_1 T_1 \parallel^2}{[\operatorname{tr}(I - \gamma_1 A_1)]^2}, \quad V_{2g} = \frac{n \parallel T_2 - A_2 T_2 \parallel^2}{[\operatorname{tr}(I - \gamma_2 A_2)]^2},$$

respectively. A_1, A_2 are the influence or hat matrices of the corresponding fitting problems and they depend on the smoothing parameters. The parameters γ_1, γ_2 are sometimes used to inflate the GCV objective function to make sure that the models are smoother [34, 3]. There are efficient algorithms available in the mgcv [40, 37, 34] package (e.g. magic function) to estimate the smoothing parameters along with the model parameters and we use them in our implementation for CMP GAM.

The inference for the spline regression coefficients in GAM is developed using a bayesian view of the smoothing process, in which the smoothing penalties are induced by improper Gaussian priors on β , γ and $\hat{\beta}$, $\hat{\gamma}$ are also the modes of the posterior densities of β , γ [33, 30, 21]. Please refer to Wood [36, 39] for more details. Based on the results from Wood [36, 38], the large sample posterior distribution for the regression spline coefficients in CMP GAM is

$$\boldsymbol{\beta}|\boldsymbol{\eta}_1, Y \sim \mathcal{N}(\hat{\boldsymbol{\beta}}, \Sigma_{\boldsymbol{\beta}}), \quad \boldsymbol{\gamma}|\boldsymbol{\eta}_2, Y \sim \mathcal{N}(\hat{\boldsymbol{\gamma}}, \Sigma_{\boldsymbol{\gamma}}),$$

where $\Sigma_{\beta} = (X^T W_1 X + \sum_{j=1}^p \eta_{1j} S_{1j})^{-1} \phi_1$, $\Sigma_{\gamma} = (Z^T W_2 Z + \sum_{j=1}^q \eta_{2j} S_{2j})^{-1} \phi_2$ with ϕ_1 , ϕ_2 as the estimated scale parameters and W_1 , W_2 are the weight matrices at the convergence of the P-IRLS algorithm.

From a practitioner point of view, sometimes it is required to check whether a particular smooth function is significant in the model or can be discarded. More formally, to test the null hypotheses that $f_j = 0$ or $m_j = 0$ for any j, Wood [39] proposed the Wald test statistic (wt_r) , which under the null hypothesis and with a large sample follows a chi-square distribution. The Wald test statistic is defined as $wt_r = \hat{f}_j^T V_{f_j}^{r-} \hat{f}_j$ where $V_{f_j}^{r-}$ is the rankr pseudo-inverse of $V_{f_j} = \boldsymbol{x}_j V_\beta \boldsymbol{x}_j^T$. The authors have suggested that naive choices of r lead to low power or an incorrect null distribution for p-values and using the effective degrees of freedom for r is a better choice. For more details please refer to Wood [39].

The inference procedure proposed for CMP GAM is based on the marginal likelihood. If needed, one could also develop the conditional inference framework by suitably adjusting the covariance matrices. However, it is not guaranteed that the rank-r pseudo inverse remains optimal. Further, it is also possible that after the correction, the covariance matrix may not be positive definite. For this reason, we only use the inference procedure based on marginal likelihood. In practice, we found that although there are some changes in p-values, the results remain same whether we use inference procedure developed based on marginal likelihood or conditional likelihood.

6. Simulation Study for the CMP Generalized Additive Model

We conducted a simulation study to evaluate the usefulness of the CMP GAM for fitting non linear terms. We consider two examples for our simulation study; in the first example we choose one fixed value for ν such as 0.5 or 2.5 and in the second example we simulate ν using a nonparametric smooth function. While the first example is considered to compare CMP GAM with other models such as NB GAM, Poisson GAM and CMP GLM (CMP Regression), the second example is considered to showcase the flexibility of CMP GAM allowing dispersion to vary non linearly across observations.

6.1. Example 1

Inspired by the four uni-variate example [11, 40], we simulated data from a CMP GAM, with sample size 500, as following:

- Simulate x_1, x_2, x_3 and x_4 from a standard uniform U[0, 1] distribution.
- Consider the functions $f_1 = \sin(\pi x_1), f_2 = \exp(x_2), f_3 = 0.02x_3^2(1 x_3) + (0.5x_3)^2(1 x_3)^3$ and $f_4 = x_4$.
- Calculate $f = af_1(x_1) + bf_2(x_2) + cf_3(x_3)$, where a, b, c are pre-specified constants.
- Set $\lambda = \exp(f)$ and simulate data for a fixed ν .

Although we simulated 4 covariates, we have used only the first three (x_1, x_2, x_3) to compute λ , which is used to simulate the dependent variable. We still use x_4 to estimate the models. This allows us to check how the proposed method deals with over-specification. In an ideal case, the estimated model would identify x_4 as non significant. We consider two different values for ν (= 0.5, 2) to capture both over dispersion and under dispersion scenarios. Unlike standard GLMs, simulating data from the CMP distribution requires a careful consideration of parameters λ, ν . The scale of the dependent variable is in the range of $\lambda^{1/\nu}$. If we choose $\nu < 1$ (over dispersion), then we must consider smaller values for λ otherwise we will generate only very large counts (e.g. 1500, 2000), resulting in a data set that is less useful for illustrating count data models. Similarly, for $\nu > 1$, we must consider larger values for λ to avoid generating only very small counts such as 0,1,2, resulting in an extremely under dispersed data set that is less useful for comparing against count models such as Poisson or Negative Binomial GAM. To avoid the aforementioned problems, we choose two sets of different values for constants a, b, c such as $\{0.2, 0.5, -0.5\}$ for $\nu = 0.5$ and $\{1, 1, 1\}$ for $\nu = 2.5$.

We also considered a Poisson GAM and a CMP GLM model for comparison with the CMP GAM. We used the *mgcv* [40] package in R to estimate the Poisson GAM and NB GAM as they are also implemented using penalized splines. Although the current implementations of both Poisson GAM and NB GAM do not use the P-IRLS algorithm by default, we specifically used P-IRLS (*performance iteration*) algorithm to estimate these models to provide a fair comparison.

For each of the two scenarios ($\nu = 0.5, 2.5$) we used 50 bootstrap replications and recorded the significance levels for each nonparametric term in the model. We also recorded and compared their AIC. While the model equation for an additive model is $y \sim s(x_1) + s(x_2) + s(x_3) + s(x_4)$, where $s(\cdot)$ is the smooth function, the model for the CMP GLM is $y \sim x_1 + x_2 + x_3 + x_4$ without any smooth functions for covariates.

The simulation results are summarized in Table 2. The top table (a) describes the results for $\nu = 0.5$ and the bottom table (b) describes the results for $\nu = 2.5$. Since we estimated 50 models for each bootstrap data, we plotted the AIC values for each data for all the models in Figure A.1 in the Appendix. From the top plot, for $\nu = 0.5$, it can be seen that AIC for CMP GAM is consistently better than the AIC for both Poisson GAM and CMP GLM. The AIC values for NB GAM are closer to the AIC values for CMP GAM which indicates that their fits are reasonably close. While both CMP and NB GAMs declared $s(x_3)$ as non significant, Poisson GAM identified it as significant at least half the times. Similarly, for the smooth term $s(x_4)$, which is not part of the true model, Poisson GAM declared it as significant at least half the times. Ideally, CMP GLM should not produce any significant coefficients because of the nonlinear terms in the true model. However, it can be observed that x_2 is significant and this is because the function $f_2 = \exp(x_2)$ is approximately equal to $1 + x_2$ (because $x_2 \in [0, 1]$).

The results from Table 2 (b) can be interpreted similarly. Not surprisingly, in terms of the model fit, both Poisson GAM and NB GAM are very close but not better than the CMP GAM which is evident from the bottom plot in Figure A.1 in the Appendix. This is because of their inability to handle under dispersion.

6.2. Example 2

We simulate data similar to the procedure in Section 6.1. The procedure is as follows:

- We consider the same functions for f_1 , f_2 and f_3 but a different function for $f_4 = 2x_4 x_4^2$.
- Calculate $f = f_2 + f_3 + 2f_4$ and set $\lambda = \exp(f)$.
- Calculate $g = f_1$ and set $\boldsymbol{\nu} = \exp(g)$.
- Simulate y_i from $CMP(\lambda_i, \nu_i)$.

Since the dispersion parameter is generated using smooth function, the observations will have different dispersions. In the simulated data, the dispersion

		#Significant (of S	50 bootstraps)		
	$(\le 0.001, \le 0.01, \le 0.5, \le 0.1, n.s)$				
	cmp - gam	poisson-gam	nb-gam	cmp - glm	
$s(x_1)$	$(46,\!4,\!0,\!0,\!0)$	$(50,\!0,\!0,\!0,\!0)$	$(42,\!8,\!0,\!0,\!0)$	$(0,\!0,\!0,\!1,\!49)$	
$s(x_2)$	$(50,\!0,\!0,\!0,\!0)$	$(50,\!0,\!0,\!0,\!0)$	$(50,\!0,\!0,\!0,\!0)$	$(50,\!0,\!0,\!0,\!0)$	
$s(x_3)$	(0,1,1,4,44)	(1,4,13,6, 26)	(0,0,3,3,44)	$(0,\!0,\!2,\!3,\!45)$	
$s(x_4)$	(0,0,1,5,44)	(2,8,8,7,25)	(0,0,1,3,46)	(0,1,3,2,44)	
		Estimation	and Fit		
	cmp - gam	poisson - gam	nb-gam	cmp - glm	
$\nu \text{ or } \theta$	-0.78 (0.18)		9.92 (2.05)	-0.81 (0.08)	
AIC	$\boldsymbol{2726.53}$	2830.73	2734.18	2757.16	
		(a) For $\nu = 0$	5		
		(a) 101 $\nu = 0.5$	5.		
		#Significant (of	50 boostraps)		
	(:	$\frac{(a) \text{ for } \nu = 0.0}{\# \text{Significant (of } \le 0.001, \le 0.01, \le 0.00, 0.0$	50 boostraps) $(0.5, \leq 0.1, n.s)$	s)	
	(j	#Significant (of $\leq 0.001, \leq 0.01, \leq poisson - gam$	50 boostraps) $50 ext{ boostraps}$ $50 ext{ constraints}$ $50 ext{ boostraps}$ $50 ext{ boostraps}$ 50	s) cmp - glm	
$s(x_1)$	(mp - gam) (50,0,0,0,0)	$# Significant (of \leq 0.001, \leq 0.01, \leq poisson - gam (33, 15, 2, 0, 0)$	50 boostraps) 50 boostraps) $(50,5, \le 0.1, n.s, nb - gam)$ (50,0,0,0,0)	cmp - glm (0,0,0,2,48)	
$ s(x_1) \\ s(x_2) $	(mp - gam) (50,0,0,0,0) (50,0,0,0,0)	$ \frac{\text{#Significant (of}}{\text{#Significant (of}} \leq 0.001, \leq 0.01, \leq 0.01, \leq 0.01, \leq 0.01, \leq 0.001, \leq 0.00, 0.0, (33, 15, 2, 0, 0), (33, 15, 2, 0, 0), (50, 0, 0, 0, 0, 0)) $	50 boostraps) 50 boostraps) (50,0,0,0,0) (50,0,0,0,0)	s) cmp - glm (0,0,0,2,48) (50,0,0,0,0)	
$ \begin{array}{c} s(x_1)\\ s(x_2)\\ s(x_3) \end{array} $	$(mp - gam) \\ (50,0,0,0,0) \\ (50,0,0,0,0) \\ (0,0,2,5,43)$	$ \begin{array}{l} \# \text{Significant (of} \\ \leq 0.001, \leq 0.01, \leq \\ poisson - gam \\ (33, 15, 2, 0, 0) \\ (50, 0, 0, 0, 0) \\ (0, 0, 0, 0, 50) \end{array} $	50 boostraps) 50 boostraps) (50,0,0,0,0,0) (50,0,0,0,0,0) (0,0,2,5,43)	s) cmp - glm (0,0,0,2,48) (50,0,0,0,0) (0,0,2,2,46)	
$s(x_1) \\ s(x_2) \\ s(x_3) \\ s(x_4)$	$(mp - gam) \\ (50,0,0,0,0) \\ (50,0,0,0,0) \\ (0,0,2,5,43) \\ (0,0,0,8,42)$	$ \frac{\text{(a) 101 V} = 0.4}{\#\text{Significant (of}} \leq 0.001, \leq 0.01, \leq 0.01, \leq 0.01, \leq 0.01, \leq 0.01, \leq 0.00, 0.0, (33, 15, 2, 0, 0), (33, 15, 2, 0, 0), (50, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0)) $	$50 \text{ boostraps}) \\ 50 \text{ boostraps}) \\ 50.5, \le 0.1, n.s \\ nb - gam \\ (50, 0, 0, 0, 0) \\ (50, 0, 0, 0, 0) \\ (0, 0, 2, 5, 43) \\ (0, 0, 1, 8, 41) \end{cases}$	s) cmp - glm (0,0,0,2,48) (50,0,0,0,0) (0,0,2,2,46) (0,0,2,2,46)	
$s(x_1) \\ s(x_2) \\ s(x_3) \\ s(x_4)$	$(mp - gam) \\ (50,0,0,0,0) \\ (50,0,0,0,0) \\ (0,0,2,5,43) \\ (0,0,0,8,42) \\ (0,0,0,8,42) \\ (0,0,0,0,0,0) \\ (0,0,0,0,0,0) \\ (0,0,0,0,0,0) \\ (0,0,0,0,0,0) \\ (0,0,0,0,0,0) \\ (0,0,0,0,0,0) \\ (0,0,0,0,0,0) \\ (0,0,0,0,0,0) \\ (0,0,0,0,0,0) \\ (0,0,0,0,0,0) \\ (0,0,0,0,0,0) \\ (0,0,0,0,0,0) \\ (0,0,0,0,0,0) \\ (0,0,0,0,0,0) \\ (0,0,0,0,0,0) \\ (0,0,0,0,0,0,0) \\ (0,0,0,0,0,0,0) \\ (0,0,0,0,0,0,0) \\ (0,0,0,0,0,0,0) \\ (0,0,0,0,0,0,0) \\ (0,0,0,0,0,0,0,0,0) \\ (0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,$	$\frac{(a) 1017 = 0.0}{\# \text{Significant (of}} \leq 0.001, \leq 0.01, \leq 0.01, \leq 0.01, \leq 0.01, \leq 0.00, 0.0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0$	5. 50 boostraps) $50.5, \le 0.1, n.s$ nb - gam (50,0,0,0,0) (50,0,0,0,0) (0,0,2,5,43) (0,0,1,8,41) and Fit	$\begin{array}{c} s) \\ cmp - glm \\ (0,0,0,2,48) \\ ({\bf 50},0,0,0,0) \\ (0,0,2,2,46) \\ (0,0,2,2,46) \end{array}$	
$s(x_1) \ s(x_2) \ s(x_3) \ s(x_4)$	(mp - gam) = (50,0,0,0,0) = (50,0,0,0,0) = (0,0,2,5,43) = (0,0,0,8,42) = cmp - gam	$\frac{\#\text{Significant (of}}{\#\text{Significant (of}} \le 0.001, \le 0.01, \le 0.01, \le 0.01, \le 0.01, \le 0.01, \le 0.00, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, $	50 boostraps) 50 boostraps) (50,0,0,0,0) (50,0,0,0,0) (0,0,2,5,43) (0,0,1,8,41) and Fit nb - gam	$cmp - glm \\ (0,0,0,2,48) \\ (50,0,0,0,0) \\ (0,0,2,2,46) \\ (0,0,2,2,46) \\ cmp - glm$	
$s(x_1)$ $s(x_2)$ $s(x_3)$ $s(x_4)$ $\nu \text{ or } \theta$	(mp - gam) $(50,0,0,0,0)$ $(50,0,0,0,0)$ $(0,0,2,5,43)$ $(0,0,0,8,42)$ $cmp - gam$ $0.84(0.09)$	$\frac{(a) 10177 = 0.0}{\#\text{Significant (of}} \leq 0.001, \leq 0.01, \leq 0.01, \leq 0.01, \leq 0.01, \leq 0.01, \leq 0.00, 0.0, 0.0, (33, 15, 2, 0, 0), (50, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0$	50 boostraps) 50 boostraps) (50,0,0,0,0) (50,0,0,0,0) (0,0,2,5,43) (0,0,1,8,41) and Fit nb - gam 10000(0)	$cmp - glm \\ (0,0,0,2,48) \\ (50,0,0,0,0) \\ (0,0,2,2,46) \\ (0,0,2,2,46) \\ cmp - glm \\ 0.77(0.05)$	
$s(x_1)$ $s(x_2)$ $s(x_3)$ $s(x_4)$ $\nu \text{ or } \theta$ AIC	(mp - gam) $(50,0,0,0,0)$ $(50,0,0,0,0)$ $(0,0,2,5,43)$ $(0,0,0,8,42)$ $cmp - gam$ $0.84(0.09)$ 1438.74	$\frac{(a) 161 \ \nu = 0.1}{\# \text{Significant (of}} \leq 0.001, \leq 0.01, \leq 0.01, \leq 0.01, \leq 0.01, \leq 0.01, \leq 0.00, 0.00, (33, 15, 2, 0, 0), (50, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0, 0), (0, 0, 0, 0), (0, 0, 0, 0), (0, 0, 0, 0), (0, 0, 0, 0), (0, 0, 0, 0), (0, 0, 0, 0), (0, 0, 0)$	5. 50 boostraps) $(50,5) \le 0.1, n.s, nb - gam$ (50,0,0,0,0) (50,0,0,0,0) (0,0,2,5,43) (0,0,1,8,41) and Fit nb - gam 10000(0) 1568.13	$cmp - glm \\ (0,0,0,2,48) \\ (50,0,0,0,0) \\ (0,0,2,2,46) \\ (0,0,2,2,46) \\ cmp - glm \\ 0.77 (0.05) \\ 1484.84$	
$s(x_1)$ $s(x_2)$ $s(x_3)$ $s(x_4)$ $\nu \text{ or } \theta$ AIC	(mp - gam) $(50,0,0,0,0)$ $(50,0,0,0,0)$ $(0,0,2,5,43)$ $(0,0,0,8,42)$ $cmp - gam$ $0.84(0.09)$ 1438.74	$\frac{(a) 1017 = 0.0}{\# \text{Significant (of}} \leq 0.001, \leq 0.01, \leq 0.01, \leq 0.01, \leq 0.01, \leq 0.01, \leq 0.00, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, $	50 boostraps) 50 boostraps) $(50,0,5, \le 0.1, n.s, nb - gam)$ (50,0,0,0,0) (0,0,2,5,43) (0,0,1,8,41) and Fit nb - gam 10000(0) 1568.13	$\begin{array}{c} cmp - glm \\ (0,0,0,2,48) \\ ({\bf 50},0,0,0,0) \\ (0,0,2,2,46) \\ (0,0,2,2,46) \\ cmp - glm \\ {\bf 0.77}(0.05) \\ {\bf 1484.84} \end{array}$	

Table 2: Comparison of coefficient significance level and fit among CMP GAM, Poisson GAM, NB GAM and CMP GLM with 50 bootstrap replications. The model for gams is $y \sim s(x_1) + s(x_2) + s(x_3) + s(x_4)$ and for the regression $y \sim x_1 + x_2 + x_3 + x_4$. For the Estimation and Fit results the numbers in parenthesis are standard errors.

	#Significant (of 50 bootstraps) ($\leq 0.001, \leq 0.01, \leq 0.5, \leq 0.1, n.s$)		
	cmp-gam		
	m_1	m_2	
$\boldsymbol{\lambda}$ model:			
$s(x_4)$	$(50,\!0,\!0,\!0,\!0)$	$(50,\!0,\!0,\!0,\!0)$	
$s(x_2)$	$(50,\!0,\!0,\!0,\!0)$	$(50,\!0,\!0,\!0,\!0)$	
$s(x_3)$	(0,1,3,1,45)	(1,4,9,5,31)	
u model:			
γ_0	$(50,\!0,\!0,\!0,\!0)$	$(50,\!0,\!0,\!0,\!0)$	
$s(x_1)$ or γ_1	$(50,\!0,\!0,\!0,\!0)$	$(2,\!1,\!8,\!2,\!37)$	
AIC	1968.62	2891.60	

Table 3: Comparison of coefficient significance levels and fit between CMP GAMs with two different models for $\boldsymbol{\nu}$ based on 50 bootstrap replications. The models for $\boldsymbol{\nu}$ are $m_1 = -s(x_1); m_2 = -x_1$. The model for $\boldsymbol{\lambda}$ is $-s(x_4) + s(x_2) + s(x_3)$ for both m_1 and m_2 .

 (ν_i) varies from 0.7 to 2.5. We now use the data and fit two different models. We estimate ν non parametrically in the first model and parametrically in the second model.

Table 3 contains the significance levels results from the above mentioned two models using 50 bootstrap replications. As expected modeling ν non parametrically yields better fit. As seen from the results in Section 6.1, the parametric model is not able to capture the underlying nonlinear function for ν .

7. Example: Modeling Data from Bike Sharing Systems

To illustrate the use of the CMP GAM on real data, we model data from a bike sharing application. Bike sharing systems are a new generation of traditional bike rental services where the entire process that includes membership, bike rental and return has become automated. Through these systems, users are able to easily rent a bike from a particular station and return it to another location. There is a need for bike sharing programs to effectively understand the factors that influence demand so that they can better maintain inventory, schedule repairs, and manage resources. We therefore use a GAM model in this context.

Data collected by bike sharing systems typically include information on each trip taken (time stamps and locations of rental and return) and sometimes also information on the rider. Several datasets from real bike sharing systems are publicly available. We use the data available from Fanaee-T and Gama [6] on rides in Washington, DC between 2011-2012. The data is available in two formats: daily and hourly number of rentals. We chose the hourly data as it is more complex and better illustrates the new models that we introduce. The data includes information about the number of rides by casual users and registered users for every hour during the years 2011 and 2012. In addition, it also includes external information such whether the hour is on a weekday, a working day or a holiday, the weather situation (clear/cloudy/rainy), temperature, and wind speed. These external factors are considered detrimental to the demand for bikes. For a full list of attributes and their descriptions please refer to the Table A.1 in the Appendix.

We considered the following model for both counts of casual and registered users:

$$\ln(\boldsymbol{\lambda}) = \beta_0 + \beta_1 hour + \beta_2 holiday + \beta_3 weekday + \beta_4 weathersit + s(atemp) + s(hum) + s(windspeed) + s(day)$$
$$\ln(\boldsymbol{\nu}) = \gamma_0.$$

Since the attributes *atemp* and *temp* are highly correlated, we included only *atemp* in the model. We kept control variables such as *holiday* and *weathersit* as parametric terms and included other continuous variables of interest such as *hum* and *windspeed* as nonparametric components.

For this study we have only considered the January 2012 data. The same analysis can be repeated for every month or for every season. Since we have two dependent variables of interest we fit two models to this data. The first model is for the number of hourly rentals for registered users and the second model is for the number of hourly rentals for casual users. For the sake of comparison we also fit a Poisson GAM and NB GAM using the *mgcv* [40] package. As in the simulation study in Section 6, we made sure that both the Poisson and NB GAMs are estimated via the *performance iteration* algorithm.

The coefficient significance results are described in Table 4. For brevity we did not include the coefficient significance results for the parametric com-

ponents. We only reported the results for the non parametric components and the findings from these are similar to the findings from the parametric part. From the registered users results in Table 4, it can be observed that while in CMP GAM all the smooth variables are significant, in CMP and NB GAM the smooth variable *windspeed* is not significant. Since the data exhibits high levels of dispersion, which is evident from γ_0 (< 0 over dispersion; > 0 under dispersion), the significance results from NB GAM and CMP are similar. Based on the *AIC* values, not surprisingly, both CMP and NB GAM fit better than Poisson GAM.

]	Registered Users	
	cmp - gam	poisson-gam	nb-gam
	(edf)	(edf)	(edf)
s(day)	5.63^{**}	7.91***	3.40^{**}
s(atemp)	6.85^{***}	8.86***	7.11^{***}
s(hum)	7.96^{**}	8.97^{***}	8.31^{***}
s(windspeed)	2.31	8.90***	1.00
$\gamma_0 \text{ or } \theta$	-3.03^{***}		3.82^{***}
AIC	7413.55	18639.87	7563.97
RMSE	49.41	49.46	59.10
$ \hline \ \ \ \ \ \ \ \ \ \ \ \ \$	0.01, *p < 0.05		
		Casual Users	
	cmp - gam	poisson-gam	nb-gam
	(edf)	(edf)	(edf)
s(day)	8.90***	8.96***	8.82***
s(atemp)	1.00^{***}	5.12^{***}	1.00^{***}
s(hum)	8.05**	8.81***	1.78^{***}
s(windspeed)	1.68	6.11^{***}	1.00
$\gamma_0 \text{ or } \theta$	-1.36^{***}		3.20***
AIC	3990.84	4713.76	4044.07
RMSE	6.59	6.91	9.59

***p < 0.001, **p < 0.01, *p < 0.05

Table 4: Comparison of CMP GAM, Poisson GAM and NB GAM in terms of coefficient significance and fit for # rentals for both registered (top) and casual (bottom) users in January 2012.

Similar results are seen in Table 4 for the casual users. In Poisson GAM all the smooth variables are significant whereas in CMP GAM and NB GAM



Figure 5: Partial plots for the smooth variables for CMP GAM. The dependent variable is # rentals for registered users in January 2012.



Figure 6: Partial plots for the smooth variables for CMP GAM. The dependent variable is #rentals for casual users in January 2012.

the smooth variable *windspeed* is not significant. Similar to the registered users, the data is over dispersed. These results therefore also highlight potential inference errors when fitting a Poisson GAM to data with excessive dispersion.

To provide more meaningful interpretations we use partial plots. The partial plots for the CMP GAM for the registered users are shown in Figure 5. The smooth variables *day* and *atemp* exhibit an increasing trend while *humidity* exhibits a decreasing trend. This is expected because when there is high humidity people may not show much interest in riding bikes. Further, in a winter month like January, high temperature (or sunny day) encourages people to ride bikes. The insignificance of *windspeed* is evident from the plot as the smooth curve is close to zero throughout the range. To draw more meaningful interpretations we would need more domain knowledge.

For casual users one can draw interpretations from the partial plots in Figure 6. The results are similar to Figure 5 except for the smooth variable day, which shows a cyclical pattern that might indicate the high demand for bikes during weekends.

Finally, from an actual model fit perspective, both CMP and NB GAMs perform reasonably well. Figure A.2 in the Appendix compares fitted values and residuals from CMP, Poisson and NB GAMs. Since the data is hourly, we plotted the fitted values for every hour. We joined the data points to provide better visualization. We see that for both registered and casual users, the CMP and Poisson GAMs fitted values are close to the actual values while the NB GAM fitted values are not. The identical performance of the CMP and Poisson GAMs in terms of fitted values is expected as they differ only in terms of standard errors rather than point predictions.

In summary, CMP GAM can be a valuable addition for modeling count data. Despite its computational complexity, CMP GAM is very flexible as it can handle both over dispersion and under dispersion which existing methods fail to handle. Although the bike sharing data did not exhibit under dispersion, there are plenty of data sets that do. Moreover, when the researcher does not know the dispersion type (over or under) prior to modeling, CMP GAM is a safe option.

8. Conclusions and Future Directions

We introduced a flexible estimation framework for estimating a CMP GLM model that is based on the IRLS approach. This framework allows the CMP distribution to join other existing GLMs where IRLS is used for efficient estimation as well as for various modeling enhancements. This framework can be further developed to extend methods such as the lasso.

While the IRLS algorithm for CMP GLM is computationally intensive compared to an ordinary Poisson model, the computation time can be reduced by suitably parallelizing some of the computations such as the calculation of cumulants. Such parallel computing will be beneficial especially with large samples.

In this work we explored fitting additive models with penalized splines. We considered the *performance iteration* method to fit the model as it is based on the P-IRLS algorithm. One possible extension is to develop the *outer iteration* method using the Newton algorithm. The numerical derivatives required for the Newton algorithm are computationally slower and not very stable, thereby, requiring new, efficient implementation.

Acknowledgements

The authors were partially funded by research grant 105-2410-H-007-034-MY3 by the Ministry of Science and Technology in Taiwan. The authors would like to thank the AE and the two referees for their numerous suggestions that substantially improved the content and the presentation of the paper. The authors also gratefully acknowledge Prof. Li-Shan Huang and Prof. Satish Iyengar for their valuable feedback and suggestions on earlier versions of this manuscript.

References

- Abramowitz, M., Stegun, I. A., 1966. Handbook of mathematical functions. Vol. 55. Applied mathematics series.
- [2] Buja, A., Hastie, T. J., Tibshirani, R., 1989. Linear smoothers and additive models. The Annals of Statistics, 453–510.
- [3] Chambers, J. M., 1998. Programming with data: A guide to the S language. Springer Science & Business Media.
- [4] Daly, F., Gaunt, R. E., 2016. The Conway-Maxwell-Poisson distribution: distributional theory and approximation. ALEA, Lat. Am. J. Probab. Math. Stat. 13, 635–658.

- [5] Dominici, F., McDermott, A., Zeger, S. L., Samet, J. M., 2002. On the use of generalized additive models in time-series studies of air pollution and health. American journal of epidemiology 156 (3), 193–203.
- [6] Fanaee-T, H., Gama, J., 2014. Event labeling combining ensemble detectors and background knowledge. Progress in Artificial Intelligence 2, 113–127.
- [7] Gaunt, R. E., Iyengar, S., Olde Daalhuis, A. B., Simsek, B., 2016. An asymptotic expansion for the normalizing constant of the Conway-Maxwell-Poisson distribution. ArXiv:1612.06618v1.
- [8] Gillispie, S. B., Green, C. G., 2015. Approximating the Conway-Maxwell-Poisson distribution normalization constant. Statistics 49 (5), 1062–1073.
- [9] Green, P. J., 1984. Iteratively reweighted least squares for maximum likelihood estimation, and some robust and resistant alternatives. Journal of the Royal Statistical Society. Series B (Methodological), 149–192.
- [10] Gu, C., 2013. Smoothing spline ANOVA models. Vol. 297. Springer Science & Business Media.
- [11] Gu, C., Wahba, G., 1991. Minimizing gcv/gml scores with multiple smoothing parameters via the newton method. SIAM Journal on Scientific and Statistical Computing 12 (2), 383–398.
- [12] Gupta, R. C., Sim, S. Z., Ong, S. H., 2014. Analysis of discrete data by Conway-Maxwell-Poisson distribution. AStA Advances in Statistical Analysis 98 (4), 327–343.
- [13] Hastie, T. J., Tibshirani, R., 1986. Generalized additive models. Statistical science, 297–310.
- [14] Hastie, T. J., Tibshirani, R., 1990. Generalized additive models. Vol. 43. CRC Press.
- [15] Jensen, S. T., Johansen, S., Lauritzen, S. L., 1991. Globally convergent algorithms for maximizing a likelihood function. Biometrika 78 (4), 867– 877.

- [16] Keener, R. W., 2006. Statistical theory: notes for a course in theoretical statistics. Springer.
- [17] Lehmann, E. L., Casella, G., 2006. Theory of point estimation. Springer Science & Business Media.
- [18] Lu, T., Shiou, S., 2002. Inverses of 2× 2 block matrices. Computers & Mathematics with Applications 43 (1-2), 119–129.
- [19] Marschner, I. C., 2011. glm2: fitting generalized linear models with convergence problems. The R journal 3 (2), 12–15.
- [20] McCullagh, P., Nelder, J. A., 1989. Generalized linear models. Vol. 37. CRC press.
- [21] Nychka, D., 1988. Bayesian confidence intervals for smoothing splines. Journal of the American Statistical Association 83 (404), 1134–1143.
- [22] Olver, F. W., 2014. Asymptotics and special functions. Academic press.
- [23] Pollock, J., 2014. CompGLM: Conway-Maxwell-Poisson GLM and distribution functions. R package version 1.0. URL https://CRAN.R-project.org/package=CompGLM
- [24] Ramsay, J. O., 2006. Functional data analysis. Wiley Online Library.
- [25] Ruppert, D., Wand, M. P., Carroll, R. J., 2003. Semiparametric regression. No. 12. Cambridge university press.
- [26] Sellers, K. F., Lotze, T., Raim, A., 2017. COMPoissonReg: Conway-Maxwell-Poisson (COM-Poisson) Regression. R package version 0.4.1. URL https://github.com/lotze/COMPoissonReg
- [27] Sellers, K. F., Shmueli, G., 2010. A flexible regression model for count data. Annals of Applied Statistics 4 (2), 943–961.
- [28] Sellers, K. F., Shmueli, G., 2013. Data dispersion: Now you see it now you don't. Communications in Statistics-Theory and Methods 42 (17), 3134–3147.

- [29] Shmueli, G., Minka, T. P., Kadane, J. B., Borle, S., Boatwright, P., 2005. A useful distribution for fitting discrete data: revival of the Conway-Maxwell-Poisson distribution. Journal of the Royal Statistical Society: Series C (Applied Statistics) 54 (1), 127–142.
- [30] Silverman, B. W., 1985. Some aspects of the spline smoothing approach to non-parametric regression curve fitting. Journal of the Royal Statistical Society. Series B (Methodological), 1–52.
- [31] Steutel, F., 1985. Log-concave and log-convex distributions. Wiley StatsRef: Statistics Reference Online.
- [32] Stieb, D. M., Judek, S., Burnett, R. T., 2003. Meta-analysis of timeseries studies of air pollution and mortality: update in relation to the use of generalized additive models. Journal of the Air & Waste Management Association 53 (3), 258–261.
- [33] Wahba, G., 1983. Bayesian" confidence intervals" for the cross-validated smoothing spline. Journal of the Royal Statistical Society. Series B (Methodological), 133–150.
- [34] Wood, S. N., 2004. Stable and efficient multiple smoothing parameter estimation for generalized additive models. Journal of the American Statistical Association 99 (467), 673–686.
- [35] Wood, S. N., 2006. Generalized Additive Models: an introduction with R. CRC press.
- [36] Wood, S. N., 2006. On confidence intervals for generalized additive models based on penalized regression splines. Australian & New Zealand Journal of Statistics 48 (4), 445–464.
- [37] Wood, S. N., 2008. Fast stable direct fitting and smoothness selection for generalized additive models. Journal of the Royal Statistical Society: Series B (Statistical Methodology) 70 (3), 495–518.
- [38] Wood, S. N., 2011. Fast stable restricted maximum likelihood and marginal likelihood estimation of semiparametric generalized linear models. Journal of the Royal Statistical Society: Series B (Statistical Methodology) 73 (1), 3–36.

- [39] Wood, S. N., 2012. On p-values for smooth components of an extended generalized additive model. Biometrika 100 (1), 221–228.
- [40] Wood, S. N., 2017. mgcv. R package version 1.8-18. URL https://CRAN.R-project.org/package=mgcv
- [41] Yee, T. W., 2007. Vector generalized linear and additive models. Springer.

Appendix

R Package

We created an R-package (cmp) with all the methods developed in the paper. The package is available on github and can be installed by running the following R code:

The full list of attributes and their descriptions for the Bikesharing data

Name	Description
dteday	date
season	season (1:spring, 2:summer, 3:fall, 4:winter)
yr	year $(0: 2011, 1:2012)$
mnth	month $(1 \text{ to } 12)$
hr	hour $(0 \text{ to } 23)$
holiday	weather the day is holiday or not (extracted from
	http://dchr.dc.gov/page/holiday-schedule)
weekday	day of the week
working day	if day is neither weekend nor holiday is 1, otherwise is 0.
weather sit	1 = Clear, Few clouds, Partly cloudy
	2 = Mist + Cloudy, Mist + Broken clouds, Mist + Few clouds, Mist
	3 = Light Snow, Light Rain + Thunderstorm + Scattered clouds
	4 = Heavy Rain + Ice Pallets + Thunderstorm + Mist, Snow + Fog
temp	Normalized temperature in Celsius. The values are divided to $41 \pmod{41}$
atemp	Normalized feeling temperature in Celsius. The values are divided to 50 (max)
hum	Normalized humidity. The values are divided to $100 \pmod{2}$
windspeed	Normalized wind speed. The values are divided to $67 \pmod{2}$
casual	count of casual users
registered	count of registered users
cnt	count of total rental bikes including both casual and registered

Table A.1: Full list of attributes and their description for the bike sharing data



Figure A.1: Comparison of the AIC values from CMP GAM, Poisson GAM, NB GAM and CMP Regression for the 50 bootstrap replications. The X-axis denotes the bootstrap sample. Top chart: over dispersion $\nu = 0.5$; Bottom chart: under dispersion $\nu = 2.5$.



Figure A.2: Comparison of the fitted values and residuals from CMP and Poisson additive models with actual values. (—True values, —CMP GAM , —Poisson GAM and —NB GAM).