
Sparse Estimation of Multivariate Poisson Log-Normal Models from Count Data

Hao Wu

Department of Electrical
and Computer Engineering
Virginia Tech
Arlington, VA 22203, USA
wuhao723@vt.edu

Xinwei Deng

Department of Statistics
Virginia Tech
Blacksburg, VA 24061, USA
xdeng@vt.edu

Naren Ramakrishnan

Department of Computer Science
Virginia Tech
Arlington, VA 22203, USA
naren@vt.edu

Abstract

Modeling data with multivariate count responses is a challenging problem due to the discrete nature of the responses. Existing methods for univariate count responses cannot be easily extended to the multivariate case since the dependency among multiple responses needs to be properly accommodated. In this paper, we propose a multivariate Poisson log-normal regression model for multivariate data with count responses. By simultaneously estimating the regression coefficients and inverse covariance matrix over the latent variables with an efficient Monte Carlo EM algorithm, the proposed regression model takes advantages of association among multiple count responses to improve the model prediction performance. Simulation studies and applications to real world data are conducted to systematically evaluate the performance of the proposed method in comparison with conventional methods.

1 Introduction

In this decade of data science, multivariate response observations are routine in numerous disciplines. To model such datasets, multivariate regression and multi-task learning models are common techniques to study and investigate the relationships between $q \geq 2$ responses and p predictors. The former class of methods, e.g. [25, 23, 22] and [27], estimates the $p \times q$ regression coefficients as well as recover the correlation structures among response variables using regularization. The latter class of methods focuses on learning the shared features [13, 15, 14, 18] or common underlying structure(s) among multiple tasks [2, 20, 6, 31, 3] using regression approaches and enforcing regularization controls over the coefficient matrix. However, all such multivariate regression or multi-task learning models discussed above deal with continuous responses, none of them handle count data.

When responses are count variables, the Poisson model is a natural approach to model them, e.g., in domains such as influenza case count modeling [28], traffic accident analysis [24, 8] and consumer services [26]. However, Poisson regression models proposed in these works are either univariate or inferred via Bayesian approaches and no sparsity or feature selection is typically enforced over the coefficients. When count responses are multivariate, it is challenging to quantify the association among them due to the discrete nature of the data. One important approach is to model each dimension of count variables as the sum of independent Poisson variables with some common Poisson variables capturing dependencies [19]. A drawback of this method is that it can only model positive correlations. Recent literature [30, 16] models multivariate count data with novel Poisson graphical models which can handle both positive and negative dependencies. However, these works do not consider multivariate count data in the context of regression.

To consider a joint model for data with multivariate count responses, it is important to properly exploit the hidden associations among the count responses. One way to consider the joint model of multivariate count responses is via penalty-based model selection from the perspective of parameter regularization. The key idea is to allow the count responses to be independent of each other, while

the regression coefficients are required to obey a certain common sparse structure. Hence the joint modeling is enabled because of the joint estimation of regression coefficients through appropriate penalties. Such a modeling strategy leads to an explicit loss function with tractable computational characteristics. However, this method overlooks the essential correlation among multiple count responses, which could result in poor prediction performance. There are also several recent papers that develop models of multivariate count data from the lens of conditional dependency. But these method typically are restricted to approximated likelihood functions under the framework of generalized liner models.

In this work, we propose a novel multivariate Poisson log-normal model for data with multiple count responses. The motivation to adopt the log-normal model is to borrow strength from regression under the multivariate normal assumption, which can simultaneously estimate regression coefficients and covariance structure. For the proposed model, the logarithm of the Poisson rate parameters is modeled as multivariate normal with a sparse inverse covariance matrix, which combines the strengths of sparse regression and graphical modeling to improve prediction performance. Thus, this approach can fully exploit the conditional dependency among multiple count responses. Estimating such model is non-trivial since it is intractable to derive an explicit analytical solution. Thus, to facilitate the estimation of model parameters, we develop an Monte Carlo EM algorithm which allows to iteratively estimate the regression coefficients using the Lasso penalty and the inverse covariance matrix by a graphical Lasso approach. By applying the proposed model to synthetic data and a real world influenza-like illness dataset, we demonstrate the effectiveness of the proposed method when modeling multivariate data with count responses.

It is worth pointing out that the proposed method is not restricted to adopt the Lasso penalty for regression parameters. It can be easily extended to other penalties such as the adaptive Lasso, group Lasso or fused Lasso. While covariance matrix estimation and inverse covariance matrix estimation have attracted significant attention in the literature [11, 25, 23], here we use this idea in the context of a multivariate regression for count data. Thus inverse covariance matrix estimation is conducted here to improve prediction performance, not just as an unsupervised procedure. One may call such a strategy *supervised covariance estimation*, which has not been widely studied in the literature. One exception is the multivariate regression for continuous responses [25, 29]. Therefore, to the best of our knowledge, our proposed method is a first to incorporate covariance matrix estimation into a multivariate regression model of count responses.

2 Multivariate Poisson Log-Normal model

In this section, we formally specify the Multivariate Poisson Log-Normal (MVPLN) model, and propose a Monte Carlo Expectation-Maximization (MCEM) algorithm for parameter estimation in detail.

2.1 The proposed model

Consider the multivariate random variable $\mathbf{Y} = \{\mathcal{Y}^{(1)}, \mathcal{Y}^{(2)}, \dots, \mathcal{Y}^{(q)}\}^T \in \mathcal{Z}_+^q$, where the superscript T denotes the transpose, and \mathcal{Z}_+ represents the set of all positive integers. For count data, it is reasonable to make the assumption that \mathbf{Y} follows the multivariate Poisson distribution. Without loss of generality, let's assume that each dimension of \mathbf{Y} , say $\mathcal{Y}^{(i)}$, follows the univariate Poisson distribution with parameter $\theta^{(i)}$, and is conditional independent of other dimensions given $\theta^{(i)}$. That is:

$$\mathcal{Y}^{(i)} \sim \text{Poisson}(\theta^{(i)}), \theta^{(i)} \in \mathcal{R}_+, \forall i = 1, 2, \dots, q \quad (1)$$

Let $\mathbf{x} = \{x^{(1)}, x^{(2)}, \dots, x^{(p)}\}^T \in \mathcal{R}^p$ denote the predictor vector. In order to establish relationship between \mathbf{Y} and \mathbf{x} , we consider the following regression model:

$$\begin{aligned} \boldsymbol{\theta} &= \exp(\mathbf{B}^T \mathbf{x} + \boldsymbol{\varepsilon}) \\ \boldsymbol{\varepsilon} &\sim N(0, \boldsymbol{\Sigma}) \end{aligned} \quad (2)$$

where \mathbf{B} is a $p \times q$ coefficient matrix, and $\boldsymbol{\Sigma}$ is the $q \times q$ covariance matrix which captures the covariance structure of variable $\boldsymbol{\theta} = \{\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(q)}\}^T$ given \mathbf{x} . Through the variable $\boldsymbol{\theta}$, we

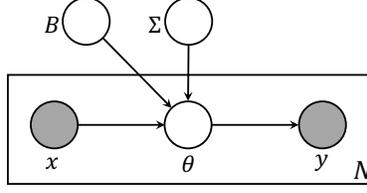


Figure 1: The plate notation of the proposed MVPLN model.

model the covariance structure of the count variable \mathcal{Y} indirectly. Fig. 1 shows the plate notation of the proposed MVPLN model.

Given n observations of the predictor $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]^T$ and corresponding responses $\mathbf{Y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n]^T$, the log-likelihood of the MVPLN model is:

$$\mathcal{L}(\mathbf{B}, \mathbf{\Sigma}) = \sum_{j=1}^n \log p(\mathcal{Y} = \mathbf{y}_j | \mathbf{x}_j), \quad (3)$$

where

$$p(\mathcal{Y} = \mathbf{y} | \mathbf{x}) = \int_{\boldsymbol{\theta}} p(\mathcal{Y} = \mathbf{y}, \boldsymbol{\theta} | \mathbf{x}) d\boldsymbol{\theta} = \int_{\boldsymbol{\theta}} p(\mathcal{Y} = \mathbf{y} | \boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathbf{x}) d\boldsymbol{\theta} \quad (4)$$

Here, $p(\mathcal{Y} = \mathbf{y} | \boldsymbol{\theta})$ and $p(\boldsymbol{\theta} | \mathbf{x})$ follow multivariate Poisson distribution and multivariate log-normal distribution as derived in Section A of Supplementary Material, respectively. To jointly infer the sparse estimations of coefficient matrix \mathbf{B} and covariance matrix $\mathbf{\Sigma}$, we adopt the regularized negative log-likelihood function with l_1 penalties as our loss function. To be specific, the loss function could be written as:

$$\mathcal{L}_p(\mathbf{B}, \mathbf{\Sigma}) = -\mathcal{L}(\mathbf{B}, \mathbf{\Sigma}) + \lambda_1 \|\mathbf{B}\|_1 + \lambda_2 \|\mathbf{\Sigma}^{-1}\|_1, \quad (5)$$

where $\|\cdot\|_1$ denote the l_1 matrix norm, and $\lambda_1 > 0, \lambda_2 > 0$ are two tuning parameters.

For convenience, we use the following notation to present the proposed MVPLN model in the rest of the paper. Normal lower case letters, e.g. x and y , represent scalars. While, bold lower case letters, e.g. \mathbf{x} and \mathbf{y} , are used to represent column vectors, and bold upper case letters in the calligraphic font, e.g. \mathcal{X} and \mathcal{Y} , denote random column vectors. Let letters with superscript in parentheses, e.g. $x^{(i)}$, denote the i^{th} component of the corresponding vector \mathbf{x} . Matrices are represented by bold upper case letters in normal font, e.g. \mathbf{X} and \mathbf{Y} . Letters in lower case with two subscripts, e.g. $x_{i,j}$, denote the (i, j) entry of the corresponding matrix \mathbf{X} .

2.2 Monte Carlo EM algorithm for parameter estimation

In order to obtain the estimations of MVPLN model parameters \mathbf{B} and $\mathbf{\Sigma}$, we could simply solve the following optimization problem:

$$\hat{\mathbf{B}}, \hat{\mathbf{\Sigma}} = \underset{\mathbf{B}, \mathbf{\Sigma}}{\operatorname{argmin}} \mathcal{L}_p(\mathbf{B}, \mathbf{\Sigma}). \quad (6)$$

However, it's difficult to directly minimize the objective function defined above due to the complicated integral in Equation (4). Thus, we turn to an iterative approach for the solution. We treat $\boldsymbol{\theta}$ as latent random variables, and apply the EM algorithm to obtain the maximum likelihood parameter estimation (MLE). However, we cannot derive the analytical form of the expected log-likelihood of the model due to the integral in Equation (4). Here, we adopt a Monte Carlo variant of the EM algorithm for an approximate solution.

2.2.1 Monte Carlo (MC) E-step

In the MC E-step of iteration $t + 1$, instead of trying to derive the close form of the conditional probability distribution of $\boldsymbol{\theta}_j$, we draw m random samples of $\boldsymbol{\theta}_j$, say $\boldsymbol{\Theta}_j = [\boldsymbol{\theta}_j^{(1)}, \boldsymbol{\theta}_j^{(2)}, \dots, \boldsymbol{\theta}_j^{(m)}]^T$, from $p(\boldsymbol{\theta}_j | \mathcal{Y} = \mathbf{y}_j, \mathbf{x}_j; \mathbf{B}^{(t)}, \mathbf{\Sigma}^{(t)})$, and approximate the expected log-likelihood function with:

$$\tilde{Q}(\mathbf{B}, \mathbf{\Sigma} | \mathbf{B}^{(t)}, \mathbf{\Sigma}^{(t)}) = \sum_{j=1}^n \frac{1}{m} \sum_{\tau=1}^m \log p(\mathcal{Y} = \mathbf{y}_j, \boldsymbol{\theta}_j^{(\tau)} | \mathbf{x}_j; \mathbf{B}^{(t)}, \mathbf{\Sigma}^{(t)}). \quad (7)$$

Drawing random samples of θ_j can be achieved with the Metropolis Hasting algorithm. In order to reduce the burn-in period of the Metropolis Hasting algorithm, we adopt the tailored normal distribution [7] as our proposal distribution. Since $p(\theta_j \mid \mathcal{Y} = \mathbf{y}_j, \mathbf{x}_j; \mathbf{B}^{(t)}, \Sigma^{(t)}) \propto p(\mathcal{Y} = \mathbf{y}_j, \theta_j \mid \mathbf{x}_j; \mathbf{B}^{(t)}, \Sigma^{(t)})$, if we let $f(\theta_j) = p(\mathcal{Y} = \mathbf{y}_j, \theta_j \mid \mathbf{x}_j; \mathbf{B}^{(t)}, \Sigma^{(t)})$, the initial value $\theta_j^{(0)}$ of the location parameter for the tailored normal distribution should be the mode of $f(\theta_j)$, and the covariance matrix is $\tau(-\mathbf{H}(\theta_j^{(0)}))^{-1}$, where $\mathbf{H}(\theta_j^{(0)})$ denotes the Hessian matrix of $\log f(\theta_j)$ at $\theta_j^{(0)}$, and τ is a tuning parameter. Considering the performance issue, we adopt a linear approximation approach with the first order Taylor expansion to solve $\theta_j^{(0)}$. In this case, the approximate analytical solution of $\theta_j^{(0)}$ is $\theta_j^{(0)} = e^{\kappa_j}$ where

$$\kappa_j = \left(\text{diag} \left(e^{\kappa_j^{(0)}} \right) + \Sigma^{(t)-1} \right)^{-1} \left(\mathbf{y}_j - \mathbf{1} + \Sigma^{(t)-1} \mathbf{B}^{(t)T} \mathbf{x}_j + \text{diag} \left(e^{\kappa_j^{(0)}} \right) \kappa_j^{(0)} - e^{\kappa_j^{(0)}} \right).$$

Here, $\kappa_j^{(0)} = \log \mathbf{y}_j$. In the case that the covariance matrix $\tau(-\mathbf{H}(\theta_j^{(0)}))^{-1}$ is not positive semidefinite, the nearest positive semidefinite matrix to $\tau(-\mathbf{H}(\theta_j^{(0)}))^{-1}$ is used to replace $\tau(-\mathbf{H}(\theta_j^{(0)}))^{-1}$ [17]. The details for Metropolis Hasting algorithm and the derivation of the tailored normal distribution as the proposal distribution are provided in Section B of the Supplementary Material.

2.2.2 M-step: maximize approximate penalized expected log-likelihood

If we let $\Omega = \Sigma^{-1}$ and $\varphi_{\tau,j} = (\log \theta_j^{(\tau)} - \mathbf{B}^T \mathbf{x}_j)$, with the Monte Carlo approximation of the expected log-likelihood in the MC E-step, the optimization problem we need to solve in the M-step of the MCEM algorithm can be reformulated as:

$$\mathbf{B}^{(t+1)}, \Sigma^{(t+1)} = \underset{\mathbf{B}, \Omega}{\text{argmin}} \left\{ \frac{1}{mn} \text{tr} \left(\Phi^T \Phi \Omega \right) - \log |\Omega| + \lambda_1 \|\mathbf{B}\|_1 + \lambda_2 \|\Omega\|_1 \right\} \quad (8)$$

where $\Phi = [\varphi_{1,1}, \varphi_{2,1}, \dots, \varphi_{m,1}, \varphi_{1,2}, \varphi_{2,2}, \dots, \varphi_{m,2}, \dots, \varphi_{m,n}]^T$. The optimization problem defined in Equation (8) is not convex. However, it is convex w.r.t. either \mathbf{B} or Ω with the other fixed [25]. Thus, we present an iterative algorithm that optimizes the objective function in Equation (8) alternatively w.r.t. \mathbf{B} and Ω .

With \mathbf{B} fixed at \mathbf{B}_0 , the optimization problem in Equation (8) yields:

$$\Omega(\mathbf{B}_0) = \underset{\Omega}{\text{argmin}} \left\{ \frac{1}{mn} \text{tr} \left(\Phi^T \Phi \Omega \right) - \log |\Omega| + \lambda_2 \|\Omega\|_1 \right\}, \quad (9)$$

which is the similar problem studied in [11]. We solve this problem with the graphical lasso approach.

When Ω is fixed at Ω_0 , we have the following optimization problem:

$$\mathbf{B}(\Omega_0) = \underset{\mathbf{B}}{\text{argmin}} \left\{ \frac{1}{mn} \text{tr} \left(\Phi^T \Phi \Omega_0 \right) + \lambda_1 \|\mathbf{B}\|_1 \right\}, \quad (10)$$

which is similar to the problem solved by Lasso, and we could adopt the cyclical coordinate descent algorithm [10] to obtain the estimation of \mathbf{B} . However, considering the computational burden already brought in by the MCMC approximation in the MC E-step, we solve the optimization problem in Equation (10) approximately where the l_1 matrix norm $\|\mathbf{B}\|_1$ is replaced by its quadratic approximation $\text{tr}(\mathbf{B}'^T \mathbf{B}')$ where $\mathbf{B}' = \mathbf{B} \circ (1/\sqrt{|\hat{\mathbf{B}}|})$. Here, \circ denotes the Hadamard (element-wise) product, $\hat{\mathbf{B}}$ denotes the current estimation of \mathbf{B} , and $1/\sqrt{|\hat{\mathbf{B}}|}$ represents the matrix that each entry is the inverse of the square root of the absolute value of the corresponding entry in $\hat{\mathbf{B}}$. With such approximation, we would get the analytical solution to the optimization problem in Equation (10) as

$$\text{vec}(\mathbf{B}) = \left[\Omega_0 \otimes \mathbf{S} + \text{diag} \left(\text{vec} \left(\frac{\lambda_1 mn}{|\hat{\mathbf{B}}|} \right) \right) \right]^{-1} \text{vec}(\mathbf{H}). \quad (11)$$

Here, $\text{vec}(\cdot)$ represents the vectorization operation over the matrix, and the two auxiliary matrices \mathbf{H} and \mathbf{S} are:

$$\mathbf{H} = \left(\sum_{j=1}^n \mathbf{X}_j^T (\log \Theta_j) \right) \Omega_0, \quad \mathbf{S} = \sum_{j=1}^n \mathbf{X}_j^T \mathbf{X}_j,$$

where \mathbf{X}_j is a $m \times p$ matrix with each row being \mathbf{x}_j for all $j = 1, 2, \dots, n$. The estimated coefficient matrix \mathbf{B} can be obtained by reorganizing the $\text{vec}(\mathbf{B})$ in Equation (11). By solving the optimization problem in Equation (9) and (10) alternatively until convergence, we will obtain the MLE of the coefficient matrix \mathbf{B} and inverse covariance matrix Ω . The detailed derivation of the algorithm for M-step is provided in Section C of the Supplementary Material.

2.3 Selection of tuning parameters

To determine the optimal values of the tuning parameters λ_1 and λ_2 , we adopt the extended Bayesian Information Criterion (EBIC) approach proposed in [5] and extended to Gaussian Graphical Models in [9]. Assume $\mathbf{B}_{\lambda_1, \lambda_2}$ and $\Omega_{\lambda_1, \lambda_2}$ denote the MLE of the model parameter \mathbf{B} and Ω with regularization parameters λ_1 and λ_2 . The EBIC value for this model is given by the following equation:

$$\begin{aligned} \text{EBIC}_\gamma(\lambda_1, \lambda_2) = & -2\tilde{Q}(\mathbf{B}_{\lambda_1, \lambda_2}, \Omega_{\lambda_1, \lambda_2}) + [v(\mathbf{B}_{\lambda_1, \lambda_2}) + v(\Omega_{\lambda_1, \lambda_2})] \log n + 2\gamma v(\mathbf{B}_{\lambda_1, \lambda_2}) \log(pq) \\ & + 4\gamma v(\Omega_{\lambda_1, \lambda_2}) \log q, \end{aligned} \quad (12)$$

where $\tilde{Q}(\mathbf{B}_{\lambda_1, \lambda_2}, \Omega_{\lambda_1, \lambda_2})$ is the approximate expected log-likelihood in Equation (7), $v(\mathbf{B}_{\lambda_1, \lambda_2})$ and $v(\Omega_{\lambda_1, \lambda_2})$ denote the number of non-zero entries in $\mathbf{B}_{\lambda_1, \lambda_2}$ and $\Omega_{\lambda_1, \lambda_2}$, respectively, and n is the number of training observations. With EBIC, the optimal values for λ_1 and λ_2 are selected by

$$(\hat{\lambda}_1, \hat{\lambda}_2) = \underset{\lambda_1, \lambda_2}{\text{argmin}} \text{EBIC}_\gamma(\lambda_1, \lambda_2).$$

3 Experiments and results

3.1 Simulation study

In our simulation study, we compare the proposed MVPLN model with the separate univariate Lasso regularized Poisson regression model (GLMNET model) (e.g., as implemented in the R `glmnet` package [12]). The regularized univariate Poisson regression is applied to each response dimension, and a Bayesian Information Criterion (BIC) is used to select regularization parameters in order to make a fair comparison. The simulation data are generated with the following approach. Each data observation in the $n \times p$ predictor matrix \mathbf{X} is independently sampled from a multivariate normal distribution $N(\boldsymbol{\mu}_X, \sigma_X \mathbf{I})$, where the location parameter $\boldsymbol{\mu}_X$ is sampled from a uniform distribution $Unif(\boldsymbol{\mu}_{\min}, \boldsymbol{\mu}_{\max})$. The corresponding observations in the $n \times q$ response matrix \mathbf{Y} are generated following the definition of the MVPLN model in Equation (1) and (2). In order to enforce sparsity, a fixed number of zeros are randomly placed into each column of the coefficient matrix \mathbf{B} . The other non-zero entries of \mathbf{B} are independently sampled from a univariate normal distribution $N(\mu_B, \sigma_B)$. Regarding the inverse covariance matrix $\Omega = \Sigma^{-1}$ for ϵ , we consider four scenarios: (1). Random Ω , where the inverse covariance matrix is generated by $\Omega = \Psi^T \Psi$ to ensure the positive semidefinite property. Each entry in Ψ is independently sampled from a uniform distribution $Unif(-1, 1)$; (2). Banded Ω , where the sparsity is enforced by the modified Cholesky decomposition [21]: $\Omega = \mathbf{T}^T \mathbf{D}^{-1} \mathbf{T}$. Here, \mathbf{T} is a lower triangular matrix with 1's on the diagonal, and \mathbf{D} is a diagonal matrix. The non-zero off diagonal elements in \mathbf{T} and diagonal elements in \mathbf{D} are independently sampled from uniform distribution $Unif(-1, 1)$ and $Unif(0, 1)$ respectively; (3). sparse Ω , where the Ω matrix is generated by performing some random row and column permutations over the banded Ω matrix; (4). Diagonal Ω , where the diagonal elements are sampled independently from standard uniform distribution. In order to make sure that the elements in the response matrix \mathbf{Y} are within the reasonable range, we scale the matrix Σ to make the largest element equal to ψ . By tuning the synthetic data generation parameters $\boldsymbol{\mu}_{\min}$, $\boldsymbol{\mu}_{\max}$, σ_X , μ_B , σ_B , and ψ , we could adjust the range and variations in the generated response matrix \mathbf{Y} .

In our experiments, we fix the number of observations in the training data at $n = 50$, and the number of observations in the test data at 20. We consider two scenarios: (1) the dimension of predictors

Table 1: Estimation errors w.r.t. \mathbf{B} and $\mathbf{\Omega}$. The standard errors are shown in the parentheses.

$\mathbf{\Omega}$	ψ	$l(\mathbf{B}, \hat{\mathbf{B}})$				$l(\mathbf{\Omega}, \hat{\mathbf{\Omega}})$			
		$p < n$		$p > n$		$p < n$		$p > n$	
		GLMNET	MVPLN	GLMNET	MVPLN	GLMNET	MVPLN	GLMNET	MVPLN
Random	0.4	2.25607 (0.04547)	1.19936 (0.01277)	1.61016 (0.01830)	1.44383 (0.01076)	NA	0.99550 (0.00131)	NA	0.99595 (0.00100)
	1.0	4.35649 (0.09258)	1.70326 (0.03620)	2.41644 (0.03039)	1.74861 (0.02796)	NA	0.99033 (0.00153)	NA	0.99151 (0.00200)
	1.6	5.37513 (0.12519)	1.80392 (0.03618)	2.87839 (0.04629)	1.94325 (0.02844)	NA	0.98928 (0.00211)	NA	0.98561 (0.00452)
	2.2	6.32172 (0.17932)	1.99852 (0.04246)	3.21878 (0.05822)	2.12487 (0.04339)	NA	0.99214 (0.00126)	NA	0.98343 (0.00328)
Banded	0.4	2.12650 (0.05377)	1.16671 (0.01882)	1.49028 (0.01747)	1.38619 (0.01133)	NA	0.98029 (0.00204)	NA	0.98500 (0.00148)
	1.0	3.57945 (0.10031)	1.59062 (0.04760)	2.13400 (0.03313)	1.59255 (0.02344)	NA	0.95796 (0.00508)	NA	0.94881 (0.00563)
	1.6	4.41182 (0.13408)	1.80692 (0.06746)	2.59768 (0.05930)	1.78361 (0.02981)	NA	0.93159 (0.00811)	NA	0.92380 (0.00874)
	2.2	5.21359 (0.18171)	2.04397 (0.07308)	2.84779 (0.07824)	2.01992 (0.05492)	NA	0.93695 (0.00681)	NA	0.90552 (0.00838)
Sparse	0.4	1.98327 (0.06026)	1.11950 (0.01556)	1.52410 (0.02270)	1.40847 (0.01107)	NA	0.98277 (0.00259)	NA	0.98205 (0.00201)
	1.0	3.43339 (0.11127)	1.50384 (0.04915)	2.13721 (0.03966)	1.60572 (0.02315)	NA	0.95978 (0.00597)	NA	0.96085 (0.00425)
	1.6	4.69189 (0.15989)	1.88319 (0.07134)	2.54446 (0.05723)	1.76144 (0.02705)	NA	0.92684 (0.00880)	NA	0.92349 (0.00852)
	2.2	5.09710 (0.21733)	2.12963 (0.07617)	2.74681 (0.08444)	1.91288 (0.04085)	NA	0.96626 (0.01344)	NA	0.90581 (0.00957)
Diagonal	0.4	1.86103 (0.05898)	1.10292 (0.01452)	1.43937 (0.01870)	1.34607 (0.01274)	NA	0.96841 (0.00324)	NA	0.97068 (0.00413)
	1.0	3.29868 (0.09724)	1.53224 (0.04655)	2.01567 (0.04295)	1.56539 (0.02745)	NA	0.88673 (0.01313)	NA	0.89628 (0.01510)
	1.6	4.33160 (0.13345)	1.84269 (0.06302)	2.39551 (0.05794)	1.70712 (0.04889)	NA	0.81895 (0.01851)	NA	0.84071 (0.02020)
	2.2	5.00582 (0.23160)	1.95903 (0.08481)	2.56122 (0.08119)	1.76716 (0.03930)	NA	0.88405 (0.02031)	NA	0.81663 (0.02034)

is less than the number of observations in training data ($p < n$); (2) the dimension of predictors is greater than or equal to the number of observations in training data ($p \geq n$). We let $p = 30, q = 5$ for the case $p < n$, and $p = 70, q = 5$ for the case $p \geq n$. For each parameter setting, the simulation is repeated for 60 times, and the reported results are averaged across the 60 replications to alleviate the randomness.

3.1.1 Estimation accuracy

To measure model estimation accuracy w.r.t. \mathbf{B} and $\mathbf{\Omega}$, we report the estimation errors by computing the distance between \mathbf{B} and $\hat{\mathbf{B}}$ (or $\mathbf{\Omega} = \mathbf{\Sigma}^{-1}$ and $\hat{\mathbf{\Omega}} = \hat{\mathbf{\Sigma}}^{-1}$) using the normalized matrix Frobenius norm:

$$l(\mathbf{B}, \hat{\mathbf{B}}) = \frac{\|\mathbf{B} - \hat{\mathbf{B}}\|_F}{\|\mathbf{B}\|_F}$$

Here, \mathbf{B} denotes the true value of coefficient matrix and $\hat{\mathbf{B}}$ represents the estimation given by the MVPLN or GLMNET model. Table 1 shows the estimation errors of coefficient matrix \mathbf{B} and inverse covariance matrix $\mathbf{\Omega}$ in various parameter settings. Since the GLMNET model cannot infer the inverse covariance matrix, we omit the corresponding results here. We can see that the proposed MVPLN model consistently outperforms the GLMNET model in all parameter settings, especially when the variation in the simulated data is large (ψ is large). Such promising results demonstrate that the proposed MVPLN model leverages the dependency structures between the multi-dimensional count responses to improve the estimation accuracy.

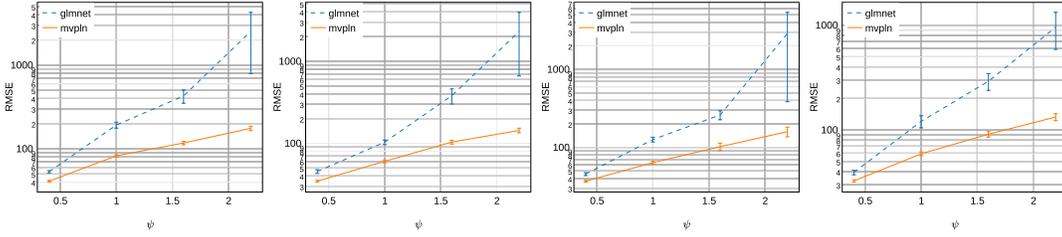


Figure 2: Average rMSE across response dimensions over test data when Ω is random, sparse, banded and diagonal (from left to right), and $p < n$. The vertical error bars indicate the standard deviation, and the Y-axis is in log scale.

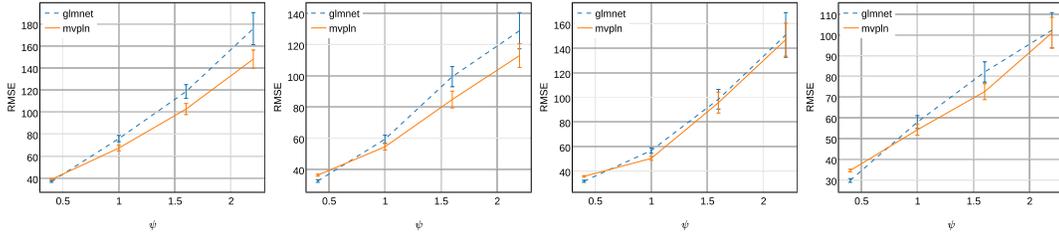


Figure 3: Average rMSE across response dimensions over test data when Ω is random, sparse, banded and diagonal (from left to right), and $p \geq n$. The vertical error bars indicate the standard deviation.

3.1.2 Prediction errors

To evaluate the prediction performance of the proposed model, we report the average root-mean-square error (rMSE) across all the response dimensions over the test data. Figure 2 and 3 show the average rMSE for the cases when $p < n$ and $p \geq n$ respectively. These figures show that when the variations in the simulated data are small (ψ is small), the prediction performances of the proposed MVPLN model and GLMNET model are comparable. As the variations in the data increase, the prediction performance of the proposed MVPLN model becomes better than GLMNET model. This demonstrates that by incorporating the dependency structures between the count responses, the proposed MVPLN model improves its prediction performance. However, when the variations in the data are small, it is difficult for the MVPLN model to take the advantage of inverse covariance matrix estimation. On the contrary, approximating the log-likelihood with MCMC techniques would impose negative effects on the model estimation and prediction accuracy. This is why we observe that when ψ is small, the proposed MVPLN model sometimes does not perform as well as the GLMNET model in term of rMSE.

3.1.3 Model convergence

Another aspect we would like to emphasize here is the model convergence performance. During the experiments over the simulated data, we notice that the GLMNET model will not always converge in some parameter settings, especially when ψ is large. As a result, no parameter estimations are given by the GLMNET model. Figure 4 shows the convergence rate (the fraction of experiment replications that converge and produce valid model estimation) over the simulated data for various parameter settings. We can see from the figure that the larger variations (larger ψ) in the data, the more frequently the GLMNET model fails to give a valid model estimation. On the other hand, the proposed MVPLN model consistently produces the valid model estimation in all of the scenarios. Such results demonstrate that the proposed MVPLN model is more robust to the variations in the underlying multivariate data with count responses.

3.2 Modeling influenza-like illness case counts

We apply the proposed MVPLN model to a real influenza-like illness (ILI) dataset for two Latin American countries, Brazil and Chile, each with four types of ILI diseases (FLUAH3, FLUB,

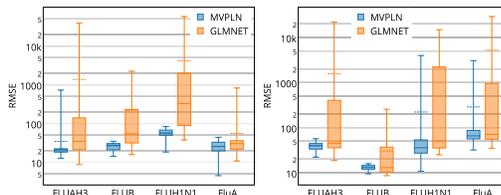
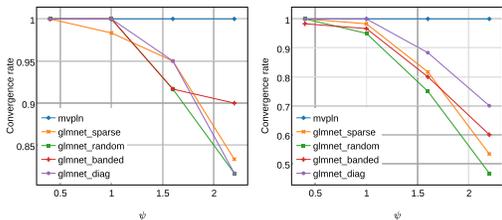


Figure 4: Convergence rate of GLMNET and MV-PLN models when $p < n$ (left) and $p > n$ (right). Since MVPLN model always converges, we use a dash single line to represent these four scenarios.

Figure 5: rMSE box plot of MVPLN and GLMNET models on the real ILI dataset for the countries of Brazil (left) and Chile (right). The dash lines indicate the mean of the rMSE.

FLUH1N1 and FLUA). The data were collected from WHO FluNet [1] from May 1st, 2012 to Dec. 27, 2014 ($n = 139$ weeks), which serves as the multivariate responses of the dataset. The predictors of this ILI dataset are the weekly counts of 108 ILI related keywords collected from the Twitter users of Brazil and Chile during the same period. Before applying the proposed MVPLN model, we preprocessed the ILI dataset with the following approach. We first clustered the 108 ILI related keywords into 20 clusters based on their weekly counts during the selected period using the k-means algorithm. Then for each cluster, we aggregated the weekly counts together for the keywords that belong to this cluster, and finally, we scaled the aggregated keyword counts for each cluster so that it has zero mean and unit standard deviation.

It should also be noticed that although this ILI dataset is time-indexed, we chose to model it as merely a multivariate dataset in our first study here since the proposed MVPLN model is not specially designed to model time series datasets. We use 70% of the preprocessed ILI dataset as the training set and the rest (30%) as the test set. We apply the proposed MVPLN model over the training set, and compute the rMSE of the test set as the criterion for the prediction performance of the model. As a comparison, we also apply the GLMNET model to the same ILI dataset, and compare the rMSE with the proposed MVPLN model. We repeat this experiment for 60 independent runs, and for each run, we shuffle the ILI dataset and re-split the training set and test set. Figure 5 shows the rMSE box plots of the proposed MVPLN model and the GLMNET model for Brazil and Chile after removing some extreme outliers. As we can see from the box plots, although the proposed MVPLN model generates slightly large rMSE over the test set for some response dimensions occasionally, in general, the rMSEs of the MVPLN model are much smaller and have less variation when compared to the GLMNET model for both Brazil and Chile, which indicates that the proposed MVPLN model is better and more stable in term of the prediction performance over the real dataset with count responses. Such results demonstrate that by leveraging the covariance structures between multiple count responses, the proposed MVPLN model improves the prediction performance. However, we also notice that for some flu types, the proposed MVPLN model sometimes generate a large rMSE value, e.g. FLUAH3 in the Brazil dataset, FLUH1N1 and FluA in the Chile dataset. The potential reason for this is likely that the data shuffling procedure happens to place most of the large-response data instances into the model training set, which could mislead the model estimation and result in an overestimation over the test set.

4 Conclusion

In this paper, we have proposed and formulated a multivariate Poisson log-normal model for datasets with count responses. By developing an MCEM algorithm, we accomplish simultaneous sparse estimations of the regression coefficients and of the inverse covariance matrix of the model. Results of simulation studies on synthetic data and an application to a real ILI dataset demonstrate that the proposed MVPLN model achieves better estimation and prediction performance versus a classical Lasso regularized Poisson regression model. Additional interesting future work for the proposed model are being conducted on the following lines. (1) Asymptotic properties of the proposed model are being further investigated; (2) instead of using MCMC techniques, we aim to develop a better approximation algorithm, e.g. using variational inference [4]; (3) we aim to develop variants of the proposed model to better deal with count data with over-dispersion and zero-inflation.

Appendix

A Distribution of multivariate count responses

Given the multivariate count response \mathcal{Y} and the predictor \mathbf{x} , with the conditional independence assumption, the probability mass function for the multivariate Poisson random variable \mathcal{Y} is

$$p(\mathcal{Y} = \mathbf{y} \mid \boldsymbol{\theta}) = \prod_{i=1}^q p(\mathcal{Y}^{(i)} = y^{(i)} \mid \theta^{(i)}) = \prod_{i=1}^q \frac{(\theta^{(i)})^{y^{(i)}} \exp(-\theta^{(i)})}{y^{(i)}!} \quad (13)$$

From the specification of the MVPLN model, since $\boldsymbol{\varepsilon} \sim N(0, \boldsymbol{\Sigma})$, if we let $\boldsymbol{\gamma} = \mathbf{B}^T \mathbf{x} + \boldsymbol{\varepsilon}$, we know that $\boldsymbol{\gamma}$ follows the multivariate normal distribution $N(\mathbf{B}^T \mathbf{x}, \boldsymbol{\Sigma})$ with density function:

$$p(\boldsymbol{\gamma} \mid \mathbf{x}) = \frac{1}{(2\pi)^{q/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2}(\boldsymbol{\gamma} - \mathbf{B}^T \mathbf{x})^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{\gamma} - \mathbf{B}^T \mathbf{x})\right)$$

Since $\boldsymbol{\theta} = \exp(\boldsymbol{\gamma}) = \exp(\mathbf{B}^T \mathbf{x} + \boldsymbol{\varepsilon})$, $\boldsymbol{\theta} \mid \mathbf{x}$ follows the multivariate log-normal distribution, and we can derive that the density function of $\boldsymbol{\theta} \mid \mathbf{x}$ is:

$$p(\boldsymbol{\theta} \mid \mathbf{x}) = p_{\boldsymbol{\gamma}}(\log(\boldsymbol{\theta}) \mid \mathbf{x}) \left| \text{diag}\left(\frac{1}{\theta^{(i)}}\right) \right| = \frac{\exp\left(-\frac{1}{2}\left(\log \boldsymbol{\theta} - \mathbf{B}^T \mathbf{x}\right)^T \boldsymbol{\Sigma}^{-1}\left(\log \boldsymbol{\theta} - \mathbf{B}^T \mathbf{x}\right)\right)}{(2\pi)^{q/2} |\boldsymbol{\Sigma}|^{1/2} \prod_{i=1}^q \theta^{(i)}}. \quad (14)$$

Thus, the probability mass function for $\mathcal{Y} \mid \mathbf{x}$ is:

$$p(\mathcal{Y} = \mathbf{y} \mid \mathbf{x}) = \int_{\boldsymbol{\theta}} p(\mathcal{Y} = \mathbf{y}, \boldsymbol{\theta} \mid \mathbf{x}) d\boldsymbol{\theta} = \int_{\boldsymbol{\theta}} p(\mathcal{Y} = \mathbf{y} \mid \boldsymbol{\theta}) p(\boldsymbol{\theta} \mid \mathbf{x}) d\boldsymbol{\theta},$$

where $p(\mathcal{Y} = \mathbf{y} \mid \boldsymbol{\theta})$ and $p(\boldsymbol{\theta} \mid \mathbf{x})$ are specified in Equation (13) and (14), respectively.

B Monte Carlo E-step in MCEM algorithm

B.1 Metropolis Hasting algorithm for sampling θ_j

Suppose in the MC E-step of iteration $t + 1$, the current estimations of the model parameters are $\mathbf{B}^{(t)}$ and $\boldsymbol{\Sigma}^{(t)}$. Thus, the conditional distribution of the latent variable $\boldsymbol{\theta}$ given $\mathbf{x}, \mathbf{y}, \mathbf{B}^{(t)}$ and $\boldsymbol{\Sigma}^{(t)}$ is:

$$p(\boldsymbol{\theta} \mid \mathcal{Y} = \mathbf{y}, \mathbf{x}; \mathbf{B}^{(t)}, \boldsymbol{\Sigma}^{(t)}) = \frac{p(\mathcal{Y} = \mathbf{y}, \boldsymbol{\theta} \mid \mathbf{x}; \mathbf{B}^{(t)}, \boldsymbol{\Sigma}^{(t)})}{p(\mathcal{Y} = \mathbf{y} \mid \mathbf{x}; \mathbf{B}^{(t)}, \boldsymbol{\Sigma}^{(t)})}. \quad (15)$$

Then, the expected log-likelihood of the model under $p(\boldsymbol{\theta} \mid \mathcal{Y} = \mathbf{y}, \mathbf{x}; \mathbf{B}^{(t)}, \boldsymbol{\Sigma}^{(t)})$ would be:

$$\begin{aligned} Q(\mathbf{B}, \boldsymbol{\Sigma} \mid \mathbf{B}^{(t)}, \boldsymbol{\Sigma}^{(t)}) &= E_{p(\boldsymbol{\theta} \mid \mathcal{Y} = \mathbf{y}, \mathbf{x})}[\mathcal{L}(\mathbf{B}, \boldsymbol{\Sigma})] \\ &= \sum_{j=1}^n E_{p(\boldsymbol{\theta}_j \mid \mathcal{Y} = \mathbf{y}_j, \mathbf{x}_j)}[\log p(\mathcal{Y} = \mathbf{y}_j, \boldsymbol{\theta}_j \mid \mathbf{x}_j; \mathbf{B}, \boldsymbol{\Sigma})]. \end{aligned} \quad (16)$$

In order to compute the approximate expected log-likelihood, we adopt the MCMC technique to sample the $\boldsymbol{\theta}_j$ from $p(\boldsymbol{\theta}_j \mid \mathcal{Y} = \mathbf{y}_j, \mathbf{x}_j; \mathbf{B}^{(t)}, \boldsymbol{\Sigma}^{(t)})$. Since $\mathbf{y}_j, \mathbf{x}_j, \mathbf{B}^{(t)}$ and $\boldsymbol{\Sigma}^{(t)}$ are all known values, which makes $p(\mathcal{Y} = \mathbf{y}_j \mid \mathbf{x}_j; \mathbf{B}^{(t)}, \boldsymbol{\Sigma}^{(t)})$ a constant. In this case, Equation (15) yields

$$p(\boldsymbol{\theta}_j \mid \mathcal{Y} = \mathbf{y}_j, \mathbf{x}_j; \mathbf{B}^{(t)}, \boldsymbol{\Sigma}^{(t)}) \propto p(\mathcal{Y} = \mathbf{y}_j, \boldsymbol{\theta}_j \mid \mathbf{x}_j; \mathbf{B}^{(t)}, \boldsymbol{\Sigma}^{(t)}).$$

Let $f(\boldsymbol{\theta}_j) = p(\mathcal{Y} = \mathbf{y}_j, \boldsymbol{\theta}_j \mid \mathbf{x}_j; \mathbf{B}^{(t)}, \boldsymbol{\Sigma}^{(t)})$ and $g(\boldsymbol{\theta}^* \mid \boldsymbol{\theta})$ be the density function of the proposal distribution. Algorithm 1 illustrates the Metropolis Hasting algorithm for sampling $\boldsymbol{\theta}_j$ from $p(\boldsymbol{\theta}_j \mid \mathcal{Y} = \mathbf{y}_j, \mathbf{x}_j; \mathbf{B}^{(t)}, \boldsymbol{\Sigma}^{(t)})$.

Algorithm 1: Metropolis Hasting algorithm for sampling θ_j

input : $\mathbf{y}_j, \mathbf{x}_j, \mathbf{B}^{(t)}$ and $\Sigma^{(t)}$.

output : m samples $\Theta_j = \{\theta_j^{(1)}, \theta_j^{(2)}, \dots, \theta_j^{(m)}\}^T$.

```

1 Choose  $\theta_j^{(0)}$  as initial value, and let  $\tau \leftarrow 1$ ;
2 while  $|\Theta_j| < m$  do
3   Draw a candidate  $\theta_j^*$  from  $g(\theta_j^* | \theta_j^{(\tau-1)})$ ;
4    $\alpha \leftarrow \min\left(\frac{f(\theta_j^*)/g(\theta_j^*|\theta_j^{(\tau-1)})}{f(\theta_j^{(\tau-1)})/g(\theta_j^{(\tau-1)}|\theta_j^*)}, 1\right)$ ;
5   Accept  $\theta_j^{(*)}$  as  $\theta_j^{(\tau)}$  with probability  $\alpha$ ;
6   if  $\theta_j^{(*)}$  is accepted then
7      $\Theta_j \leftarrow \Theta_j \cup \{\theta_j^{(\tau)}\}$ ;
8      $\tau \leftarrow \tau + 1$ ;
9   end
10 end
11 return  $\Theta_j$ ;

```

B.2 Derivation of the tailored normal distribution as proposal distribution

To find the mode of $p(\theta_j | \mathcal{Y} = \mathbf{y}_j, \mathbf{x}_j; \mathbf{B}^{(t)}, \Sigma^{(t)})$, we need to solve the following optimization problem:

$$\theta_j^{(0)} = \underset{\theta_j}{\operatorname{argmax}}\{\log f(\theta_j)\},$$

let $F(\theta_j) = \log f(\theta_j) = \log(p(\mathcal{Y} = \mathbf{y}_j | \theta_j, \mathbf{x}_j; \mathbf{B}^{(t)}, \Sigma^{(t)})p(\theta_j | \mathbf{x}_j; \mathbf{B}^{(t)}, \Sigma^{(t)}))$. By combining Equation (13) and (14) together, we can derive that:

$$F(\theta_j) = (\mathbf{y}_j - \mathbf{1})^T \log \theta_j - \mathbf{1}^T \theta_j - \frac{1}{2} \left(\log \theta_j - \mathbf{B}^{(t)T} \mathbf{x}_j \right)^T \Sigma^{(t)-1} \left(\log \theta_j - \mathbf{B}^{(t)T} \mathbf{x}_j \right) + C, \quad (17)$$

where $\mathbf{1}$ denotes a column vector of 1s, and C represents the sum of all the constants in $\log f(\theta_j)$. Then, the first order and second order derivatives of $F(\theta_j)$ w.r.t. θ_j are

$$\begin{aligned} \nabla F(\theta_j) &= \frac{dF(\theta_j)}{d\theta_j} = \operatorname{diag}\left(\frac{1}{\theta_j^{(i)}}\right) \left[(\mathbf{y}_j - \mathbf{1}) - \Sigma^{(t)-1} \left(\log \theta_j - \mathbf{B}^{(t)T} \mathbf{x}_j \right) \right] - \mathbf{1} \quad (18) \\ \mathbf{H}(\theta_j) &= \operatorname{diag}\left(-\frac{y_j^{(i)} - 1}{\theta_j^{(i)2}}\right) + \operatorname{diag}\left(-\frac{1}{\theta_j^{(i)2}}\right) \operatorname{diag}\left(\Sigma^{(t)-1} \left(\log \theta_j - \mathbf{B}^{(t)T} \mathbf{x}_j \right)\right) \\ &\quad + \operatorname{diag}\left(\frac{1}{\theta_j^{(i)}}\right) \Sigma^{(t)-1} \operatorname{diag}\left(\frac{1}{\theta_j^{(i)}}\right) \quad (19) \end{aligned}$$

Let $\nabla F(\theta_j) = 0$, and we could get that the initial value $\theta_j^{(0)}$ of the location parameter for the tailored normal distribution is the solution to the following equation:

$$\theta_j + \Sigma^{(t)-1} \log \theta_j = \mathbf{y}_j - \mathbf{1} + \Sigma^{(t)-1} \mathbf{B}^{(t)T} \mathbf{x}_j \quad (20)$$

which can be solved by any numerical root discovering algorithms. However, taking performance issue into account, we let $\kappa_j = \log \theta_j$, and adopt a linear approximation to e^{κ_j} with its first order Taylor expansion at $\kappa_j^{(0)} = \log \mathbf{y}_j$. In this case, Equation (20) becomes:

$$e^{\kappa_j^{(0)}} + \operatorname{diag}\left(e^{\kappa_j^{(0)}}\right) \left(\kappa_j - \kappa_j^{(0)} \right) + \Sigma^{(t)-1} \kappa_j = \mathbf{y}_j - \mathbf{1} + \Sigma^{(t)-1} \mathbf{B}^{(t)T} \mathbf{x}_j. \quad (21)$$

Solving Equation (21) for κ_j , the location parameter (mean) θ_j of the tailored normal distribution is given by $\theta_j^{(0)} = e^{\kappa_j}$ where

$$\kappa_j = \left(\text{diag} \left(e^{\kappa_j^{(0)}} \right) + \Sigma^{(t)-1} \right)^{-1} \left(\mathbf{y}_j - \mathbf{1} + \Sigma^{(t)-1} \mathbf{B}^{(t)T} \mathbf{x}_j + \text{diag} \left(e^{\kappa_j^{(0)}} \right) \kappa_j^{(0)} - e^{\kappa_j^{(0)}} \right),$$

and the covariance matrix is given by $\tau(-\mathbf{H}(\theta_j^{(0)}))^{-1}$. In the case that the covariance matrix $\tau(-\mathbf{H}(\theta_j^{(0)}))^{-1}$ is not positive semidefinite, the nearest positive semidefinite matrix to $\tau(-\mathbf{H}(\theta_j^{(0)}))^{-1}$ is used to replace $\tau(-\mathbf{H}(\theta_j^{(0)}))^{-1}$.

C M-step in the MCEM algorithm

C.1 The optimization problem in M-step

The joint distribution of $(\mathcal{Y} = \mathbf{y}_j, \theta_j^{(\tau)})$ given $\mathbf{x}_j, \mathbf{B}^{(t)}$ and $\Sigma^{(t)}$ is:

$$p(\mathcal{Y} = \mathbf{y}_j, \theta_j^{(\tau)} | \mathbf{x}_j; \mathbf{B}^{(t)}, \Sigma^{(t)}) = p(\mathcal{Y} = \mathbf{y}_j | \theta_j^{(\tau)}) p(\theta_j^{(\tau)} | \mathbf{x}_j; \mathbf{B}^{(t)}, \Sigma^{(t)})$$

where $p(\mathcal{Y} = \mathbf{y}_j | \theta_j^{(\tau)})$ and $p(\theta_j^{(\tau)} | \mathbf{x}_j; \mathbf{B}^{(t)}, \Sigma^{(t)})$ are given by Equation (13) and (14) respectively. Let $\Omega = \Sigma^{-1}$ and $\varphi_{\tau,j} = (\log \theta_j^{(\tau)} - \mathbf{B}^T \mathbf{x}_j)$. Combining the approximated expected log-likelihood we derived in the MC E-step in Section 2.2.1 (Equation (7) in the paper), we can reformulate $\tilde{Q}(\mathbf{B}, \Sigma | \mathbf{B}^{(t)}, \Sigma^{(t)})$ as:

$$\begin{aligned} \tilde{Q}(\mathbf{B}, \Sigma | \mathbf{B}^{(t)}, \Sigma^{(t)}) &= -\frac{1}{n} \sum_{j=1}^n \frac{1}{m} \sum_{\tau=1}^m \left[\left(\log \theta_j^{(\tau)} - \mathbf{B}^T \mathbf{x}_j \right)^T \Omega \left(\log \theta_j^{(\tau)} - \mathbf{B}^T \mathbf{x}_j \right) - \log |\Omega| \right] \\ &= -\frac{1}{mn} \text{tr} \left(\Phi^T \Phi \Omega \right) + \log |\Omega|. \end{aligned} \quad (22)$$

Then, the optimization problem we need to solve in the M-step is:

$$\begin{aligned} \mathbf{B}^{(t+1)}, \Sigma^{(t+1)} &= \underset{\mathbf{B}, \Sigma}{\text{argmin}} \left\{ -\tilde{Q}(\mathbf{B}, \Sigma | \mathbf{B}^{(t)}, \Sigma^{(t)}) + \lambda_1 \|\mathbf{B}\|_1 + \lambda_2 \|\Sigma^{-1}\|_1 \right\} \\ &= \underset{\mathbf{B}, \Omega}{\text{argmin}} \left\{ \frac{1}{mn} \text{tr} \left(\Phi^T \Phi \Omega \right) - \log |\Omega| + \lambda_1 \|\mathbf{B}\|_1 + \lambda_2 \|\Omega\|_1 \right\} \end{aligned} \quad (23)$$

C.2 Approach to solve \mathbf{B} approximately when Ω fixed

When Ω is fixed at Ω_0 , we have the following convex optimization problem:

$$\mathbf{B}(\Omega_0) = \underset{\mathbf{B}}{\text{argmin}} \left\{ \frac{1}{mn} \text{tr} \left(\Phi^T \Phi \Omega_0 \right) + \lambda_1 \|\mathbf{B}\|_1 \right\}. \quad (24)$$

The l_1 matrix norm penalty in Equation (24) could be approximated with the following approach

$$\|\mathbf{B}\|_1 \approx \text{tr} \left(\mathbf{B}'^T \mathbf{B}' \right), \quad \text{where } \mathbf{B}' = \mathbf{B} \circ \frac{1}{\sqrt{|\hat{\mathbf{B}}|}}.$$

Here, \circ denotes the Hadamard (element-wise) product. If we write Φ into the following block matrix

$$\Phi = \begin{bmatrix} \log \Theta_1 - \mathbf{X}_1 \mathbf{B} \\ \log \Theta_2 - \mathbf{X}_2 \mathbf{B} \\ \vdots \\ \log \Theta_n - \mathbf{X}_n \mathbf{B} \end{bmatrix},$$

where \mathbf{X}_j is $m \times p$ matrix with each row being \mathbf{x}_j for all $j = 1, 2, \dots, n$, the objective function of the optimization problem in (24) can be written as:

$$\eta(\mathbf{B}) = \lambda_1 \text{tr} \left(\mathbf{B}'^T \mathbf{B}' \right) + \frac{1}{mn} \sum_{j=1}^n \text{tr} \left((\log \Theta_j - \mathbf{X}_j \mathbf{B})^T (\log \Theta_j - \mathbf{X}_j \mathbf{B}) \Omega_0 \right). \quad (25)$$

Algorithm 2: M-step of the MCEM algorithm

input : $\mathbf{X}, \{\Theta_j\}, \Omega_0, \mathbf{B}_0, \lambda_1$ and λ_2 .**output** : MLE of \mathbf{B} and Ω .

```
1  $t \leftarrow -1$ ;  
2 repeat  
3    $t \leftarrow t + 1$ ;  
4    $\Phi \leftarrow \begin{bmatrix} \log \Theta_1 - \mathbf{X}_1 \mathbf{B}^{(t)} \\ \log \Theta_2 - \mathbf{X}_2 \mathbf{B}^{(t)} \\ \vdots \\ \log \Theta_n - \mathbf{X}_n \mathbf{B}^{(t)} \end{bmatrix}$ ;  
5    $\Omega^{(t+1)} \leftarrow \text{Graphical\_Lasso}(\Phi, \lambda_2)$ ;  
6    $\mathbf{S} \leftarrow \sum_{j=1}^n \mathbf{X}_j^T \mathbf{X}_j$ ;  
7    $\mathbf{H} \leftarrow \sum_{j=1}^n \mathbf{X}_j^T (\log \Theta_j) \Omega^{(t+1)}$ ;  
8    $\mathbf{B}^{(t+1)} \leftarrow \left[ \Omega^{(t+1)} \otimes \mathbf{S} + \text{diag} \left( \text{vec} \left( \frac{\lambda_1 mn}{|\hat{\mathbf{B}}|} \right) \right) \right]^{-1} \text{vec}(\mathbf{H})$ ;  
9 until convergence;  
10 return  $(\mathbf{B}^{(t+1)}, \Omega^{(t+1)})$ ;
```

Taking the first order derivative of $\eta(\mathbf{B})$ w.r.t. \mathbf{B} and setting it to zero, we have

$$\left(\sum_{j=1}^n \mathbf{X}_j^T \mathbf{X}_j \right) \mathbf{B} \Omega_0 + \mathbf{B} \circ \frac{\lambda_1 mn}{|\hat{\mathbf{B}}|} = \left(\sum_{j=1}^n \mathbf{X}_j^T (\log \Theta_j) \right) \Omega_0 \quad (26)$$

If we let $\left(\sum_{j=1}^n \mathbf{X}_j^T (\log \Theta_j) \right) \Omega_0 = \mathbf{H}$ and $\sum_{j=1}^n \mathbf{X}_j^T \mathbf{X}_j = \mathbf{S}$, and apply the matrix vectorization operator $\text{vec}(\cdot)$ to both sides of Equation (26), we have:

$$\left(\Omega_0^T \otimes \mathbf{S} \right) \text{vec}(\mathbf{B}) + \text{vec} \left(\frac{\lambda_1 mn}{|\hat{\mathbf{B}}|} \right) \circ \text{vec}(\mathbf{B}) = \text{vec}(\mathbf{H}).$$

Here, \otimes represents the Kronecker product. By pulling $\text{vec}(\mathbf{B})$ out from the left hand side of the above equation, we can get:

$$\left[\Omega_0 \otimes \mathbf{S} + \text{diag} \left(\text{vec} \left(\frac{\lambda_1 mn}{|\mathbf{B}_{est}|} \right) \right) \right] \text{vec}(\mathbf{B}) = \text{vec}(\mathbf{H}).$$

Thus, the solution to the optimization problem in Equation (24) is

$$\text{vec}(\mathbf{B}) = \left[\Omega_0 \otimes \mathbf{S} + \text{diag} \left(\text{vec} \left(\frac{\lambda_1 mn}{|\hat{\mathbf{B}}|} \right) \right) \right]^{-1} \text{vec}(\mathbf{H}), \quad (27)$$

and the estimated coefficient matrix \mathbf{B} can be obtained by reorganizing the $\text{vec}(\mathbf{B})$ in the above equation.

C.3 Algorithm pseudo code for M-step

By solving \mathbf{B} and Ω alternatively with the other fixed at the value of the last estimate until convergence, we will obtain the MLE of the coefficient matrix \mathbf{B} and inverse covariance matrix Ω for the current iteration of MCEM algorithm. Algorithm 2 summarizes the M-step of the MCEM algorithm.

References

- [1] WHO FluNet, 2015. http://www.who.int/influenza/gisrs_laboratory/flunet/en/.

- [2] A. Argyriou, T. Evgeniou, and M. Pontil. Multi-task feature learning. In *NIPS*, pages 41–48, 2007.
- [3] A. Argyriou, C. A. Micchelli, M. Pontil, and Y. Ying. A spectral regularization framework for multi-task structure learning. In *NIPS*, 2007.
- [4] D. M. Blei, A. Kucubelbir, and J. D. McAuliffe. Variational inference: A review for statisticians, 2016. <https://arxiv.org/abs/1601.00670>.
- [5] J. Chen and Z. Chen. Extended Bayesian information criteria for model selection with large model spaces. *Biometrika*, 95(3):759–771, 2008.
- [6] J. Chen, J. Zhou, and J. Ye. Integrating low-rank and group-sparse structures for robust multi-task learning. In *KDD '11*, pages 42–50, 2011.
- [7] S. Chib, E. Greenberg, and R. Winkelmann. Posterior simulation and bayes factors in panel count data models. *Journal of Econometrics*, 86(1):33 – 54, 1998.
- [8] K. El-Basyouny and T. Sayed. Collision prediction models using multivariate poisson-lognormal regression. *Accident Analysis and Prevention*, 41(4):820 – 828, 2009.
- [9] R. Foygel and M. Drton. Extended bayesian information criteria for gaussian graphical models. In *NIPS*, pages 604–612, 2010.
- [10] J. Friedman, T. Hastie, H. Höfling, and R. Tibshirani. Pathwise coordinate optimization. *The Annals of Applied Statistics*, 1(2):302–332, 2007.
- [11] J. Friedman, T. Hastie, and R. Tibshirani. Sparse inverse covariance estimation with the graphical lasso. *Biostatistics*, 9(3):432–441, July 2008.
- [12] J. Friedman, T. Hastie, N. Simon, and R. Tibshirani. Lasso and elastic-net regularized generalized linear models, 2014. URL <http://cran.r-project.org/web/packages/glmnet/glmnet.pdf>.
- [13] P. Gong, J. Ye, and C. shui Zhang. Multi-stage multi-task feature learning. In *NIPS*, 2012.
- [14] P. Gong, J. Ye, and C. Zhang. Robust multi-task feature learning. In *KDD '12*, pages 895–903, 2012.
- [15] P. Gong, J. Zhou, W. Fan, and J. Ye. Efficient multi-task feature learning with calibration. In *KDD '14*, pages 761–770, 2014.
- [16] F. Hadiji, A. Molina, S. Natarajan, and K. Kersting. Poisson dependency networks: Gradient boosted models for multivariate count data. *Machine Learning*, 100(2):477–507, 2015.
- [17] N. J. Higham. Computing the nearest correlation matrix — a problem from finance. *IMA Journal of Numer. Anal.*, 22(3):329–343, 2002.
- [18] A. Jalali, S. Sanghavi, C. Ruan, and P. K. Ravikumar. A dirty model for multi-task learning. In *NIPS*, pages 964–972, 2010.
- [19] D. Karlis. An em algorithm for multivariate poisson distribution and related models. *Journal of Applied Statistics*, 30(1):63–77, 2003.
- [20] A. Kumar and H. Daumé III. Learning task grouping and overlap in multi-task learning. In *ICML '12*, 2012.
- [21] E. Levina, A. Rothman, J. Zhu, et al. Sparse estimation of large covariance matrices via a nested lasso penalty. *The Annals of Applied Statistics*, 2(1):245–263, 2008.
- [22] H. Liu, L. Wang, and T. Zhao. Multivariate regression with calibration. In *NIPS*, pages 127–135, 2014.
- [23] A. C. Lozano, H. Jiang, and X. Deng. Robust sparse estimation of multiresponse regression and inverse covariance matrix via the l2 distance. In *KDD '13*, pages 293–301, 2013.
- [24] J. Ma, K. M. Kockelman, and P. Damien. A multivariate Poisson-lognormal regression model for prediction of crash counts by severity, using Bayesian methods. *Accident Analysis and Prevention*, 40:964–975, 2008.
- [25] A. J. Rothman, E. Levina, and J. Zhu. Sparse multivariate regression with covariance estimation. *Journal of Computational and Graphical Statistics*, 19(4):947–962, 2010.
- [26] H. Wang, M. U. Kalwani, and T. Akçura. A bayesian multivariate poisson regression model of cross-category store brand purchasing behavior. *Journal of Retailing & Consumer Services*, 14(6):369–382, 2007.
- [27] W. Wang, Y. Liang, and E. P. Xing. Block regularized lasso for multivariate multi-response linear regression. In *AISTATS*, pages 608–617, 2013.
- [28] Z. Wang, P. Chakraborty, S. R. Mekaru, J. S. Brownstein, J. Ye, and N. Ramakrishnan. Dynamic poisson autoregression for influenza-like-illness case count prediction. In *KDD '15*, pages 1285–1294, 2015.
- [29] M. Wytock and J. Z. Kolter. Sparse gaussian conditional random fields: Algorithms, theory, and application to energy forecasting. In *ICML '13*, pages 1265–1273, 2013.
- [30] E. Yang, P. K. Ravikumar, G. I. Allen, and Z. Liu. On poisson graphical models. In *NIPS '13*. 2013.
- [31] S. Yu, V. Tresp, and K. Yu. Robust multi-task learning with t-processes. In *ICML '07*, 2007.