

Learning the Structure of Mixed Graphical Models

Jason D. Lee & Trevor J. Hastie

To cite this article: Jason D. Lee & Trevor J. Hastie (2015) Learning the Structure of Mixed Graphical Models, Journal of Computational and Graphical Statistics, 24:1, 230-253, DOI: [10.1080/10618600.2014.900500](https://doi.org/10.1080/10618600.2014.900500)

To link to this article: <https://doi.org/10.1080/10618600.2014.900500>



View supplementary material [↗](#)



Published online: 31 Mar 2015.



Submit your article to this journal [↗](#)



Article views: 2526



View related articles [↗](#)



View Crossmark data [↗](#)



Citing articles: 34 View citing articles [↗](#)

Learning the Structure of Mixed Graphical Models

Jason D. LEE and Trevor J. HASTIE

We consider the problem of learning the structure of a pairwise graphical model over continuous and discrete variables. We present a new pairwise model for graphical models with both continuous and discrete variables that is amenable to structure learning. In previous work, authors have considered structure learning of Gaussian graphical models and structure learning of discrete models. Our approach is a natural generalization of these two lines of work to the mixed case. The penalization scheme involves a novel symmetric use of the group-lasso norm and follows naturally from a particular parameterization of the model. Supplementary materials for this article are available online.

Key Words: Structure learning; Lasso; Group lasso.

1. INTRODUCTION

Many authors have considered the problem of learning the edge structure and parameters of sparse undirected graphical models. We will focus on using the l_1 regularizer to promote sparsity. This line of work has taken two separate paths: one for learning continuous valued data and one for learning discrete valued data. However, typical data sources contain both continuous and discrete variables: population survey data, genomics data, URL-click pairs, etc. For genomics data, in addition to the gene expression values, we have attributes attached to each sample such as gender, age, ethnicity, etc. In this work, we consider learning mixed models with both continuous Gaussian variables and discrete categorical variables.

For only continuous variables, previous work assumes a multivariate Gaussian (Gaussian graphical) model with mean 0 and inverse covariance Θ . Θ is then estimated via the graphical lasso by minimizing the regularized negative log-likelihood $\ell(\Theta) + \lambda \|\Theta\|_1$. Several efficient methods for solving this can be found in Friedman, Hastie, and Tibshirani (2008); Banerjee, El Ghaoui, and d'Aspremont (2008). Because the graphical lasso problem is computationally challenging, several authors considered methods related to the pseudo-likelihood (PL) and nodewise regression (Meinshausen and Bühlmann 2006; Peng et al. 2009; Friedman, Hastie, and Tibshirani 2010). For discrete models, previous work focuses

Jason D. Lee, Institute of Computational and Mathematical Engineering, Stanford University (E-mail: jdl17@stanford.edu). Trevor J. Hastie, Department of Statistics, Stanford University (E-mail: hastie@stanford.edu). An earlier version of this article was published at AISTATS 2013. This version is significantly expanded with new experimental results, comparisons, and theoretical results.

© 2015 *American Statistical Association, Institute of Mathematical Statistics, and Interface Foundation of North America*

Journal of Computational and Graphical Statistics, Volume 24, Number 1, Pages 230–253

DOI: [10.1080/10618600.2014.900500](https://doi.org/10.1080/10618600.2014.900500)

Color versions of one or more of the figures in the article can be found online at www.tandfonline.com/r/jcgs.

on estimating a pairwise Markov random field of the form $p(y) \propto \exp \sum_{r \leq j} \phi_{rj}(y_r, y_j)$, where ϕ_{rj} are pairwise potentials. The maximum likelihood problem is intractable for models with a moderate to large number of variables (high-dimensional) because it requires evaluating the partition function and its derivatives. Again previous work has focused on the PL approach (Lee, Ganapathi, and Koller 2006; Schmidt et al. 2008; Höfling and Tibshirani 2009; Guo et al. 2010; Schmidt 2010; Ravikumar, Wainwright, and Lafferty 2010; Jalali et al. 2011).

Our main contribution here is to propose a model that connects the discrete and continuous models previously discussed. The conditional distributions of this model are two widely adopted and well understood models: multiclass logistic regression and Gaussian linear regression. In addition, in the case of only discrete variables, our model is a pairwise Markov random field; in the case of only continuous variables, it is a Gaussian graphical model. Our proposed model leads to a natural scheme for structure learning that generalizes the graphical Lasso. Here, the parameters occur as singletons, vectors or blocks, which we penalize using group-lasso norms, in a way that respects the symmetry in the model. Since each parameter block is of different size, we also derive a calibrated weighting scheme to penalize each edge fairly. We also discuss a conditional model (conditional random field) that allows the output variables to be mixed, which can be viewed as a multivariate response regression with mixed output variables. Similar ideas have been used to learn the covariance structure in multivariate response regression with continuous output variables (Witten and Tibshirani 2009; Kim, Sohn, and Xing 2009; Rothman, Levina, and Zhu 2010).

In Section 2, we introduce our new mixed graphical model and discuss previous approaches to modeling mixed data. Section 3 discusses the pseudo-likelihood (PL) approach to parameter estimation and connections to generalized linear models. Section 4 discusses a natural method to perform structure learning in the mixed model. Section 5 presents the calibrated regularization scheme, Section 6 discusses the consistency of the estimation procedures, and Section 7 discusses two methods for solving the optimization problem. Finally, Section 8 discusses a conditional random field extension and Section 9 presents empirical results on a census population survey dataset and synthetic experiments.

2. MIXED GRAPHICAL MODEL

We propose a pairwise graphical model on continuous and discrete variables. The model is a pairwise Markov random field with density $p(x, y; \Theta)$ proportional to

$$\exp \left(\sum_{s=1}^p \sum_{t=1}^p -\frac{1}{2} \beta_{st} x_s x_t + \sum_{s=1}^p \alpha_s x_s + \sum_{s=1}^p \sum_{j=1}^q \rho_{sj}(y_j) x_s + \sum_{j=1}^q \sum_{r=1}^q \phi_{rj}(y_r, y_j) \right). \quad (1)$$

Here, x_s denotes the s th of p continuous variables, and y_j the j th of q discrete variables. The joint model is parameterized by $\Theta = [\{\beta_{st}\}, \{\alpha_s\}, \{\rho_{sj}\}, \{\phi_{rj}\}]$. The discrete y_r takes on L_r states. The model parameters are β_{st} continuous-continuous edge potential, α_s continuous node potential, $\rho_{sj}(y_j)$ continuous-discrete edge potential, and $\phi_{rj}(y_r, y_j)$ discrete-discrete edge potential. $\rho_{sj}(y_j)$ is a function taking L_j values $\rho_{sj}(1), \dots, \rho_{sj}(L_j)$. Similarly, $\phi_{rj}(y_r, y_j)$ is a bivariate function taking on $L_r \times L_j$ values. Later, we will think of $\rho_{sj}(y_j)$ as a vector of length L_j and $\phi_{rj}(y_r, y_j)$ as a matrix of size $L_r \times L_j$.

The two most important features of this model are:

1. the conditional distributions are given by Gaussian linear regression and multiclass logistic regressions;
2. the model simplifies to a multivariate Gaussian in the case of only continuous variables and simplifies to the usual discrete pairwise Markov random field in the case of only discrete variables.

The conditional distributions of a graphical model are of critical importance. The absence of an edge corresponds to two variables being conditionally independent. The conditional independence can be read off from the conditional distribution of a variable on all others. For example, in the multivariate Gaussian model, x_s is conditionally independent of x_t iff the partial correlation coefficient is 0. The partial correlation coefficient is also the regression coefficient of x_t in the linear regression of x_s on all other variables. Thus, the conditional independence structure is captured by the conditional distributions via the regression coefficient of a variable on all others. Our mixed model has the desirable property that the two type of conditional distributions are simple Gaussian linear regressions and multiclass logistic regressions. This follows from the pairwise property in the joint distribution. In more detail:

1. The conditional distribution of y_r given the rest is multinomial, with probabilities defined by a multiclass logistic regression where the covariates are the other variables x_s and $y_{\setminus r}$ (denoted collectively by z in the right-hand side):

$$p(y_r = k | y_{\setminus r}, x; \Theta) = \frac{\exp(\omega_k^T z)}{\sum_{l=1}^{L_r} \exp(\omega_l^T z)} = \frac{\exp(\omega_{0k} + \sum_j \omega_{kj} z_j)}{\sum_{l=1}^{L_r} \exp(\omega_{0l} + \sum_j \omega_{lj} z_j)}. \quad (2)$$

Here we use a simplified notation, which we make explicit in Section 3.1. The discrete variables are represented as dummy variables for each state, for example, $z_j = \mathbb{1}[y_u = k]$, and for continuous variables $z_s = x_s$.

2. The conditional distribution of x_s given the rest is Gaussian, with a mean function defined by a linear regression with predictors $x_{\setminus s}$ and y_r .

$$E(x_s | x_{\setminus s}, y_r; \Theta) = \omega^T z = \omega_0 + \sum_j z_j \omega_j \quad (3)$$

$$p(x_s | x_{\setminus s}, y_r; \Theta) = \frac{1}{\sqrt{2\pi}\sigma_s} \exp\left(-\frac{1}{2\sigma_s^2}(x_s - \omega^T z)^2\right).$$

As before, the discrete variables are represented as dummy variables for each state $z_j = \mathbb{1}[y_u = k]$ and for continuous variables $z_s = x_s$.

The exact form of the conditional distributions (2) and (3) are given in (11) and (10) in Section 3.1, where the regression parameters ω_j are defined in terms of the parameters Θ .

The second important aspect of the mixed model is the two special cases of only continuous and only discrete variables.

1. Continuous variables only. The pairwise mixed model reduces to the familiar multivariate Gaussian parameterized by the symmetric positive-definite inverse covariance matrix $B = \{\beta_{st}\}$ and mean $\mu = B^{-1}\alpha$,

$$p(x) \propto \exp\left(-\frac{1}{2}(x - B^{-1}\alpha)^T B(x - B^{-1}\alpha)\right).$$

2. Discrete variables only. The pairwise mixed model reduces to a pairwise discrete (second-order interaction) Markov random field,

$$p(y) \propto \exp\left(\sum_{j=1}^q \sum_{r=1}^q \phi_{rj}(y_r, y_j)\right).$$

Although these are the most important aspects, we can characterize the joint distribution further. The conditional distribution of the continuous variables given the discrete follow a multivariate Gaussian distribution, $p(x|y) = \mathcal{N}(\mu(y), B^{-1})$. Each of these Gaussian distributions share the same inverse covariance matrix B but differ in the mean parameter, since all the parameters are pairwise. By standard multivariate Gaussian calculations,

$$p(x|y) = \mathcal{N}(B^{-1}\gamma(y), B^{-1}) \quad (4)$$

$$\{\gamma(y)\}_s = \alpha_s + \sum_j \rho_{sj}(y_j) \quad (5)$$

$$p(y) \propto \exp\left(\sum_{j=1}^q \sum_{r=1}^j \phi_{rj}(y_r, y_j) + \frac{1}{2}\gamma(y)^T B^{-1}\gamma(y)\right). \quad (6)$$

Thus, we see that the continuous variables conditioned on the discrete are multivariate Gaussian with common covariance, but with means that depend on the value of the discrete variables. The means depend additively on the values of the discrete variables since $\{\gamma(y)\}_s = \sum_{j=1}^r \rho_{sj}(y_j)$. The marginal $p(y)$ has a known form, so for models with few number of discrete variables we can sample efficiently.

2.1 RELATED WORK ON MIXED GRAPHICAL MODELS

Lauritzen (1996) proposed a type of mixed graphical model, with the property that conditioned on discrete variables, $p(x|y) = \mathcal{N}(\mu(y), \Sigma(y))$. The homogeneous mixed graphical model enforces common covariance, $\Sigma(y) \equiv \Sigma$. Thus, our proposed model is a special case of Lauritzen's mixed model with the following assumptions: common covariance, additive mean assumptions and the marginal $p(y)$ factorizes as a pairwise discrete Markov random field. With these three assumptions, the full model simplifies to the mixed pairwise model presented. Although the full model is more general, the number of parameters scales exponentially with the number of discrete variables, and the conditional distributions are

not as convenient. For each state of the discrete variables there is a mean and covariance. Consider an example with q binary variables and p continuous variables; the full model requires estimates of 2^q mean vectors and covariance matrices in p dimensions. Even if the homogeneous constraint is imposed on Lauritzen's model, there are still 2^q mean vectors for the case of binary discrete variables. The full mixed model is very complex and cannot be easily estimated from data without some additional assumptions. In comparison, the mixed pairwise model has number of parameters $O((p + q)^2)$ and allows for a natural regularization scheme which makes it appropriate for high-dimensional data.

An alternative to the regularization approach that we take in this article, is the limited-order correlation hypothesis testing method by Tur and Castelo (2012). The authors developed a hypothesis test via likelihood ratios for conditional independence. However, they restricted to the case where the discrete variables are marginally independent so the maximum likelihood estimates are well-defined for $p > n$.

There is a line of work regarding parameter estimation in undirected mixed models that are decomposable: any path between two discrete variables cannot contain only continuous variables. These models allow for fast exact maximum likelihood estimation through node-wise regressions, but are only applicable when the structure is known and $n > p$ (Edwards 2000). There is also related work on parameter learning in directed mixed graphical models. Since our primary goal is to learn the graph structure, we forgo exact parameter estimation and use the PL. Similar to the exact maximum likelihood in decomposable models, the PL can be interpreted as node-wise regressions that enforce symmetry.

After we proposed our model,¹ the independent work of Cheng, Levina, and Zhu (2013) appeared, which considers a more complicated mixed graphical model. Their model includes higher order interaction terms by allowing the covariance of the continuous variables to be a function of the categorical variables y , which results in a larger model similar to Lauritzen's model. This results in a model with $O(p^2q + q^2)$ parameters, as opposed to $O(p^2 + q^2)$ in our proposed pairwise model. We believe that in the high-dimensional setting where data sparsity is an issue a simpler model is an advantage.

To our knowledge, this work is the first to consider convex optimization procedures for learning the edge structure in mixed graphical models.

3. PARAMETER ESTIMATION: MAXIMUM LIKELIHOOD AND PL

Given samples $(x_i, y_i)_{i=1}^n$, we want to find the maximum likelihood estimate of Θ . This can be done by minimizing the negative log-likelihood of the samples:

$$\ell(\Theta) = - \sum_{i=1}^n \log p(x_i, y_i; \Theta), \quad (7)$$

¹Cheng, Levina, and Zhu (2013), <http://arxiv.org/abs/1304.2810>, appeared on arXiv 11 months after our original article was put on arXiv, <http://arxiv.org/abs/1205.5012>.

where

$$\begin{aligned} \log p(x, y; \Theta) = & \sum_{s=1}^p \sum_{t=1}^p -\frac{1}{2} \beta_{st} x_s x_t + \sum_{s=1}^p \alpha_s x_s + \sum_{s=1}^p \sum_{j=1}^q \rho_{sj}(y_j) x_s \\ & + \sum_{j=1}^q \sum_{r=1}^j \phi_{rj}(y_r, y_j) - \log Z(\Theta). \end{aligned} \quad (8)$$

The negative log-likelihood is convex, so standard gradient-descent algorithms can be used for computing the maximum likelihood estimates. The major obstacle here is $Z(\Theta)$, which involves a high-dimensional integral. Since the pairwise mixed model includes both the discrete and continuous models as special cases, maximum likelihood estimation is at least as difficult as the two special cases, the first of which is a well-known computationally intractable problem. We defer the discussion of maximum likelihood estimation to the supplementary material.

3.1 PSEUDO-LIKELIHOOD

The PL method (Besag 1975) is a computationally efficient and consistent estimator formed by products of all the conditional distributions:

$$\tilde{\ell}(\Theta|x, y) = - \sum_{s=1}^p \log p(x_s|x_{\setminus s}, y; \Theta) - \sum_{r=1}^q \log p(y_r|x, y_{\setminus r}; \Theta). \quad (9)$$

The conditional distributions $p(x_s|x_{\setminus s}, y; \theta)$ and $p(y_r = k|y_{\setminus r}, x; \theta)$ take on the familiar form of linear Gaussian and (multiclass) logistic regression, as we pointed out in (2) and (3). Here, are the details:

- The conditional distribution of a continuous variable x_s is Gaussian with a linear regression model for the mean, and unknown variance.

$$p(x_s|x_{\setminus s}, y; \Theta) = \frac{\sqrt{\beta_{ss}}}{\sqrt{2\pi}} \exp \left(-\frac{\beta_{ss}}{2} \left(\frac{\alpha_s + \sum_j \rho_{sj}(y_j) - \sum_{t \neq s} \beta_{st} x_t}{\beta_{ss}} - x_s \right)^2 \right) \quad (10)$$

- The conditional distribution of a discrete variable y_r with L_r states is a multinomial distribution, as used in (multiclass) logistic regression. Whenever a discrete variable is a predictor, each of its levels contribute an additive effect; continuous variables contribute linear effects.

$$p(y_r|y_{\setminus r}, x; \Theta) = \frac{\exp \left(\sum_s \rho_{sr}(y_r) x_s + \phi_{rr}(y_r, y_r) + \sum_{j \neq r} \phi_{rj}(y_r, y_j) \right)}{\sum_{l=1}^{L_r} \exp \left(\sum_s \rho_{sr}(l) x_s + \phi_{rr}(l, l) + \sum_{j \neq r} \phi_{rj}(l, y_j) \right)}. \quad (11)$$

Taking the negative log of both gives us

$$-\log p(x_s | x_{\setminus s}, y; \Theta) = -\frac{1}{2} \log \beta_{ss} + \frac{\beta_{ss}}{2} \left(\frac{\alpha_s}{\beta_{ss}} + \sum_j \frac{\rho_{sj}(y_j)}{\beta_{ss}} - \sum_{t \neq s} \frac{\beta_{st}}{\beta_{ss}} x_t - x_s \right)^2 \quad (12)$$

$$-\log p(y_r | y_{\setminus r}, x; \Theta) = -\log \frac{\exp \left(\sum_s \rho_{sr}(y_r) x_s + \phi_{rr}(y_r, y_r) + \sum_{j \neq r} \phi_{rj}(y_r, y_j) \right)}{\sum_{l=1}^{L_r} \exp \left(\sum_s \rho_{sr}(l) x_s + \phi_{rr}(l, l) + \sum_{j \neq r} \phi_{rj}(l, y_j) \right)}. \quad (13)$$

A generic parameter block, θ_{uv} , corresponding to an edge (u, v) appears twice in the PL, once for each of the conditional distributions $p(z_u | z_v)$ and $p(z_v | z_u)$.

Proposition 1. The negative log-PL in (9) is jointly convex in all the parameters $\{\beta_{ss}, \beta_{st}, \alpha_s, \phi_{rj}, \rho_{sj}\}$ over the region $\beta_{ss} > 0$.

We prove Proposition 1 in the supplementary materials.

3.2 SEPARATE NODE-WISE REGRESSION

A simple approach to parameter estimation is via separate node-wise regressions; a generalized linear model is used to estimate $p(z_s | z_{\setminus s})$ for each s . Separate regressions were used by Meinshausen and Bühlmann (2006) for the Gaussian graphical model and by Ravikumar, Wainwright, and Lafferty (2010) for the Ising model. The method can be thought of as an asymmetric form of the PL since the PL enforces that the parameters are shared across the conditionals. Thus, the number of parameters estimated in the separate regression is approximately double that of the PL, so we expect that the PL outperforms at low sample sizes and low regularization regimes. The node-wise regression was used as our baseline method since it is straightforward to extend it to the mixed model. As we predicted, the PL or joint procedure outperforms separate regressions; see top left box of Figures 5 and 6. Liu and Ihler (2012, 2011) confirmed that the separate regressions are outperformed by PL in numerous synthetic settings.

Concurrent work of Yang et al. (2012, 2013) extends the separate node-wise regression model from the special cases of Gaussian and categorical regressions to generalized linear models, where the univariate conditional distribution of each node $p(x_s | x_{\setminus s})$ is specified by a generalized linear model (e.g., Poisson, categorical, Gaussian). By specifying the conditional distributions, Besag (1974) showed that the joint distribution is also specified. Thus, another way to justify our mixed model is to define the conditionals of a continuous variable as Gaussian linear regression and the conditionals of a categorical variable as multiple logistic regression and use the results in Besag (1974) to arrive at the joint distribution in (1). However, the neighborhood selection algorithm in Yang et al. (2012, 2013) is restricted to models of the form $p(x) \propto \exp \left(\sum_s \theta_s x_s + \sum_{s,t} \theta_{st} x_s x_t + \sum_s C(x_s) \right)$. In particular, this procedure cannot be applied to edge selection in our pairwise mixed model in (1) or the categorical model in (2) with greater than two states. Our baseline method of

separate regressions is closely related to the neighborhood selection algorithm they proposed; the baseline can be considered as a generalization of Yang et al. (2012, 2013) to allow for more general pairwise interactions with the appropriate regularization to select edges. Unfortunately, the theoretical results in Yang et al. (2012, 2013) do not apply to the baseline nodewise regression method, nor the joint PL.

4. CONDITIONAL INDEPENDENCE AND PENALTY TERMS

In this section, we show how to incorporate edge selection into the maximum likelihood or PL procedures. In the graphical representation of probability distributions, the absence of an edge $e = (u, v)$ corresponds to a conditional independency statement that variables x_u and x_v are conditionally independent given all other variables (Koller and Friedman 2009). We would like to maximize the likelihood subject to a penalization on the number of edges since this results in a sparse graphical model. In the pairwise mixed model, there are three type of edges

1. β_{st} is a scalar that corresponds to an edge from x_s to x_t . $\beta_{st} = 0$ implies x_s and x_t are conditionally independent given all other variables. This parameter is in two conditional distributions, corresponding to either x_s or x_t is the response variable, $p(x_s|x_{\setminus s}, y; \Theta)$ and $p(x_t|x_{\setminus t}, y; \Theta)$.
2. ρ_{sj} is a vector of length L_j . If $\rho_{sj}(y_j) = 0$ for all values of y_j , then y_j and x_s are conditionally independent given all other variables. This parameter is in two conditional distributions, corresponding to either x_s or y_j being the response variable: $p(x_s|x_{\setminus s}, y; \Theta)$ and $p(y_j|x, y_{\setminus j}; \Theta)$.
3. ϕ_{rj} is a matrix of size $L_r \times L_j$. If $\phi_{rj}(y_r, y_j) = 0$ for all values of y_r and y_j , then y_r and y_j are conditionally independent given all other variables. This parameter is in two conditional distributions, corresponding to either y_r or y_j being the response variable, $p(y_r|x, y_{\setminus r}; \Theta)$ and $p(y_j|x, y_{\setminus j}; \Theta)$.

For conditional independencies that involve discrete variables, the absence of that edge requires that the entire matrix ϕ_{rj} or vector ρ_{sj} is 0.² The form of the pairwise mixed model motivates the following regularized optimization problem

$$\underset{\Theta}{\text{minimize}} \ell_{\lambda}(\Theta) = \ell(\Theta) + \lambda \left(\sum_{s < t} \mathbb{1}[\beta_{st} \neq 0] + \sum_{sj} \mathbb{1}[\rho_{sj} \neq 0] + \sum_{r < j} \mathbb{1}[\phi_{rj} \neq 0] \right). \quad (14)$$

See Figure 1 for a visualization of the edges. All parameters that correspond to the same edge are grouped in the same indicator function. This problem is nonconvex, so we replace the l_0 sparsity and group sparsity penalties with the appropriate convex relaxations. For

²If $\rho_{sj}(y_j) = \text{constant}$, then x_s and y_j are also conditionally independent. However, the unpenalized term α will absorb the constant, so the estimated $\rho_{sj}(y_j)$ will never be constant for $\lambda > 0$.

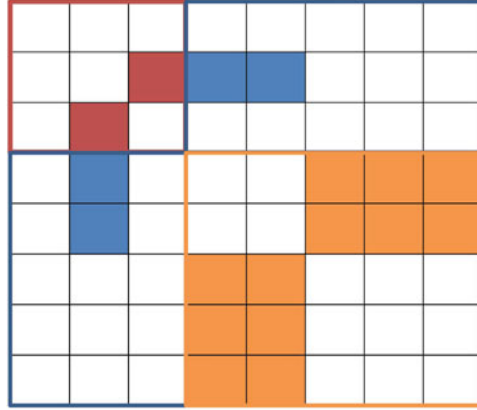


Figure 1. Symmetric matrix represents the parameters Θ of the model. This example has $p = 3$, $q = 2$, $L_1 = 2$, and $L_2 = 3$. The red square corresponds to the continuous graphical model coefficients B and the solid red square is the scalar β_{st} . The blue square corresponds to the coefficients ρ_{sj} and the solid blue square is a vector of parameters $\rho_{sj}(\cdot)$. The orange square corresponds to the coefficients ϕ_{rj} and the solid orange square is a matrix of parameters $\phi_{rj}(\cdot, \cdot)$. The matrix is symmetric, so each parameter block appears in two of the conditional probability regressions.

scalars, we use the absolute value (l_1 norm), for vectors we use the l_2 norm, and for matrices we use the Frobenius norm. This choice corresponds to the standard relaxation from group l_0 to group l_1/l_2 (group lasso) norm (Yuan and Lin 2006; Bach et al. 2011),

$$\underset{\Theta}{\text{minimize}} \quad \ell_{\lambda}(\Theta) = \ell(\Theta) + \lambda \left(\sum_{s=1}^p \sum_{t=1}^{s-1} |\beta_{st}| + \sum_{s=1}^p \sum_{j=1}^q \|\rho_{sj}\|_2 + \sum_{j=1}^q \sum_{r=1}^{j-1} \|\phi_{rj}\|_F \right). \quad (15)$$

5. CALIBRATED REGULARIZERS

In (15) each of the group penalties are treated as equals, irrespective of the size of the group. We suggest a calibration or weighting scheme to balance the load in a more equitable way. We introduce weights for each group of parameters and show how to choose the weights such that each parameter set is treated equally under p_F , the fully factorized independence model.³

$$\underset{\Theta}{\text{minimize}} \quad \ell(\Theta) + \lambda \left(\sum_{t=1}^p \sum_{s=t+1}^p w_{st} |\beta_{st}| + \sum_{s=1}^p \sum_{j=1}^q w_{sj} \|\rho_{sj}\|_2 + \sum_{j=1}^q \sum_{r=1}^{j-1} w_{rj} \|\phi_{rj}\|_F \right). \quad (16)$$

Based on the KKT conditions (Friedman et al. 2007), the parameter group θ_g is nonzero if

$$\left\| \frac{\partial \ell}{\partial \theta_g} \right\| > \lambda w_g$$

³Under the independence model p_F is fully factorized $p(x, y) = \prod_{s=1}^p p(x_s) \prod_{r=1}^q p(y_r)$.

where θ_g and w_g represents one of the parameter groups and its corresponding weight.

Now $\frac{\partial \ell}{\partial \theta_g}$ can be viewed as a generalized residual, and for different groups these are different dimensions—for example, scalar/vector/matrix. So, even under the independence model (when all terms should be zero), one might expect some terms $\left\| \frac{\partial \ell}{\partial \theta_g} \right\|$ to have a better random chance of being nonzero (e.g., those of bigger dimensions).

Thus for all parameters to be on equal footing, we would like to choose the weights w_g such that

$$w_g = E_{p_F} \left\| \frac{\partial \ell}{\partial \theta_g} \right\| \times \text{constant}, \quad (17)$$

where p_F is the fully factorized (independence) model. We will refer to these as the exact weights. We do not have a closed form expression for computing them, but they can be easily estimated by a simple null-model simulation. We also propose an approximation that can be computed exactly.

It is straightforward to compute $E_{p_F} \left\| \frac{\partial \ell}{\partial \theta_g} \right\|^2$ in closed form, which leads to the approximate weights

$$w_g \propto \sqrt{E_{p_F} \left\| \frac{\partial \ell}{\partial \theta_g} \right\|^2}. \quad (18)$$

In the supplementary material, we show that for the three types of edges this leads to the expressions

$$\begin{aligned} w_{st} &= \sigma_s \sigma_t, \\ w_{sj} &= \sigma_s \sqrt{\sum_a p_a (1 - p_a)}, \end{aligned}$$

and

$$w_{rj} = \sqrt{\sum_a p_a (1 - p_a) \sum_b q_b (1 - q_b)}, \quad (19)$$

where σ_s is the standard deviation of the continuous variable x_s , $p_a = \Pr(y_r = a)$ and $q_b = \Pr(y_j = b)$. For all three types of parameters, the weight has the form of $w_{uv}^2 = \text{tr}(\text{cov}(z_u))\text{tr}(\text{cov}(z_v))$, where z represents a generic variable and $\text{cov}(z)$ is the variance-covariance matrix of z .

We conducted a small simulation study to show that calibration is needed. Consider a model with four independent variables: two continuous with variance 10 and 1, and two discrete variables with 10 and 2 levels.

There are six candidate edges in this model and from row 1 of [Table 1](#) we can see the sizes of the gradients are different. In fact, the ratio of the largest gradient to the smallest gradient is greater than 4. The edges ρ_{11} and ρ_{12} involving the first continuous variable with variance 10 have larger edge weights than the corresponding edges, ρ_{21} and ρ_{22} involving the second continuous variable with variance 1. Similarly, the edges involving the first

Table 1. Penalty weights in a six-edge model. Row 1 shows the exact weights w_g computed via (17) using Monte Carlo simulation

	$\left\ \frac{\partial \ell}{\partial \phi_{12}} \right\ _F$	$\left\ \frac{\partial \ell}{\partial \rho_{11}} \right\ _2$	$\left\ \frac{\partial \ell}{\partial \rho_{21}} \right\ _2$	$\left\ \frac{\partial \ell}{\partial \rho_{12}} \right\ _2$	$\left\ \frac{\partial \ell}{\partial \rho_{22}} \right\ _2$	$\left \frac{\partial \ell}{\partial \beta_{12}} \right $
Exact weights w_g (17)	0.18	0.63	0.19	0.47	0.15	0.53
Approximate weights w_g (19)	0.13	0.59	0.18	0.44	0.13	0.62

NOTE: Row 2 shows the approximate weights computed via (19). Note that the weights are far from uniform, and the approximate weights are close to the exact weights.

discrete variable with 10 levels are larger than the edges involving the second discrete variable with two levels. This reflects our intuition that larger variance and longer vectors will have larger norm.

The approximate weights from Equation (19) are a very good approximation to $\|\nabla \ell\|$. Since the weights are only defined up to a proportionality constant, the cosine similarity is an appropriate measure of the quality of approximation. For this simulation, the cosine similarity is

$$\text{sim}(w, \|\nabla \ell\|) = 0.993,$$

which is extremely close to 1.

Using the weights from Table 1, we conducted a second simulation to record which edge would enter first when the four variables are independent. The results are shown in Table 2. Both exact and approximate calibration perform much better than no calibration, but neither deliver the ideal 1/6th probabilities one might desire in a situation like this. Of course, calibrating to the expectations of the gradient does not guarantee equal entry probability, but it brings them much closer.

The exact weights do not have simple closed-form expressions, but they can be easily computed via Monte Carlo. This can be done by simulating independent Gaussians and multinomials with the appropriate marginal variance σ_s and marginal probabilities p_a , then approximating the expectation in (17) by an average. The computational cost of this procedure is negligible compared to fitting the mixed model, so in practice either the exact or approximate weights can be used.

Table 2. Fraction of times an edge is the first selected by the group lasso regularizer, based on 1000 simulation runs

	ϕ_{12}	ρ_{11}	ρ_{21}	ρ_{12}	ρ_{22}	β_{12}
No calibration $w_g = 1$	0.000	0.487	0.000	0.163	0.000	0.350
Exact w_g (17)	0.101	0.092	0.097	0.249	0.227	0.234
Approximate w_g (19)	0.144	0.138	0.134	0.196	0.190	0.198

NOTE: Ideally each edge should be first 1/6th of the time (0.167), with a standard error of 0.012. The group lasso with equal weights (first row) is highly unbalanced. Using the exact weights from (17) is quite good (second row), while the approximate weighing scheme of (19) (third row) appears to perform the best.

6. MODEL SELECTION CONSISTENCY

In this section, we study the model selection consistency, whether the correct edge set is selected and the parameter estimates are close to the truth, of the PL and maximum likelihood estimators. Consistency can be established using the framework first developed in Ravikumar, Wainwright, and Lafferty (2010) and later extended to general m -estimators by Lee, Sun, and Taylor (2013). Instead of stating the full results and proofs, we will illustrate the type of theorems that can be shown and defer the rigorous statements to the supplementary material.

First, we define some notation. Recall that Θ is the vector of parameters being estimated $\{\beta_{ss}, \beta_{st}, \alpha_s, \phi_{rj}, \rho_{sj}\}$, Θ^* be the true parameters that estimated the model, and $Q = \nabla^2 \ell(\Theta^*)$. Both maximum likelihood and PL estimation procedures can be written as a convex optimization problem of the form

$$\text{minimize } \ell(\Theta) + \lambda \sum_{g \in G} \|\Theta_g\|_2, \quad (20)$$

where $\ell(\theta) = \{\ell_{\text{ML}}, \ell_{\text{PL}}\}$ is one of the two log-likelihoods. The regularizer

$$\sum_{g \in G} \|\Theta_g\| = \lambda \left(\sum_{s=1}^p \sum_{t=1}^{s-1} |\beta_{st}| + \sum_{s=1}^p \sum_{j=1}^q \|\rho_{sj}\|_2 + \sum_{j=1}^q \sum_{r=1}^{j-1} \|\phi_{rj}\|_F \right).$$

The set G indexes the edges β_{st} , ρ_{sj} , and ϕ_{rj} , and Θ_g is one of the three types of edges. Let A and I represent the active and inactive groups in Θ , so $\Theta_g^* \neq 0$ for any $g \in A$ and $\Theta_g^* = 0$ for any $g \in I$.

Let $\hat{\Theta}$ be the minimizer to Equation (20). Then $\hat{\Theta}$ satisfies,

1. $\|\hat{\Theta} - \Theta^*\|_2 \leq C \sqrt{\frac{|A| \log |G|}{n}}$
2. $\hat{\Theta}_g = 0$ for $g \in I$.

The exact statement of the theorem is given in the supplementary material.

7. OPTIMIZATION ALGORITHMS

In this section, we discuss two algorithms for solving (15): the proximal gradient and the proximal Newton methods. This is a convex optimization problem that decomposes into the form $f(x) + g(x)$, where f is smooth and convex and g is convex but possibly nonsmooth. In our case, f is the negative log-likelihood or negative log-PL and g are the group sparsity penalties.

Block coordinate descent is a frequently used method when the nonsmooth function g is the l_1 or group l_1 . It is especially easy to apply when the function f is quadratic, since each block coordinate update can be solved in closed form for many different nonsmooth g (Friedman et al. 2007). The smooth f in our particular case is not quadratic, so each block update cannot be solved in closed form. However, in certain problems (sparse in-

verse covariance), the update can be approximately solved by using an appropriate inner optimization routine (Friedman, Hastie, and Tibshirani 2008).

7.1 PROXIMAL GRADIENT

Problems of this form are well-suited for the proximal gradient and accelerated proximal gradient algorithms as long as the proximal operator of g can be computed (Beck and Teboulle 2010; Combettes and Pesquet 2011)

$$\text{prox}_t(x) = \underset{u}{\operatorname{argmin}} \frac{1}{2t} \|x - u\|^2 + g(u). \quad (21)$$

For the sum of l_2 group sparsity penalties considered, the proximal operator takes the familiar form of soft-thresholding and group soft-thresholding (Bach et al. 2011). Since the groups are nonoverlapping, the proximal operator simplifies to scalar soft-thresholding for β_{sj} and group soft-thresholding for ρ_{sj} and ϕ_{rj} .

The class of proximal gradient and accelerated proximal gradient algorithms is directly applicable to our problem. These algorithms work by solving a first-order model at the current iterate x_k

$$\underset{u}{\operatorname{argmin}} f(x_k) + \nabla f(x_k)^T (u - x_k) + \frac{1}{2t} \|u - x_k\|^2 + g(u) \quad (22)$$

$$= \underset{u}{\operatorname{argmin}} \frac{1}{2t} \|u - (x_k - t \nabla f(x_k))\|^2 + g(u) \quad (23)$$

$$= \text{prox}_t(x_k - t \nabla f(x_k)). \quad (24)$$

The proximal gradient iteration is given by $x_{k+1} = \text{prox}_t(x_k - t \nabla f(x_k))$, where t is determined by line search. The theoretical convergence rates and properties of the proximal gradient algorithm and its accelerated variants are well-established (Beck and Teboulle 2010). The accelerated proximal gradient method achieves linear convergence rate of $O(c^k)$ when the objective is strongly convex and the sublinear rate $O(1/k^2)$ for nonstrongly convex problems.

The TFOCS framework (Becker, Candès, and Grant 2011) is a package that allows us to experiment with six different variants of the accelerated proximal gradient algorithm. The TFOCS authors found that the Auslender-Teboulle algorithm exhibited less oscillatory behavior, and proximal gradient experiments in the next section were done using the Auslender-Teboulle implementation in TFOCS.

7.2 PROXIMAL NEWTON ALGORITHMS

The class of proximal Newton algorithms is a second-order analog of the proximal gradient algorithms with a quadratic convergence rate (Schmidt 2010; Schmidt, Kim, and Sra 2011; Lee, Sun, and Saunders 2012). It attempts to incorporate second-order information about the smooth function f into the model function. At each iteration, it minimizes a quadratic model centered at x_k

$$\underset{u}{\operatorname{argmin}} f(x_k) + \nabla f(x_k)^T (u - x_k) + \frac{1}{2t} (u - x_k)^T H(u - x_k) + g(u) \quad (25)$$

$$= \operatorname{argmin}_u \frac{1}{2t} (u - x_k + t H^{-1} \nabla f(x_k))^T H (u - x_k + t H^{-1} \nabla f(x_k)) + g(u) \quad (26)$$

$$= \operatorname{argmin}_u \frac{1}{2t} \|u - (x_k - t H^{-1} \nabla f(x_k))\|_H^2 + g(u) \quad (27)$$

$$:= H \operatorname{prox}_t (x_k - t H^{-1} \nabla f(x_k)) \text{ where } H = \nabla^2 f(x_k). \quad (28)$$

Algorithm 1 Proximal Newton

repeat

Solve subproblem $p_k = H \operatorname{prox}_t (x_k - t H_k^{-1} \nabla f(x_k)) - x_k$ using TFOCS.

Find t to satisfy Armijo line search condition with parameter α

$$f(x_k + t p_k) + g(x_k + t p_k) \leq f(x_k) + g(x_k) - \frac{t\alpha}{2} \|p_k\|^2$$

Set $x_{k+1} = x_k + t p_k$

$k = k + 1$

until $\frac{\|x_k - x_{k+1}\|}{\|x_k\|} < \text{tol}$

The $H \operatorname{prox}$ operator is analogous to the proximal operator, but in the $\|\cdot\|_H$ -norm. It simplifies to the proximal operator if $H = I$, but in the general case of positive definite H there is no closed-form solution for many common nonsmooth $g(x)$ (including l_1 and group l_1). However, if the proximal operator of g is available, each of these subproblems can be solved efficiently with proximal gradient. In the case of separable g , coordinate descent is also applicable. Fast methods for solving the subproblem $H \operatorname{prox}_t (x_k - t H^{-1} \nabla f(x_k))$ include coordinate descent methods, proximal gradient methods, or Barzilai-Borwein (Friedman et al. 2007; Wright, Nowak, and Figueiredo 2009; Beck and Teboulle 2010; Combettes and Pesquet 2011). The proximal Newton framework allows us to bootstrap many previously developed solvers to the case of arbitrary loss function f .

Theoretical analysis in Lee, Sun, and Saunders (2012) suggests that proximal Newton methods generally require fewer outer iterations (evaluations of $H \operatorname{prox}$) than first-order methods while providing higher accuracy because they incorporate second-order information. We have confirmed empirically that the proximal Newton methods are faster when n is very large or the gradient is expensive to compute (e.g., maximum likelihood estimation). Since the objective is quadratic, coordinate descent is also applicable to the subproblems. The Hessian matrix H can be replaced by a quasi-newton approximation such as BFGS/L-BFGS/SR1. In our implementation, we use the PNOPT implementation (Lee, Sun, and Saunders 2012).

7.3 PATH ALGORITHM

Frequently in machine learning and statistics, the regularization parameter λ is heavily dependent on the dataset. Here, λ is generally chosen via cross-validation or holdout set performance, so it is convenient to provide solutions over an interval of $[\lambda_{\min}, \lambda_{\max}]$. We

start the algorithm at $\lambda_1 = \lambda_{\max}$ and solve, using the previous solution as warm start, for $\lambda_2 > \dots > \lambda_{\min}$. We find that this reduces the cost of fitting an entire path of solutions (see [Figure 4](#)). λ_{\max} can be chosen as the smallest value such that all parameters are 0 by using the KKT equations (Friedman et al. [2007](#)).

8. CONDITIONAL MODEL

In addition to the variables, we would like to model, there are often additional features or covariates that affect the dependence structure of the variables. For example in genomic data, in addition to expression values, we have attributes associated to each subject such as gender, age and ethnicity. These additional attributes affect the dependence of the expression values, so we can build a conditional model that uses the additional attributes as features. In this section, we show how to augment the pairwise mixed model with features.

Conditional models only model the conditional distribution $p(z|f)$, as opposed to the joint distribution $p(z, f)$, where z are the variables of interest to the prediction task and f are features. These models are frequently used in practice (Lafferty, McCallum, and Pereira [2001](#)). In addition to observing x and y , we observe features f and we build a graphical model for the conditional distribution $p(x, y|f)$. Consider a full pairwise model $p(x, y, f)$ of the form (1). We then choose to only model the joint distribution over only the variables x and y to give us $p(x, y|f)$ which is of the form

$$p(x, y|f; \Theta) = \frac{1}{Z(\Theta|f)} \exp \left(\sum_{s=1}^p \sum_{t=1}^p -\frac{1}{2} \beta_{st} x_s x_t + \sum_{s=1}^p \alpha_s x_s + \sum_{s=1}^p \sum_{j=1}^q \rho_{sj}(y_j) x_s \right. \\ \left. + \sum_{j=1}^q \sum_{r=1}^j \phi_{rj}(y_r, y_j) + \sum_{l=1}^F \sum_{s=1}^p \gamma_{ls} x_s f_l + \sum_{l=1}^F \sum_{r=1}^q \eta_{lr}(y_r) f_l \right). \quad (29)$$

We can also consider a more general model where each pairwise edge potential depends on the features

$$p(x, y|f; \Theta) = \frac{1}{Z(\Theta|f)} \exp \left(\sum_{s=1}^p \sum_{t=1}^p -\frac{1}{2} \beta_{st}(f) x_s x_t + \sum_{s=1}^p \alpha_s(f) x_s \right. \\ \left. + \sum_{s=1}^p \sum_{j=1}^q \rho_{sj}(y_j, f) x_s + \sum_{j=1}^q \sum_{r=1}^j \phi_{rj}(y_r, y_j, f) \right) \quad (30)$$

(29) is a special case of this where only the node potentials depend on features and the pairwise potentials are independent of feature values. The specific parameterized form we consider is $\phi_{rj}(y_r, y_j, f) \equiv \phi_{rj}(y_r, y_j)$ for $r \neq j$, $\rho_{sj}(y_j, f) \equiv \rho_{sj}(y_j)$, and $\beta_{st}(f) = \beta_{st}$. The node potentials depend linearly on the feature values, $\alpha_s(f) = \alpha_s + \sum_{l=1}^F \gamma_{ls} x_s f_l$, and $\phi_{rr}(y_r, y_r, f) = \phi_{rr}(y_r, y_r) + \sum_l \eta_{lr}(y_r) f_l$.

9. EXPERIMENTAL RESULTS

We present experimental results on synthetic data, survey data, and on a conditional model.

9.1 SYNTHETIC EXPERIMENTS

In the synthetic experiment, the training points are sampled from a true model with 10 continuous variables and 10 binary variables. The edge structure is shown in Figure 2(a). The λ is chosen proportional to $\sqrt{\frac{\log(p+q)}{n}}$ as suggested by the theoretical results in Section 6. We experimented with 3 values $\lambda = \{1, 5, 10\} \sqrt{\frac{\log(p+q)}{n}}$ and chose $\lambda = 5 \sqrt{\frac{\log(p+q)}{n}}$ so that the true edge set was recovered by the algorithm for the sample size $n = 2000$. We see from the experimental results that recovery of the correct edge set undergoes a sharp phase transition, as expected. With $n = 1000$ samples, the PL is recovering the correct edge set with probability nearly 1. The maximum likelihood was performed using an exact evaluation of the gradient and log-partition. The poor performance of the maximum likelihood estimator is explained by the maximum likelihood objective violating the irrepresentable condition; a similar example was discussed by Ravikumar, Wainwright, and Lafferty (2010, Section 3.1.1), where the maximum likelihood is not irrepresentable, yet the neighborhood selection procedure is. The phase transition experiments were done using the proximal Newton algorithm discussed in Section 7.2.

We also run the proximal Newton algorithm for a sequence of instances with $p = q = 10, 50, 100, 500, 1000$ and $n = 500$. The largest instance has 2000 variables and takes 12.5 hr to complete. The timing results are summarized in Figure 3.

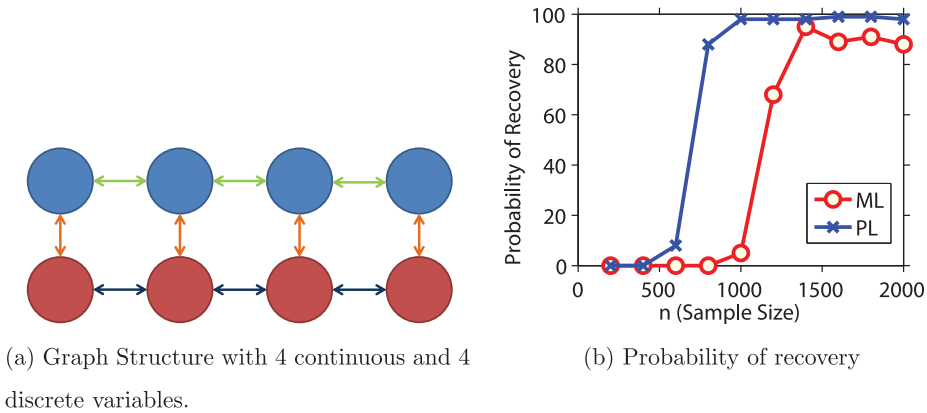


Figure 2. (a) shows the graph used in the synthetic experiments for $p = q = 4$; the experiment actually used $p=10$ and $q=10$. Blue nodes are continuous variables, red nodes are binary variables and the orange, green and dark blue lines represent the three types of edges. (b) is a plot of the probability of correct edge recovery, meaning every true edge is selected and no nonedge is selected, at a given sample size using maximum likelihood and PL. Results are averaged over 100 trials. (a) Graph Structure with four continuous and four discrete variables. (b) Probability of recovery.

$p + q$	Time per iteration (sec)	Total time (min)	Number of iterations
20	0.13	0.003	13
100	4.39	1.32	18
200	18.44	6.45	21
1000	245.34	139	34
2000	1025.6	752	44

Figure 3. Timing experiments for various instances of the graph in Figure 2(a). The number of variables range from 20 to 2000 with $n = 500$.

9.2 SURVEY EXPERIMENTS

The census survey dataset we consider consists of 11 variables, of which two are continuous and nine are discrete: age (continuous), log-wage (continuous), year (7 states), sex (2 states), marital status (5 states), race (4 states), education level (5 states), geographic region (9 states), job class (2 states), health (2 states), and health insurance (2 states). The dataset was assembled by Steve Miller of OpenBI.com from the March 2011 Supplement to Current Population Survey data. All the evaluations are done using a holdout test set of size 100,000 for the survey experiments. The regularization parameter λ is varied over the interval $[5 \times 10^{-5}, 0.7]$ at 50 points equispaced on log-scale for all experiments. In practice, λ can be chosen to minimize the holdout log-PL.

9.2.1 Model Selection. In Figure 4, we study the model selection performance of learning a graphical model over the 11 variables under different training samples sizes. We see that as the sample size increases, the optimal model is increasingly dense, and less regularization is needed.

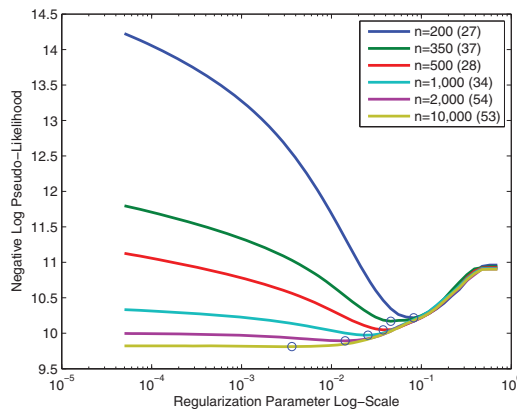


Figure 4. Model selection under different training set sizes. Circle denotes the lowest test set negative log-PL and the number in parentheses is the number of edges in that model at the lowest test negative log-PL. The saturated model has 55 edges.

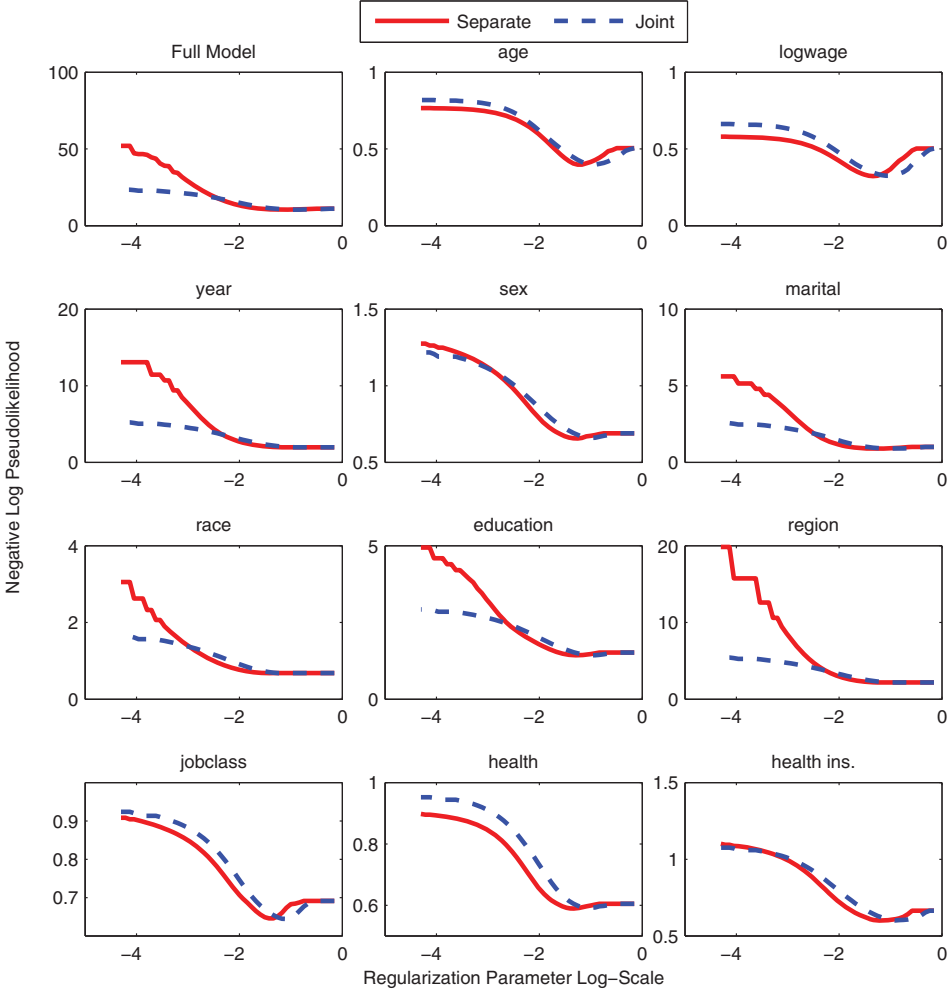


Figure 5. Separate regression versus pseudo-likelihood $n = 100$. The y-axis is the appropriate regression loss for the response variable. For low levels of regularization and at small training sizes, the pseudo-likelihood seems to overfit less; this may be due to a global regularization effect from fitting the joint distribution as opposed to separate regressions.

9.2.2 Comparing Against Separate Regressions. A sensible baseline method to compare against is a separate regression algorithm. This algorithm fits a linear Gaussian or (multiclass) logistic regression of each variable conditioned on the rest. We can evaluate the performance of the PL by evaluating $-\log p(x_s | x_{\setminus s}, y)$ for linear regression and $-\log p(y_r | y_{\setminus r}, x)$ for (multiclass) logistic regression. Since regression is directly optimizing this loss function, it is expected to do better. The PL objective is similar, but has half the number of parameters as the separate regressions since the coefficients are shared between two of the conditional likelihoods. From Figures 5 and 6, we can see that the PL performs very similarly to the separate regressions and sometimes even outperforms regression. The benefit of the PL is that we have learned parameters of the joint distribution

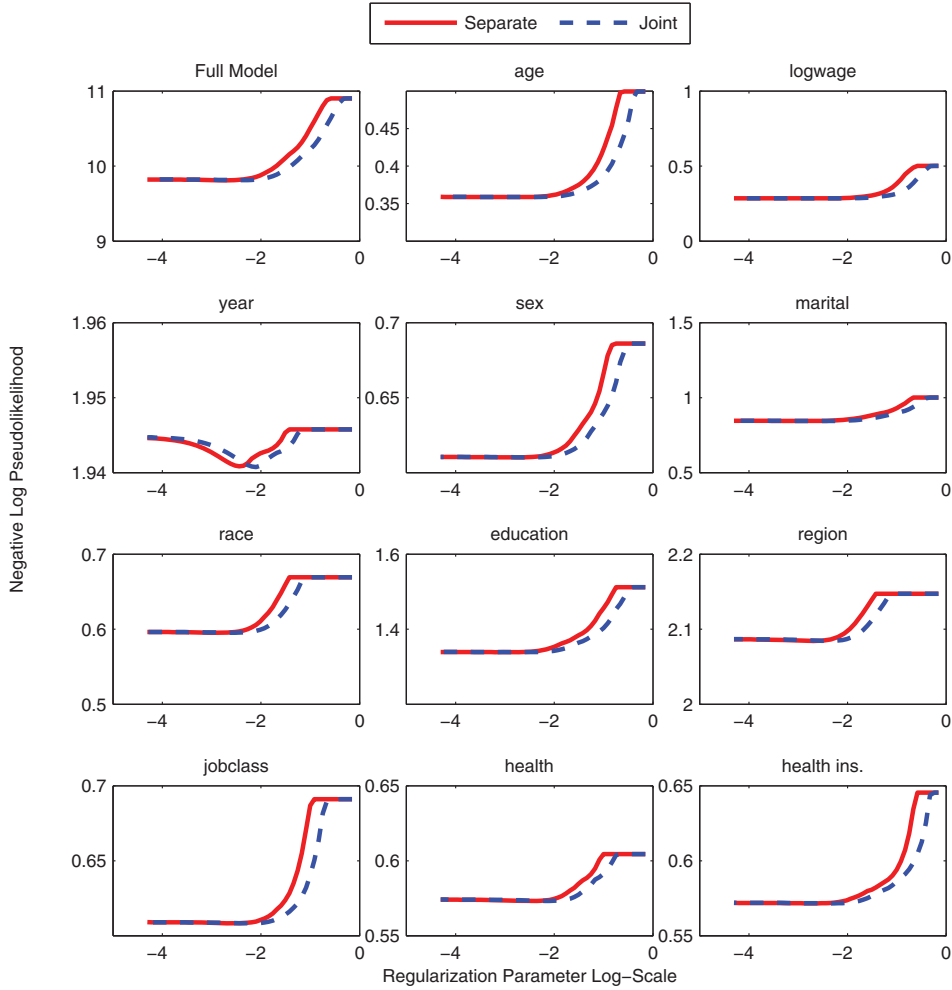


Figure 6. Separate regression versus pseudo-likelihood $n=10,000$. y-axis is the appropriate regression loss for the response variable. At large sample sizes, separate regressions and pseudo-likelihood perform very similarly. This is expected since this is nearing the asymptotic regime.

$p(x, y)$ and not just of the conditionals $p(x_s|y, x_{\setminus s})$. On the test dataset, we can compute quantities such as conditionals over arbitrary sets of variables $p(y_A, x_B|y_{A^c}, x_{B^c})$ and marginals $p(x_A, y_B)$ (Koller and Friedman 2009). This would not be possible using the separate regressions.

9.2.3 Conditional Model. Using the conditional model (29), we model only the three variables logwage, education(5), and jobclass(2). The other 8 variables are only used as features. The conditional model is then trained using the PL. We compare against the generative model that learns a joint distribution on all 11 variables. From Figure 7, we see that the conditional model outperforms the generative model, except at small sample sizes. This is expected since the conditional distribution models less variables. At very small sample sizes and small λ , the generative model outperforms the conditional model. This is

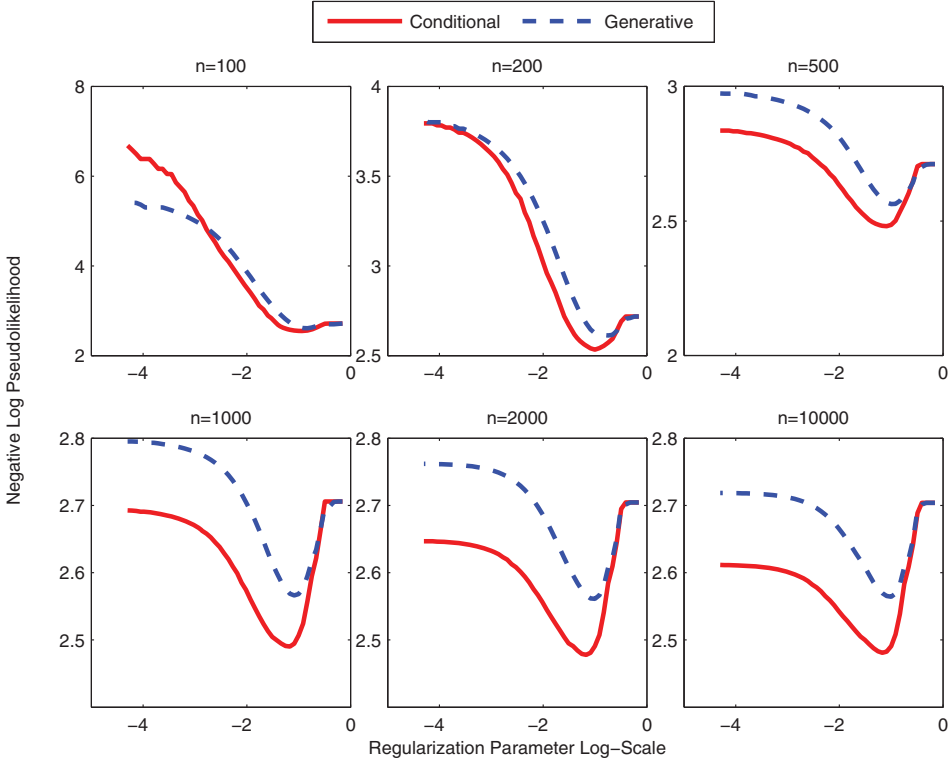


Figure 7. Conditional model versus generative model at various sample sizes. y-axis is test set performance is evaluated on negative log-PL of the conditional model. The conditional model outperforms the full generative model at except the smallest sample size $n = 100$.

likely because generative models converge faster (with less samples) than discriminative models to its optimum.

9.2.4 Maximum Likelihood versus PL. The maximum likelihood estimates are computable for very small models such as the conditional model previously studied. The PL was originally motivated as an approximation to the likelihood that is computationally tractable. We compare the maximum likelihood and maximum PL on two different evaluation criteria: the negative log-likelihood and negative log-PL. In Figure 8, we find that the PL outperforms maximum likelihood under both the negative log-likelihood and negative log-PL. We would expect that the PL trained model does better on the PL evaluation and maximum likelihood trained model does better on the likelihood evaluation. However, we found that the PL trained model outperformed the maximum likelihood trained model on both evaluation criteria. Although asymptotic theory suggests that maximum likelihood is more efficient than the PL, this analysis is applicable because of the finite sample regime and misspecified model. See Liang and Jordan (2008) for asymptotic analysis of PL and maximum likelihood under a well-specified model. We also observed the PL slightly outperforming the maximum likelihood in the synthetic experiment of Figure 2(b).

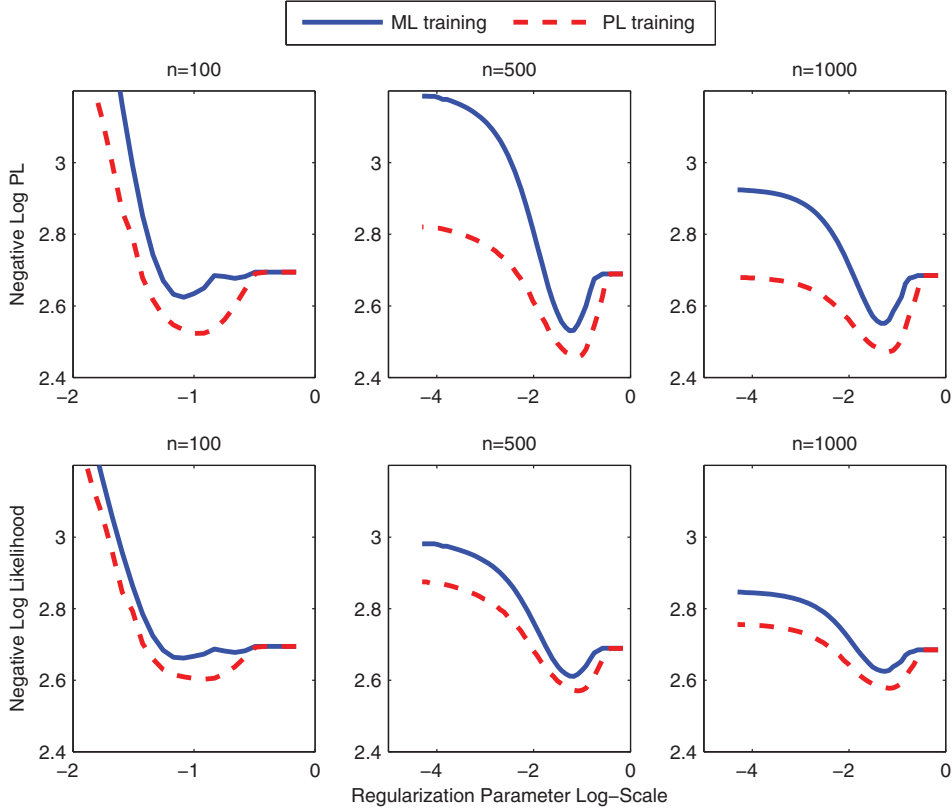


Figure 8. Maximum likelihood versus pseudo-likelihood (PL). The y-axis for top row is the negative log-PL. The y-axis for bottom row is the negative log-likelihood. PL outperforms maximum likelihood across all the experiments.

10. CONCLUSION

This work proposes a new pairwise mixed graphical model, which combines the Gaussian graphical model and discrete graphical model. Due to the introduction of discrete variables, the maximum likelihood estimator is computationally intractable, so we investigated the PL estimator. To learn the structure of this model, we use the appropriate group sparsity penalties with a calibrated weighing scheme. Model selection consistency results are shown for the mixed model using the maximum likelihood and PL estimators. The extension to a conditional model is discussed, since these are frequently used in practice.

We proposed two efficient algorithms for the purpose of estimating the parameters of this model, the proximal Newton and the proximal gradient algorithms. The proximal Newton algorithm is shown to scale to graphical models with 2000 variables on a standard desktop. The model is evaluated on synthetic and the current population survey data, which demonstrates the PL performs well compared to maximum likelihood and nodewise regression.

For future work, it would be interesting to incorporate other discrete variables such as poisson or binomial variables and non-Gaussian continuous variables. This would broaden

the scope of applications that mixed models could be used for. Our work is a first step in that direction.

SUPPLEMENTARY MATERIALS

Code. MATLAB Code that implements structure learning for the mixed graphical model.

Supplementary Material. Technical appendices on sampling from the mixed model, maximum likelihood, calibration weights, model selection consistency, and convexity.

ACKNOWLEDGMENTS

We thank Percy Liang and Rahul Mazumder for helpful discussions. The work on consistency follows from a collaboration with Yuekai Sun and Jonathan Taylor. Jason Lee is supported by the Department of Defense (DoD) through the National Defense Science & Engineering Graduate Fellowship (NDSEG) Program, National Science Foundation Graduate Research Fellowship Program, and the Stanford Graduate Fellowship. Trevor Hastie was partially supported by grant DMS-1007719 from the National Science Foundation, and grant RO1-EB001988-15 from the National Institutes of Health.

[Received August 2013. Revised February 2014.]

REFERENCES

- Bach, F., Jenatton, R., Mairal, J., and Obozinski, G. (2011), “Optimization With Sparsity-Inducing Penalties,” *Foundations and Trends in Machine Learning*, 4, 1–106. Available at <http://dx.doi.org/10.1561/22000000015>. [238,242]
- Banerjee, O., El Ghaoui, L., and d’Aspremont, A. (2008), “Model Selection Through Sparse Maximum Likelihood Estimation for Multivariate Gaussian or Binary Data,” *The Journal of Machine Learning Research*, 9, 485–516. [230]
- Beck, A., and Teboulle, M. (2010), “Gradient-Based Algorithms With Applications to Signal Recovery Problems,” in *Convex Optimization in Signal Processing and Communications*, eds. D. Palomar and Y. Eldar, 42–88, Cambridge: Cambridge University Press. [242,243]
- Becker, S. R., Candès, E. J., and Grant, M. C. (2011), “Templates for Convex Cone Problems With Applications to Sparse Signal Recovery,” *Mathematical Programming Computation*, 3, 165–218. [242]
- Besag, J. (1974), “Spatial Interaction and the Statistical Analysis of Lattice Systems,” *Journal of the Royal Statistical Society, Series B*, 36, 192–236. [236]
- (1975), “Statistical Analysis of Non-Lattice Data,” *The Statistician*, 179–195. [235]
- Cheng, J., Levina, E., and Zhu, J. (2013), “High-Dimensional Mixed Graphical Models,” *arXiv preprint arXiv:1304.2810*. [234]
- Combettes, P. L., and Pesquet, J. C. (2011), “Proximal Splitting Methods in Signal Processing,” in *Fixed-Point Algorithms for Inverse Problems in Science and Engineering*, eds. H. H. Bauschke et al., New York: Springer, pp. 185–212. [242,243]
- Edwards, D. (2000), *Introduction to Graphical Modelling*, New York: Springer. [234]
- Friedman, J., Hastie, T., Höfling, H., and Tibshirani, R. (2007), “Pathwise Coordinate Optimization,” *The Annals of Applied Statistics*, 1, 302–332. [238,241,243]
- Friedman, J., Hastie, T., and Tibshirani, R. (2008), “Sparse Inverse Covariance Estimation With the Graphical Lasso,” *Biostatistics*, 9, 432–441. [230,242]

- (2010), “Applications of the Lasso and Grouped Lasso to the Estimation of Sparse Graphical Models,” Technical Report, Stanford University. [230]
- Guo, J., Levina, E., Michailidis, G., and Zhu, J. (2010), “Joint Structure Estimation for Categorical Markov Networks,” available at <http://www.stat.lsa.umich.edu/~elevina>. [231]
- Höfling, H., and Tibshirani, R. (2009), “Estimation of Sparse Binary Pairwise Markov Networks Using Pseudo-Likelihoods,” *The Journal of Machine Learning Research*, 10, 883–906. [231]
- Jalali, A., Ravikumar, P., Vasuki, V., Sanghavi, S., ECE, UT, and CS, UT (2011), “On Learning Discrete Graphical Models Using Group-Sparse Regularization,” in *Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS)*. [231]
- Kim, S., Sohn, K.-A., and Xing, E. P. (2009), “A Multivariate Regression Approach to Association Analysis of a Quantitative Trait Network,” *Bioinformatics*, 25, i204–i212. [231]
- Koller, D., and Friedman, N. (2009), *Probabilistic Graphical Models: Principles and Techniques*, Cambridge MA: The MIT Press. [237,248]
- Lafferty, J., McCallum, A., and Pereira, F. C. N. (2001), “Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data,” Departmental Papers (CIS), Paper 159, available at http://repository.upenn.edu/cis_papers/159 [244]
- Lauritzen, S. L. (1996), *Graphical Models* (Vol. 17), Oxford: Clarendon Press. [233]
- Lee, J. D., Sun, Y., and Saunders, M. A. (2012), “Proximal Newton-Type Methods for Minimizing Convex Objective Functions in Composite Form,” *arXiv preprint arXiv:1206.1623*. [242,243]
- Lee, J. D., Sun, Y., and Taylor, J. (2013), “On Model Selection Consistency of M-estimators With Geometrically Decomposable Penalties,” *arXiv preprint arXiv:1305.7477*. [241]
- Lee, S. I., Ganapathi, V., and Koller, D. (2006), “Efficient Structure Learning of Markov Networks Using L1 regularization,” in *Advances in Neural Information Processing Systems*, 817–827. [231]
- Liang, P., and Jordan, M. I. (2008), “An Asymptotic Analysis of Generative, Discriminative, and Pseudolikelihood Estimators,” in *Proceedings of the 25th International Conference on Machine Learning*, ACM, pp. 584–591. [249]
- Liu, Q., and Ihler, A. (2011), “Learning Scale Free Networks by Reweighted l1 Regularization,” in *Proceedings of the 14th International Conference on Artificial Intelligence and Statistics (AISTATS)*. [236]
- (2012), “Distributed Parameter Estimation via Pseudo-Likelihood,” in *Proceedings of the International Conference on Machine Learning (ICML)*. [236]
- Meinshausen, N., and Bühlmann, P. (2006), “High-Dimensional Graphs and Variable Selection With the Lasso,” *The Annals of Statistics*, 34, 1436–1462. [230,236]
- Peng, J., Wang, P., Zhou, N., and Zhu, J. (2009), “Partial Correlation Estimation by Joint Sparse Regression Models,” *Journal of the American Statistical Association*, 104, 735–746. [230]
- Ravikumar, P., Wainwright, M. J., and Lafferty, J. D. (2010), “High-Dimensional Ising Model Selection Using l1-Regularized Logistic Regression,” *The Annals of Statistics*, 38, 1287–1319. [231,236,241,245]
- Rothman, A. J., Levina, E., and Zhu, J. (2010), “Sparse Multivariate Regression With Covariance Estimation,” *Journal of Computational and Graphical Statistics*, 19, 947–962. [231]
- Schmidt, M. (2010), “Graphical Model Structure Learning With l1-Regularization,” PhD thesis, University of British Columbia. [231,242]
- Schmidt, M., Kim, D., and Sra, S. (2011), “Projected Newton-Type Methods in Machine Learning,” in *Optimization for Machine Learning*, Cambridge, MA: MIT Press. [242]
- Schmidt, M., Murphy, K., Fung, G., and Rosales, R. (2008), “Structure Learning in Random Fields for Heart Motion Abnormality Detection,” *CVPR, IEEE Computer Society*. [231]
- Tur, I., and Castelo, R. (2012), “Learning Mixed Graphical Models From Data With p Larger Than n ,” *arXiv preprint arXiv:1202.3765*. [234]
- Witten, D. M., and Tibshirani, R. (2009), “Covariance-Regularized Regression and Classification for High Dimensional Problems,” *Journal of the Royal Statistical Society, Series B*, 71, 615–636. [231]

- Wright, S. J., Nowak, R. D., and Figueiredo, M. A. T. (2009), “Sparse Reconstruction by Separable Approximation,” *IEEE Transactions on Signal Processing*, 57, 2479–2493. [243]
- Yang, E., Ravikumar, P., Allen, G. I., and Liu, Z. (2012), “Graphical Models via Generalized Linear Models,” in *Advances in Neural Information Processing Systems* 25, 1367–1375. [236]
- (2013), “On Graphical Models via Univariate Exponential Family Distributions,” *arXiv preprint arXiv:1301.4183*. [236]
- Yuan, M., and Lin, Y. (2006), “Model Selection and Estimation in Regression With Grouped Variables,” *Journal of the Royal Statistical Society, Series B*, 68, 49–67. [238]