
Learning Higher-Order Graph Structure with Features by Structure Penalty

Shilin Ding^{1*}, Grace Wahba^{1,2,3*}, and Xiaojin Zhu^{2*}

Department of {¹Statistics, ²Computer Sciences, ³Biostatistics and Medical Informatics}
University of Wisconsin-Madison, WI 53705
{sding, wahba}@stat.wisc.edu, jerryzhu@cs.wisc.edu

Abstract

In discrete undirected graphical models, the conditional independence of node labels Y is specified by the graph structure. We study the case where there is another input random vector X (e.g. observed features) such that the distribution $P(Y | X)$ is determined by functions of X that characterize the (higher-order) interactions among the Y 's. The main contribution of this paper is to learn the graph structure and the functions conditioned on X at the same time. We prove that discrete undirected graphical models with feature X are equivalent to multivariate discrete models. The reparameterization of the potential functions in graphical models by conditional log odds ratios of the latter offers advantages in representation of the conditional independence structure. The functional spaces can be flexibly determined by kernels. Additionally, we impose a Structure Lasso (SLasso) penalty on groups of functions to learn the graph structure. These groups with overlaps are designed to enforce hierarchical function selection. In this way, we are able to shrink higher order interactions to obtain a sparse graph structure.

1 Introduction

In undirected graphical models (UGMs), a graph is defined as $G = (V, E)$, where $V = \{1, \dots, K\}$ is the set of nodes and $E \subset V \times V$ is the set of edges between the nodes. The graph structure specifies the conditional independence among nodes. Much prior work has focused on graphical model structure learning without conditioning on X . For instance, Meinshausen and Bühlmann [1] and Peng *et al.* [2] studied sparse covariance estimation of Gaussian Markov Random Fields. The covariance matrix fully determines the dependence structure in the Gaussian distribution. But it is not the case for non-elliptical distributions, such as the discrete UGMs. Ravikumar *et al.* [3] and Höfling and Tibshirani [4] studied variable selection of Ising models based on l_1 penalty. Ising models are special cases of discrete UGMs with (usually) only pairwise interactions, and without features. We focused on discrete UGMs with both higher order interactions and features. It is important to note that the graph structure may change conditioned on different X 's, thus our approach may lead to better estimates and interpretation.

In addressing the problem of structure learning with features, Liu *et al.* [5] assumed Gaussian distributed Y given X , and they partitioned the space of X into bins. Schmidt *et al.* [6] proposed a framework to jointly learn pairwise CRFs and parameters with block- l_1 regularization. Bradley and Guestrin [7] learned tree CRF that recovers a max spanning tree of a complete graph based on heuristic pairwise link scores. These methods utilize only pairwise information to scale to large graphs. The closest work is Schmidt and Murphy [8], which examined the higher-order graphical structure

*SD wishes to acknowledge the valuable comments from Stephen J. Wright and Sijian Wang. Research of SD and GW is supported in part by NIH Grant EY09946, NSF Grant DMS-0906818 and ONR Grant N0014-09-1-0655. Research of XZ is supported in part by NSF IIS-0953219, IIS-0916038.

learning problem without considering features. They used an active set method to learn higher order interactions in a greedy manner. Their model is over-parameterized, and the hierarchical assumption is sufficient but not necessary for conditional independence in the graph.

To the best of our knowledge, no previous work addressed the issue of graph structure learning of all orders while conditioning on input features. Our contributions include a reparameterization of UGMs with bivariate outcomes into multivariate Bernoulli (MVB) models. The set of conditional log odds ratios in MVB models are complete to represent the effects of features on responses and their interactions at all levels. The sparsity in the set of functions are sufficient and necessary for the conditional independence in the graph, i.e., two nodes are conditionally independent iff the pairwise interaction is constant zero; and the higher order interaction among a subset of nodes means none of the variables is separable from the others in the joint distribution.

To obtain a sparse graph structure, we impose Structure Lasso (SLasso) penalty on groups of functions with overlaps. SLasso can be viewed as group lasso with overlaps. Group lasso [9] leads to selection of variables in groups. Jacob *et al.* [10] considered the penalty on groups with arbitrary overlaps. Zhao *et al.* [11] set up the general framework for hierarchical variable selection with overlapping groups, which we adopt here for the functions. Our groups are designed to shrink higher order interactions similar to hierarchical inclusion restriction in Schimdt and Murphy [8]. We give a proximal linearization algorithm that efficiently learns the complete model. Global convergence is guaranteed [12]. We then propose a greedy search algorithm to scale our method up to large graphs as the number of parameters grows exponentially.

2 Conditional Independence in Discrete Undirected Graphical Models

In this section, we first discuss the relationship between the multivariate Bernoulli (MVB) model and the UGM whose nodes are binary, i.e. $Y_i = 0$ or 1. At the end, we will give the representation of the general discrete UGM where Y_i takes value in $\{0, \dots, m-1\}$. In UGMs, the distribution of multivariate discrete random variables Y_1, \dots, Y_K given X is:

$$P(Y_1 = y_1, \dots, Y_K = y_K | X) = \frac{1}{Z(X)} \prod_{C \in \mathcal{C}} \Phi_C(y_C; X) \quad (1)$$

where $Z(X)$ is the normalization factor. The distribution is factorized according to the cliques in the graph. A clique $C \subseteq \Omega = \{1, \dots, K\}$ is the set of nodes that are fully connected. $\Phi_C(y_C; X)$ is the potential function on C , indexed by $y_C = (y_i)_{i \in C}$. This factorization follows from the Markov property: any two nodes not in a clique are conditionally independent given others [13]. So \mathcal{C} does not have to comply with the graph structure, as long as it is sufficient. For example, the most general choice for any given graph is $\mathcal{C} = \{\Omega\}$. See Theorem 2.1 and Example 2.1 for details.

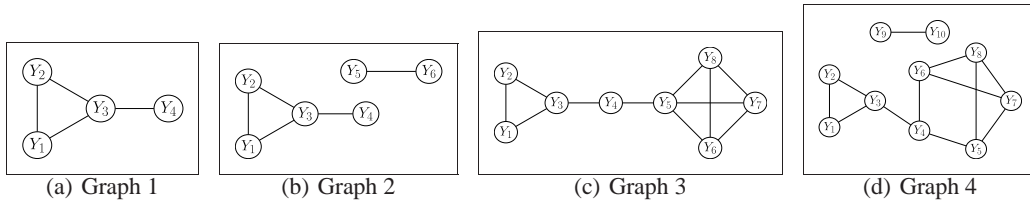


Figure 1: Graphical model examples.

Given the graph structure, the potential functions characterize the distribution on the graph. But if the graph is unknown in advance, estimating the potential functions on all possible cliques tends to be over-parameterized [8]. Furthermore, $\log \Phi_C(y_C; X) = 0$ is sufficient for the conditional independence among the nodes but not necessary (see Example 2.1). To avoid these problems, we introduce the MVB model that is equivalent to (1) with binary nodes, i.e. $Y_i = 0$ or 1. The MVB distribution is:

$$\begin{aligned} P(Y_1 = y_1, \dots, Y_K = y_k | X = x) &= \exp \left\{ \sum_{\omega \in \Psi_K} y^\omega f^\omega - b(f) \right\} \\ &= \exp \left\{ y_1 f^1(x) + \dots + y_K f^K(x) + \dots + y_1 y_2 f^{1,2}(x) + \dots + y_1 \dots y_K f^{1, \dots, K}(x) - b(f) \right\} \end{aligned} \quad (2)$$

Here, we use the following notations. Let $\overline{\Psi}_K$ be the power set of $\Omega = \{1, \dots, K\}$, and use $\Psi_K = \overline{\Psi}_K - \{\emptyset\}$ to index the $2^K - 1$ f^ω 's in (2). Let ω denotes a set in Ψ_K , define $\mathcal{Y} = (y^1, \dots, y^\omega, \dots, y^\Omega)$ be the augmented response with $y^\omega = \prod_{i \in \omega} y_i$. And $f = (f^1, \dots, f^\omega, \dots, f^\Omega)$ is the vector of conditional log odds ratios [14]. We assume f^ω is in a Reproducing Kernel Hilbert Space (RKHS) \mathcal{H}^ω with kernel K^ω [15]. For example, in our simulation we choose f^ω to be B-spline (see supplementary material). We focus on estimating the set of $f^\omega(x)$ with feature x where the sparsity in the set specifies the graph structure.

We present the following lemma and theorem which show the equivalence between UGM and MVB:

Lemma 2.1. *In a MVB model, define the odd-even partition of the power set of ω as: $\Psi_{odd}^\omega = \{\kappa \subseteq \omega \mid |\kappa| = |\omega| - k, \text{ where } k \text{ is odd}\}$, and $\Psi_{even}^\omega = \{\kappa \subseteq \omega \mid |\kappa| = |\omega| - k, \text{ where } k \text{ is even}\}$. Note $|\Psi_{odd}^\omega| = |\Psi_{even}^\omega| = 2^{|\omega|-1}$. The following property holds:*

$$f^\omega = \log \frac{\prod_{\kappa \in \Psi_{even}^\omega} P(Y_i = 1, i \in \kappa; Y_j = 0, j \in \Omega \setminus \kappa | X)}{\prod_{\kappa \in \Psi_{odd}^\omega} P(Y_i = 1, i \in \kappa; Y_j = 0, j \in \Omega \setminus \kappa | X)}, \quad b(f) = \log \frac{Z(x)}{\prod_{C \in \mathcal{C}} \Phi_C(0; x)} \quad (3)$$

Theorem 2.1. *A UGM of the general form (1) with binary nodes is equivalent to a MVB model of (2). In addition, the following are equivalent: 1) There is no $|C|$ -order interaction in $\{Y_i, i \in C\}$; 2) There is no clique $C \in \Psi_K$ in the graph; 3) $f^\omega = 0$ for all ω such that $C \subseteq \omega$.*

A proof is given in Appendix. It states that there is a clique C in the graph, iff there is $\omega \supseteq C$, $f^\omega \neq 0$ in MVB model. The advantage of modeling by MVB is that the sparsity in f^ω 's is sufficient and necessary for the conditional independence in the graph, thus fully specifying the graph structure. Specifically, Y_i, Y_j are conditionally independent iff $f^\omega = 0, \omega \supseteq \{i, j\}$. This showed the interaction is non-zero iff all the nodes involved are not conditionally independent.

Example 2.1. *When $K = 2, \Omega = \{1, 2\}, \mathcal{C} = \{\emptyset\}$, denote $\Phi_\Omega(Y_1 = 1, Y_2 = 1; X)$ as Φ_{11} for simplicity, then $P(Y_1 = 1, Y_2 = 1 | X) = \frac{1}{2} \Phi_{11}$. Define $\Phi_{10}, \Phi_{01}, \Phi_{00}$ similarly, then the distribution with UGM parameterization is determined. The relation between UGM and MVB is*

$$f^1 = \log \frac{\Phi_{10}}{\Phi_{00}}, \quad f^2 = \log \frac{\Phi_{01}}{\Phi_{00}}, \quad f^{1,2} = \log \frac{\Phi_{11} \cdot \Phi_{00}}{\Phi_{01} \cdot \Phi_{10}}$$

Note, the independence between Y_1 and Y_2 implies: $f^{1,2} = 0$ or $\Phi_{11} \cdot \Phi_{00} = \Phi_{01} \cdot \Phi_{10}$. Therefore, $f^{1,2}$ being zero in MVB model is sufficient and necessary for the conditional independence in the model. On the other hand, $\log \Phi_C = 0$ is a sufficient condition but not necessary.

The distribution of a general discrete UGM where $Y_k \in \{0, \dots, m-1\}$ can be extended from (2).

Lemma 2.2. *Let $V = \{1, \dots, m-1\}, y_\omega = (y_i)_{i \in \omega}$, then*

$$P(Y_1 = y_1, \dots, Y_K = y_K | X) = \exp \left\{ \sum_{\omega=1}^{\Omega} \sum_{v \in V^{|\omega|}} I(y_\omega = v) f_v^\omega - b(f) \right\} \quad (4)$$

where I is an indicator function and V^n is the tensor product of n V 's. Each f^ω is a $|V|^{|\omega|}$ vector.

3 Structure Penalty

In many applications, the assumption is that the graph has very few large cliques. Similar to the hierarchical inclusion restriction in Schmidt and Murphy [8], we will include a higher order interaction only when all its subsets are included. Our model is very flexible in that $f^\omega(x)$ can be in an arbitrary RKHS.

Let $y(i) = (y_1(i), \dots, y_K(i)), x(i) = (x_1(i), \dots, x_p(i))$ be the i th data point. There are $|\Psi_K| = 2^K - 1$ functions in total. We first consider learning the full model when K is small, and later propose a greedy search algorithm to scale to large graphs. The penalized log likelihood model is:

$$\min I_\lambda(f) = L(f) + \lambda J(f) = \sum_{i=1}^n \left(-\mathcal{Y}(i)^T f(x(i)) + b(f) \right) + \lambda J(f) \quad (5)$$

where $L(f)$ is the negative log likelihood and $J(\cdot)$ is the structure penalty. The hierarchical assumption is that if there is no interaction on clique C , then all f^ω should be zero, for $\omega \supseteq C$. The penalty is designed to shrink such f^ω toward zero. We consider the Structure Lasso (SLasso) penalty guided by the lattice in Figure 2. The lattice T has $2^K - 1$ nodes: $1, \dots, \omega, \dots, \Omega$. There is an edge from ω_1 to ω_2 if and only if $\omega_1 \subset \omega_2$ and $|\omega_1| + 1 = |\omega_2|$. Jenatton *et al.* [16] discussed how to define the groups to achieve different nonzero patterns.

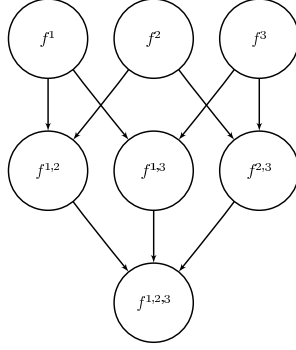


Figure 2: Hierarchical lattice for penalty

Let $T_v = \{\omega \in \Psi_K | v \subseteq \omega\}$ be the subgraph rooted at v in T , including all the descendants of v . Denote $f^{T_v} = (f^\omega)_{\omega \in T_v}$. All the functions are categorized into groups with overlaps as (T_1, \dots, T_Ω) . The SLasso penalty on the group T_v is: $J(f^{T_v}) = p_v \sqrt{\sum_{\omega \in T_v} \|f^\omega\|_{\mathcal{H}^\omega}^2}$ where p_v is the weight for the penalty on T_v , empirically chosen as $\frac{1}{|T_v|}$. Then, the objective is:

$$\min_f I_\lambda(f) = L(f) + \lambda \sum_v p_v \sqrt{\sum_{\omega \in T_v} \|f^\omega\|_{\mathcal{H}^\omega}^2} \quad (6)$$

The following theorem shows that by minimizing the objective (6), f^{ω_1} will enter the model before f^{ω_2} if $\omega_1 \subset \omega_2$. That is to say, if f^{ω_1} is zero, there will be no higher order interactions on ω_2 . It is an extension of Theorem 1 in Zhao *et al.* [11] and the proof is given in Appendix.

Theorem 3.1. *Objective (6) is convex, thus the minimal is attainable. Let $\omega_1, \omega_2 \in \Psi_K$ and $\omega_1 \subset \omega_2$. If \hat{f} is the minimizer of (6) given the observations, that is, $0 \in \partial I_\lambda(\hat{f})$ which is the subgradient of I_λ at \hat{f} , then $\hat{f}^{\omega_2} = 0$ almost surely if $\hat{f}^{\omega_1} = 0$.*

Example 3.1. *If $K = 3$, $f = (f^1, f^2, f^3, f^{1,2}, f^{1,3}, f^{2,3}, f^{1,2,3})$. The group at node 1 in Figure 2 is $f^{T_1} = (f^1, f^{1,2}, f^{1,3}, f^{1,2,3})$ and $J(f^{T_1}) = p_1 \sqrt{\|f^1\|^2 + \|f^{1,2}\|^2 + \|f^{1,3}\|^2 + \|f^{1,2,3}\|^2}$.*

4 Parameter Estimation

In this section, we discuss parameter estimation where the ω th function space is linear as $\mathcal{H}^\omega = \{1\} \oplus \mathcal{H}_1^\omega$ for simplicity. $\{1\}$ refers to the constant function space, and \mathcal{H}_1^ω is a RKHS with a linear kernel. The functions in \mathcal{H}^ω have the form $f^\omega(x) = c_0^\omega + \sum_{j=1}^p c_j^\omega x_j$. Its norm is $\|f^\omega\|_{\mathcal{H}^\omega} = \|c^\omega\|$, where $\|\cdot\|$ stands for Euclidean l_2 norm. Here, we denote $c^\omega = (c_0^\omega, \dots, c_p^\omega)^T \in \mathbb{R}^{p+1}$ as a vector of length $p+1$ and $c = (c^\omega)_{\omega \in \Psi_K} \in \mathbb{R}^{\bar{p}}$ is the concatenated vector of all parameters of length $\bar{p} = (p+1) \cdot |\Psi_K|$. Let $c^{T_v} = (c^\omega)_{\omega \in T_v}$ be a $(p+1) \cdot |T_v|$ vector, then the objective (6) is now:

$$\min_c I_\lambda(c) = L(c) + \lambda \sum_v p_v \|c^{T_v}\| \quad (7)$$

4.1 Estimating the complete model on small graphs

Many applications do not involve a large amount of responses, so it is desirable to learn the complete model when the graph is small for consistency reasons. We propose a method to optimize (7) of the

Algorithm 1 Proximal Linearization Algorithm

Input: $c_0, \alpha_0, \zeta > 1, tol > 0$
repeat
 Choose $\alpha_k \in [\alpha_{min}, \alpha_{max}]$
 Solve Eq (8) for $d_k = c - c_k$
 while $\delta_k = I_\lambda(c_k) - I_\lambda(c_k + d_k) < \|d_k\|^3$ **do**
 // Insufficient decrease
 Set $\alpha_k = \max(\alpha_{min}, \zeta \alpha_k)$
 Solve Eq (8) for d_k
 end while
 Set $\alpha_{k+1} = \alpha_k / \zeta$
 Set $c_{k+1} = c_k + d_k$
until $\delta_k < tol$

complete model with all interaction levels by iteratively solving the following proximal linearization problem as discussed in Wright [12]:

$$\min_c L_k + \nabla L_k^T(c - c_k) + \frac{\alpha_k}{2} \|c - c_k\|^2 + \lambda J(c) \quad (8)$$

where $L_k = L(c_k)$, and α_k is a positive scalar chosen adaptively at k th step. With slight abuse of notation, we denote c_k as the value of c at k th step. Algorithm 1 summarized the framework of solving (7). Following the analysis in Wright [12], we can ensure that the proximal linearization algorithm will converge for the negative log-likelihood loss function with the SLasso penalty.

However, solving group lasso with overlaps is not trivial due to the non-smoothness at the singular point. In recent years, several papers have addressed this problem. Jacob *et al.* [10] duplicated the design matrix columns that appear in group overlaps, then solved the problem as group lasso without overlaps. Kim and Xing [17] reparameterized the group norm with additional dummy variables. They alternatively optimized the model parameters and the dummy ones at each step. It is efficient for the quadratic loss function on Gaussian data, but might not scale well in our case. Instead, we solve (8) by its smooth and convex dual problem [18]. The details are in the supplementary material.

4.2 Estimating large graphs

The above algorithm is efficient on small graphs ($K < 20$). It usually terminates within 20 iterations in our experiments. However, the issue of estimating a complete model is the exponential number of f^ω 's and the same amount of groups involved in objective (7). It is intractable when the graph becomes large. The hierarchical assumption and the SLasso penalty lend themselves naturally to a greedy search algorithm:

1. Start from the set of main effects as $A_0 = \{f^1, \dots, f^K\}$.
2. In step i , remove the nodes that are not in A_i from the lattice in Figure 2. Obtain a sparse estimation of the functions in A_i by algorithm (1). Denote the resulting sparse set A'_i .
3. Let $A_{i+1} = A'_i$. Keep adding a higher order interaction into A_{i+1} if all its subsets of interactions are included in A'_i . And also add this node into the lattice in Figure 2.

Iterate step 2 and 3 until convergence. The algorithm is similar to the active set method in Schmidt and Murphy [8]. It has multiple runs of algorithm (1) to enforce the hierarchical assumption. It is not guaranteed to converge to the global optimum. Nonetheless, our empirical experiments show its ability to scale to large graphs.

5 Experiments

5.1 Toy Data

In the simulation, we create 6 toy graphs. The first four graphs are depicted in Figure 1. Graph 5 has 100 nodes where the first 8 nodes have the same structure as in Figure 1(c) and the others are independent. Graph 6 also has 100 nodes where the first 10 nodes have the same connection as in Figure 1(d) and the others are independent. We generate 100 datasets for each structure to evaluate

the performance. The sample size of each dataset is 1000. Here is how the first data set is generated: The length of the feature vector, p , is set to 5 in our experiment, i.e. $X = (X_1, \dots, X_5)$. Each $f^\omega(x) = c_0^\omega + \sum_{j=1}^5 g_j^\omega(x_j)$ where $g_j^\omega(x_j) = \sum_{k=1}^D c_{jk}^\omega B_k(x_j)$ is spanned by the B-spline basis functions $\{B_k(\cdot)\}_{k=1, \dots, D}$ (see the supplementary material), where D is chosen to be 5. The true set of the model parameters, c_{jk}^ω , is uniformly sampled from $\{-5, -4, \dots, 5\}$. We set the intercepts c_0^ω in main effects to 1, and those in second or higher order interactions to 2. The features, X_j , are i.i.d uniform on $[-1, 1]$. Then, Y is sampled according to the probability in equation (2).

We use GACV (generalized approximate cross validation) and BGACV (B-type GACV) [19] to choose the regularization parameter λ for the complete model (graphs 1-4). We call these variants of SLasso Complete-GACV and Complete-BGACV. We use AIC for greedy search (Greedy-AIC) in graphs 5 and 6 due to computational consideration. The range of λ is chosen according to Koh *et al.* [20]. The details of the tuning methods are discussed in the supplementary material. The R package, BMN, is used as a baseline [4].

Table 1: Number of true positive and false positive functions

Graph	Method	$f^{1,2}$	$f^{1,3}$	$f^{2,3}$	$f^{3,4}$	$f^{1,2,3}$	$f^{5,7,8}$	$f^{5,6,7,8}$	FP
1	BMN	60	76	70	60	0	-	-	162
	Complete-GACV	100	100	100	94	84	-	-	136
	Complete-BGACV	86	83	83	72	14	-	-	11
2	BMN	44	50	38	58	0	-	-	412
	Complete-GACV	100	99	100	99	83	-	-	341
	Complete-BGACV	88	91	88	78	33	-	-	64
3	BMN	72	64	60	60	0	0	0	830
	Complete-GACV	91	87	81	92	62	71	33	412
	Complete-BGACV	36	22	23	93	0	39	0	162
4	BMN	48	34	37	29	0	0	-	774
	Complete-GACV	92	98	94	90	54	45	-	693
	Complete-BGACV	68	68	71	62	0	0	-	144
5	BMN	38	28	26	22	0	0	0	9476
	Greedy-AIC	99	99	98	97	22	21	0	1997
6	BMN	28	26	14	26	0	0	-	9672
	Greedy-AIC	100	100	100	99	24	15	-	3458

In Table 1, we count, for each function f^ω , the number of runs out of 100 where f^ω is recovered ($\|c^\omega\| \neq 0$). If a recovered function is in the true model, it is considered a true positive, otherwise a false positive. The main effects are always detected correctly, thus are not listed in the table. SLasso is more effective compared to BMN which only considers pairwise interactions.

In Figure 3, we show the learning results in terms of true positive rate (TPR) as sample size increases from 100 to 1000. The experimental setting is the same as before. The TPRs improve with increasing sample size. GACV achieves better TPR, but higher FPR compared to BGACV. Our method outperforms BMN in all six graphs.

5.2 Case Study: Census Bureau County Data

We use the county data from U.S. Census Bureau¹ to validate our method. We remove the counties that have missing values and obtain 2668 entries in total. The outcomes of this study are summarized in Table 2. ‘‘Vote’’ [21] is coded as 1 if the Republican candidate won in the 2004 presidential election. To dichotomize the remaining outcomes, the national mean is selected as a threshold. The data is standardized to mean 0 and variance 1. The following features are included: Housing unit change in percent from 2000-2006, percent of ethnic groups, percent foreign born, percent people over 65, percent people under 18, percent people with a high school education, percent people with a bachelors degree; birth rate, death rate, per capita government expenditure in dollars. By adjusting λ , we observe new interactions enter the model. The graph structure of $\lambda = 0.1559$ is

¹<http://www.census.gov/statab/www/ccdb.html>

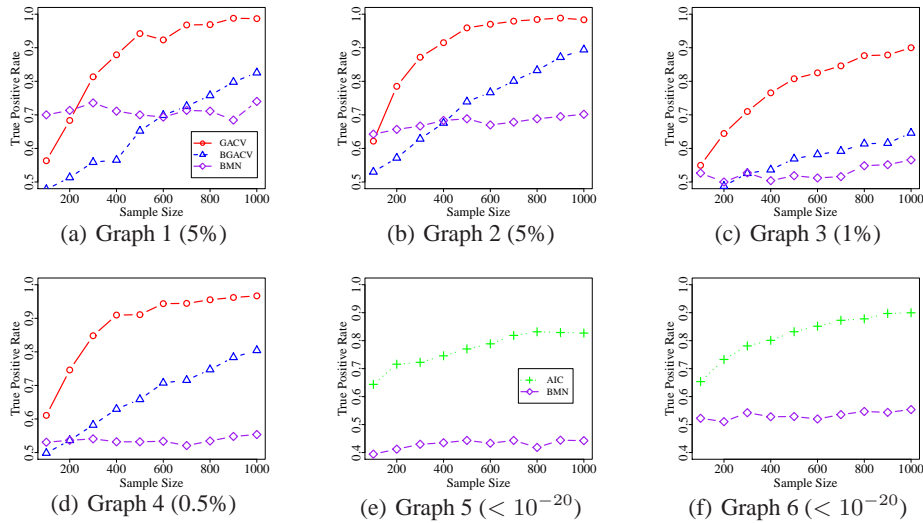


Figure 3: The True Positive Rate (TPR) of graph structure learning methods with increasing sample size. The percentage in the bracket is the upper bound of False Positive Rate (FPR) in each experiment. BMN always has larger FPR compared to SLasso.

Table 2: Selected response variables

Response	Description	Positive%
Vote	2004 votes for Republican presidential candidate	81.11
Poverty	Poverty Rate	52.70
VCrime	Violent Crime Rate, eg. murder, robbery	23.09
PCrime	Property Crime Rate, eg. burglary	6.82
URate	Unemployment Rate	51.35
PChange	Population change in percent from 2000 to 2006	64.96

shown in Figure 4(a). The results of BMN (the tuning parameter is 0.015) is in Figure 4(b). The unemployment rate plays an important role as a hub as discovered by SLasso, but not by BMN.

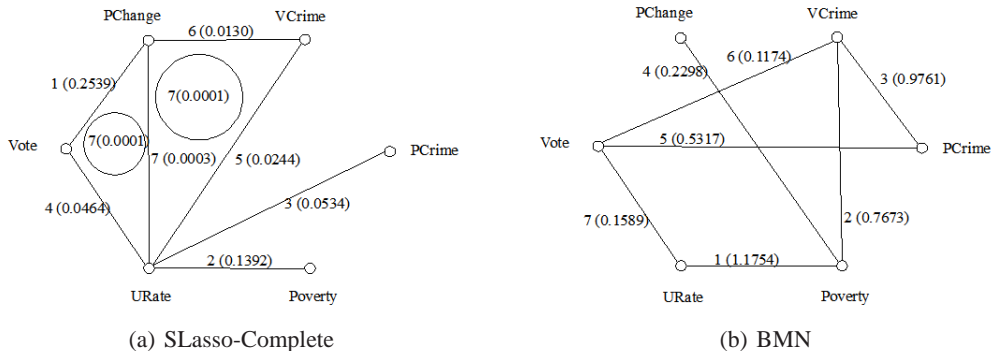


Figure 4: Interactions of response variables in the Census Bureau data. The first number on the edge is the order at which the link is recovered. The number in bracket is the function norm on the clique and the absolute value of the elements in the concentration matrix, respectively. We note SLasso discovers at 7th step two third-order interactions which are displayed by two circles in (a).

We analyze the link between “Vote” and “PChange”. Though the marginal correlation between them (without X) is only 0.0389, which is the second lowest absolute pairwise correlation, the

link is firstly recovered by S_Lasso. It has been suggested that there is indeed a connection². This shows that after taking features into account, the dependence structure of response variables may change and hidden relations could be discovered. The main factors in this case are “percentage of housing unit change” (X_1) and “population percentage of people over 65” (X_2). The part of the fitted model shown below suggests that as housing units increase, the counties are more likely to have both positive results for “Vote” and “PChange”. But this tendency will be counteracted by the increase of people over 65: the responses are less likely to take both positive values.

$$\begin{aligned}\hat{f}^{Vote} &= 0.2913 \cdot X_1 + 0.3475 \cdot X_2 + \dots \\ \hat{f}^{PChange} &= 1.4726 \cdot X_1 - 0.3709 \cdot X_2 + \dots \\ \hat{f}^{Vote, PChange} &= 0.1358 \cdot X_1 - 0.0458 \cdot X_2 + \dots\end{aligned}$$

6 Conclusions

Our S_Lasso method can learn the graph structure that is specified by the conditional log odds ratios conditioned on input features X , which allows the graphical model depending on features. The modeling interprets well, since $f^\omega = 0$ iff there is no such clique. An efficient algorithm is given to estimate the complete model. A greedy approach is applied when the graph is large. S_Lasso can be extended to model a general discrete UGM, where Y_k takes value in $\{0, \dots, m-1\}$. Also, there exist rich selections of the function forms, which makes the model more flexible and powerful, though modification is needed in solving the proximal subproblem for non-parametric families.

A Proof

A.1 Proof of Theorem 2.1

Proof. Given UGM (1), the corresponding parameterization in MVB model is shown in (3) of Lemma 2.1. Conversely, given the MVB model of (2), the cliques can be determined by the nonzero f^ω : clique C exists if $C = \omega$ and $f^\omega \neq 0$. Then the maximal cliques can be inferred from the graph structure. And suppose they are C_1, \dots, C_m . Let $\omega_i = C_i$, for $i = 1, \dots, m$, and $\kappa_1 = \emptyset$, $\kappa_i = C_i \cap (C_{i-1} \cup \dots \cup C_1)$, $i = 2, \dots, m$. Then the parameterization is:

$$\Phi_{C_i}(y_{C_i}; x) = \exp(S^{\omega_i}(y; x) - S^{\kappa_i}(y; x)) \quad \text{and} \quad Z(x) = \exp(b(f)) \quad (9)$$

where $S^\omega(y; x) = \sum_{\kappa \subseteq \omega} y^\kappa f^\kappa(x)$. Thus, UGM (1) with bivariate nodes is equivalent to MVB (2).

In the latter part of the theorem, $1 \Rightarrow 2$ and $3 \Rightarrow 1$ follow naturally from the Markov property of graphical models. To show $2 \Rightarrow 3$, let y_C^ω be a realization of y_C such that $y_C^\omega = (y_i^\omega)_{i \in C}$ where $y_i^\omega = 1$ if $i \in \omega$ and $y_i^\omega = 0$ otherwise. Notice that whenever $\kappa \cap C = \kappa' \cap C$, we have $y_C^\kappa = y_C^{\kappa'}$. For any possible $v = \kappa \cap C$, $\kappa' \in \{\kappa | \kappa = v \cup u, \text{ s.t. } u \subseteq \omega - v\}$ will satisfy the condition: $\kappa' \cap C = v$. There are $2^{|\omega - v|}$ such κ' in total due to the choice of u . Also, they appear in the nominator and denominator of equation (3) equally. So, for any $C \in \mathcal{C}$,

$$\prod_{\kappa \in \Psi_{\text{even}}^\omega} \Phi_C(y_C^\kappa; x) = \prod_{\kappa \in \Psi_{\text{odd}}^\omega} \Phi_C(y_C^\kappa; x) \quad (10)$$

It follows that $f^\omega = 0$ by (3). \square

A.2 Proof of Theorem 3.1

Proof. We give the proof for the linear case. The convexity of I_λ is easy to check, since L and $J(f^{T_v})$ are all convex in c . Suppose there is some $\omega_2 \supset \omega_1$ s.t. $\hat{c}^{\omega_2} \neq 0$ and $\hat{c}^{\omega_1} = 0$, by the groups constructed through Figure 2, $\|\hat{c}^{T_v}\| = \|(\hat{c}^\omega)_{v \subseteq \omega}\| \neq 0$ for all $v \subseteq \omega_1$. So the partial derivative of the objective (7) with respect to c^{ω_1} at \hat{c}^{ω_1} is

$$\left. \frac{\partial L}{\partial c^{\omega_1}} \right|_{c^{\omega_1} = \hat{c}^{\omega_1}} + \lambda \sum_{v \subseteq \omega_1} p_v \frac{\hat{c}^{\omega_1}}{\|\hat{c}^{T_v}\|} = 0 \quad (11)$$

Thus, the probability of $\{\hat{c}^{\omega_2} \neq 0\}$ equals to the probability of $\{\left. \frac{\partial L}{\partial c^{\omega_1}} \right|_{c^{\omega_1} = \hat{c}^{\omega_1}} = 0\}$, which is 0. \square

²<http://www.ipsos-mori.com/researchpublications/researcharchive/2545/Analysis-Population-change-turnout-the-election.aspx>

B Multivariate Bernoulli model

The multivariate Bernoulli (MVB) model of K random variables has $2^K - 1$ natural parameters [22]. Given the predictive variable X , these parameters are functions of X , called conditional log odds ratios. From the distribution of the MVB, $f^\omega(X)$ can be written as:

$$f^\omega(x) = \log OR(Y_i, i \in \omega | Y_j = 0, j \notin \omega; X = x) \quad (12)$$

Here, the odds ratios are calculated recursively

$$OR(Y_i | X = x) = \frac{P(Y_i = 1 | X = x)}{1 - P(Y_i = 1 | X = x)}, \quad (13)$$

$$OR(Y_i, i \in \omega \cup \{k\} | X = x) = \frac{OR(Y_i, i \in \omega | Y_k = 1, X = x)}{OR(Y_i, i \in \omega | Y_k = 0, X = x)}, \text{ with } k \notin \omega \quad (14)$$

The following two notations are useful in optimization and parameter tuning:

$$S^\omega(y; x) = \sum_{\kappa \subset \omega} y^\kappa f^\kappa(x); \quad S^\omega(x) = \sum_{\kappa \subset \omega} f^\kappa(x); \quad (15)$$

It follows from the definition of the conditional log odds ratio in (12) that

$$\exp(S^\omega(x)) = \frac{P(Y_i = 1, i \in \omega, \text{ and } Y_j = 0, j \in \Omega - \omega | X = x)}{P(Y_i = 0, i \in \Omega | X = x)} \quad (16)$$

Then the normalization factor is:

$$\exp(b(f(x))) = 1 + \sum_{\omega \in \Psi_K} \exp(S^\omega(x)) \quad (17)$$

In practice, the $\exp(b(f(x)))$ is calculated by the junction tree algorithm to avoid enumerating 2^K possible values of Y , which is intractable in large graphs.

C Dual of the proximal linearization problem

To solve the following objective of the proximal linearization problem

$$\min_c L_k + \nabla L_k^T (c - c_k) + \frac{\alpha_k}{2} \|c - c_k\|^2 + \lambda J(c) \quad (18)$$

we solve its dual problem as suggested in Liu and Ye [18]. Let $Z = \{v \in \Psi_K | \|c^{T_v}\| = 0\}$, and $\bar{Z} = \Psi_K - Z$ be the complement. Define $s_v, v \in \Psi_K$ as:

$$s_v \in \mathbb{S}_v = \{s = (s^\omega)_{\omega \in \Psi_K} \mid s \in \mathbb{R}^{\bar{p}}, \|s\| \leq \lambda p_v, s^\omega = 0 \text{ if } \omega \in T_v\} \quad (19)$$

then the subgradient of (8) is:

$$\nabla L + \alpha_k (c - c_k) + \sum_{v \in Z} s_v + \sum_{u \in \bar{Z}} r_u \quad (20)$$

where s_v is the subgradient of $\lambda p_v \|c^{T_v}\|$ for $v \in Z$ and r_u is the subgradient for $u \in \bar{Z}$:

$$r_u = \arg \max_{s_u} \langle s_u, c \rangle, \text{ for } u \in \bar{Z} \quad (21)$$

The subgradient s_v is in a unit ball of certain subspace of $\mathbb{R}^{\bar{p}}$. These subspaces are not perpendicular to each other. Thus, s_v 's are not separable, and closed form solution of (8) cannot be obtained. We solve the proximal subproblem (8) by its smoothing and convex dual problem. Note (8) is equivalent to:

$$\min_{c \in \mathbb{R}^{\bar{p}}} \max_{S \in \mathbb{S}} \phi(c, S) = \nabla L_k^T (c - c_k) + \frac{\alpha_k}{2} \|c - c_k\|^2 + \sum_{v \in \Omega} \langle s_v, c \rangle \quad (22)$$

where S is a $\bar{p} \times |\Psi_K|$ matrix whose columns are s_v . $\mathbb{S} = \{S \mid S = (s_1, \dots, s_v, \dots, s_\Omega), s_v \in \mathbb{S}_v \text{ for } v \in \Psi_K\}$ is the feasible region of S . Since $\phi(\cdot, S)$ is lower semicontinuous and $\phi(c, \cdot)$ is

upper semicontinuous, there exists a saddle point and the max and min are exchangeable. The solution of minimizing $\phi(c, S)$ is:

$$\tilde{c} = \arg \min_c \phi(c, S) = c_k - \frac{1}{\alpha_k} \nabla L_k - \frac{1}{\alpha_k} \sum_v s_v \quad (23)$$

Substitute \tilde{c} back into (22), we have the dual problem of (8) as:

$$\max_{S \in \mathbb{S}} \eta(S) = -\frac{1}{2} \left\| \sum_v s_v \right\|^2 + (\alpha_k c_k - \nabla L_k)^T \sum_v s_v \quad (24)$$

Following the proof in Liu and Ye [18], we can show that $\eta(S)$ is convex and Lipschitz continuous. The differential is $\alpha_k \tilde{c} e^T$ where $e \in \mathbb{R}^{\tilde{p}}$ is a vector of ones. Hence, (24) can be solved by existing gradient methods.

D B-spline

Given m knots, $t_0 \leq t_1 \leq \dots \leq t_{m-1}$, the B-spline basis functions of degree d are defined recursively [23]:

$$b_{k,0} = \begin{cases} 1; & \text{if } t_k \leq t < t_{k+1} \\ 0; & \text{otherwise} \end{cases}, \text{ for } k = 0, \dots, m-2$$

$$b_{k,l} = \frac{t - t_k}{t_{k+l} - t_k} b_{k,l-1}(t) + \frac{t_{k+l+1} - t}{t_{k+l+1} - t_{k+1}} b_{k+1,l-1}(t), \text{ for } k = 0, \dots, m-d-2; l = 0, \dots, d$$

Let $B_k(\cdot) = b_{k,d}(\cdot)$, then $\{B_k, k = 0, \dots, m-d-2\}$ are $m-d-1$ basis functions, which span the functional space \mathcal{B} . The B-spline curve in \mathcal{B} is:

$$g(t) = \sum_{k=0}^{m-d-2} c_k B_k(t) \quad (25)$$

where c_k 's are the control points to be estimated. In our simulation studies, c_k 's are assumed to be one dimensional scalars for simplicity.

We let each $f^\omega(x)$ where $x = (x_1, \dots, x_p)'$ be in $\mathcal{B}_0 \oplus \mathcal{B}_1 \oplus \dots \oplus \mathcal{B}_p$. Here, \mathcal{B}_0 is a space of constant functions and $\mathcal{B}_j, j = 1 \dots, p$ is a B-spline functional space on domain $x_j \in \mathcal{X}_j$. Therefore,

$$f^\omega(x) = c_0^\omega + \sum_{j=1}^p g_j(x_j) \quad (26)$$

where $g_j(x_j) \in \mathcal{B}_j$ are defined in (25).

E Tuning

For i -th data point $(y(i), x(i))$, denote $S_i^\omega = S^\omega(x(i))$, then the normalization factor of the i -th data is $b_i = b(f(x(i))) = \log(1 + \sum_\omega \exp S_i^\omega)$. The mean of the augmented response $\mathcal{Y}(i)$ in the MVB model is:

$$\mu(i) = E[\mathcal{Y}(i)|x(i), f] = (\mu^1(i), \dots, \mu^\kappa(i), \dots, \mu^\Omega(i)) \quad (27)$$

$$\text{where } \mu^\kappa(i) = \frac{\partial b_i}{\partial f^\kappa} = \frac{\sum_{\omega \in T_\kappa} \exp S_i^\omega}{\exp b_i} \quad (28)$$

The $|\Psi_K| \times |\Psi_K|$ covariance matrix of the augmented response is:

$$W(i) = \text{var}(\mathcal{Y}(i)|x(i), f) \quad (29)$$

where the (α, β) -th element of $W(i)$ is:

$$W_{\alpha,\beta}(i) = \frac{\partial^2 b_i}{\partial f^\alpha \partial f^\beta} = \frac{\sum_{\omega \in T_\alpha \cap T_\beta} \exp S_i^\omega}{\exp b_i} - \mu^\alpha(i) \cdot \mu^\beta(i) \quad (30)$$

Let R_v be a $\tilde{p} \times \tilde{p}$ diagonal matrix whose (i, i) -th element is 1 if $c_i \neq 0$. Then, the v -th group penalty $J(f^{T_v})$ can be written as:

$$J(f^{T_v}) = p_v \sqrt{\sum_{\omega \in T_v} \|f^\omega\|^2} = p_v \|R_v c\| \quad (31)$$

Note R_v is symmetric and $R_v \cdot R_v = R_v$, direct calculation yields the derivative and Hessian of the penalty term:

$$\frac{\partial J}{\partial c} = \sum_{v: R_v c \neq 0} p_v \frac{R_v c}{\|R_v c\|} \quad (32)$$

$$\frac{\partial^2 J}{\partial c \partial c^T} = \sum_{v: R_v c \neq 0} p_v J_v = \sum_{v: R_v c \neq 0} p_v \frac{R_v (\|R_v c\|^2 I - c \cdot c^T) R_v}{\|R_v c\|^3} \quad (33)$$

where $J_v \doteq (R_v (\|R_v c\|^2 I - c \cdot c^T) R_v) / \|R_v c\|^3$. Denote the grand design matrix as:

$$D = (D(1)^T \quad \dots \quad D(n)^T)^T \quad (34)$$

$$\text{where } D(i) = \begin{pmatrix} x(i)^T & 0 & \dots & 0 \\ 0 & x(i)^T & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & x(i)^T \end{pmatrix} \quad (35)$$

Suppose there are N non-zero elements of c at location $\{a_1, \dots, a_N\}$. Let \tilde{D} be the matrix composed by the a_1, \dots, a_N th column of D . Then, the Hessian matrix of I_λ is:

$$\frac{\partial^2 I_\lambda}{\partial c \partial c^T} = \frac{\partial^2 L}{\partial c \partial c^T} + \lambda \frac{\partial^2 J}{\partial c \partial c^T} = \tilde{D}^T W \tilde{D} + \lambda \sum_{v: R_v c \neq 0} p_v J_v \quad (36)$$

Let H be the $n|\Psi_K| \times n|\Psi_K|$ influence matrix that implies

$$f_{\lambda, \epsilon} - f_\lambda \approx H \epsilon \quad (37)$$

where ϵ is a small perturbation on \mathcal{Y} ; $f_\lambda = D c_\lambda$ is the estimated function value with tuning parameter λ ; and $f_{\lambda, \epsilon}$ is the estimated function value with the perturbation. Then, the analysis of the first order Taylor expansion of $\frac{\partial I_\lambda}{\partial c}(\mathcal{Y} + \epsilon, c_{\lambda, \epsilon})$ leads to the formulation of H as follows (refer to Xiang and Wahba [24] and Ma [19] Chapter 3 for more details)

$$H = \tilde{D} \left(\frac{\partial^2 I_\lambda}{\partial c \partial c^T} \right)^{-1} \tilde{D}^T = \tilde{D} \left(\tilde{D}^T W \tilde{D} + \lambda \sum_{v: R_v c \neq 0} p_v J_v \right)^{-1} \tilde{D}^T \quad (38)$$

The (i, j) -th $q \times q$ submatrix of H is

$$H(i, j) = \tilde{D}(i)^T \left(\frac{\partial^2 I_\lambda}{\partial c \partial c^T} \right)^{-1} \tilde{D}(j) \quad (39)$$

Let $Q(i) = I - H(i, i)W(i)$ for $i = 1, \dots, n$, define the generalized average matrix, denoted as \bar{Q} , of $\{Q(i), i = 1, \dots, n\}$ as follows

$$\bar{Q} = (\delta - \gamma) I_{q \times q} + \gamma \cdot e e^T = \begin{pmatrix} \delta & \gamma & \dots & \gamma \\ \gamma & \delta & \dots & \gamma \\ \vdots & \vdots & \ddots & \vdots \\ \gamma & \gamma & \dots & \delta \end{pmatrix} \quad (40)$$

where e is the unit vector of length q and

$$\delta = \frac{1}{nq \sum_{i=1}^n \text{tr}(Q(i))}, \quad \gamma = \frac{1}{nq(q-1)} [e^T Q(i) e - \text{tr}(Q(i))] \quad (41)$$

Let \bar{H} be the generalized average of $\{H(i, i), i = 1, \dots, n\}$, the GACV score is

$$GACV(\lambda) = OBS(\lambda) + \frac{1}{n} \sum_{i=1}^n \mathcal{Y}(i)^T \bar{Q}^{-1} \bar{H} (\mathcal{Y}(i) - \mu(i)) \quad (42)$$

where

$$OBS(\lambda) = \frac{1}{n} \left[-\mathcal{Y}(i)^T f_\lambda(x(i)) + b(f_\lambda(x(i))) \right] \quad (43)$$

is the observed log-likelihood.

The degrees of freedom of multivariate Bernoulli data is generally difficult to obtain. But we can have a good approximation from GACV [25] as

$$\hat{df}(\lambda) = \sum_{i=1}^n \mathcal{Y}(i)^T \bar{Q}^{-1} \bar{H} (\mathcal{Y}(i) - \mu(i)) \quad (44)$$

So the BGACV score can be defined as

$$BGACV(\lambda) = OBS(\lambda) + \frac{1}{n} \frac{\log n}{2} \sum_{i=1}^n \mathcal{Y}(i)^T \bar{Q}^{-1} \bar{H} (\mathcal{Y}(i) - \mu(i)) \quad (45)$$

For the model selection criteria AIC, the degree of freedom is approximated by the number of non-zero c_{jk} 's in the group penalty.

References

- [1] N. Meinshausen and P. Bühlmann. High-dimensional graphs and variable selection with the lasso. *The Annals of Statistics*, 34(3):1436–1462, 2006.
- [2] J. Peng, P. Wang, N. Zhou, and J. Zhu. Partial correlation estimation by joint sparse regression models. *Journal of the American Statistical Association*, 104(486):735–746, 2009.
- [3] P. Ravikumar, M.J. Wainwright, and J. Lafferty. High-dimensional Ising model selection using l_1 -regularized logistic regression. *Annals of Statistics*, 38(3):1287–1319, 2010.
- [4] H. Höfling and R. Tibshirani. Estimation of sparse binary pairwise markov networks using pseudo-likelihoods. *The Journal of Machine Learning Research*, 10:883–906, 2009.
- [5] Han Liu, Xi Chen, John Lafferty, and Larry Wasserman. Graph-valued regression. In J. Lafferty, C. K. I. Williams, J. Shawe-Taylor, R.S. Zemel, and A. Culotta, editors, *Advances in Neural Information Processing Systems 23*, pages 1423–1431. 2010.
- [6] M. Schmidt, K. Murphy, G. Fung, and R. Rosales. Structure learning in random fields for heart motion abnormality detection. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8, 2008.
- [7] J.K. Bradley and C. Guestrin. Learning tree conditional random fields. In *Proceedings of the 27th International Conference on Machine Learning*, pages 127–134, 2010.
- [8] M. Schmidt and K. Murphy. Convex structure learning in log-linear models: Beyond pairwise potentials. In *Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2010.
- [9] M. Yuan and Y. Lin. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 68(1):49–67, 2006.
- [10] L. Jacob, G. Obozinski, and J.P. Vert. Group Lasso with overlap and graph Lasso. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pages 433–440, 2009.
- [11] P. Zhao, G. Rocha, and B. Yu. The composite absolute penalties family for grouped and hierarchical variable selection. *Annals of Statistics*, 37(6A):3468–3497, 2009.
- [12] S.J. Wright. Accelerated block-coordinate relaxation for regularized optimization. Technical report, Department of Computer Science, University of Wisconsin-Madison, 2010.
- [13] M.J. Wainwright and M.I. Jordan. Graphical models, exponential families, and variational inference. *Foundations and Trends® in Machine Learning*, 1:1–305, 2008.
- [14] F. Gao, G. Wahba, R. Klein, and B. Klein. Smoothing Spline ANOVA for multivariate Bernoulli observations, with application to ophthalmology data. *Journal of the American Statistical Association*, 96(453):127, 2001.
- [15] G. Wahba. *Spline Models for Observational Data*. Society for Industrial Mathematics, 1990.
- [16] R. Jenatton, J.Y. Audibert, and F. Bach. Structured variable selection with sparsity-inducing norms. *arXiv:0904.3523*, 2009.
- [17] S. Kim and E.P. Xing. Tree-guided group lasso for multi-task regression with structured sparsity. In *Proceedings of 27th International Conference on Machine Learning*, pages 543–550, Haifa, Israel, 2010.
- [18] J. Liu and J. Ye. Fast overlapping group lasso. *arXiv:1009.0306v1*, 2010.
- [19] Xiwen Ma. *Penalized Regression in Reproducing Kernel Hilbert Spaces With Randomized Covariate Data*. PhD thesis, Department of Statistics, University of Wisconsin-Madison, 2010.
- [20] K. Koh, S.J. Kim, and S. Boyd. An interior-point method for large-scale l_1 -regularized logistic regression. *Journal of Machine Learning Research*, 8(8):1519–1555, 2007.
- [21] R.M. Scammon, A.V. McGilivray, and R. Cook. *America Votes 26: 2003-2004, Election Returns By State*. CQ Press, 2005.
- [22] J. Whittaker. *Graphical models in applied multivariate statistics*. Wiley (Chichester England and New York), 1990.
- [23] C. De Boor. *A practical guide to splines*. Applied Mathematical Sciences, 1978.
- [24] D. Xiang and G. Wahba. A generalized approximate cross validation for smoothing splines with non-Gaussian data. *Statistica Sinica*, 6:675–692, 1996.
- [25] W. Shi, G. Wahba, S. Wright, K. Lee, R. Klein, and B. Klein. LASSO-Patternsearch algorithm with application to ophthalmology and genomic data. *Statistics and its Interface*, 1(1):137, 2008.