

# Sparse regression with clustered predictors

Aiden Kenny, Danielle Solomon, and Kumer Das  
Department of Mathematics  
Lamar University

January 24, 2020

## Abstract

We investigate the potential effectiveness of using clustering algorithms to generate a grouping structure for high-dimensional data sets. Using various regularization techniques, we seek to determine if the generated groups are truly relevant to the response and if the accuracy and interpretability of the models can be improved. We apply the clustered group structure to two real-world data sets.

## 1 Introduction

The idea of using data and information to train models that are both accurate and interpretable has been around for decades. One desires to build a model based on the *predictors* that is both accurate and interpretable; we want our models to correctly predict the outcome and we want to know which predictors are responsible. However, in the age of big data it is becoming increasingly common that a data set is *high-dimensional*, meaning the number of predictors  $p$  vastly exceeds the number of observations  $n$ . In this setting, many longstanding statistical modeling techniques, such as linear and logistic regression, no longer suffice. Regularization is a popular technique that imposes a penalty on the original model; in some cases the models are *sparse*, meaning they are very interpretable.

It is sometimes the case that the predictors of a model belong to some kind of pre-defined group, and the response is now based on these groups, as opposed to the individual predictors. More advanced regularization methods have been developed to accommodate for group structure, and assuming that the groups are well-represented, can greatly improve the accuracy and interpretability of the models. Unfortunately, while the response could truly be dependent on the group structure, the actual grouping structure is unknown beforehand. In this situation, one would desire to properly identify the grouping structure and build a model based on the result.

The goal of this paper is to investigate the effect that clustering can have on regularized models. We seek to answer two questions:

1. Can clustering algorithms be used to properly identify a grouping structure in a data set?
2. Can grouping the predictors using the clustering information improve the accuracy and interpretability of the model?

In Section 2 we provide a brief overview of the various regularization techniques and clustering algorithms we used in our study. Sections 3 and 4 investigate the effect of clustering predictors on two real-world genomic data sets, and we close with a discussion in Section 5.

## 2 Methodology

### 2.1 Logistic regression

In many situations, the response variable of a data set is categorical in nature, and we wish to assign an observation to one of the response variables given its inputs, a process known as classification. We seek to model the *probability* that an observation falls into a given class. It is often the case where the response belongs to one of two classes coded as  $\mathcal{G} = \{0, 1\}$ . In this binary setting, one popular approach to modeling the probabilities is *logistic regression*.

Suppose we have  $n$  observations and  $p$  predictors stored in a data matrix  $\mathbf{X} = \{x_{i,j}\}$  for  $i = 1, \dots, n$  and  $j = 1, \dots, p$ , along with a response vector  $\mathbf{y} = (y_1, \dots, y_n)$ , where  $y_i \in \{0, 1\}$ . If  $p(\mathbf{x}_i) = \mathbb{P}(Y = 1 \mid X = \mathbf{x}_i)$ , where  $\mathbf{x}_i = (x_{i,1}, \dots, x_{i,p})$  is the  $i$ th observation in  $\mathbf{X}$ , then the probability is modeled (as the log-odds) by

$$\log \left( \frac{p(\mathbf{x}_i)}{1 - p(\mathbf{x}_i)} \right) = \beta_0 + \mathbf{x}_i^T \boldsymbol{\beta}. \quad (1)$$

From this, the estimated response  $\hat{y}_i$  is 1 if  $p(\mathbf{x}_i) \geq 0.5$  and 0 otherwise. The coefficients  $\beta_0$  and  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)$  are estimated from the data by minimizing the negative log-likelihood function

$$L(\beta_0, \boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n \left[ \log \left( 1 + e^{\beta_0 + \mathbf{x}_i^T \boldsymbol{\beta}} \right) - y_i (\beta_0 + \mathbf{x}_i^T \boldsymbol{\beta}) \right]. \quad (2)$$

### 2.2 Basic regularization

In general, a regularized linear model seeks to minimize a penalized version of (2) of the form

$$Q(\beta_0, \boldsymbol{\beta}) = L(\beta_0, \boldsymbol{\beta}) + \lambda P(\boldsymbol{\beta}),$$

where  $P(\boldsymbol{\beta})$  is some type of penalty imposed on the coefficient vector  $\boldsymbol{\beta}$ . The tuning parameter  $\lambda \geq 0$  effectively controls the severity of the penalty; as the value of  $\lambda$  increases, more shrinkage is imposed on the coefficients.

Various regularization methods have been introduced throughout the years using different penalty functions, with each method shrinking the coefficients in a different way. *Ridge regression* (Hoerl & Kennard, 1970) imposes a squared  $\ell_2$  norm on  $\boldsymbol{\beta}$ , and seeks to minimize

$$Q(\beta_0, \boldsymbol{\beta}) = L(\beta_0, \boldsymbol{\beta}) + \frac{\lambda}{2} \|\boldsymbol{\beta}\|_2^2. \quad (3)$$

The  $\ell_2$  norm causes continuous shrinkage of the estimated coefficients. A major drawback to ridge regression is that it produces dense models, i.e. models where  $\beta_j \neq 0$  for all  $j$ , an undesirable characteristic for an interpretable model.

An alternative similar to ridge regression is the *lasso* (Tibshirani, 1996), which minimizes

$$Q(\beta_0, \boldsymbol{\beta}) = L(\beta_0, \boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\|_1. \quad (4)$$

Here, an  $\ell_1$  norm is imposed on  $\boldsymbol{\beta}$ , as opposed to a squared  $\ell_2$  norm. Unlike ridge regression, the lasso is able to perform variable selection, forcibly setting many estimated coefficients to zero, producing sparse models. The resulting sparsity of the model often makes the lasso more preferable than ridge regression in the high-dimensional setting. Unfortunately, the lasso has several caveats as well; in the high-dimensional setting the lasso will select at most  $n$  predictors, and if several predictors are highly-correlated, the lasso will select only one and force the others to zero.

A generalization to ridge regression and the lasso, which attempts to combine the benefits while negating the drawbacks, is the *elastic net*<sup>1</sup> (Zou & Hastie, 2005), which minimizes

$$Q(\beta_0, \beta) = L(\beta_0, \beta) + \lambda \left[ \alpha \|\beta\|_1 + \frac{1 - \alpha}{2} \|\beta\|_2^2 \right]. \quad (5)$$

This penalty is a linear combination of (3) and (4), and the mixing parameter  $\alpha \in [0, 1]$  is used to determine how much of each type of penalty is imposed on the model;  $\alpha = 0$  corresponds to ridge regression, while  $\alpha = 1$  gives the lasso.

### 2.3 The group setting

Much work has been done to develop penalties that exploit pre-determined group structure. Suppose that the predictors of  $\mathbf{X}$  are split into  $K$  non-overlapping groups, with  $S_k$  denoting the size of the  $k$ th group. For  $k = 1, \dots, K$ , let  $\mathbf{X}_k \in \mathbb{R}^{n \times S_k}$  denote the data matrix with the predictors in group  $k$ , and let  $\beta_k = (\beta_{k,1}, \dots, \beta_{k,S_k})$  be the sub-vector of  $\beta$  corresponding to the  $k$ th group.

The *group lasso* (“gLasso”) (Yuan & Lin, 2006) imposes an  $\ell_2$  norm on each of the coefficient sub-vectors; it minimizes

$$Q(\beta_0, \beta) = L(\beta_0, \beta) + \lambda \sum_{k=1}^K \sqrt{S_k} \|\beta_k\|_2. \quad (6)$$

The group lasso was later extended to logistic regression by Meier, Van De Geer, and Bühlmann (2008). The  $\ell_2$  penalties on each of the coefficient sub-vectors creates sparsity among the different groups while performing ridge shrinkage within each group. As a result, the group lasso unfortunately only induces sparsity at the group level, and if a group is determined to be significant, *all* of the group’s predictors will be nonzero.

Both Yuan and Lin (2006) and Meier et al. (2008) assume that the data is orthonormal within each group, i.e.  $\mathbf{X}_k^T \mathbf{X}_k = \mathbf{I}$  for all  $k$ . This is almost never the case in practice, so one would want to orthonormalize each  $\mathbf{X}_k$  before minimizing (6). However, as Simon and Tibshirani (2012) show, this actually changes the penalty to

$$Q(\beta_0, \beta) = L(\beta_0, \beta) + \lambda \sum_{k=1}^K \sqrt{S_k} \|\mathbf{X}_k \beta_k\|_2. \quad (7)$$

This alternative penalty is theoretically and computationally superior (P. Breheny & Huang, 2015) to (6), so for the rest of this paper we refer to (7) when speaking about the group lasso.

There are several methods used in practice to induce sparsity both within and among groups, a feature known as *bi-variate selection*. One method is to combine the lasso and the group lasso penalties as a linear combination, similar to the elastic net; the resulting penalty is known as the *sparse group lasso* (Simon, Friedman, Hastie, & Tibshirani, 2013), and minimizes

$$Q(\beta_0, \beta) = L(\beta_0, \beta) + \lambda \left[ \alpha \|\beta\|_1 + (1 - \alpha) \sum_{k=1}^K \sqrt{S_k} \|\beta_k\|_2 \right]. \quad (8)$$

With this penalty, sparsity is induced at the group level, and elastic net-type shrinkage is imposed within each group. Unfortunately, unlike the group lasso, there is no way to orthonormalize each group without corrupting the within-group sparsity effect, making any implementation of the sparse group lasso difficult compared to other penalties.

<sup>1</sup>Zou and Hastie (2005) call this penalty the *naïve* elastic net penalty, and suggest that scaling the estimated coefficients up by a factor of  $1 + \lambda(1 - \alpha)$  improves prediction accuracy. However, in their paper describing the implementation of the elastic net, Friedman, Hastie, and Tibshirani (2010) abandon this distinction.

An alternative method for performing bi-variate selection is the *composite minimax concave penalty*<sup>2</sup> (“cMCP”) (P. Breheny & Huang, 2009), which minimizes

$$Q(\beta_0, \beta) = L(\beta_0, \beta) + \sum_{k=1}^K f_{\lambda, \Gamma_k} \left( \sum_{s=1}^{S_k} f_{\lambda, \gamma}(|\beta_{k,s}|) \right). \quad (9)$$

Here  $f_{\lambda, \gamma}(\cdot)$  is the minimax concave penalty (Zhang, 2007), given by

$$f_{\lambda, \gamma}(\phi) = \begin{cases} \lambda\phi - \frac{\phi^2}{2\gamma}, & \text{if } \phi \leq \gamma\lambda \\ \frac{1}{2}\gamma\lambda^2, & \text{if } \phi > \gamma\lambda \end{cases}. \quad (10)$$

The intuition behind (10) is to counter the aggressive shrinkage that the lasso imposes on large coefficients. The parameter  $\gamma > 0$  controls the “range” of this counter, and the penalty becomes the lasso as  $\gamma \rightarrow \infty$ . We can see that (9) is effectively applying the penalty twice, once to induce sparsity within each group and then again to induce sparsity among the groups. The outer parameter is set to  $\Gamma_k = \frac{1}{2}S_k\gamma\lambda$ , while the inner penalty  $\gamma$  is the same throughout.

## 2.4 Regularization based on principal components

Let  $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T$  be the singular value decomposition of the data matrix, and let  $m = \text{rank}(\mathbf{X})$ . The principal axes, or right singular vectors, are given by the columns of  $\mathbf{V} \in \mathbb{R}^{p \times m}$ , and  $\mathbf{d} = (d_1, \dots, d_m)$  are the singular values such that  $d_1 \geq \dots \geq d_m > 0$ .  $\mathbf{D} \in \mathbb{R}^{m \times m}$  is a diagonal matrix whose diagonal entries are the elements of  $\mathbf{d}$ .

Principal components lasso (“pcLasso”) (Tay, Friedman, & Tibshirani, 2018) minimizes

$$Q(\beta_0, \beta) = L(\beta_0, \beta) + \lambda\|\beta\|_1 + \frac{\theta}{2}\beta^T \left( \mathbf{V}\mathbf{D}_{d_1^2-d_j^2} \mathbf{V}^T \right) \beta, \quad (11)$$

where  $\lambda$  and  $\theta$  are two separate tuning parameters. The diagonal matrix  $\mathbf{D}_{d_1^2-d_j^2} \in \mathbb{R}^{m \times m}$  has diagonal inputs that are functions of the singular values of  $\mathbf{X}$ , and is given by

$$\mathbf{D}_{d_1^2-d_j^2} = \text{diag}(d_1^2 - d_1^2, d_1^2 - d_2^2, \dots, d_1^2 - d_m^2). \quad (12)$$

This “pcLasso penalty” has the result of imposing less shrinkage in the direction of the leading principal axis and more severe shrinkage in the directions of subsequent principal axes. In other words,  $\beta$  is biased in the direction of the leading principal axis. The presence of the  $\ell_1$  norm allows pcLasso to simultaneously perform feature selection.

pcLasso can also be modified to exploit group structure. Let  $\mathbf{X}_k = \mathbf{U}_k\mathbf{D}_k\mathbf{V}_k^T$  be the singular value decomposition for the  $k$ th group matrix, and let  $m_k = \text{rank}(\mathbf{X}_k)$ . Then the columns of  $\mathbf{V}_k$  and  $\mathbf{d}_k = (d_{k,1}, \dots, d_{k,m_k})$  are the principal axes and singular values of  $\mathbf{X}_k$ , respectively. In this setting, pcLasso seeks to minimize

$$Q(\beta_0, \beta) = L(\beta_0, \beta) + \lambda\|\beta\|_1 + \frac{\theta}{2} \sum_{k=1}^K \sqrt{S_k} \beta_k^T \left( \mathbf{V}_k \mathbf{D}_{d_{k,1}^2-d_{k,j}^2} \mathbf{V}_k^T \right) \beta_k. \quad (13)$$

Similar to (12), the matrix  $\mathbf{D}_{d_{k,1}^2-d_{k,j}^2} \in \mathbb{R}^{m_k \times m_k}$  is given by

$$\mathbf{D}_{d_{k,1}^2-d_{k,j}^2} = \text{diag}(d_{k,1}^2 - d_{k,1}^2, d_{k,1}^2 - d_{k,2}^2, \dots, d_{k,1}^2 - d_{k,m_k}^2). \quad (14)$$

We now see that pcLasso biases each coefficient sub-vector  $\beta_k$  in the direction of that group’s leading principal axis, all while producing sparse models. Unlike the group lasso and cMCP, pcLasso does not require each group matrix  $\mathbf{X}_k$  to be orthonormal.

<sup>2</sup>P. Breheny and Huang (2009) originally denote cMCP as the *group* MCP. To avoid confusion, Huang, Breheny, and Ma (2012) recommend denoting (9) as the *composite* MCP.

## 2.5 Clustering methods

In general, a clustering algorithm seeks to partition the predictors of a data set into different sub-groups based on some dissimilarity measure; ideally, the dissimilarity will be low for predictors within the same cluster and high for predictors in separate clusters. Various clustering algorithms and dissimilarity measures exist that seek to achieve this goal, so for this report we only focus on two simple algorithms.

*K-means clustering* (MacQueen et al., 1967) clusters the predictors into  $K$  non-overlapping groups based on their Euclidean distance in the observation space. Let  $C_k$  be the set of all predictors that belong to group  $k$ , for  $k = 1, \dots, K$ , and let  $S_k$  denote the number of predictors in group  $k$ .  $K$ -means clustering seeks to minimize the total within-cluster variation, given by

$$W_K = \sum_{k=1}^K \sum_{j \in C_k} \|\mathbf{x}_j - \bar{\mathbf{x}}_k\|_2^2, \quad (15)$$

where  $\mathbf{x}_j$  is the  $j$ th predictor and  $\bar{\mathbf{x}}_k = \frac{1}{S_k} \sum_{j \in C_k} \mathbf{x}_j$  is the  $k$ th group centroid.

One drawback of  $K$ -means clustering is that the number of clusters  $K$  must be supplied before  $W_K$  can be minimized. Given that we have no information about the group structure beforehand, we desire some type of measurement to determine the optimal number of clusters. The GAP statistic (Tibshirani, Walther, & Hastie, 2001), defined as

$$\text{Gap}(m) = \mathbb{E}[\log(W_m)] - \log(W_m), \quad (16)$$

attempts to choose an optimal  $K$  from the rate that  $W_K$  decreases. Given a maximum amount of clusters  $M$ , the gap statistic is estimated using  $B$  Monte Carlo random samples; the optimal number of clusters  $K$  is chosen when  $\text{Gap}(K) \geq \text{Gap}(K+1) - \delta(K+1)$ , where

$$\delta(m) = \text{sd}_m \sqrt{1 + \frac{1}{B}}$$

and  $\text{sd}_m$  is the standard deviation of the  $B$  estimated  $W_m$ 's.

Another appealing clustering algorithm is *hierarchical clustering*, which groups the predictors into nested clusters. While Euclidean distance could be used as the choice of dissimilarity, we chose to use correlation instead to investigate how this choice effects the resulting group structure. Unlike  $K$ -means clustering, there is no criteria that could be used to choose the optimal number of clusters. Fortunately, given that the clusters are nested, they can be represented in a dendrogram where the number of clusters can be chosen by the user.

## 3 Colon data set

The colon data set, originally introduced by Alon et al. (1999), contains the gene expressions of 2,000 genes for 62 different tissue samples, i.e.  $n = 62$  and  $p = 2,000$ . Of the 62 tissue samples, 40 of the samples tested positive for colon cancer, while 22 tested negative.

### 3.1 Clustering information

Figure 1 shows the clustering information for the colon data set. The top-left panel measures the gap statistics for  $m = 1, \dots, 20$ , and chooses  $K = 9$  as the optimal number of clusters using  $K$ -means clustering. The top-right panel plots the predictors against the first two columns of  $\mathbf{Z} = \mathbf{V}\mathbf{D}$  (where  $\mathbf{V}\mathbf{D}\mathbf{U}^T = \mathbf{X}^T$  is the SVD of the transposed data matrix, so the columns of  $\mathbf{Z}$  are the principal components of  $\mathbf{X}^T$ ), along with the labeled groups that each predictor belongs to.

The bottom panel displays the corresponding dendrogram using correlation as the dissimilarity measure for hierarchical clustering. We decided that  $K = 7$  was a reasonable cut-off for this dendrogram.

### 3.2 Results

For all of the following methods, we randomly split the data set into a training set and a test set, both with 31 observations each. Each model was fit on the training set, and its performance was measured on the test set. Models that have additional parameters besides  $\lambda$  were fit using a grid of values of said parameter (e.g.  $\alpha$  for the elastic net), and the model with the lowest deviance was chosen.

We first fit the following regularized models on the colon data set *without* any grouping structure:

1. The lasso: was fit using `glmnet` version 2.0.16.
2. The elastic net: was fit for  $\alpha = 0.95, 0.8, 0.6, 0.4, 0.2, 0.05$ , using `glmnet` version 2.0.16.
3. pcLasso: was fit for  $\text{rat}^3 = 0.95, 0.9, 0.75, 0.5, 0.25, 0.1$ , using `pcLasso` version 1.1.

Next, we fit the following models on the colon data set using the grouping structure obtained from  $K$ -means clustering:

1. gLasso: was fit using `grpreg` version 3.2.1.
2. sgLasso: was fit for  $\alpha = 0.95, 0.8, 0.6, 0.4, 0.2, 0.05$ , using `SGL` version 1.2.
3. cMCP: was fit for  $\gamma = 30$ , using `grpreg` version 3.2.1.
4. pcLasso: was fit for  $\text{rat} = 0.95, 0.9, 0.75, 0.5, 0.25, 0.1$ , using `pcLasso` version 1.1.

The `clusGap` function from `cluster` version 2.1.0 was used to calculate the gap statistic, and the groups were clustered using the `kmeans` function. Finally, the process above was repeated using the grouping structure from hierarchical clustering, which was obtained using the `hclust` function.

The results from the various models are presented in Table 1. Included in the table are the values of the optimized parameters, the cross-validation deviance, the number of missclassifications on the test set, the number of nonzero coefficients in the final model (including the intercept  $\beta_0$ ), the number of significant groups in the final model (if a single group contained a nonzero coefficient, then it is considered significant), and the area under the curve (AUC) measurements. The corresponding ROC curves for each model have been printed in Figure 2.

We can see that pcLasso using hierarchical clustering performs the best in terms of missclassifications, having only incorrectly predicting 3 of the 31 test observations. When looking at the missclassifications for all of the models, we can see that the lasso, the elastic net, and pcLasso without clustering all perform about the same.

With  $K$ -means clustering, gLasso missclassifies 12 of the 31 test observations, which is worse than a null model (the test set had 11 observations that tested negative for cancer and 20 that tested positive). sgLasso performs slightly worse, and cMCP and pcLasso perform about the same as the non-clustered models. Looking at the number of significant groups, we can see that pcLasso does not perform bi-variate selection at all, since all nine groups are represented in the final model. This, along with the poor performance of gLasso and sgLasso, indicate that the grouping structure obtained using  $K$ -means clustering is not sufficient.

A similar situation occurs with hierarchical clustering. This time we see that sgLasso missclassifies the most test observations, and gLasso also performs worse than the non-grouped models. Interestingly, the trained cMCP model using hierarchical clustering is identical in size to the model using  $K$ -means clustering. Finally, pcLasso here both missclassifies the least amount

---

<sup>3</sup>As opposed to testing over a grid of values for  $\theta$ , [Tay et al. \(2018\)](#) suggest specifying a value of `rat` instead. More details can be found in their paper.

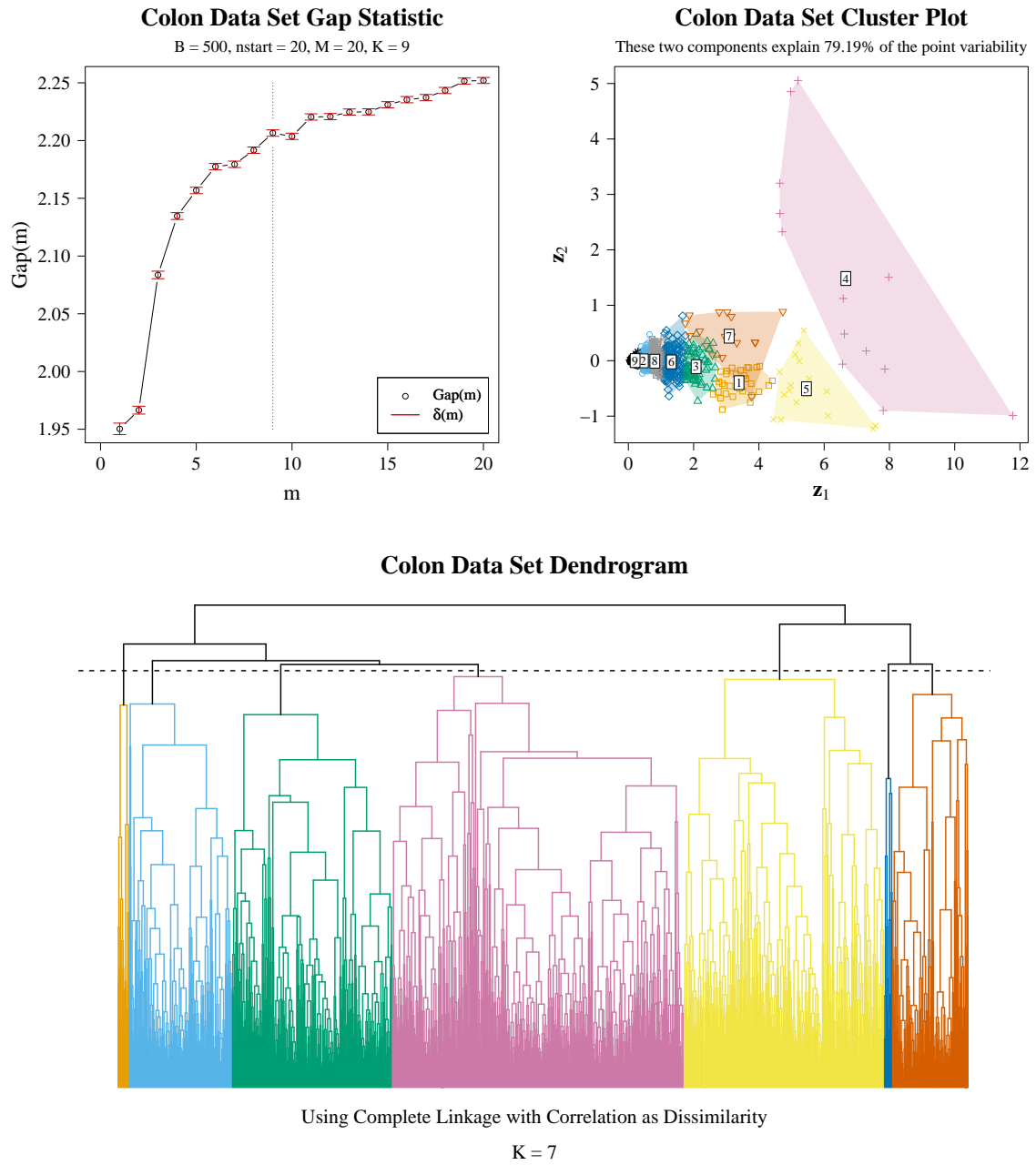


Figure 1: *Clustering information for the colon data set.*

	No Clustering			K-means Clustering				Hierarchical Clustering			
	—			K = 9				K = 7			
	Lasso	Elastic Net	pcLasso	gLasso	sgLasso	cMCP	pcLasso	gLasso	sgLasso	cMCP	pcLasso
Parameters	$\lambda = 0.0642$	$\lambda = 0.0853$	$\lambda = 8.75$	$\lambda = 0.00562$	$\lambda = 0.0215$	$\lambda = 0.0838$	$\lambda = 35.32$	$\lambda = 0.0621$	$\lambda = 0.0194$	$\lambda = 0.0813$	$\lambda = 163.94$
	—	$\alpha = 0.95$	<b>rat</b> = 0.95	—	$\alpha = 0.4$	$\gamma = 30$	<b>rat</b> = 0.25	—	$\alpha = 0.05$	$\gamma = 30$	<b>rat</b> = 0.95
Deviance	0.938	0.945	0.561	0.490	0.853	1.01	0.525	0.987	0.863	0.995	0.548
Misclass.	6/31	6/31	5/31	12/31	7/31	6/31	5/31	8/31	11/31	6/31	3/31
Sig. Coef.	16	19	30	49	29	11	30	20	461	11	7
Sig. Groups	—	—	—	3	2	4	9	1	1	4	6
AUC	0.936	0.945	0.850	0.695	0.745	0.932	0.891	0.623	0.814	0.932	0.818

Table 1: *The performance of various models on the colon data set.*

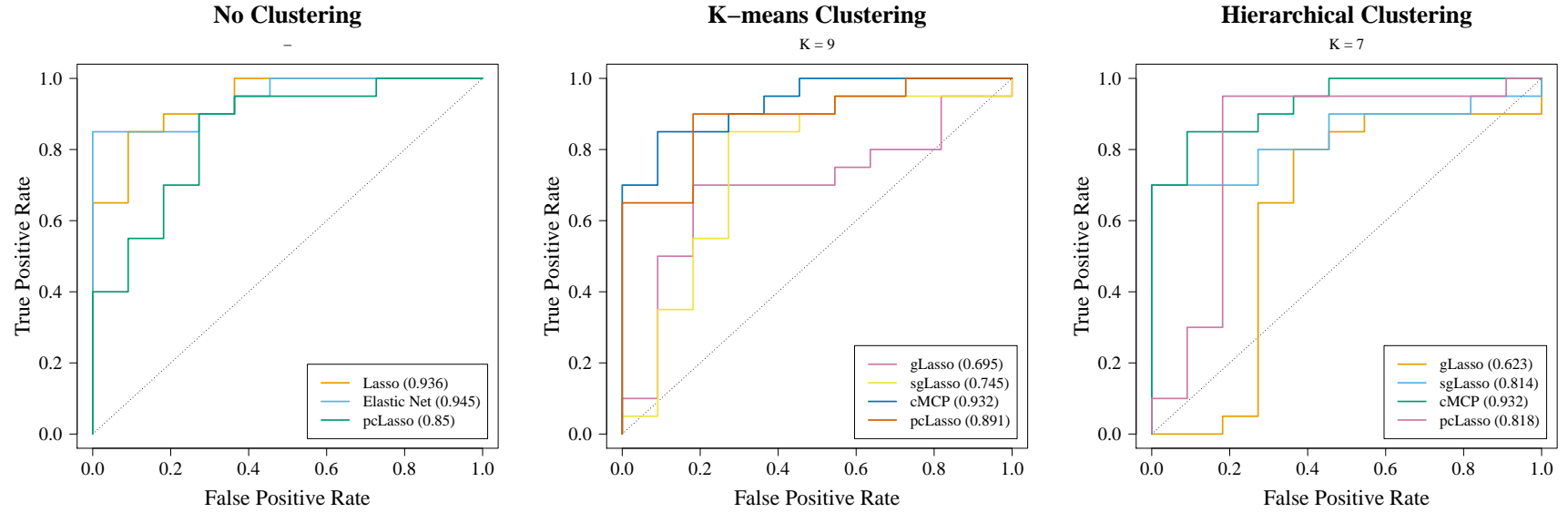


Figure 2: *The ROC curves for the colon data set.*



of test observations and has the least number of significant coefficients. There are six significant coefficients (keep in mind the table includes the intercept term) as well as six non-zero coefficients, so again pcLasso does not induce shrinkage at the group level. It is worth noting that pcLasso using hierarchical clustering also has less significant coefficients than all of the non-clustered models, so it is both more accurate and more interpretable.

## 4 Leukemia data set

The leukemia data set, from [Golub et al. \(1999\)](#), contains the gene expressions of 7,128 genes for 72 different patients ( $n = 72$  and  $p = 7,128$ ). The response is the type of leukemia each patient has; 47 were diagnosed with acute lymphoblastic leukemia (ALL) while 25 were diagnosed with acute myeloid leukemia (AML).

### 4.1 Clustering information

The clustering information for the leukemia data set has been printed in Figure 3. The gap statistics for  $m = 1, \dots, 20$  have been printed in the top left panel, and  $K = 19$  is chosen as the optimal number of clusters. The top right panels plots the predictors against the first two principal components of the observation space, with each group labeled. Finally, the bottom panel shows the leukemia dendrogram; we chose  $K = 5$  at the optimal number of clusters.

### 4.2 Results

The same general set up for the colon data set was applied to the leukemia data set; the results have been printed in Table 2 and the ROC curves have been printed in Figure 4. As with the colon data set, it seems that the clustering algorithms are unable to identify a sufficient grouping structure.

For the non-clustered models, both the elastic net and pcLasso perform the same in terms of missclassifications. However, the elastic net has a markedly lower deviance, less significant coefficients, and a higher AUC, making it a decisively better model in this situation.

Both gLasso and sgLasso perform extremely poorly for both  $K$ -means and hierarchical clustering; in fact, gLasso actually fits a null model in both cases. cMCP again fits models with similar sizes in both cases, and while the missclassification rate is much better than gLasso and sgLasso, it still performs worse than the non-clustered models. As with the colon data set, clustered pcLasso is the overall winner, having only missclassified two of the 36 test observations in both cases. In addition, we see that pcLasso does not induce shrinkage at the group level. And while both models perform the same in prediction accuracy, pcLasso using hierarchical clustering has a slightly lower deviance, significantly less significant predictors, and a slightly higher AUC, so it can be considered the best model for all of the clustered predictors.

Unlike the colon data set, however, there is some type of trade-off that one would have to consider when choosing an overall best model. While pcLasso with hierarchical clustering has the lowest number of missclassifications, it has 45 significant predictors. On the other hand, the lasso and the elastic net have 13 and 27, respectively, so even though their missclassification rates are slightly higher, they are more interpretable models. It is entirely reasonable for one to choose a model with one additional incorrect prediction if it is easier to interpret and explain.

## 5 Discussion

With both data sets, we saw that the group lasso and the sparse group lasso performed very poorly relative to the other methods, especially for the leukemia data set. In addition, pcLasso

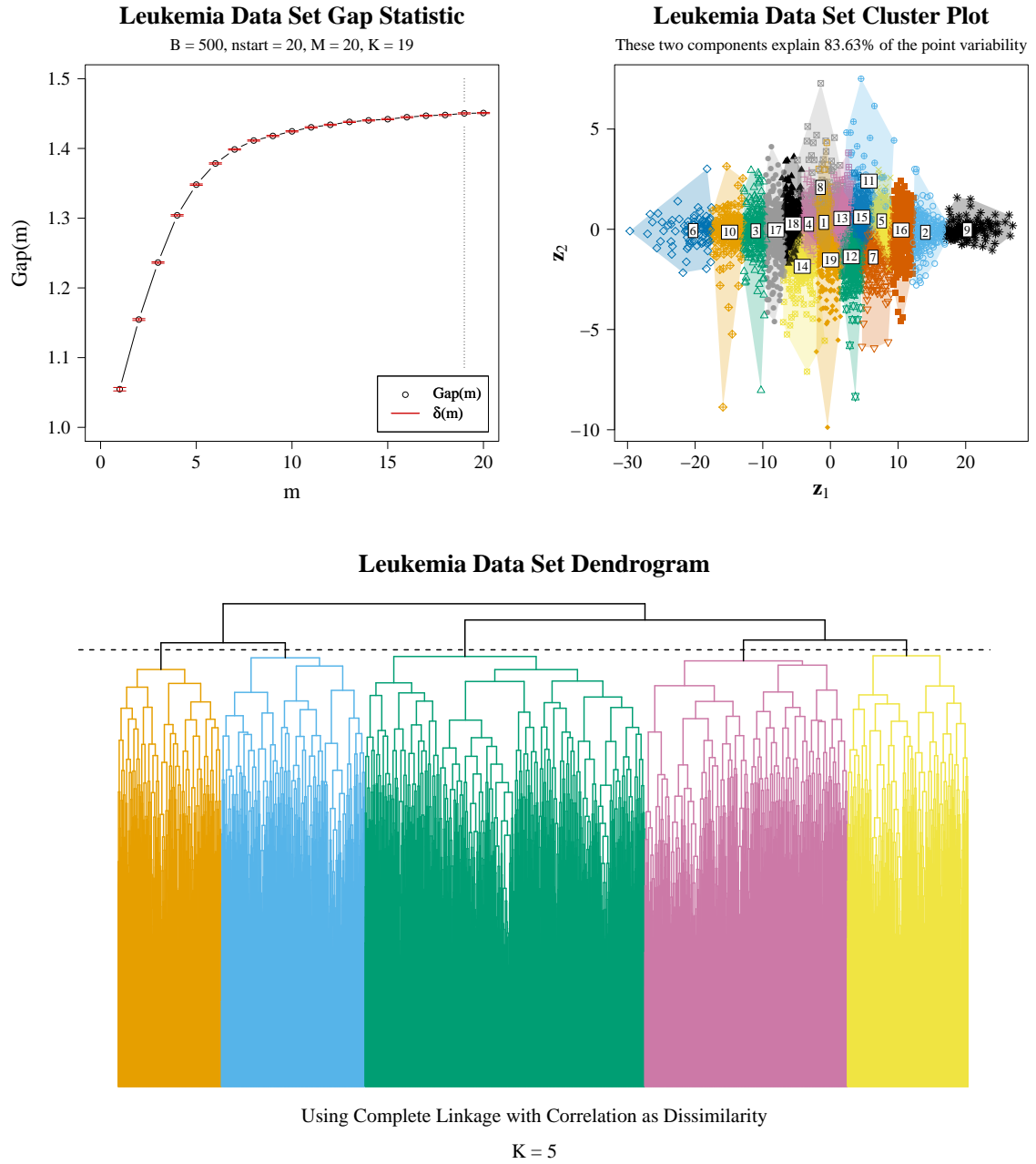
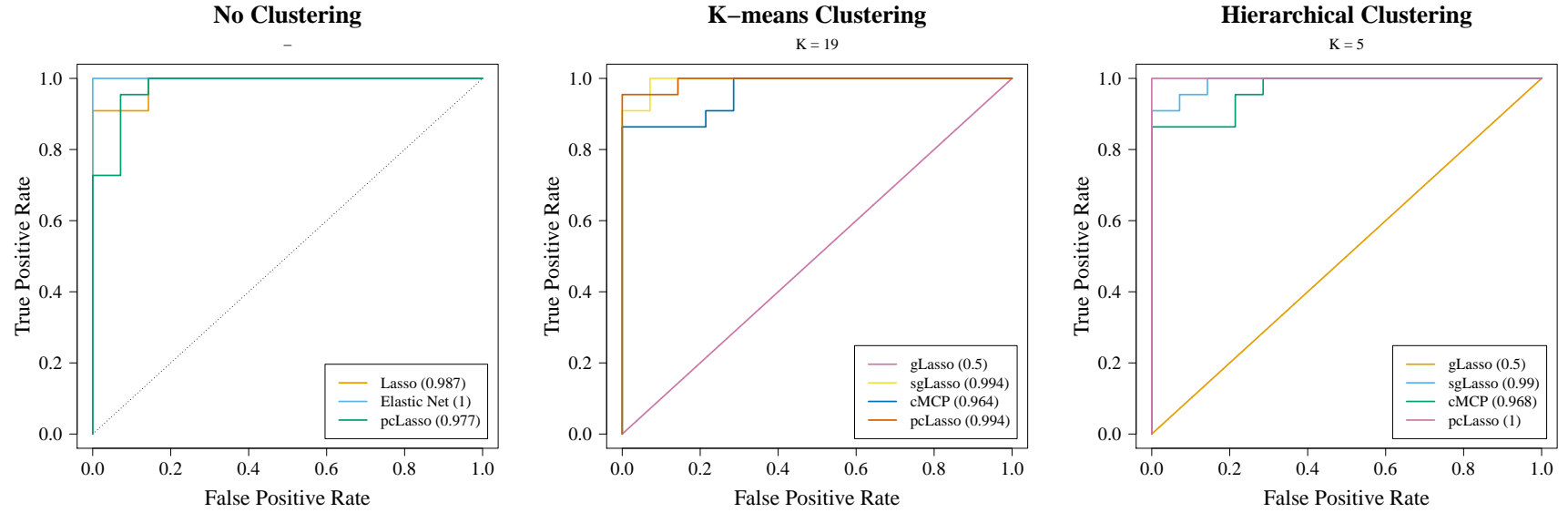


Figure 3: *Clustering information for the leukemia data set.*

	No Clustering			<i>K</i> -means Clustering				Hierarchical Clustering			
	–			<i>K</i> = 19				<i>K</i> = 5			
	Lasso	Elastic Net	pcLasso	gLasso	sgLasso	cMCP	pcLasso	gLasso	sgLasso	cMCP	pcLasso
Parameters	$\lambda = 0.00407$	$\lambda = 0.00509$	$\lambda = 0.00548$	$\lambda = 0.0779$	$\lambda = 0.0457$	$\lambda = 0.111$	$\lambda = 0.00548$	$\lambda = 0.0779$	$\lambda = 0.0236$	$\lambda = 0.111$	$\lambda = 0.00548$
	–	$\alpha = 0.8$	<b>rat</b> = 0.95	–	$\alpha = 0.95$	$\gamma = 30$	<b>rat</b> = 0.9	–	$\alpha = 0.4$	$\gamma = 30$	<b>rat</b> = 0.95
Deviance	0.241	0.0834	0.465	1.240	0.730	0.671	0.442	1.240	0.731	0.671	0.439
Misclass.	5/36	3/36	3/36	14/36	13/36	6/36	2/36	14/36	14/36	6/36	2/36
Sig. Coef.	14	28	41	1	6	6	76	1	772	7	46
Sig. Groups	–	–	–	0	1	2	19	0	1	2	5
AUC	0.987	1.000	0.977	0.500	0.994	0.964	0.994	0.500	0.990	0.968	1.000

Table 2: *The performance of various models on the leukemia data set.*Figure 4: *The ROC curves for the leukemia data set.*

did not perform bi-variate selection. Both of these facts are indications that both  $K$ -means clustering and hierarchical clustering were *ineffective* in properly identifying a grouping structure that was relevant to the response. However, we also observed that pcLasso with clustered predictors had a slightly lower number of missclassifications than the non-clustered models, showing that even though the groups were not relevant to the response, the models can potentially become more accurate and interpretable.

There are several avenues that one could follow for further research:

- Looking at the top right panels of Figures 1 and 3, we can see that the clusters are not well-separated at all. Despite this, the gap statistic split the predictors into a rather large number of clusters, especially for the leukemia data set, which may be a result of the large number of predictors. Perhaps another measure to determine the optimal number of clusters can be used. It is also worth mentioning that the computation time when using the `clusGap` function was on the magnitude of *hours*, making it non-practical for larger data sets.
- The two clustering algorithms used were chosen because of their simplicity and reputation, but one could argue that their poor performance is entirely expected. Much work has been done in the field of unsupervised learning, and one could use a different clustering algorithm to generate the group structure.
- There are many more grouped regularization models that can be employed that perform bi-variate selection, such as the group exponential lasso (“GEL”) (P. J. Breheny, 2015) or the group bridge (Huang, Ma, Xie, & Zhang, 2009), that can be efficiently implemented using the `grpreg` package. Each of these methods have their own interesting properties, and one could investigate if their use can improve over the non-clustered models.

**Acknowledgements:** We would like to thank J. Kenneth Tay for his helpful comments about pcLasso and undergraduate research in general. We would also like to thank Patrick Breheny for his helpful comments about his `grpreg` package. The authors were supported by NSF Award 1757717.

## References

- Alon, U., Barkai, N., Notterman, D. A., Gish, K., Ybarra, S., Mack, D., & Levine, A. J. (1999). Broad patterns of gene expression revealed by clustering analysis of tumor and normal colon tissues probed by oligonucleotide arrays. *Proceedings of the National Academy of Sciences*, 96(12), 6745–6750.
- Breheny, P., & Huang, J. (2009). Penalized methods for bi-level variable selection. *Statistics and its interface*, 2(3), 369.
- Breheny, P., & Huang, J. (2015). Group descent algorithms for nonconvex penalized linear and logistic regression models with grouped predictors. *Statistics and computing*, 25(2), 173–187.
- Breheny, P. J. (2015). The group exponential lasso for bi-level variable selection. *Biometrics*, 71 3, 731–40.
- Friedman, J., Hastie, T., & Tibshirani, R. (2010). Regularization paths for generalized linear models via coordinate descent. *Journal of statistical software*, 33(1), 1.
- Golub, T. R., Slonim, D. K., Tamayo, P., Huard, C., Gaasenbeek, M., Mesirov, J. P., ... others (1999). Molecular classification of cancer: class discovery and class prediction by gene expression monitoring. *science*, 286(5439), 531–537.
- Hoerl, A. E., & Kennard, R. W. (1970). Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, 12(1), 55–67.

- Huang, J., Breheny, P., & Ma, S. (2012). A selective review of group selection in high-dimensional models. *Statistical science: a review journal of the Institute of Mathematical Statistics*, 27(4).
- Huang, J., Ma, S., Xie, H., & Zhang, C.-H. (2009). A group bridge approach for variable selection. *Biometrika*, 96(2), 339–355.
- MacQueen, J., et al. (1967). Some methods for classification and analysis of multivariate observations. In *Proceedings of the fifth berkeley symposium on mathematical statistics and probability* (Vol. 1, pp. 281–297).
- Meier, L., Van De Geer, S., & Bühlmann, P. (2008). The group lasso for logistic regression. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 70(1), 53–71.
- Simon, N., Friedman, J., Hastie, T., & Tibshirani, R. (2013). A sparse-group lasso. *Journal of Computational and Graphical Statistics*, 22(2), 231–245.
- Simon, N., & Tibshirani, R. (2012). Standardization and the group lasso penalty. *Statistica Sinica*, 22(3), 983.
- Tay, J. K., Friedman, J., & Tibshirani, R. (2018, Oct). Principal component-guided sparse regression. *arXiv e-prints*, arXiv:1810.04651.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1), 267–288.
- Tibshirani, R., Walther, G., & Hastie, T. (2001). Estimating the number of clusters in a data set via the gap statistic. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 63(2), 411–423.
- Yuan, M., & Lin, Y. (2006). Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 68(1), 49–67.
- Zhang, C. H. (2007). Penalized linear unbiased selection. *Department of Statistics and Bioinformatics, Rutgers University*, 3.
- Zou, H., & Hastie, T. (2005). Regularization and variable selection via the elastic net. *Journal of the royal statistical society: series B (statistical methodology)*, 67(2), 301–320.