

Covariance-regularized regression and classification for high-dimensional problems

Daniela M. Witten†

*Department of Statistics, Stanford University, 390 Serra Mall, Stanford CA 94305, USA.
E-mail: dwitten@stanford.edu*

Robert Tibshirani

Departments of Statistics and Health Research & Policy, Stanford University, 390 Serra Mall, Stanford CA 94305, USA. E-mail: tibs@stat.stanford.edu

Summary. We propose covariance-regularized regression, a family of methods for prediction in high-dimensional settings that uses a shrunken estimate of the inverse covariance matrix of the features in order to achieve superior prediction. An estimate of the inverse covariance matrix is obtained by maximizing the log likelihood of the data, under a multivariate normal model, subject to a penalty; it is then used to estimate coefficients for the regression of the response onto the features. We show that ridge regression, the lasso, and the elastic net are special cases of covariance-regularized regression, and we demonstrate that certain previously unexplored forms of covariance-regularized regression can outperform existing methods in a range of situations. The covariance-regularized regression framework is extended to generalized linear models and linear discriminant analysis, and is used to analyze gene expression data sets with multiple class and survival outcomes.

Keywords: regression, classification, $n \ll p$, covariance regularization, variable selection

1. Introduction

In high-dimensional regression problems, where p , the number of features, is nearly as large as, or larger than, n , the number of observations, ordinary least squares regression does not provide a satisfactory solution. A remedy for the shortcomings of least squares is to modify the sum of squared errors criterion used to estimate the regression coefficients, using penalties that are based on the magnitudes of the coefficients:

$$\hat{\beta} = \arg \min_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda_1 \|\beta\|^{p_1} + \lambda_2 \|\beta\|^{p_2} \quad (1)$$

(Here, the notation $\|\beta\|^s$ is used to indicate $\sum_{i=1}^p |\beta_i|^s$.) Many popular regularization methods fall into this framework. For instance, when $\lambda_2 = 0$, $p_1 = 0$ gives best subset selection, $p_1 = 2$ gives ridge regression (Hoerl & Kennard 1970), and $p_1 = 1$ gives the lasso (Tibshirani 1996). More generally, for $\lambda_2 = 0$ and $p_1 \geq 0$, the above equation defines the bridge estimators (Frank & Friedman 1993). Equation 1 defines the naive elastic net in the case that $p_1 = 1$ and $p_2 = 2$ (Zou & Hastie 2005). In this paper, we present a new approach to regularizing linear regression that involves applying a penalty not to the sum of squared errors, but rather to the log likelihood of the data under a multivariate normal model.

†Corresponding author.

The least squares solution is $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$. In multivariate normal theory, the entries of $(\mathbf{X}^T \mathbf{X})^{-1}$ that equal zero correspond to pairs of variables that have no sample partial correlation; in other words, pairs of variables that are conditionally independent, given all of the other features in the data. Non-zero entries of $(\mathbf{X}^T \mathbf{X})^{-1}$ correspond to non-zero partial correlations. One way to perform regularization of least squares regression is to shrink the matrix $(\mathbf{X}^T \mathbf{X})^{-1}$; in fact, this is done by ridge regression, since the ridge solution can be written as $\hat{\beta}_{ridge} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$. Here, we propose a more general approach to shrinkage of the inverse covariance matrix. Our method involves estimating a regularized inverse covariance matrix by maximizing the log likelihood of the data under a multivariate normal model, subject to a constraint on the elements of the inverse covariance matrix. In doing this, we attempt to distinguish between variables that truly are partially correlated with each other and variables that in fact have zero partial correlation. We then use this regularized inverse covariance matrix in order to obtain regularized regression coefficients. We call the class of regression methods defined by this procedure the **Scout**.

In Section 2, we present the Scout criteria and explain the method in greater detail. We also discuss connections between the Scout and pre-existing regression methods. In particular, we show that ridge regression, the lasso, and the elastic net are special cases of the Scout. In addition, we present some specific members of the Scout class that perform well relative to pre-existing methods in a variety of situations. In Sections 3, 4, and 5, we demonstrate the use of these methods in regression, classification, and generalized linear model settings on simulated data and on a number of gene expression data sets.

2. The Scout Method

2.1. The General Scout Family

Let $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_p)$ denote an $n \times p$ matrix of data, where n is the number of observations and p the number of features. Let \mathbf{y} denote a vector of length n , containing a response value for each observation. Assume that the columns of \mathbf{X} are standardized, and that \mathbf{y} is centered. We can create a matrix $\tilde{\mathbf{X}} = (\mathbf{X} \ \mathbf{y})$, which has dimension $n \times (p + 1)$. If we assume that $\tilde{\mathbf{X}} \sim \mathbf{N}(\mathbf{0}, \Sigma)$, then we can find the maximum likelihood estimator of the population inverse covariance matrix Σ^{-1} by maximizing

$$\log(\det \Sigma^{-1}) - \text{tr}(\mathbf{S} \Sigma^{-1}) \quad (2)$$

where $\mathbf{S} = \begin{pmatrix} \mathbf{S}_{xx} & \mathbf{S}_{xy} \\ \mathbf{S}_{xy}^T & S_{yy} \end{pmatrix}$ is the empirical covariance matrix of $\tilde{\mathbf{X}}$. Assume for a moment that \mathbf{S} is invertible. Then, the maximum likelihood estimator for Σ^{-1} is \mathbf{S}^{-1} (we use the fact that $\frac{d}{d\mathbf{W}} \log \det \mathbf{W} = \mathbf{W}^{-1}$ for a symmetric positive definite matrix \mathbf{W}). Let $\Theta = \begin{pmatrix} \Theta_{xx} & \Theta_{xy} \\ \Theta_{xy}^T & \Theta_{yy} \end{pmatrix}$ denote a symmetric estimate of Σ^{-1} . The problem of regressing \mathbf{y} onto \mathbf{X} is closely related to the problem of estimating Σ^{-1} , since the least squares coefficients for the regression equal $-\frac{\Theta_{xy}}{\Theta_{yy}}$ for $\Theta = \mathbf{S}^{-1}$ (this follows from the partitioned inverse formula; see e.g. Mardia et al. (1979) page 459). If $p > n$, then some type of regularization is needed in order to estimate the regression coefficients, since \mathbf{S} is not invertible. Even if $p < n$, we may want to shrink the least squares coefficients in some way in order to achieve superior prediction. The connection between estimation

of Θ and estimation of the least squares coefficients suggests the possibility that rather than shrinking the coefficients β by applying a penalty to the sum of squared errors for the regression of \mathbf{y} onto \mathbf{X} , as is done for example in ridge regression or the lasso, we can obtain shrunken β estimates through maximization of the penalized log likelihood of the data.

To do this, one could estimate Σ^{-1} as Θ that maximizes

$$\log(\det \Theta) - \text{tr}(\mathbf{S}\Theta) - J(\Theta) \quad (3)$$

where $J(\Theta)$ is a penalty function. For example, $J(\Theta) = \|\Theta\|^p$ denotes the sum of absolute values of the elements of Θ if $p = 1$, and it denotes the sum of squared elements of Θ if $p = 2$. Our regression coefficients would then be given by the formula $\beta = -\frac{\Theta_{xy}}{\Theta_{yy}}$. However, recall that the ij element of Θ is zero if and only if the partial correlation of $\tilde{\mathbf{x}}_i$ with $\tilde{\mathbf{x}}_j$ (conditional on all of the other variables in $\tilde{\mathbf{X}}$) is zero. (This follows from the definition of the partial correlation, and again from the partitioned inverse formula.) Note that \mathbf{y} is included in $\tilde{\mathbf{X}}$. So it does not make sense to regularize the elements of Θ as presented above, because we really care about the partial correlations of pairs of variables given the other variables, as opposed to the partial correlations of pairs of variables given the other variables and the response.

For these reasons, rather than obtaining an estimate of Σ^{-1} by maximizing the penalized log likelihood in Equation 3, we estimate it via a two-stage maximization, given in the following algorithm:

The Scout Procedure for General Penalty Functions

1. Compute $\hat{\Theta}_{\mathbf{xx}}$, which maximizes

$$\log(\det \Theta_{\mathbf{xx}}) - \text{tr}(\mathbf{S}_{\mathbf{xx}} \Theta_{\mathbf{xx}}) - J_1(\Theta_{\mathbf{xx}}). \quad (4)$$

2. Compute $\hat{\Theta}$, which maximizes

$$\log(\det \Theta) - \text{tr}(\mathbf{S}\Theta) - J_2(\Theta), \quad (5)$$

where the top left $p \times p$ submatrix of $\hat{\Theta}$ is constrained to equal $\hat{\Theta}_{\mathbf{xx}}$, the solution to Step 1.

3. Compute $\hat{\beta}$, defined by $\hat{\beta} = -\frac{\hat{\Theta}_{xy}}{\hat{\Theta}_{yy}}$.

4. Compute $\hat{\beta}^* = c\hat{\beta}$, where c is the coefficient for the regression of \mathbf{y} onto $\mathbf{X}\hat{\beta}$.

$\hat{\beta}^*$ denotes the regularized coefficients obtained using this new method. Step 1 of the Scout Procedure involves obtaining shrunken estimates of $(\Sigma_{\mathbf{xx}})^{-1}$ in order to smooth our estimates of which variables are conditionally independent. Step 2 involves obtaining shrunken estimates of Σ^{-1} , conditional on $(\Sigma^{-1})_{\mathbf{xx}} = \hat{\Theta}_{\mathbf{xx}}$, the estimate obtained in Step 1. Thus, we obtain regularized estimates of which predictors are dependent on \mathbf{y} , given all of the other predictors. The scaling in the last step is performed because it has been found, empirically, to improve performance.

By penalizing the entries of the inverse covariance matrix of the predictors in Step 1 of the Scout Procedure, we are attempting to distinguish between pairs of variables that

Table 1. Special cases of the Scout.

| $J_1(\Theta_{\mathbf{xx}})$ | $J_2(\Theta)$ | <i>Method</i> |
|-----------------------------------|----------------|------------------|
| 0 | 0 | Least Squares |
| $\text{tr}(\Theta_{\mathbf{xx}})$ | 0 | Ridge Regression |
| $\text{tr}(\Theta_{\mathbf{xx}})$ | $\ \Theta\ ^1$ | Elastic Net |
| 0 | $\ \Theta\ ^1$ | Lasso |
| 0 | $\ \Theta\ ^2$ | Ridge Regression |

truly are conditionally dependent, and pairs of variables that appear to be conditionally dependent due only to chance. We are searching, or **scouting**, for variables that truly are correlated with each other, conditional on all of the other variables. Our hope is that sets of variables that truly are conditionally dependent will also be related to the response. In the context of a microarray experiment, where the variables are genes and the response is some clinical outcome, this assumption is reasonable: we seek genes that are part of a pathway related to the response. One expects that such genes will also be conditionally dependent. In Step 2, we shrink our estimates of the partial correlation between each predictor and the response, given the shrunken partial correlations between the predictors that we estimated in Step 1. In contrast to ordinary least squares regression, which uses the inverse of the empirical covariance matrix to compute regression coefficients, we jointly model the relationship that the p predictors have with each other and with the response in order to obtain shrunken regression coefficients.

We define the **Scout family** of estimated coefficients for the regression of \mathbf{y} onto \mathbf{X} as the solutions $\hat{\beta}^*$ obtained in Step 4 of the Scout Procedure. We refer to the penalized log likelihoods in Steps 1 and 2 of the Scout Procedure as the first and second **Scout criteria**.

In the rest of the paper, when we discuss properties of the Scout, for ease of notation we will ignore the scale factor in Step 4 of the Scout Procedure. For instance, if we claim that two procedures yield the same regression coefficients, we more specifically mean that the regression coefficients are the same up to scaling by a constant factor.

Least squares, the elastic net, the lasso, and ridge regression result from the Scout Procedure with appropriate choices of J_1 and J_2 (up to a scaling by a constant). Details are in Table 1. The first two results can be shown directly by differentiating the Scout criteria, and the others follow from Equation 11 in Section 2.4.

2.2. L_p Penalties

Throughout the remainder of this paper, with the exception of Section 3.2, we will exclusively be interested in the case that $J_1(\Theta) = \lambda_1 \|\Theta\|^{p_1}$ and $J_2(\Theta) = \frac{\lambda_2}{2} \|\Theta\|^{p_2}$, where the norm is taken elementwise over the entries of Θ , and where $\lambda_1, \lambda_2 \geq 0$. For ease of notation, $\text{Scout}(p_1, p_2)$ will refer to the solution to the Scout criterion with J_1 and J_2 as just mentioned. If $\lambda_2 = 0$, then this will be indicated by $\text{Scout}(p_1, \cdot)$, and if $\lambda_1 = 0$, then this will be indicated by $\text{Scout}(\cdot, p_2)$. Therefore, in the rest of this paper, the Scout Procedure will be as follows:

The Scout Procedure with L_p Penalties

1. Compute $\hat{\Theta}_{\mathbf{xx}}$, which maximizes

$$\log(\det \Theta_{\mathbf{xx}}) - \text{tr}(\mathbf{S}_{\mathbf{xx}} \Theta_{\mathbf{xx}}) - \lambda_1 \|\Theta_{\mathbf{xx}}\|^{p_1}. \quad (6)$$

2. Compute $\hat{\Theta}$, which maximizes

$$\log(\det \Theta) - \text{tr}(\mathbf{S}\Theta) - \frac{\lambda_2}{2} \|\Theta\|^{p_2}, \quad (7)$$

where the top left $p \times p$ submatrix of $\hat{\Theta}$ is constrained to equal $\hat{\Theta}_{\mathbf{xx}}$, the solution to Step 1. Note that because of this constraint, the penalty really is only being applied to the last row and column of $\hat{\Theta}$.

3. Compute $\hat{\beta}$, defined by $\hat{\beta} = -\frac{\hat{\Theta}_{\mathbf{xy}}}{\hat{\Theta}_{\mathbf{yy}}}$.

4. Compute $\hat{\beta}^* = c\hat{\beta}$, where c is the coefficient for the regression of \mathbf{y} onto $\mathbf{X}\hat{\beta}$.

2.3. Simple Example

Here, we present a toy example in which $n = 20$ observations on $p = 19$ variables are generated under the model $\mathbf{y} = \mathbf{X}\beta + \epsilon$, where $\beta_j = j$ for $j \leq 10$ and $\beta_j = 0$ for $j > 10$, and where $\epsilon_i \sim N(0, 25)$. In addition, the first 10 variables have correlation 0.5 with each other; the rest are uncorrelated.

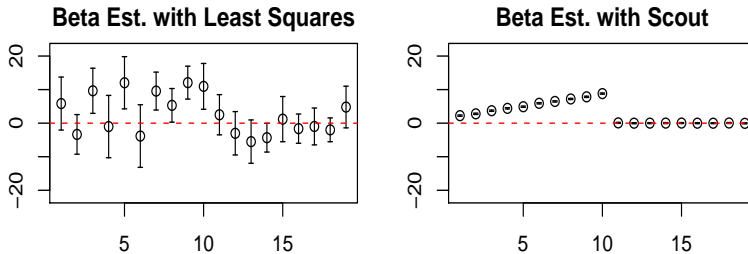


Fig. 1. Data were generated under a simple model. Average coefficient estimates (over 500 repetitions) and 95% confidence intervals are shown. Panels show $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ (left) and $\text{Scout}(1, \cdot)$ (right). The red dashed line indicates $y = 0$.

Figure 1 shows the average over 500 simulations of the least squares regression coefficients and the $\text{Scout}(1, \cdot)$ regression estimate. It is not surprising that least squares performs poorly in this situation, since n is barely larger than p . $\text{Scout}(1, \cdot)$ performs quite well; though it results in coefficient estimates that are slightly biased, they have much lower variance. This simple example demonstrates that benefits can result from the use of a shrunk estimate of the inverse covariance matrix.

2.4. Maximization of the Scout Criteria with L_p Penalties

If $\lambda_1 = 0$, then the maximum of the first Scout criterion is given by $(\mathbf{S}_{\mathbf{xx}})^{-1}$ (if $\mathbf{S}_{\mathbf{xx}}$ is invertible). In the case that $\lambda_1 > 0$ and $p_1 = 1$, maximization of the first Scout criterion has been studied extensively; see e.g. Meinshausen & Bühlmann (2006). The solution can be found via the “graphical lasso”, an efficient algorithm given by Banerjee et al. (2008) and Friedman et al. (2007) that involves iteratively regressing one row of the

estimated covariance matrix onto the others, subject to an L_1 constraint, in order to update the estimate for that row.

If $\lambda_1 > 0$ and $p_1 = 2$, the solution to Step 1 of the Scout Procedure is even easier. We want to find $\Theta_{\mathbf{xx}}$ that maximizes

$$\log(\det \Theta_{\mathbf{xx}}) - \text{tr}(\mathbf{S}_{\mathbf{xx}} \Theta_{\mathbf{xx}}) - \lambda \|\Theta_{\mathbf{xx}}\|^2. \quad (8)$$

Differentiating with respect to $\Theta_{\mathbf{xx}}$, we see that the maximum solves

$$\Theta_{\mathbf{xx}}^{-1} - 2\lambda \Theta_{\mathbf{xx}} = \mathbf{S}_{\mathbf{xx}}. \quad (9)$$

This equation implies that $\Theta_{\mathbf{xx}}$ and $\mathbf{S}_{\mathbf{xx}}$ share the same eigenvectors. Letting θ_i denote the i^{th} eigenvalue of $\Theta_{\mathbf{xx}}$ and letting s_i denote the i^{th} eigenvalue of $\mathbf{S}_{\mathbf{xx}}$, it is clear that

$$\frac{1}{\theta_i} - 2\lambda \theta_i = s_i. \quad (10)$$

We can easily solve for θ_i , and can therefore solve the first Scout criterion exactly in the case $p_1 = 2$, in essentially just the computational cost of obtaining the eigenvalues of $\mathbf{S}_{\mathbf{xx}}$.

It turns out that if $p_2 = 1$ or $p_2 = 2$, then it is not necessary to maximize the second Scout criterion directly, as there is an easier alternative:

Claim 1. *For $p_2 \in \{1, 2\}$, the solution to Step 3 of the Scout Procedure is equal to the solution to the following, up to scaling by a constant:*

$$\hat{\beta} = \arg \min_{\beta} \{ \beta^T \hat{\Sigma}_{\mathbf{xx}} \beta - 2\mathbf{S}_{\mathbf{xy}}^T \beta + \lambda_2 \|\beta\|^{p_2} \} \quad (11)$$

where $\hat{\Sigma}_{\mathbf{xx}}$ is the inverse of the solution to Step 1 of the Scout Procedure.

(The proof of Claim 1 is in Section 8.1.1 in the Appendix.) Therefore, we can replace Steps 2 and 3 of the Scout Procedure with an L_{p_2} regression. It is trivial to show that if $\lambda_2 = 0$ in the Scout Procedure, then the Scout solution is given by $\hat{\beta} = (\hat{\Sigma}_{\mathbf{xx}})^{-1} \mathbf{S}_{\mathbf{xy}}$. It also follows that if $\lambda_1 = 0$, then the cases $\lambda_2 = 0$, $p_2 = 1$, and $p_2 = 2$ correspond to ordinary least squares regression (if the empirical covariance matrix is invertible), the lasso, and ridge regression, respectively.

In addition, we will show in Section 2.5.1 that if $p_1 = 2$ and $p_2 = 1$, then the Scout can be re-written as an elastic net problem with slightly different data; therefore, fast algorithms for solving the elastic net (Friedman et al. 2008) can be used to solve *Scout*(2, 1). The methods for maximizing the Scout criteria are summarized in Table 2.

We compared computation times for *Scout*(2, \cdot), *Scout*(1, \cdot), *Scout*(2, 1), and *Scout*(1, 1) on an example with $n = 100$, $\lambda_1 = \lambda_2 = 0.2$, and \mathbf{X} dense. All timings were carried out on a Intel Xeon 2.80 GHz processor. Table 3 shows the number of CPU seconds required for each of these methods for a range of values of p (the number of features). For all methods, after the Scout coefficients have been estimated for a given set of parameter values, estimation for different parameter values is faster because an approximate estimate of the inverse covariance matrix is available for use as an initial value (when $p_1 = 1$) or because the eigen decomposition has already been computed (when $p_1 = 2$).

Scout(p_1, p_2) involves the use of two tuning parameters, λ_1 and λ_2 ; in practice, these are chosen by cross-validating over a grid of (λ_1, λ_2) values. In Section 2.7, we

Table 2. Maximization of the Scout criteria: special cases.

| | $\lambda_2 = 0$ | $p_2 = 1$ | $p_2 = 2$ |
|-----------------|--------------------|------------------------------|---------------------------------|
| $\lambda_1 = 0$ | Least Squares | L_1 Regression | L_2 Regression |
| $p_1 = 1$ | Graphical Lasso | Graphical Lasso + L_1 Reg. | Graphical Lasso + L_2 Reg. |
| $p_1 = 2$ | Eigenvalue Problem | Elastic Net | Eigenvalue Problem + L_2 Reg. |

Table 3. Timing comparisons for maximization of the Scout criteria, in CPU seconds. $\lambda_1 = \lambda_2 = 0.2$, $n = 100$, \mathbf{X} dense, and p is the number of features.

| p | $Scout(1, \cdot)$ | $Scout(1, 1)$ | $Scout(2, \cdot)$ | $Scout(2, 1)$ |
|------|-------------------|---------------|-------------------|---------------|
| 500 | 1.685 | 1.700 | 0.034 | 0.072 |
| 1000 | 22.432 | 22.504 | 0.083 | 0.239 |
| 2000 | 241.289 | 241.483 | 0.260 | 0.466 |

present a Bayesian connection to the first scout criterion. An editor suggested that as an alternative to cross-validating over λ_1 , one could instead draw from the posterior distribution of $\Theta_{\mathbf{xx}}$.

2.5. Properties of the Scout

In this section, for ease of notation, we will consider an equivalent form of the Scout Procedure obtained by replacing $\mathbf{S}_{\mathbf{xx}}$ with $\mathbf{X}^T \mathbf{X}$ and $\mathbf{S}_{\mathbf{xy}}$ with $\mathbf{X}^T \mathbf{y}$.

2.5.1. Similarities between Scout, Ridge Regression, and the Elastic Net

Let $\mathbf{U}_{n \times p} \mathbf{D}_{p \times p} \mathbf{V}_{p \times p}^T$ denote the singular value decomposition of \mathbf{X} with d_i the i^{th} diagonal element of \mathbf{D} and $d_1 \geq d_2 \geq \dots \geq d_r > d_{r+1} = \dots = d_p = 0$, where $r = \text{rank}(\mathbf{X}) \leq \min(n, p)$. Consider $Scout(2, p_2)$. As previously discussed, the first step in the Scout Procedure corresponds to finding Θ that solves

$$\Theta^{-1} - 2\lambda_1 \Theta = \mathbf{X}^T \mathbf{X}. \quad (12)$$

Since Θ and $\mathbf{X}^T \mathbf{X}$ therefore share the same eigenvectors, it follows that $\Theta^{-1} = \mathbf{V}(\mathbf{D}^2 + \tilde{\mathbf{D}}^2)\mathbf{V}^T$ where $\tilde{\mathbf{D}}^2$ is a $p \times p$ diagonal matrix with i^{th} diagonal entry equal to $\frac{1}{2}(-d_i^2 + \sqrt{d_i^4 + 8\lambda_1})$. It is not difficult to see that ridge regression, $Scout(2, \cdot)$, and $Scout(2, 2)$ result in similar regression coefficients:

$$\begin{aligned} \hat{\beta}_{rr} &= (\mathbf{V}(\mathbf{D}^2 + c\mathbf{I})\mathbf{V}^T)^{-1} \mathbf{X}^T \mathbf{y} \\ \hat{\beta}_{Scout(2, \cdot)} &= (\mathbf{V}(\mathbf{D}^2 + \tilde{\mathbf{D}}^2)\mathbf{V}^T)^{-1} \mathbf{X}^T \mathbf{y} \\ \hat{\beta}_{Scout(2, 2)} &= (\mathbf{V}(\mathbf{D}^2 + \tilde{\mathbf{D}}^2 + \lambda_2 \mathbf{I})\mathbf{V}^T)^{-1} \mathbf{X}^T \mathbf{y} \end{aligned} \quad (13)$$

Therefore, while ridge regression simply adds a constant to the diagonal elements of \mathbf{D} in the least squares solution, $Scout(2, \cdot)$ instead adds a function that is monotone decreasing in the value of the diagonal element. (The consequences of this alternative shrinkage are explored under a latent variable model in Section 2.6). $Scout(2, 2)$ is a compromise between $Scout(2, \cdot)$ and ridge regression.

In addition, we note that the solutions to the naive elastic net and *Scout*(2,1) are quite similar to each other:

$$\begin{aligned}
\hat{\beta}_{enet} &= \arg \min_{\beta} \beta^T \mathbf{V}(\mathbf{D}^2 + c\mathbf{I})\mathbf{V}^T \beta - 2\beta^T \mathbf{X}^T \mathbf{y} + \lambda_2 \|\beta\|^1 \\
&= \arg \min_{\beta} \beta^T \mathbf{X}^T \mathbf{X} \beta - 2\beta^T \mathbf{X}^T \mathbf{y} + \lambda_2 \|\beta\|^1 + c\|\beta\|^2 \\
&= \arg \min_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda_2 \|\beta\|^1 + c\|\beta\|^2 \\
\hat{\beta}_{Scout(2,1)} &= \arg \min_{\beta} \beta^T \mathbf{V}(\mathbf{D}^2 + \tilde{\mathbf{D}}^2)\mathbf{V}^T \beta - 2\beta^T \mathbf{X}^T \mathbf{y} + \lambda_2 \|\beta\|^1 \\
&= \arg \min_{\beta} \beta^T \mathbf{V}(\frac{1}{2}\mathbf{D}^2 + \frac{1}{2}\tilde{\mathbf{D}}^2)\mathbf{V}^T \beta - 2\beta^T \mathbf{X}^T \mathbf{y} + \lambda_2 \|\beta\|^1 + \sqrt{2\lambda_1} \|\beta\|^2 \\
&= \arg \min_{\beta} \|\mathbf{y}^* - \mathbf{X}^* \beta\|^2 + \lambda_2 \|\beta\|^1 + \sqrt{2\lambda_1} \|\beta\|^2 \tag{14}
\end{aligned}$$

where $\tilde{\mathbf{D}}^2$ is the diagonal matrix with elements $\sqrt{d_i^4 + 8\lambda_1} - \sqrt{8\lambda_1}$, and where $\mathbf{X}^* = \frac{1}{\sqrt{2}} \begin{pmatrix} \mathbf{X} \\ \tilde{\mathbf{D}}\mathbf{V}^T \end{pmatrix}$, $\mathbf{y}^* = \begin{pmatrix} \sqrt{2}\mathbf{y} \\ 0 \end{pmatrix}$. From Equation 14, it is clear that *Scout*(2,1) solutions can be obtained using software for the elastic net on data \mathbf{X}^* and \mathbf{y}^* . In addition, given the similarity between the elastic net and *Scout*(2,1) solutions, it is not surprising that *Scout*(2,1) shares some of the elastic net's desirable properties, as is shown in Section 2.5.2.

2.5.2. Variable Grouping Effect

Zou & Hastie (2005) show that unlike the lasso, the elastic net and ridge regression have a variable grouping effect: correlated variables result in similar coefficients. The same is true of *Scout*(2,1):

Claim 2. *Assume that the predictors are standardized and that \mathbf{y} is centered. Let ρ denote the correlation between \mathbf{x}_i and \mathbf{x}_j , and let $\hat{\beta}$ denote the solution to *Scout*(2,1). If $\hat{\beta}_i \hat{\beta}_j \neq 0$, then the following holds:*

$$|\hat{\beta}_i - \hat{\beta}_j| \leq \sqrt{\frac{2(1-\rho)}{\lambda_1}} \|\mathbf{y}\| \tag{15}$$

The proof of Claim 2 is in Section 8.1.2 in the Appendix. Similar results hold for *Scout*(2,·) and *Scout*(2,2), without the requirement that $\hat{\beta}_i \hat{\beta}_j \neq 0$.

2.5.3. Connections to Regression with Orthogonal Features

Assume that the features are standardized, and consider the Scout criterion with $p_1 = 1$. For λ_1 sufficiently large, the solution $\hat{\Theta}_{\mathbf{xx}}$ to the first Scout criterion (Equation 6) is a diagonal matrix with diagonal elements $\frac{1}{\lambda_1 + \mathbf{x}_i^T \mathbf{x}_i}$. (More specifically, if $\lambda_1 \geq |\mathbf{x}_i^T \mathbf{x}_j|$ for all $i \neq j$, then the Scout criterion with $p_1 = 1$ results in a diagonal matrix; see Banerjee et al. (2008) Theorem 4). Thus, if $\hat{\beta}_i$ is the i^{th} component of the *Scout*(1,·) solution, then $\hat{\beta}_i = \frac{\mathbf{x}_i^T \mathbf{y}}{\lambda_1 + 1}$. If $\lambda_2 > 0$, then the resulting Scout solutions with $p_2 = 1$ are given by a variation of the univariate soft thresholding formula for L_1 regression:

$$\hat{\beta}_i = \frac{1}{\lambda_1 + 1} \text{sgn}(\mathbf{x}_i^T \mathbf{y}) \max(0, |\mathbf{x}_i^T \mathbf{y}| - \frac{\lambda_2}{2}) \quad (16)$$

Similarly, if $p_2 = 2$, the resulting Scout solutions are given by the following formula:

$$\hat{\beta} = (1 + \lambda_1 + \lambda_2)^{-1} \mathbf{X}^T \mathbf{y} \quad (17)$$

Therefore, as the parameter λ_1 is increased, the solutions that are obtained range (up to a scaling) from the ordinary L_{p_2} multivariate regression solution to the regularized regression solution for orthonormal features.

2.6. An Underlying Latent Variable Model

Let \mathbf{X} be a $n \times p$ matrix of n observations on p variables, and \mathbf{y} a $n \times 1$ vector of response values. Suppose that \mathbf{X} and \mathbf{y} are generated under the following latent variable model:

$$\begin{aligned} \mathbf{X} &= d_1 \mathbf{u}_1 \mathbf{v}_1^T + d_2 \mathbf{u}_2 \mathbf{v}_2^T \\ d_1, d_2 &> 0 \\ \mathbf{y} &= \mathbf{u}_1 + \epsilon \\ \text{Var}(\epsilon) &= \sigma^2 \mathbf{I} \\ \text{E}(\epsilon) &= 0 \end{aligned} \quad (18)$$

$$\text{E}(\epsilon) = 0 \quad (19)$$

where \mathbf{u}_i and \mathbf{v}_i are the singular vectors of \mathbf{X} , and ϵ is a $n \times 1$ vector of noise.

Claim 3. *Under this model, if $d_1 > d_2$ and the tuning parameters for ridge and Scout(2, ·) are chosen so that the resulting estimators have the same amount of bias, then the estimator given by Scout(2, ·) will have lower variance.*

The proof of Claim 3 is given in Section 8.1.3 of the Appendix. Note that if \mathbf{v}_1 and \mathbf{v}_2 are sparse with non-overlapping regions of non-sparsity, then the model yields a block diagonal covariance matrix with two blocks, where one of the blocks of correlated features is associated with the outcome. In the case of gene expression data, these blocks could represent gene pathways, one of which is responsible for, and has expression that is correlated with, the outcome. Claim 3 shows that if the signal associated with the relevant gene pathway is sufficiently large, then Scout(2, ·) will provide a benefit over ridge.

2.7. Bayesian connection to the first scout criterion

Consider the model $\mathbf{y} = \mathbf{X}\beta + \epsilon$, where ϵ_i are independent Gaussian random variables. It is well-known that ridge regression, the lasso, and the elastic net can be viewed as the Bayes posterior modes under various priors, since they involve solving for β that minimizes a criterion of the form

$$(\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) + \lambda_1 \|\beta\|^{p_1} + \lambda_2 \|\beta\|^{p_2}. \quad (20)$$

Ridge regression corresponds to a normal prior on the elements of β , the lasso corresponds to a double-exponential prior, and the elastic net corresponds to a prior that is a combination of the two.

Similarly, we can think of the solution to the first scout criterion as the Bayes mode of the posterior distribution given by $\mathbf{X} \sim N(0, \mathbf{\Sigma})$ and a prior on the elements of $\mathbf{\Sigma}^{-1}$, such that for $i \leq j$, $(\mathbf{\Sigma}^{-1})_{ij}$ is independent and identically distributed with either a Gaussian distribution (if $p_1 = 2$) or a double-exponential distribution (if $p_1 = 1$). Formally, this would have the potential difficulty that draws from the prior distribution are not constrained to be positive semi-definite.

3. Numerical Studies: Regression via the Scout

3.1. Simulated Data

We compare the performance of ordinary least squares, the lasso, the elastic net, *Scout*(2, 1), and *Scout*(1, 1) on a suite of six simulated examples. The first four simulations are based on those used in the original elastic net paper (Zou & Hastie 2005) and the original lasso paper (Tibshirani 1996). The fifth and sixth are of our own invention. All simulations are based on the model $\mathbf{y} = \mathbf{X}\beta + \sigma\epsilon$ where $\epsilon \sim N(0, \mathbf{I})$. For each simulation, each data set consists of a small training set, a small validation set (used to select the values of the various parameters) and a large test set. We indicate the size of the training, validation, and test sets using the notation $\cdot / \cdot / \cdot$. The simulations are as follows:

1. Each data set consists of 20/20/200 observations, 8 predictors with coefficients $\beta = (3, 1.5, 0, 0, 2, 0, 0, 0)^T$, and $\sigma = 3$. $\mathbf{X} \sim N(0, \mathbf{\Sigma})$, where $\Sigma_{ij} = 0.5^{|i-j|}$.
2. This simulation is as in Simulation 1, except that $\beta_i = 0.85$ for all i .
3. Each data set consists of 100/100/400 observations and 40 predictors. $\beta_i = 0$ for $i \in 1, \dots, 10$ and for $i \in 21, \dots, 30$; for all other i , $\beta_i = 2$. We also set $\sigma = 15$. $\mathbf{X} \sim N(0, \mathbf{\Sigma})$, where $\Sigma_{ij} = 0.5$ for $i \neq j$, and $\Sigma_{ii} = 1$.
4. Each data set consists of 50/50/400 observations and 40 predictors. $\beta_i = 3$ for $i \in 1, \dots, 15$ and $\beta_i = 0$ for $i \in 16, \dots, 40$, and $\sigma = 15$. The predictors are generated as follows:

$$\begin{aligned}
 \mathbf{x}_i &= \mathbf{z}_1 + \epsilon_i^x, \mathbf{z}_1 \sim N(0, \mathbf{I}), i = 1, \dots, 5 \\
 \mathbf{x}_i &= \mathbf{z}_2 + \epsilon_i^x, \mathbf{z}_2 \sim N(0, \mathbf{I}), i = 6, \dots, 10 \\
 \mathbf{x}_i &= \mathbf{z}_3 + \epsilon_i^x, \mathbf{z}_3 \sim N(0, \mathbf{I}), i = 11, \dots, 15
 \end{aligned}
 \tag{21}$$

Also, $\mathbf{x}_i \sim N(0, \mathbf{I})$ are independent and identically distributed for $i = 16, \dots, 40$, and $\epsilon_i^x \sim N(0, 0.01\mathbf{I})$ are independent and identically distributed for $i = 1, \dots, 15$.

5. Each data set consists of 50/50/400 observations and 50 predictors; $\beta_i = 2$ for $i < 9$ and $\beta_i = 0$ for $i \geq 9$. $\sigma = 6$ and $\mathbf{X} \sim N(0, \mathbf{\Sigma})$, where $\Sigma_{ij} = 0.5 \times 1_{i,j \leq 9}$ for $i \neq j$, and $\Sigma_{ii} = 1$.
6. As in Simulation 1, but $\beta = (3, 1.5, 0, 0, 0, 0, -1, -1)^T$.

These simulations cover a variety of settings: Simulations 1, 3, 4, 5, and 6 have sparse β , Simulations 1, 2, 4, 5, and 6 have a sparse inverse covariance matrix, and Simulation 4 is characterized by groups of variables that contribute to the response.

Table 4. Mean squared error averaged over 200 simulated data sets is shown for each simulation. Standard errors are given in parentheses. For each simulation, the two methods with lowest average mean squared errors are shown in bold. Least squares was not performed for Simulation 5, because $p = n$.

| <i>Simulation</i> | <i>Least Squares</i> | <i>Lasso</i> | <i>ENet</i> | <i>Scout(1, 1)</i> | <i>Scout(2, 1)</i> |
|-------------------|----------------------|--------------|--------------------|--------------------|--------------------|
| Sim 1 | 7.72(0.46) | 2.83(0.16) | 2.28(0.13) | 2.22(0.13) | 2.29(0.13) |
| Sim 2 | 7.72(0.46) | 3.26(0.13) | 2.28(0.11) | 1.31(0.09) | 1.54(0.09) |
| Sim 3 | 158.29(3.66) | 44.07(0.80) | 30.86(0.47) | 20.44(0.25) | 18.94(0.28) |
| Sim 4 | 1094.84(44.75) | 54.79(2.30) | 25.06(1.62) | 30.21(1.61) | 28.37(1.52) |
| Sim 5 | NA | 10.91(0.38) | 2.46(0.09) | 1.62(0.09) | 2.18(0.11) |
| Sim 6 | 7.72(0.46) | 2.95(0.16) | 2.34(0.13) | 2.12(0.11) | 2.15(0.11) |

Table 5. Median L_2 distance over 200 simulated data sets is shown for each simulation; details are as in Table 4.

| <i>Simulation</i> | <i>Least Squares</i> | <i>Lasso</i> | <i>ENet</i> | <i>Scout(1, 1)</i> | <i>Scout(2, 1)</i> |
|-------------------|----------------------|--------------|-------------------|--------------------|--------------------|
| Sim 1 | 3.05(0.10) | 1.74(0.05) | 1.65(0.08) | 1.58(0.05) | 1.62(0.06) |
| Sim 2 | 3.05(0.10) | 1.95(0.02) | 1.62(0.03) | 0.90(0.03) | 1.04(0.04) |
| Sim 3 | 17.03(0.22) | 8.91(0.09) | 7.70(0.06) | 6.15(0.01) | 5.83(0.03) |
| Sim 4 | 168.40(5.13) | 17.40(0.16) | 3.85(0.13) | 5.19(2.3) | 3.80(0.14) |
| Sim 5 | NA | 3.48(0.06) | 2.08(0.06) | 1.15(0.03) | 1.55(0.05) |
| Sim 6 | 3.05(0.10) | 1.76(0.06) | 1.53(0.05) | 1.48(0.04) | 1.50(0.03) |

For each simulation, 200 data sets were generated, and the average mean squared errors (with standard errors given in parentheses) are given in Table 4. The Scout provides an improvement over the lasso in all simulations. Both Scout methods result in lower mean squared error than the elastic net in Simulations 2, 3, 5, and 6; in Simulations 1 and 4, the Scout methods are quite competitive. Table 5 shows median L_2 distances between the true and estimated coefficients for each of the models.

Though $Scout(2, 1)$ and $Scout(1, 1)$ perform well relative to the elastic net and lasso in all six simulations, neither dominates the others in all cases. For a given application, we recommend selecting a regression method based on cross-validation error (with the caveat that $Scout(1, 1)$ is slow when the number of features is very large).

3.2. Scout Using Alternative Covariance Estimators

A referee asked whether a different estimator of the inverse covariance matrix of \mathbf{X} could be used in Step 2 of the Scout Procedure, rather than the solution to the first scout criterion. A large body of literature exists on estimation of covariance and inverse covariance matrices. Examples include James & Stein (1961), Haff (1979), Dey & Srinivasan (1985), Bickel & Levina (2008), and Rothman et al. (2008). Any positive definite covariance estimate can be plugged in for $\hat{\Sigma}$ in the equation from Claim 1:

$$\hat{\beta} = \arg \min_{\beta} \{ \beta^T \hat{\Sigma} \beta - 2 \mathbf{S}_{\mathbf{xy}}^T \beta + \lambda_2 \|\beta\|^s \} \quad (22)$$

We explore that possibility here with two covariance estimates: the estimator of James & Stein (1961), and the hard-thresholding estimator of Bickel & Levina (2008). The James-Stein estimator takes the form $\hat{\Sigma}^{JS} = \mathbf{T} \mathbf{D} \mathbf{T}^T$ where \mathbf{T} is a lower triangular matrix with positive elements on the diagonal such that $\mathbf{T} \mathbf{T}^T = \mathbf{X}^T \mathbf{X}$, and \mathbf{D} is a diagonal matrix with diagonal elements $d_i = \frac{1}{n+p+1-2i}$. It is the constant risk minimax estimator under Stein's loss. $\hat{\Sigma}^{BL}$, the estimator of Bickel & Levina (2008), is obtained

Table 6. On Simulation 2, we compare two new estimators obtained by plugging in $\hat{\Sigma}^{JS}$ and $\hat{\Sigma}^{BL}$ to Equation 22. Tuning parameter values were chosen by cross-validation, and standard errors are in parentheses.

| <i>Quantity</i> | <i>Scout(JS, 1)</i> | <i>Scout(BL, 1)</i> | <i>Scout(1, 1)</i> | <i>Scout(2, 1)</i> |
|-----------------------|---------------------|---------------------|--------------------|--------------------|
| Mean MSE | 3.79(0.15) | 2.42(0.12) | 1.31(0.09) | 1.54(0.09) |
| Median L_2 Distance | 3.34(0.14) | 1.94(0.10) | 0.90(0.03) | 1.04(0.04) |

by hard-thresholding each element of the empirical covariance matrix. With $s = 1$ in Equation 22, the resulting methods (which we call $Scout(JS, 1)$ and $Scout(BL, 1)$) are compared to $Scout(2, 1)$ and $Scout(1, 1)$ on Simulation 2, described in Section 3.1. The results are shown in Table 6. In this example, $Scout(JS, 1)$ and $Scout(BL, 1)$ do not perform as well as $Scout(1, 1)$ and $Scout(2, 1)$.

3.3. Making Use of Observations without Response Values

In Step 1 of the Scout Procedure, we estimate the inverse covariance matrix based on the training set \mathbf{X} data, and in Steps 2-4, we compute a penalized least squares solution based on that estimated inverse covariance matrix and $\widehat{\text{Cov}}(\mathbf{X}, \mathbf{y})$. Step 1 of this procedure does not involve the response \mathbf{y} at all.

Now, consider a situation in which one has access to a large amount of \mathbf{X} data, but responses are known for only some of the observations. (For instance, this could be the case for a medical researcher who has clinical measurements on hundreds of cancer patients, but survival times for only dozens of patients.) More specifically, let \mathbf{X}_1 denote the observations for which there is an associated response \mathbf{y} , and let \mathbf{X}_2 denote the observations for which no response data is available. Then, one could estimate the inverse covariance matrix in Step 1 of the Scout Procedure using both \mathbf{X}_1 and \mathbf{X}_2 , and perform Step 2 using $\widehat{\text{Cov}}(\mathbf{X}_1, \mathbf{y})$. By also using \mathbf{X}_2 in Step 1, we achieve a more accurate estimate of the inverse covariance matrix than would have been possible using only \mathbf{X}_1 .

Such an approach will not provide an improvement in all cases. For instance, consider the trivial case in which the response is a linear function of the predictors, $p < n$, and there is no noise: $\mathbf{y} = \mathbf{X}_1\beta$. Then, the least squares solution, using only \mathbf{X}_1 and not \mathbf{X}_2 , is $\hat{\beta} = (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{y} = (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T \mathbf{X}_1 \beta = \beta$. In this case, it clearly is best to only use \mathbf{X}_1 in estimating the inverse covariance matrix. However, one can imagine situations in which one can use \mathbf{X}_2 to obtain a more accurate estimate of the inverse covariance matrix.

Consider a model in which a latent variable has generated some of the features, as well as the response. In particular, suppose that the data are generated as follows:

$$\begin{aligned}
 x_{ij} &= 2u_i + \epsilon_{ij}, j = 1, \dots, 5, i = 1, \dots, n \\
 x_{ij} &= \epsilon_{ij}, j = 6, \dots, 10, i = 1, \dots, n \\
 y_i &= 8u_i + 4\epsilon'_i, i = 1, \dots, n
 \end{aligned}
 \tag{23}$$

In addition, we let $\epsilon_{ij}, \epsilon'_i, u_i \sim N(0, 1)$ i.i.d. The first five variables are “signal” variables, and the rest are “noise” variables. Suppose that we have three sets of observations: a training set of size $n = 12$, for which the \mathbf{y} values are known, a test set of size $n = 200$, for which we wish to predict the \mathbf{y} values, and an additional set of size $n = 36$ observations

Table 7. Making use of observations w/o response values: Set-up.

| | <i>Sample Size</i> | <i>Response Description</i> |
|-----------------|--------------------|---------------------------------|
| Training Set | 12 | Available |
| Test Set | 200 | Unavailable - Must be predicted |
| Additional Obs. | 36 | Unavailable - Not of interest |

Table 8. Making use of observations w/o response values: Results. Standard errors are shown in parentheses. The “null model” predicts each test set outcome value using the mean of the training set outcomes.

| <i>Method</i> | <i>Mean Squared Prediction Error</i> |
|---|--------------------------------------|
| <i>Scout</i> (1, ·) w/Additional Observations | 25.65 (0.38) |
| <i>Scout</i> (1, ·) w/o Additional Observations | 29.92 (0.62) |
| Elastic Net | 32.38 (1.04) |
| Lasso | 47.24 (3.58) |
| Least Squares | 1104.9 (428.84) |
| Null Model | 79.24 (0.3) |

for which we do not know the \mathbf{y} values and do not wish to predict them. This layout is shown in Table 7.

We compare the performances of the Scout and other regression methods. The Scout method is applied in two ways: using only the training set \mathbf{X} values to estimate the inverse covariance matrix, and using also the observations without response values. All tuning parameter values are chosen by 6-fold cross-validation. The results in Table 8 are the average mean squared prediction errors obtained over 500 simulations. From the table, it is clear that both versions of Scout outperform all of the other methods. In addition, using observations that do not have response values does result in a significant improvement.

In this example, twelve labeled observations on ten variables do not suffice to reliably estimate the inverse covariance matrix. The Scout can make use of the observations that lack response values in order to improve the estimate of the inverse covariance matrix, thereby yielding superior predictions.

The use of unlabeled data for classification and regression is sometimes called *semi-supervised learning*. The use of unlabeled observations for linear discriminant analysis dates back to O’Neill (1978); other classical work in this area can be found in McLachlan (1992). It is currently an active area of research in the statistical and machine learning communities. Liang et al. (2007) presents an overview of the use of unlabeled data for prediction, as well as the underlying theory.

4. Classification via the Scout

In classification problems, linear discriminant analysis (LDA) can be used if $n > p$. However, when $p > n$, regularization of the within-class covariance matrix is necessary. Regularized linear discriminant analysis is discussed in Friedman (1989) and Guo et al. (2007). In Guo et al. (2007), the within-class covariance matrix is shrunken, as in ridge regression, by adding a multiple of the identity matrix to the empirical covariance matrix. Here, we instead estimate a shrunken within-class inverse covariance matrix by maximizing the log likelihood of the data, under a multivariate normal model, subject to an L_p penalty on its elements.

4.1. Details of Extension of Scout to Classification

Consider a classification problem with K classes; each observation belongs to some class $k \in 1, \dots, K$. Let $C(i)$ denote the class of training set observation i , which is denoted X_i . Our goal is to classify observations in an independent test set.

Let $\hat{\mu}_k$ denote the $p \times 1$ vector that contains the mean of observations in class k , and let $\mathbf{S}_{\text{wc}} = \frac{1}{n-K} \sum_{k=1}^K \sum_{i:C(i)=k} (X_i - \hat{\mu}_k)(X_i - \hat{\mu}_k)^T$ denote the estimated within-class covariance matrix (based on the training set) that is used for ordinary LDA. Then, the Scout procedure for classification is as follows:

The Scout Procedure for Classification

1. Compute the shrunken within-class inverse covariance matrix $\hat{\Sigma}_{\text{wc},\lambda}^{-1}$ as follows:

$$\hat{\Sigma}_{\text{wc},\lambda}^{-1} = \arg \max_{\Sigma^{-1}} \{ \log \det \Sigma^{-1} - \text{tr}(\mathbf{S}_{\text{wc}} \Sigma^{-1}) - \lambda \|\Sigma^{-1}\|_s \} \quad (24)$$

where λ is a shrinkage parameter.

2. Classify test set observation X to class k' if $k' = \arg \max_k \delta_k^\lambda(X)$, where

$$\delta_k^\lambda(X) = X^T \hat{\Sigma}_{\text{wc},\lambda}^{-1} \hat{\mu}_k - \frac{1}{2} \hat{\mu}_k^T \hat{\Sigma}_{\text{wc},\lambda}^{-1} \hat{\mu}_k + \log \pi_k \quad (25)$$

and π_k is the frequency of class k in the training set.

This procedure is analogous to LDA, but we have replaced \mathbf{S}_{wc} with a shrunken estimate.

This classification rule performs quite well on real microarray data (as is shown below), but has the drawback that it makes use of all of the genes. We can remedy this in one of two ways. We can apply the method described above to only the genes with highest univariate rankings on the training data; this is done in the next section. Alternatively, we can apply an L_1 penalty in estimating the quantity $\hat{\Sigma}_{\text{wc},\lambda}^{-1} \hat{\mu}_k$; note (from Equation 25) that sparsity in this quantity will result in a classification rule that is sparse in the features. Details of this second method, which is not implemented here, are given in Section 8.2 of the Appendix. We will refer to the method detailed in Equation 25 as *Scout*(s, \cdot) because the penalized log likelihood that is maximized in Equation 24 is analogous to the first Scout criterion in the regression case. The tuning parameter λ in Equations 24 and 25 can be chosen via cross-validation.

4.2. Ramaswamy Data

We assess the performance of this method on the Ramaswamy microarray data set, which is discussed in detail in Ramaswamy et al. (2001) and explored further in Zhu & Hastie (2004) and Guo et al. (2007). It consists of a training set of 144 samples and a test set of 54 samples, each of which contains measurements on 16,063 genes. The samples are classified into 14 distinct cancer types. We compare the performance of *Scout*(2, \cdot) to nearest shrunken centroids (NSC) (Tibshirani et al. (2002) and Tibshirani et al. (2003)), L_2 penalized multiclass logistic regression (Zhu & Hastie 2004), the support vector machine (SVM) with one-versus-all classification (Ramaswamy et al. 2001), regularized discriminant analysis (RDA) (Guo et al. 2007), random forests (Breiman 2001), and k-nearest neighbors (KNN). For each method, tuning parameter values were chosen by cross-validation. The results can be seen in Tables 9 and 10; *Scout*(2, \cdot) performed quite

Table 9. Methods are compared on the Ramaswamy Data. All methods were performed on the cube-rooted data, after centering and scaling each patient.

| <i>Method</i> | <i>CV Error</i> | <i>Test Error</i> | <i>Number of Genes Used</i> |
|-----------------------------|-----------------|-------------------|-----------------------------|
| NSC | 35 | 17 | 5217 |
| L_2 Penalized Multinomial | 29 | 15 | 16063 |
| SVM | 26 | 14 | 16063 |
| RDA | 27 | 11 | 9141 |
| KNN | 41 | 29 | 16063 |
| Random Forests | 40 | 24 | 16063 |
| <i>Scout</i> (2, ·) | 22 | 11 | 16063 |

Table 10. Methods are compared on the Ramaswamy Data. Methods were run on the cube-rooted data after centering and scaling the patients, using only the 4,000 genes with highest training set F-statistics.

| <i>Method</i> | <i>Test Error</i> | <i>Number of Genes Used</i> |
|-----------------------------|-------------------|-----------------------------|
| NSC | 21 | 3999 |
| L_2 Penalized Multinomial | 12 | 4000 |
| SVM | 11 | 4000 |
| RDA | 10 | 3356 |
| KNN | 17 | 4000 |
| Random Forests | 17 | 4000 |
| <i>Scout</i> (2, ·) | 7 | 4000 |

well, especially when only the 4,000 genes with highest training set F-statistics were used (Tusher et al. 2001).

5. Extension to Generalized Linear Models and the Cox Model

We have discussed the application of the Scout to classification and regression problems, and we have shown examples in which these methods perform well. In fact, the Scout can also be used in fitting generalized linear models, by replacing the iteratively reweighted least squares step with a covariance-regularized regression. In particular, we discuss the use of the Scout in the context of fitting a Cox proportional hazards model for survival data. We present an example involving four lymphoma microarray datasets in which the Scout results in improved performance relative to other methods.

5.1. Details of Extension of Scout to the Cox Model

Consider survival data of the form $(y_i, \mathbf{x}^i, \delta_i)$ for $i \in 1, \dots, n$, where δ_i is an indicator variable that equals 1 if observation i is complete and 0 if censored, and \mathbf{x}^i is a vector of predictors (x_1^i, \dots, x_p^i) for individual i . Failure times are $t_1 < t_2 < \dots < t_k$; there are d_i failures at time t_i . We wish to estimate the parameter $\beta = (\beta_1, \dots, \beta_p)^T$ in the proportional hazards model $\lambda(t|x) = \lambda_o(t)\exp(\sum_j x_j \beta_j)$. We assume that censoring is noninformative. Letting $\eta = \mathbf{X}\beta$, D the set of indices of the failures, R_r the set of indices of the individuals at risk at time t_r , and D_r the set of indices of the failures at t_r , the partial likelihood is given as follows (see e.g. Kalbfleisch & Prentice (1980)):

$$L(\beta) = \prod_{r \in D} \frac{\exp(\sum_{j \in D_r} \eta_j)}{(\sum_{j \in R_r} \exp(\eta_j))^{d_r}} \quad (26)$$

In order to fit the proportional hazards model, we must find the β that maximizes the partial likelihood above. Let $l(\beta)$ denote the log partial likelihood, $\mathbf{u} = \frac{\partial l}{\partial \eta}$, and $\mathbf{A} = -\frac{\partial^2 l}{\partial \eta \eta^T}$. The iteratively reweighted least squares algorithm that implements the Newton-Raphson method, for β_0 the value of β from the previous step, involves finding β that solves

$$\mathbf{X}^T \mathbf{A} \mathbf{X} (\beta - \beta_0) = \mathbf{X}^T \mathbf{u}. \quad (27)$$

This is equivalent to finding β that minimizes

$$\|\mathbf{y}^* - \mathbf{X}^* \beta^*\|^2 \quad (28)$$

where $\mathbf{X}^* = \mathbf{A}^{1/2} \mathbf{X}$, $\mathbf{y}^* = \mathbf{A}^{-1/2} \mathbf{u}$, $\beta^* = \beta - \beta_0$ (Green 1984).

The traditional iterative reweighted least squares algorithm involves solving the above least squares problem repeatedly, recomputing \mathbf{y}^* and \mathbf{X}^* at each step and setting β_0 equal to the solution β attained at the previous iteration. We propose to solve the above equation using the Scout, rather than by a simple linear regression. We have found empirically that good results are obtained if we initially set $\beta_0 = 0$, and then perform just one Newton-Raphson step (using the Scout). This is convenient, since for data sets with many features, solving a Scout regression can be time-consuming. Therefore, our implementation of the Scout method for survival data involves simply performing one Newton-Raphson step, beginning with $\beta_0 = 0$.

Using the notation $\Theta = \begin{pmatrix} \Theta_{\mathbf{xx}} & \Theta_{\mathbf{xy}} \\ \Theta_{\mathbf{xy}}^T & \Theta_{\mathbf{yy}} \end{pmatrix}$ and $\mathbf{S} = \begin{pmatrix} \mathbf{X}^T \mathbf{A} \mathbf{X} & \mathbf{X}^T \mathbf{u} \\ \mathbf{u}^T \mathbf{X} & \mathbf{u}^T \mathbf{A}^{-1} \mathbf{u} \end{pmatrix}$, the Scout Procedure for survival data is almost identical to the regression case, as follows:

The Scout Procedure for the Cox Model

1. Let $\hat{\Theta}_{\mathbf{xx}}$ maximize

$$\log \det \Theta_{\mathbf{xx}} - \text{tr}(\mathbf{S}_{\mathbf{xx}} \Theta_{\mathbf{xx}}) - \lambda_1 \|\Theta_{\mathbf{xx}}\|^{p_1}. \quad (29)$$

2. Let $\hat{\Theta}$ maximize

$$\log \det \Theta - \text{tr}(\mathbf{S} \Theta) - \lambda_2 \|\Theta\|^{p_2}, \quad (30)$$

where the top $p \times p$ submatrix of Θ is constrained to equal $\hat{\Theta}_{\mathbf{xx}}$, obtained in the previous step.

3. Compute $\hat{\beta} = -\frac{\hat{\Theta}_{\mathbf{xy}}}{\hat{\Theta}_{\mathbf{yy}}}$.

4. Let $\hat{\beta}^* = c \hat{\beta}$, where c is the coefficient of a Cox proportional hazards model fit to \mathbf{y} using $\mathbf{X} \hat{\beta}$ as a predictor.

$\hat{\beta}^*$ obtained in Step 4 is the vector of estimated coefficients for the Cox proportional hazards model. In the procedure above, $\lambda_1, \lambda_2 > 0$ are tuning parameters. In keeping with the notation of previous sections, we will refer to the resulting coefficient estimates as $Scout(p_1, p_2)$.

Table 11. Mean of $2(\log(L) - \log(L_o))$ on survival data. For each data set, the two highest mean values of $2(\log(L) - \log(L_o))$ are shown in bold.

| <i>Data Set</i> | L_1 Cox | SPC | Scout(1,1) | Scout(2,1) |
|-----------------|--------------------|--------------------|--------------------|--------------------|
| Hummel | 2.640(0.99) | 3.823(0.87) | 4.245(1.07) | 3.293(0.91) |
| Monti | 1.647(0.36) | 1.231(0.38) | 2.149(0.46) | 2.606(0.47) |
| Rosenwald | 4.129(0.94) | 3.542(1.17) | 3.987(0.94) | 4.930(1.47) |
| Shipp | 1.903(0.48) | 1.004(0.39) | 2.807(0.73) | 2.627(0.60) |

Table 12. Median Number of Genes Used for Survival Data.

| <i>Data Set</i> | L_1 Cox | SPC | Scout(1,1) | Scout(2,1) |
|-----------------|-----------|-----|------------|------------|
| Hummel | 14 | 33 | 78 | 13 |
| Monti | 18.5 | 17 | 801.5 | 144.5 |
| Rosenwald | 37.5 | 32 | 294 | 85 |
| Shipp | 5.5 | 10 | 4.5 | 5 |

5.2. Lymphoma Data

We illustrate the effectiveness of the Scout method on survival data using four different data sets, all involving survival times and gene expression measurements for patients with diffuse large B-cell lymphoma. The four data sets are as follows: Rosenwald et al. (2002) (“Rosenwald”), which consists of 240 patients, Shipp et al. (2002) (“Shipp”), which consists of 58 patients, Hummel et al. (2006) (“Hummel”), which consists of 81 patients, and Monti et al. (2005) (“Monti”), which consists of 129 patients. For consistency and ease of comparison, we considered only a subset of around 1482 genes that were present in all four data sets.

We randomly split each of the data sets into a training set, a validation set, and a test set of equal sizes. For each data set, we fit four models to the training set: the L_1 penalized Cox proportional hazards (“ L_1 Cox”) method of Park & Hastie (2007), the supervised principal components (SPC) method of Bair & Tibshirani (2004), *Scout*(2,1), and *Scout*(1,1). For each data set, we chose the tuning parameter values that resulted in the predictor that gave the highest log likelihood when used to fit a Cox proportional hazards model on the validation set (this predictor was $\mathbf{X}_{\text{val}}\beta_{\text{train}}$ for L_1 Cox and Scout, and it was the first supervised principal component for SPC). We tested the resulting models on the test set. The mean value of $2(\log(L) - \log(L_o))$ over ten separate training/test/validation set splits is given in Table 11, where L denotes the likelihood of the Cox proportional hazards model fit on the test set using the predictor obtained from the training set (for L_1 Cox and Scout, this was $\mathbf{X}_{\text{test}}\beta_{\text{train}}$, and for SPC, this was the first supervised principal component), and L_o denotes the likelihood of the null model. From Tables 11 and 12, it is clear that the Scout results in predictors that are on par with, if not better than, the competing methods on all four data sets.

6. Discussion

We have presented covariance-regularized regression, a class of regression procedures (the “Scout” family) obtained by estimating the inverse covariance matrix of the data by maximizing the log likelihood of the data under a multivariate normal model, subject to a penalty. We have shown that three well-known regression methods - ridge, the lasso, and the elastic net - fall into the covariance-regularized regression framework. In addition, we have explored some new methods within this framework. We have extended the covariance-regularized regression framework to classification and generalized linear

model settings, and we have demonstrated the performance of the resulting methods on a number of gene expression data sets.

A drawback of the Scout method is that when $p_1 = 1$ and the number of features is large, then maximizing the first Scout criterion can be quite slow. When more than a few thousand features are present, Scout with $p_1 = 1$ is not a viable option at present. However, Scout with $p_1 = 2$ is very fast, and we are confident that computational and algorithmic improvements will lead to increases in the number of features for which the Scout criteria can be maximized with $p_1 = 1$.

Roughly speaking, the method in this paper consists of two steps:

- (a) The features \mathbf{X} are used to obtain a regularized estimate of the inverse covariance matrix; this can be thought of as “pre-processing” the features.
- (b) The pre-processed features are combined with the outcome \mathbf{y} in order to obtain estimated regression coefficients.

In Step (a), the features are pre-processed without using the outcome \mathbf{y} . Indeed, many methods in the machine learning literature involve pre-processing the features without using the outcome. Principal components regression is a classical example of this; a more recent example with much more extensive pre-processing is in Hinton et al. (2006).

It has been shown that in order for the the lasso to exhibit model selection consistency, certain conditions on the feature matrix \mathbf{X} must be satisfied (see, for instance, the “irrepresentability condition” of Zhao & Yu (2006)). A reviewer asked whether Scout can offer a remedy in situations where these conditions are not satisfied. This is an interesting question that seems quite difficult to answer. We hope that it will be addressed in future work.

Covariance-regularized regression represents a new way to understand existing regularization methods for regression, as well as an approach to develop new regularization methods that appear to perform better in many examples.

7. Acknowledgments

We thank an editor and two reviewers for helpful comments. We thank Trevor Hastie for showing us the solution to the penalized log likelihood with an L_2 penalty. We thank both Trevor Hastie and Jerome Friedman for valuable discussions and for providing the code for the L_2 penalized multinomial and the elastic net. Daniela Witten was supported by a National Defense Science and Engineering Graduate Fellowship. Robert Tibshirani was partially supported by National Science Foundation Grant DMS-9971405 and National Institutes of Health Contract N01-HV-28183.

8. Appendix

8.1. Proofs of Claims

8.1.1. Proof of Claim 1

First, suppose that $p_2 = 1$. Consider the penalized log-likelihood

$$\log \det \Theta - \text{tr}(\mathbf{S}\Theta) - \frac{\lambda_2}{2} \|\Theta\|^1 \tag{31}$$

with $\Theta_{\mathbf{xx}}$, the top left $p \times p$ submatrix of Θ , fixed to equal the matrix that maximizes the log likelihood in Step 1 of the Scout Procedure. It is clear that if $\hat{\Theta}$ maximizes the log likelihood, then $(\hat{\Theta}^{-1})_{yy} = S_{yy} + \frac{\lambda_2}{2}$. The subgradient equation for maximization of the remaining portion of the log-likelihood is

$$0 = (\Theta^{-1})_{\mathbf{xy}} - \mathbf{S}_{\mathbf{xy}} - \frac{\lambda_2}{2}\Gamma \quad (32)$$

where $\Gamma_i = 1$ if the i^{th} element of $\Theta_{\mathbf{xy}}$ is positive, $\Gamma_i = -1$ if the i^{th} element of $\Theta_{\mathbf{xy}}$ is negative, and otherwise Γ_i is between -1 and 1 .

Let $\beta = \Theta_{\mathbf{xx}}(\Theta^{-1})_{\mathbf{xy}}$. Therefore, we equivalently wish to find β that solves

$$0 = 2(\Theta_{\mathbf{xx}})^{-1}\beta - 2\mathbf{S}_{\mathbf{xy}} - \lambda_2\Gamma. \quad (33)$$

From the partitioned inverse formula, it is clear that $\text{sgn}(\beta) = -\text{sgn}(\Theta_{\mathbf{xy}})$. Therefore, our task is equivalent to finding β which minimizes

$$\beta^T(\Theta_{\mathbf{xx}})^{-1}\beta - 2\mathbf{S}_{\mathbf{xy}}^T\beta + \lambda_2\|\beta\|^1. \quad (34)$$

Of course, this is Equation 11. It is an L_1 -penalized regression of \mathbf{y} onto \mathbf{X} , using only the inner products, with $\mathbf{S}_{\mathbf{xx}}$ replaced with $(\Theta_{\mathbf{xx}})^{-1}$. In other words, $\hat{\beta}$ that solves Equation 11 is given by $\Theta_{\mathbf{xx}}(\Theta^{-1})_{\mathbf{xy}}$, where Θ solves Step 2 of the Scout Procedure.

Now, the solution to Step 3 of the Scout Procedure is $-\frac{\Theta_{\mathbf{xy}}}{\Theta_{\mathbf{yy}}}$. By the partitioned inverse formula, $\Theta_{\mathbf{xx}}(\Theta^{-1})_{\mathbf{xy}} + \Theta_{\mathbf{xy}}(\Theta^{-1})_{\mathbf{yy}} = 0$, so $-\frac{\Theta_{\mathbf{xy}}}{\Theta_{\mathbf{yy}}} = \frac{\Theta_{\mathbf{xx}}(\Theta^{-1})_{\mathbf{xy}}}{(\Theta^{-1})_{\mathbf{yy}}\Theta_{\mathbf{yy}}} = \frac{\beta}{(\Theta^{-1})_{\mathbf{yy}}\Theta_{\mathbf{yy}}}$. In other words, the solution to Step 3 of the Scout Procedure and the solution to Equation 11 differ by a factor of $(\Theta^{-1})_{\mathbf{yy}}\Theta_{\mathbf{yy}}$. Since Step 4 of the Scout Procedure involves scaling the solution to Step 3 by a constant, it is clear that one can replace Step 3 of the Scout Procedure with the solution to Equation 11.

Now, suppose $p_2 = 2$. To find $\Theta_{\mathbf{xy}}$ that maximizes this penalized log-likelihood, we take the gradient and set it to zero:

$$0 = (\Theta^{-1})_{\mathbf{xy}} - \mathbf{S}_{\mathbf{xy}} - \frac{\lambda_2}{2}\Theta_{\mathbf{xy}} \quad (35)$$

Again, let $\beta = \Theta_{\mathbf{xx}}(\Theta^{-1})_{\mathbf{xy}}$. Therefore, we equivalently wish to find β that solves

$$0 = 2(\Theta_{\mathbf{xx}})^{-1}\beta - 2\mathbf{S}_{\mathbf{xy}} + 2\lambda_3\beta \quad (36)$$

for some new constant λ_3 , using the fact, from the partitioned inverse formula, that $-\frac{\beta}{(\Theta^{-1})_{\mathbf{yy}}} = \Theta_{\mathbf{xy}}$. The solution β minimizes

$$\beta^T(\Theta_{\mathbf{xx}})^{-1}\beta - 2\mathbf{S}_{\mathbf{xy}}^T\beta + \lambda_3\beta^T\beta.$$

Of course, this is again Equation 11. Therefore, $\hat{\beta}$ that solves Equation 11 is given (up to scaling by a constant) by $\Theta_{\mathbf{xx}}(\Theta^{-1})_{\mathbf{xy}}$, where Θ solves Step 2 of the Scout Procedure. As before, by the partitioned inverse formula, and since Step 4 of the Scout Procedure involves scaling the solution to Step 3 by a constant, it is clear that one can replace Step 3 of the Scout Procedure with the solution to Equation 11.

8.1.2. Proof of Claim 2

If $\hat{\beta}$ minimizes Equation 14, then since $\hat{\beta}_i\hat{\beta}_j \neq 0$, it follows that

$$\frac{\lambda_2}{2}(\text{sgn}(\hat{\beta}_i) - \text{sgn}(\hat{\beta}_j)) + \sqrt{2\lambda_1}(\hat{\beta}_i - \hat{\beta}_j) = (\mathbf{x}_i^* - \mathbf{x}_j^*)^T(\mathbf{y}^* - \mathbf{X}^*\hat{\beta}), \quad (37)$$

and hence that

$$\sqrt{2\lambda_1}|\hat{\beta}_i - \hat{\beta}_j| \leq |(\mathbf{x}_i^* - \mathbf{x}_j^*)^T(\mathbf{y}^* - \mathbf{X}^*\hat{\beta})|. \quad (38)$$

Note that

$$\|\mathbf{y}^* - \mathbf{X}^* \hat{\beta}\|^2 \leq \|\mathbf{y}^* - \mathbf{X}^* \beta\|^2 + \lambda_2 \|\hat{\beta}\|^1 + \sqrt{2\lambda_1} \|\hat{\beta}\|^2 \leq \|\mathbf{y}^*\|^2 = 2\|\mathbf{y}\|^2. \quad (39)$$

Therefore,

$$|\hat{\beta}_i - \hat{\beta}_j| \leq \sqrt{\frac{1}{2\lambda_1}} \|\mathbf{x}_i^* - \mathbf{x}_j^*\| \|\mathbf{y}\| \sqrt{2}. \quad (40)$$

Now, $\|\mathbf{x}_i^* - \mathbf{x}_j^*\|^2 = \frac{1}{2} \|\mathbf{x}_i - \mathbf{x}_j\|^2 + \frac{1}{2} \|(\bar{\mathbf{D}}\mathbf{V}^T)_i - (\bar{\mathbf{D}}\mathbf{V}^T)_j\|^2$. Since we assumed that the features are standardized, it follows that $\|\mathbf{x}_i^* - \mathbf{x}_j^*\|^2 = 1 - \rho + \frac{1}{2} \|(\bar{\mathbf{D}}\mathbf{V}^T)_i - (\bar{\mathbf{D}}\mathbf{V}^T)_j\|^2$ where ρ is the correlation between \mathbf{x}_i and \mathbf{x}_j . It also is easy to see that $\|(\bar{\mathbf{D}}\mathbf{V}^T)_i - (\bar{\mathbf{D}}\mathbf{V}^T)_j\|^2 \leq 1 - \rho$. Therefore, it follows that

$$|\hat{\beta}_i - \hat{\beta}_j| \leq \sqrt{\frac{2(1-\rho)}{\lambda_1}} \|\mathbf{y}\|. \quad (41)$$

8.1.3. Proof of Claim 3

Consider the latent variable model given in Section 2.6; note that under this model,

$$\mathbf{y} = \mathbf{X}\beta + \epsilon, \quad (42)$$

where $\beta = \frac{1}{d_1} \mathbf{v}_1$. In addition,

$$\mathbf{X}^T \mathbf{X} = d_1^2 \mathbf{v}_1 \mathbf{v}_1^T + d_2^2 \mathbf{v}_2 \mathbf{v}_2^T = \sum_{j=1}^p d_j^2 \mathbf{v}_j \mathbf{v}_j^T \quad (43)$$

where $d_3 = \dots = d_p = 0$ and $\mathbf{v}_1, \dots, \mathbf{v}_p$ are orthonormal. We consider two options for the regression of \mathbf{y} onto \mathbf{X} : ridge regression and *Scout*(2, \cdot). Let $\hat{\beta}^{rr}$ and $\hat{\beta}^{sc}$ denote the resulting estimates, and let λ^{rr} and λ^{sc} be the tuning parameters of the two methods, respectively.

$$\begin{aligned} \hat{\beta}^{rr} &= (\mathbf{X}^T \mathbf{X} + \lambda^{rr} \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \\ &= \left(\sum_{j=1}^p \frac{1}{d_j^2 + \lambda^{rr}} \mathbf{v}_j \mathbf{v}_j^T \right) (d_1 \mathbf{v}_1 \mathbf{u}_1^T + d_2 \mathbf{v}_2 \mathbf{u}_2^T) (\mathbf{u}_1 + \epsilon) \\ &= \frac{d_1}{d_1^2 + \lambda^{rr}} \mathbf{v}_1 + \left(\frac{d_1}{d_1^2 + \lambda^{rr}} \mathbf{v}_1 \mathbf{u}_1^T + \frac{d_2}{d_2^2 + \lambda^{rr}} \mathbf{v}_2 \mathbf{u}_2^T \right) \epsilon \end{aligned} \quad (44)$$

Similarly, the solution to *Scout*(2, \cdot) is as follows:

$$\begin{aligned} \hat{\beta}^{sc} &= \left(\sum_{j=1}^p \frac{2}{d_j^2 + \sqrt{d_j^4 + 8\lambda^{sc}}} \mathbf{v}_j \mathbf{v}_j^T \right) (d_1 \mathbf{v}_1 \mathbf{u}_1^T + d_2 \mathbf{v}_2 \mathbf{u}_2^T) (\mathbf{u}_1 + \epsilon) \\ &= \frac{2d_1}{d_1^2 + \sqrt{d_1^4 + 8\lambda^{sc}}} \mathbf{v}_1 + \left(\frac{2d_1}{d_1^2 + \sqrt{d_1^4 + 8\lambda^{sc}}} \mathbf{v}_1 \mathbf{u}_1^T + \frac{2d_2}{d_2^2 + \sqrt{d_2^4 + 8\lambda^{sc}}} \mathbf{v}_2 \mathbf{u}_2^T \right) \epsilon \end{aligned} \quad (45)$$

The biases of $\hat{\beta}^{rr}$ and $\hat{\beta}^{sc}$ are as follows:

$$\begin{aligned} E(\hat{\beta}^{rr} - \beta) &= \left(\frac{d_1}{d_1^2 + \lambda^{rr}} - \frac{1}{d_1} \right) \mathbf{v}_1 \\ E(\hat{\beta}^{sc} - \beta) &= \left(\frac{2d_1}{d_1^2 + \sqrt{d_1^4 + 8\lambda^{sc}}} - \frac{1}{d_1} \right) \mathbf{v}_1 \end{aligned} \quad (46)$$

and the variances are as follows:

$$\begin{aligned}\text{Var}(\hat{\beta}^{rr}) &= \left(\frac{d_1}{d_1^2 + \lambda^{rr}}\right)^2 \mathbf{v}_1 \mathbf{v}_1^T \sigma^2 + \left(\frac{d_2}{d_2^2 + \lambda^{rr}}\right)^2 \mathbf{v}_2 \mathbf{v}_2^T \sigma^2 \\ \text{Var}(\hat{\beta}^{sc}) &= \left(\frac{2d_1}{d_1^2 + \sqrt{d_1^4 + 8\lambda^{sc}}}\right)^2 \mathbf{v}_1 \mathbf{v}_1^T \sigma^2 + \left(\frac{2d_2}{d_2^2 + \sqrt{d_2^4 + 8\lambda^{sc}}}\right)^2 \mathbf{v}_2 \mathbf{v}_2^T \sigma^2\end{aligned}\quad (47)$$

The following relationship between λ^{rr} and λ^{sc} results in equal biases:

$$\lambda^{rr} = \frac{-d_1^2 + \sqrt{d_1^4 + 8\lambda^{sc}}}{2} \quad (48)$$

From now on, we assume that Equation 48 holds. Then, if $d_1 > d_2$, it follows that $\text{Var}(\hat{\beta}^{sc}) < \text{Var}(\hat{\beta}^{rr})$. In other words, if the portion of \mathbf{X} that is correlated with \mathbf{y} has a stronger signal than the portion that is orthogonal to \mathbf{y} , then (for a given amount of bias) *Scout*(2, \cdot) will have lower variance than ridge.

8.2. Feature Selection for Scout LDA

The method that we propose in Section 4.1 can be easily modified in order to perform built-in feature selection. Using the notation in Section 4.1, we observe that

$$\hat{\mu}_k = \arg \min_{\mu_k} \left\{ \sum_{i:C(i)=k} (X_i - \mu_k)^T \hat{\Sigma}_{\mathbf{w}\mathbf{c},\lambda}^{-1} (X_i - \mu_k) \right\} \quad (49)$$

and so we replace $\hat{\mu}_k$ in Equation 25 with

$$\hat{\mu}_k^{\lambda,\rho} = \arg \min_{\mu_k} \left\{ \sum_{i:C(i)=k} (X_i - \mu_k)^T \hat{\Sigma}_{\mathbf{w}\mathbf{c},\lambda}^{-1} (X_i - \mu_k) + \rho \|\hat{\Sigma}_{\mathbf{w}\mathbf{c},\lambda}^{-1} \mu_k\|^1 \right\}. \quad (50)$$

The above can be solved via an L_1 regression, and it gives the following classification rule for a test observation X :

$$\delta_k^{\lambda,\rho}(X) = X^T \hat{\Sigma}_{\mathbf{w}\mathbf{c},\lambda}^{-1} \hat{\mu}_k^{\lambda,\rho} - \frac{1}{2} (\hat{\mu}_k^{\lambda,\rho})^T \hat{\Sigma}_{\mathbf{w}\mathbf{c},\lambda}^{-1} \hat{\mu}_k^{\lambda,\rho} + \log \pi_k \quad (51)$$

References

- Bair, E. & Tibshirani, R. (2004), ‘Semi-supervised methods to predict patient survival from gene expression data’, *PLOS Biology* **2**, 511–522.
- Banerjee, O., El Ghaoui, L. E. & d’Aspremont, A. (2008), ‘Model selection through sparse maximum likelihood estimation for multivariate gaussian or binary data’, *Journal of Machine Learning Research* **9**, 485–516.
- Bickel, P. & Levina, E. (2008), ‘Covariance regularization by thresholding’, *Annals of Statistics* **36**, 2577–2604.
- Breiman, L. (2001), ‘Random forests’, *Machine Learning* **45**, 5–32.
- Dey, D. & Srinivasan, C. (1985), ‘Estimation of a covariance matrix under Stein’s loss’, *Annals of Statistics* **13**, 1581–1591.
- Frank, I. & Friedman, J. (1993), ‘A statistical view of some chemometrics regression tools (with discussion)’, *Technometrics* **35**(2), 109–148.

- Friedman, J. (1989), 'Regularized discriminant analysis', *Journal of the American Statistical Association* **84**, 165–175.
- Friedman, J., Hastie, T. & Tibshirani, R. (2007), 'Sparse inverse covariance estimation with the graphical lasso', *Biostatistics* **9**, 432–441.
- Friedman, J., Hastie, T. & Tibshirani, R. (2008), 'Regularization paths for generalized linear models via coordinate descent'.
- Green, P. J. (1984), 'Iteratively reweighted least squares for maximum likelihood estimation, and some robust and resistant alternatives', *Journal of the Royal Statistical Society, Series B* **46**, 149–192.
- Guo, Y., Hastie, T. & Tibshirani, R. (2007), 'Regularized linear discriminant analysis and its application in microarrays', *Biostatistics* **8**, 86–100.
- Haff, L. (1979), 'Estimation of the inverse covariance matrix: random mixtures of the inverse Wishart matrix and the identity', *Annals of Statistics* **7**, 1264–1276.
- Hinton, G., Osindero, S. & Teh, Y. (2006), 'A fast learning algorithm for deep belief nets', *Neural Computation* **18**, 1527–1553.
- Hoerl, A. E. & Kennard, R. (1970), 'Ridge regression: Biased estimation for nonorthogonal problems', *Technometrics* **12**, 55–67.
- Hummel, M., Bentink, S., Berger, H., Klappwe, W., Wessendorf, S., Barth, F. T. E., Bernd, H.-W., Cogliatti, S. B., Dierlamm, J., Feller, A. C., Hansmann, M. L., Haralambieva, E., Harder, L., Hasenclever, D., Kuhn, M., Lenze, D., Lichter, P., Martin-Subero, J. I., Moller, P., Muller-Hermelink, H.-K., Ott, G., Parwaresch, R. M., Pott, C., Rosenwald, A., Rosolowski, M., Schwaenen, C., Sturzenhockecker, B., Szczepanowski, M., Trautmann, H., Wacker, H.-H., Spang, R., Loeffler, M., Trumper, L., Stein, H. & Siebert, R. (2006), 'A biological definition of Burkitt's lymphoma from transcriptional and genomic profiling', *New England Journal of Medicine* **354**, 2419–2430.
- James, W. & Stein, C. (1961), 'Estimation with quadratic loss', *Proceedings of the Fourth Berkeley Symposium on Mathematics and Statistical Probability* **1**, 361–379.
- Kalbfleisch, J. & Prentice, R. (1980), *The statistical analysis of failure time data*, Wiley, New York.
- Liang, F., Mukherjee, S. & West, M. (2007), 'The use of unlabeled data in predictive modeling', *Statistical Science* **22**, 189–205.
- Mardia, K., Kent, J. & Bibby, J. (1979), *Multivariate Analysis*, Academic Press.
- McLachlan, G. J. (1992), *Discriminant Analysis and Statistical Pattern Recognition*, Wiley, New York.
- Meinshausen, N. & Bühlmann, P. (2006), 'High dimensional graphs and variable selection with the lasso', *Annals of Statistics* **34**, 1436–1462.
- Monti, S., Savage, K. J., Kutok, J. L., Feuerhake, F., Kurtin, P., Mihm, M., Wu, B., Pasqualucci, L., Neuberg, D., Aguiar, R. C. T., Dal Cin, P., Ladd, C., Pinkus, G. S., Salles, G., Harris, N. L., Dalla-Favera, R., Habermann, T. M., Aster, J. C., Golub, T. R. & Shipp, M. A. (2005), 'Molecular profiling of diffuse large B-cell lymphoma identifies robust subtypes including one characterized by host inflammatory response', *Blood* **105**, 1851–1861.

- O'Neill, T. (1978), 'Normal discrimination with unclassified observations', *Journal of the American Statistical Association* **73**, 821–826.
- Park, M. Y. & Hastie, T. (2007), 'An L_1 regularization path algorithm for generalized linear models', *Journal of the Royal Statistical Society Series B* **69**(4), 659–677.
- Ramaswamy, S., Tamayo, P., Rifkin, R., Mukherjee, S., Yeang, C., Angelo, M., Ladd, C., Reich, M., Latulippe, E., Mesirov, J., Poggio, T., Gerald, W., Loda, M., Lander, E. & Golub, T. (2001), 'Multiclass cancer diagnosis using tumor gene expression signature', *PNAS* **98**, 15149–15154.
- Rosenwald, A., Wright, G., Chan, W. C., Connors, J. M., Campo, E., Fisher, R. I., Gascoyne, R. D., Muller-Hermelink, H. K., Smeland, E. B. & Staudt, L. M. (2002), 'The use of molecular profiling to predict survival after chemotherapy for diffuse large B-cell lymphoma', *The New England Journal of Medicine* **346**, 1937–1947.
- Rothman, A., Levina, E. & Zhu, J. (2008), 'Sparse permutation invariant covariance estimation', *Electronic Journal of Statistics* **2**, 494–515.
- Shipp, M. A., Ross, K. N., Tamayo, P., Weng, A. P., Kutok, J. L., Aguiar, R. C., Gaasenbeek, M., Angelo, M., Reich, M., Pinkus, G. S., Ray, T. S., Koval, M. A., Last, K. W., Norton, A., Lister, T. A., Mesirov, J., Neuberg, D. S., Lander, E. S., Aster, J. C. & Golub, T. R. (2002), 'Diffuse large B-cell lymphoma outcome prediction by gene-expression profiling and supervised machine learning', *Nature Medicine* **8**, 68–74.
- Tibshirani, R. (1996), 'Regression shrinkage and selection via the lasso', *J. Royal. Statist. Soc. B.* **58**, 267–288.
- Tibshirani, R., Hastie, T., Narasimhan, B. & Chu, G. (2002), 'Diagnosis of multiple cancer types by shrunken centroids of gene expression', *Proc. Natl. Acad. Sci.* **99**, 6567–6572.
- Tibshirani, R., Hastie, T., Narasimhan, B. & Chu, G. (2003), 'Class prediction by nearest shrunken centroids, with applications to DNA microarrays', *Statistical Science* **18**, 104–117.
- Tusher, V. G., Tibshirani, R. & Chu, G. (2001), 'Significance analysis of microarrays applied to the ionizing radiation response', *Proceedings of the National Academy of Sciences of the United States of America* **98**, 5116–5121.
- Zhao, P. & Yu, B. (2006), 'On model selection consistency of lasso', *Journal of Machine Learning Research* **7**, 2541–2563.
- Zhu, J. & Hastie, T. (2004), 'Classification of gene microarrays by penalized logistic regression', *Biostatistics* **5**(2), 427–443.
- Zou, H. & Hastie, T. (2005), 'Regularization and variable selection via the elastic net', *J. Royal. Stat. Soc. B.* **67**, 301–320.