

Variable Selection Techniques

Gabriel Ackall and Connor Shrader

June 9, 2021

In mathematical modeling, especially when using linear regression, it is often important to reduce the number of predictors, or variables, used to predict an output. This can help to increase the interpretability of a model and reduce its error. Below are five methods used to select a linear model using a subset of the predictors.

1 Introduction

Linear regression is one of the simplest forms of statistical learning models. It assumes that some predictor variable y can be modeled as a linear combination of a set of predictor variables x_1, x_2, \dots, x_p . More precisely, a linear regression model assumes that

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \epsilon \quad (1)$$

where $\beta_0, \beta_1, \dots, \beta_p$ are coefficients and ϵ is a random error with mean zero and variance σ^2 . We also assume that this error has no correlation between different observations. When working with real data sets, the coefficient values are usually unknown; the goal of linear regression is to compute coefficient estimates $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$ that are good estimates for the actual coefficients.

Let $\beta = [\beta_0, \beta_1, \dots, \beta_p]^T$ be a vector of coefficient values, and let $\hat{\beta} = [\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p]^T$ represent a vector of coefficient estimates. Also, let \mathbf{X} be a $n \times (p + 1)$ matrix containing n observations in p variables. \mathbf{X} also contains a column with every entry equal to 1; the corresponding coefficient estimate will represent $\hat{\beta}_0$. Finally, let \mathbf{y} represent the vector of response values for each observation.

The most common way to compute $\hat{\beta}$ is with ordinary least squares. This method computes $\hat{\beta}$ by minimizing the residual sum of squares:

$$\hat{\beta}^{\text{OLS}} = \arg \min_{\beta} \left\{ \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right\} \quad (2)$$

$$= \arg \min_{\beta} \{ (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) \} \quad (3)$$

Ordinary least squares is favored for a few reasons. For one, the solution is very easy to compute; it is given by

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (4)$$

Another advantage of ordinary least squares comes from the Gauss-Markov Theorem. This theorem states that of all the linear unbiased estimators, ordinary least squares has the lowest variance. Alternative unbiased estimators will necessarily have a larger variance, which can lead to more inaccurate predictions.

Why, then, would we use anything other than ordinary least squares? One issue is that the matrix $\mathbf{X}^T \mathbf{X}$ is not always invertible. In fact, this matrix is never invertible in the case when $p > n$. Hence, in some situations, we cannot use ordinary least squares. Even if the matrix is invertible, other methods may be

preferable. Despite having no bias in ideal situations, ordinary least squares may still have high variance. Consequently, it may be worthwhile to sacrifice some bias in order to decrease variance.

The rest of this document describes numerous ways to modify ordinary least squares to obtain models with lower variance. Some of these algorithms also use variable selection; these methods set some coefficient estimates to zero if they do not have a strong influence on model accuracy. This not only helps reduce variance, but it also creates more interpretable models.

2 Subset Selection

First, we discuss subset selection techniques. In general, these techniques are discrete algorithms that attempt to find a good model using only a subset of the available predictors.

2.1 Best Subset Selection

Best subset selection is a method for selecting the most influential predictors that minimize error in a least squares linear regression. It does this by fitting a linear regression model to every possible combination of predictors and then choosing the best model of all the possible combinations. The best model is determined through the use of a test error estimate. The most common examples of test error estimation indicators are Akaike information criterion (AIC), Bayesian information criterion (BIC), Adjusted R^2 , and cross validation.

While best subset selection results in the best possible model given the predictors, it is very computationally expensive. As the number of predictors increases, the number of linear models that best subset selection has to fit increases exponentially. This can be seen in Table 1. Thus, for models with more than 40 predictors, this can become infeasible for most computers to compute [6]. Given that in many scenarios, especially those seen in medicine with genomic data or in scenarios where there are more predictors than data samples, there can be many thousands of predictors, and best subset selection becomes impossible.

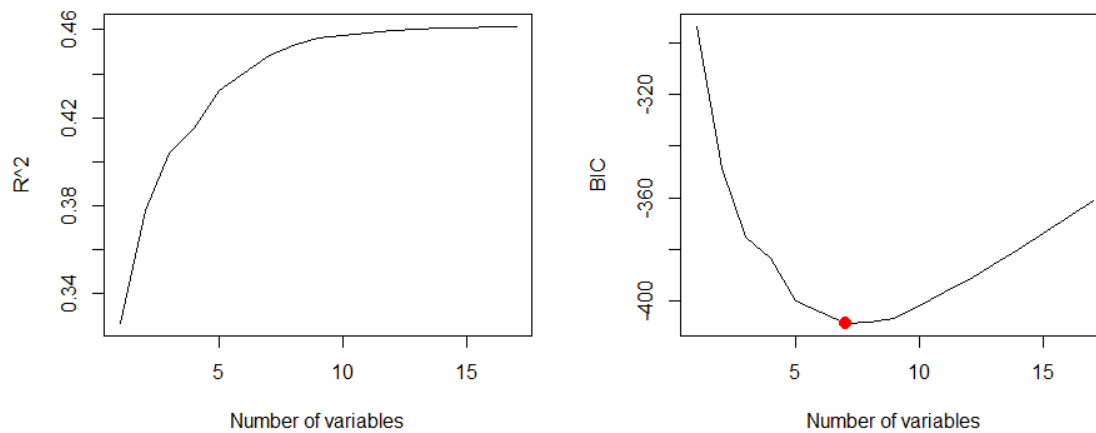
Table 1: Number of fitted models depending on number of predictors (p)

p	Fitted Models
2	$2^2 = 4$
10	$2^{10} = 1024$
100	$2^{100} > 10^{30}$
k	2^k

Figure 1 below demonstrates how the number of predictors can affect R^2 and BIC when using best subset selection. This plot was created by using the `leaps` library (which provides a function to run best subset selection) [1]. We used the `College` dataset provided by the `ISLR` library [5], and fit linear regression models using `Grad.Rate` as the response. We see that as the number of predictors increases, R^2 always increases (as expected). On the other hand, BIC is minimized with a moderate number of variables (between seven and nine). According to the BIC statistic, the best model has seven variables. The code used for this figure is in the `r` folder of the REU GitHub repository.

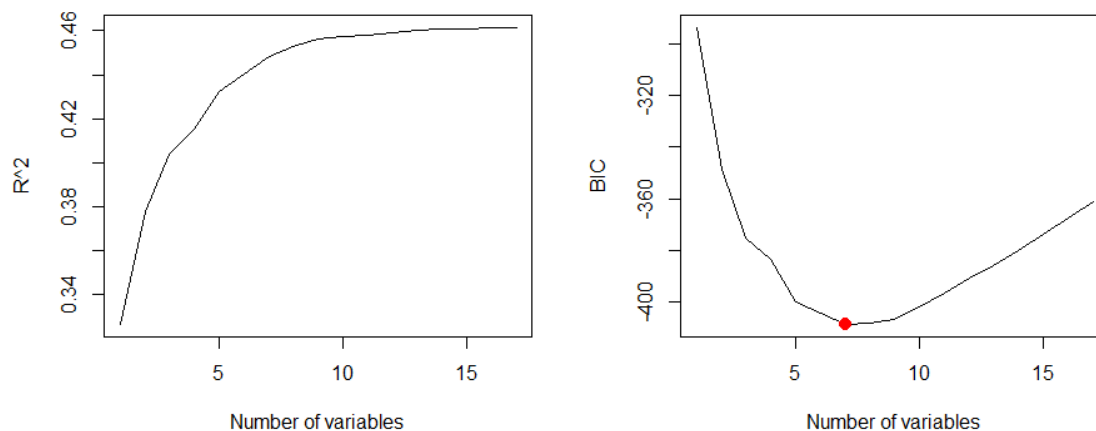
2.2 Forward Stepwise Selection

Forward stepwise selection aims to approximate the best combination of predictors in a linear regression model, but with a more computationally efficient method than best subset selection. Forward stepwise selection starts without using any predictors. It then slowly begins adding the most important predictor to the model. The predictor is chosen to minimize statistics such as p-value, AIC, BIC, or Adjusted R^2 , to name a few. This is repeated until a stopping point is reached which can be defined by p-value, AIC, BIC, and more.

Figure 1: R^2 and BIC when applying best subset selection

This process is much more computationally efficient than best subset selection, but it does not necessarily result in the best combination of parameters in the linear regression and is not guaranteed to result in the best model.

Figure 2 shows the R^2 and BIC statistics when fitting models using forward stepwise selection. Again, we predicted `Grad.Rate` using the `College` data set using the `leaps` library. The results are almost identical to what we saw for best subset selection. Even though the plots are similar, the specific model chosen by forward stepwise selection is actually different than the one found using best subset selection.

Figure 2: R^2 and BIC when applying forward stepwise selection

2.3 Backward Stepwise Selection

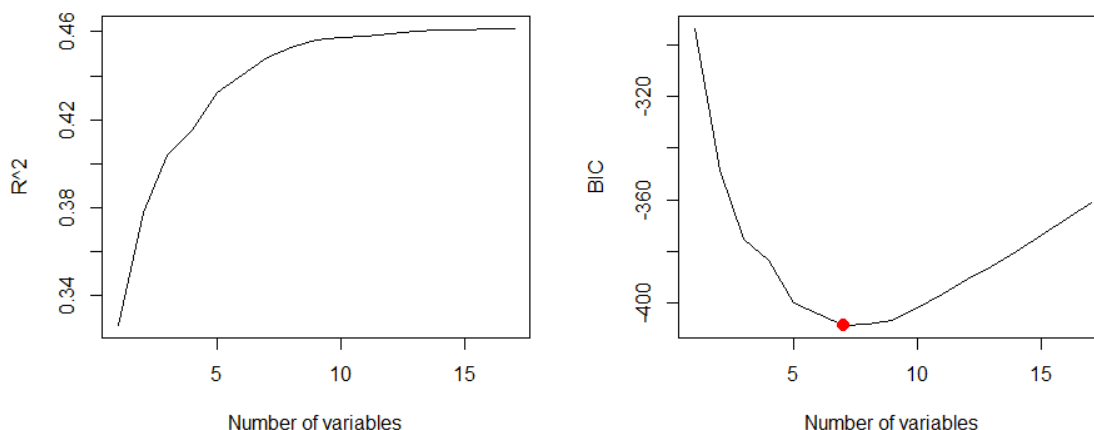
Backwards stepwise selection works very similarly to forward stepwise selection, except that it starts with every single predictor included in the least squares linear regression. Instead of adding predictors like

in forward stepwise selection, backward stepwise selection removes the least important predictor in each iteration. Similar to the forward method, the importance of a predictor can be determined by its p-value, AIC, BIC, or Adjusted R^2 . This is repeated until a pre-determined stopping point is reached.

Backward stepwise selection can often result in better models than forward stepwise selection because it is guaranteed to test all the predictors together. This is different from forward stepwise selection that can sometimes suppress predictors, especially those that are collinear. For these reasons, when its use is possible, backward stepwise selection is preferred to forward stepwise selection. However, in cases where the number of predictors are greater than the number of samples, backward stepwise selection is impossible. In these case, forward stepwise selection must be used.

Figure 3 shows R^2 and BIC after applying backward stepwise selection to the `College` data set. Again, the results are very similar to Figures 1 and 2, but the particular models chosen by the algorithm were slightly different.

Figure 3: R^2 and BIC when applying backward stepwise selection



2.4 Hybrid Stepwise Selection

One weakness of the forward stepwise and backward stepwise methods is that they are greedy algorithms; in general, they will not find the best model for a given number of predictors. One way to improve model accuracy is to use hybrid stepwise selection, which allows for both forward steps and backward steps [4].

The algorithm could start with either zero predictors or all predictors. In each iteration, the method would either add a new predictor to the model or remove a predictor that does not increase performance. Like the forward and backward stepwise selection methods, this algorithm terminates when the model cannot be improved further; measuring the accuracy of the model can be determined using the AIC or BIC.

Although this strategy is slightly more computationally expensive than forward stepwise or backward stepwise selection, a hybrid approach may improve model results while still avoiding the unrealistic runtime of best subset selection.

2.5 Forward Stagewise Selection

One last method for feature selection is called forward stagewise regression. Like forward stepwise selection, forward stagewise selection starts by fitting a model using none of the predictors. In each iteration, the

method chooses the predictor most closely correlated to the residuals of the current model, and fits a simple linear regression using the predictor against the residuals. The coefficient for this predictor in the simple model is then added to the corresponding coefficient in the other model. This process is repeated until none of the predictors are correlated with the residuals.

Note that in each iteration of this algorithm, only one of the coefficients is changed. As a result, this method has a long runtime. In the long run, forward stagewise selection is still competitive compared to the strategies previously discussed.

3 Penalized Regression

The following sections discuss several modifications to the ordinary least squares model that penalize large coefficient estimates. Recall that an ordinary least squares model assumes that some response variable y can be modeled by

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_p x_p + \epsilon \quad (5)$$

where x_1, x_2, \dots, x_p are predictors, $\beta_0, \beta_1, \dots, \beta_p$ are coefficients, and ϵ is some random error with mean 0 and variance σ^2 . An ordinary least squares model estimates the coefficients as

$$\hat{\beta}^{\text{OLS}} = \arg \min_{\beta} \left\{ \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right\} \quad (6)$$

In general, a penalized regression imposes some additional restriction that punishes large coefficient estimates. The coefficient estimates are then given by

$$\hat{\beta} = \arg \min_{\beta} \left\{ \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \sum_{j=1}^p p(\beta_j) \right\} \quad (7)$$

where $p(\beta_j)$ is some penalty function. Generally, this function should be large when the coefficient is large so that a model with smaller coefficients is favored. Penalized regression is especially useful when $p > n$ because the penalties can reduce variance and perform variable selection at the cost of some bias.

Because these models place penalties on large coefficient estimates, it is common to standardize the predictors to have a mean of 0 and a variance of 1. This ensures that none of the predictors have a disproportionate effect on the estimated coefficients.

3.1 Ridge Regression

Ridge regression helps to solve multicollinearity in predictors while also minimizing insignificant predictors. While it does not minimize these insignificant predictors completely to 0 and thus cannot be considered a variable selection method, it still proves very useful in large datasets.

Ridge regression works by minimizing Residual Sum Squared (RSS) plus a penalty as seen in Equation 8. λ is a tuning parameter and can be used to determine how much of an effect the penalty has on the regression. if $\lambda = 0$, then the regression acts exactly like ordinary least squares regression, but if $\lambda \rightarrow \infty$, then $\beta_j \rightarrow 0$ and the regression line will be a horizontal line at the intercept, β_0 .

$$\hat{\beta}^{\text{ridge}} = \arg \min_{\beta} \left\{ \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\} \quad (8)$$

An alternative way to express ridge regression is with the equation

$$\hat{\beta}^{\text{ridge}} = \arg \min_{\beta} \left\{ \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right\} \quad \text{subject to} \quad \sum_{j=1}^p \beta_j^2 \leq t \quad (9)$$

for some tuning parameter t .

3.2 Lasso Regression

The least absolute shrinkage and selection operation, often referred to as *lasso*, is a shrinkage method with a very similar form to lasso regression [7, 5, 6]. The coefficient estimates satisfy

$$\hat{\beta}^{\text{lasso}} = \arg \min_{\beta} \left\{ \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right) + \lambda \sum_{j=1}^p |\beta_j| \right\} \quad (10)$$

If $\lambda = 0$, then the lasso model is equivalent to the ordinary least squares model; if $\lambda \rightarrow \infty$, then the coefficients for all predictors will be set to 0. An equivalently way to define lasso regression is by

$$\hat{\beta}^{\text{ridge}} = \arg \min_{\beta} \left\{ \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \right\} \quad \text{subject to} \quad \sum_{j=1}^p |\beta_j| \leq t \quad (11)$$

where t is a tuning parameter.

One useful property of the lasso method is that it can perform variable selection by setting some coefficients to zero. To understand why lasso regression can perform variable selection whereas ridge regression cannot, consider Figure 4 below. This figure demonstrates the case when $p = 2$ and $t = 1$. The red diamond on the left represents the condition $|\beta_1| + |\beta_2| < 1$ for ridge regression, while the circle on the right represents the condition $\beta_1^2 + \beta_2^2 < 1$ for lasso regression. The black curves represent contours of the residual sum of squares error for values of β_1 and β_2 . The black point in the center of these curves is where the RSS is minimized; this represents the values of β_1 and β_2 that would be selected by ordinary least squares.

In the left plot, the intersection of the black curve and the red diamond represents the parameter values chosen by lasso regression; this point minimizes the RSS under the condition $|\beta_1| + |\beta_2| < 1$. Because the red region is a square, this intersection occurs at a point where $\beta_1 = 0$; hence, the lasso method removes the predictor β_1 . On the other hand, the circular shape of the constrained region for ridge regression cannot perform variable selection because the intersection does not occur at one of the axes.

The lasso method is particularly useful in the case where $p > n$ because of its ability to select variables; a model with fewer variables has less variance and is more interpretable. One major downside of lasso regression is that it does not handle multicollinearity as nicely as ridge regression. Another downside of lasso regression is that it does not have a closed-form solution, which can lead to instability in the model. This can be demonstrated by the greater spread of β coefficients compared to ridge regression, as visible in Figure 5.

3.3 Elastic Net Regression

Elastic net regression serves as a combination between ridge and lasso regression. It can handle multicollinearity as well as perform variable selection. The coefficients for elastic net regression can be determined

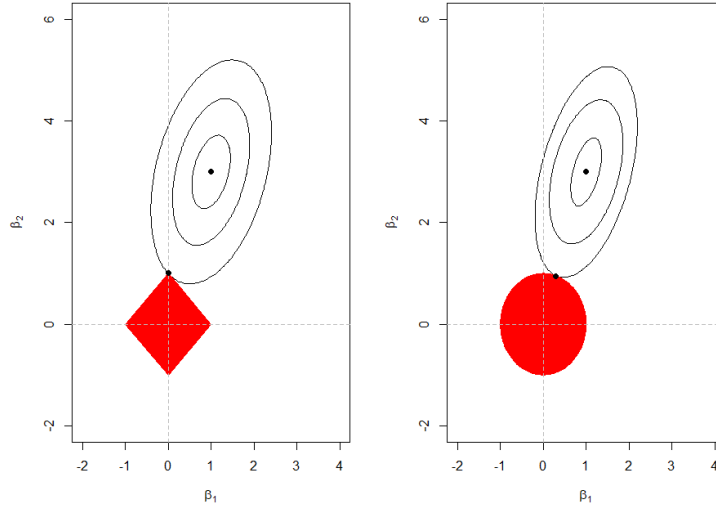
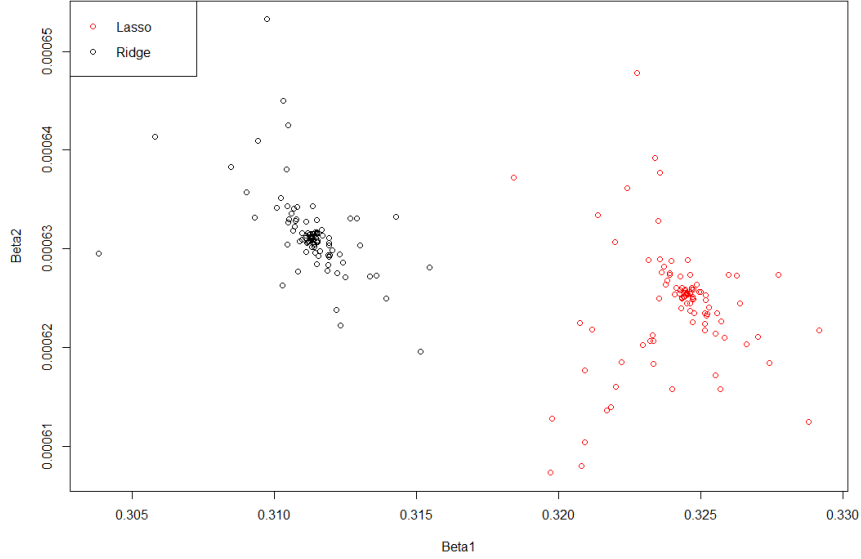
Figure 4: Error and constraint curves for the lasso and ridge models when $p = 2$.

Figure 5: Instability of Lasso and Ridge Regression to Changes in Training Data



by

$$\hat{\beta}^{\text{ENet}} = \arg \min_{\beta} \left\{ \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda_2 \sum_{j=1}^p \beta_j^2 + \lambda_1 \sum_{j=1}^p |\beta_j| \right\} \quad (12)$$

where λ_1 and λ_2 are both tuning parameters to be determined later.

An important limitation to note is that elastic net performs best when it is close to either ridge or lasso regression, meaning that either $\lambda_1 \gg \lambda_2$ or vice versa [10]. Additionally, because elastic net requires two tuning parameters, this makes it much more difficult to determine the best combination of tuning parameters to minimize error in the regression. However, this problem has been largely solved through by the LARS-EN

algorithm developed by Zou et. al. which efficiently solves for the tuning parameters.

3.4 Adaptive Lasso Regression

Normally in lasso regression, each predictor is weighted the same in the penalty function. Adaptive lasso regression is different in that a weight, \hat{w}_j is multiplied to the penalty function. The coefficients for adaptive lasso regression as designed by Zou et. al. [9] can be defined by

$$\hat{\beta}^{\text{adaptive}} = \arg \min_{\beta} \left\{ \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right) + \lambda \sum_{j=1}^p \hat{w}_j |\beta_j| \right\} \quad (13)$$

where λ is a tuning parameter to be determined later and \hat{w}_j is defined as $\frac{1}{|\hat{\beta}|^\gamma}$ with γ being a chosen parameter greater than 0.

Because of the weight that is implemented in adaptive lasso regression, zero-coefficients have a weight that is inflated up to infinity, and thus are punished much more harshly than large coefficients whose weight is much smaller in comparison. This is a similar rationale to SCAD and helps to reduce some of the bias from lasso regression. Bridge regression is the general form of lasso regression from which adaptive lasso originates from. When $\gamma < 1$, bridge regression as shown in Equation 14 is not continuous, which results in model prediction instability. However, adaptive lasso regression is completely continuous and thus has much more consistent coefficients when fitted.

$$\hat{\beta}^{\text{lasso}} = \arg \min_{\beta} \left\{ \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right) + \lambda \sum_{j=1}^p |\beta_j|^\gamma \right\} \quad (14)$$

3.5 Smoothly Clipped Absolute Deviation Regression

One major flaw of the lasso method is that the penalty punishes large coefficients, even if those coefficients should be large. One way to modify the lasso method is to use the *smoothly clipped absolute deviation* (SCAD) penalty [3]. The goal of this method is to punish large coefficients less severely, which can help mitigate some of the bias introduced by the lasso method.

$$\hat{\beta}^{\text{SCAD}} = \arg \min_{\beta} \left\{ \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right) + \lambda \sum_{j=1}^p J_a(\beta_j, \lambda) \right\} \quad (15)$$

Here, $J_a(\beta, \lambda)$ is a penalty function that satisfies

$$\frac{dJ_a(\beta, \lambda)}{d\beta} = \lambda \cdot \text{sign}(\beta) \left[I(|\beta| < \lambda) + \frac{(a\lambda - |\beta|)_+}{(a-1)\lambda} I(|\beta| > \lambda) \right] \quad (16)$$

where $\lambda \geq 0$ and $a \geq 2$ are tuning parameters. An equivalent way to write this is

$$\frac{dJ_a(\beta, \lambda)}{d\beta} = \begin{cases} \lambda, & |\beta| \leq \lambda \\ \frac{a\lambda - |\beta|}{a-1}, & \lambda < |\beta| < a\lambda \\ 0, & a\lambda < |\beta| \end{cases} \quad (17)$$

This penalty function does not punish coefficients with large magnitude as heavily as the lasso method. In fact, if the magnitude of a coefficient is larger than $a\lambda$, then the penalty becomes constant. See Figure 6a for a plot of the SCAD penalty as a function of the coefficient value.

Integrating with respect to β [2], we see that

$$J_a(\beta, \lambda) = \begin{cases} \lambda|\beta|, & |\beta| \leq \lambda \\ \frac{2a\lambda|\beta| - \beta^2 - \lambda^2}{2(a-1)}, & \lambda < |\beta| < a\lambda \\ \frac{\lambda^2(a+1)}{2}, & a\lambda < |\beta| \end{cases} \quad (18)$$

3.6 Minimax Concave Penalty Regression

The minimax concave penalty (MCP) method is very similar to smoothly clipped absolute deviation [8, 2]. Both methods are used to avoid the high bias caused by the lasso method. MCP uses a penalty function that satisfies

$$\frac{dJ_a(\beta, \lambda)}{d\beta} = \begin{cases} \text{sign}(\beta) \left(\lambda - \frac{|\beta|}{a} \right), & |\beta| \leq a\lambda \\ 0, & a\lambda < |\beta| \end{cases} \quad (19)$$

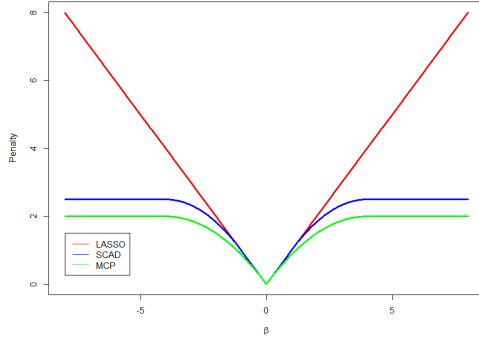
where $\lambda \geq 0$ and $a > 0$ are tuning parameters. Integrating [2], we see that

$$J_a(\beta, \lambda) = \begin{cases} \lambda|\beta| - \frac{\beta^2}{2a}, & |\beta| \leq a\lambda \\ \frac{1}{2}a\lambda^2, & a\lambda < |\beta| \end{cases} \quad (20)$$

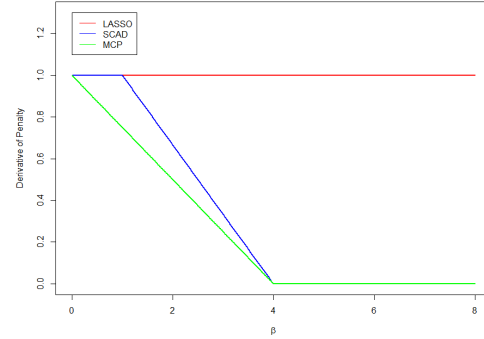
Figure 6 below shows the penalty functions (and their derivatives) for LASSO, SCAD, and MCP as a function of a coefficient value β . We see that LASSO applies a much stronger penalty to large coefficients than SCAD or MCP. Also, note that SCAD starts with a derivative equal to that of the lasso for small values of β ; on the other hand, the derivative of the penalty function for MCP starts decreasing immediately.

Figure 6: Penalty functions for LASSO, SCAD, and MCP, as well as their derivatives. These plots use $\lambda = 2$ and $a = 3$.

(a) Penalty functions for LASSO, SCAD, and MCP

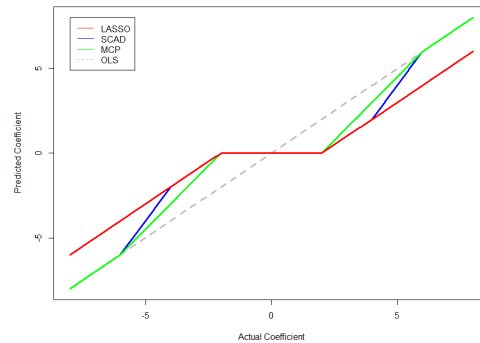


(b) Derivatives of the penalty functions for LASSO, SCAD, and MCP



Now, consider the case where $p = 1$ (there is only one predictor). Figure 7 shows the solutions given when using LASSO, SCAD, and MCP on such a model. The x -axis gives the actual coefficient for the single variable, and the y -axis represents the coefficient estimate produced using each of the algorithms. We used the particular values $\lambda = 2$ and $a = 3$. The gray line is the identity function, which also equals the solution obtained using ordinary least squares.

We see that all three methods set the predicted value to zero when the actual coefficient is small. Also, note that the LASSO method is always off from the identity function when the coefficient is large. On the other hand, SCAD and MCP merge with the identity function when the coefficient is sufficiently large. This shows that both SCAD and MCP can avoid the high bias that LASSO introduces.

Figure 7: Solutions for LASSO, SCAD, and MCP for a single predictor when $\lambda = 2$, and $a = 3$.

References

- [1] Thomas Lumley based on Fortran code by Alan Miller. *leaps: Regression Subset Selection*, 2020. R package version 3.1.
- [2] Breheny. Adaptive lasso, mcp, and scad. URL: <https://myweb.uiowa.edu/pbreheny/7600/s16/notes/2-29.pdf>, 2016.
- [3] Jianqing Fan and Runze Li. Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American statistical Association*, 96(456):1348–1360, 2001.
- [4] Jerome Friedman, Trevor Hastie, Robert Tibshirani, et al. *The elements of statistical learning*, volume 1. Springer series in statistics New York, 2001.
- [5] Gareth James, Daniela Witten, Trevor Hastie, and Rob Tibshirani. *ISLR: Data for an Introduction to Statistical Learning with Applications in R*, 2017. R package version 1.2.
- [6] Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani. *An introduction to statistical learning*, volume 112. Springer, 2013.
- [7] Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.
- [8] Cun-Hui Zhang et al. Nearly unbiased variable selection under minimax concave penalty. *The Annals of statistics*, 38(2):894–942, 2010.
- [9] Hui Zou. The adaptive lasso and its oracle properties. *Journal of the American statistical association*, 101(476):1418–1429, 2006.
- [10] Hui Zou and Trevor Hastie. Regularization and variable selection via the elastic net. *Journal of the royal statistical society: series B (statistical methodology)*, 67(2):301–320, 2005.