ECON3389 Machine Learning in Economics

Module 4 Feature Selection in Linear Models

Alberto Cappello

Department of Economics, Boston College

Fall 2024

Overview

Agenda:

- Ridge regression.
- Lasso regression.

Readings:

• ISLR Chapter 6, sections 6.1 and 6.2

Shrinkage Methods

• The subset selection methods fit a linear model that contains only a subset of the predictors. This is equivalent to setting the coefficients on excluded predictors to zero prior to running the estimation algorithm.

Shrinkage Methods

- The subset selection methods fit a linear model that contains only a subset of the predictors. This is equivalent to setting the coefficients on excluded predictors to zero prior to running the estimation algorithm.
- As an alternative, one can fit a model containing all p predictors using a technique that regularizes
 the coefficient estimates, or equivalently, that shrinks the coefficient estimates towards zero as part
 of the estimation algorithm.
- It may not be immediately obvious why such a constraint should improve the fit or if the algorithm will work in the first place, but it turns out that shrinking the coefficient estimates can significantly reduce their variance at a cost of a minor increase in bias.

Shrinkage Methods

- The subset selection methods fit a linear model that contains only a subset of the predictors. This is equivalent to setting the coefficients on excluded predictors to zero prior to running the estimation algorithm.
- As an alternative, one can fit a model containing all *p* predictors using a technique that *regularizes* the coefficient estimates, or equivalently, that *shrinks* the coefficient estimates towards zero as part of the estimation algorithm.
- It may not be immediately obvious why such a constraint should improve the fit or if the algorithm will work in the first place, but it turns out that shrinking the coefficient estimates can significantly reduce their variance at a cost of a minor increase in bias.
- Two most common shrinkage/regularization methods are ridge regression and lasso regression.

• Standard least squares regression fits the model by picking values of $\beta_0, \beta_1, \dots, \beta_p$ that minimize

$$RSS = \sum_{i=1}^{n} \left(y_i - \beta_0 - \sum_{j=1}^{p} \beta_j x_{ij} \right)^2$$

• Standard least squares regression fits the model by picking values of $\beta_0, \beta_1, \dots, \beta_p$ that minimize

$$RSS = \sum_{i=1}^{n} \left(y_i - \beta_0 - \sum_{j=1}^{p} \beta_j x_{ij} \right)^2$$

• Ridge regression instead picks coefficient values $\widehat{\beta}^R$ that minimize

$$\sum_{i=1}^{n} \left(y_i - \beta_0 - \sum_{j=1}^{p} \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^{p} \beta_j^2 = RSS + \lambda \sum_{j=1}^{p} \beta_j^2$$

• Parameter λ is the *tuning parameter*, to be determined separately.

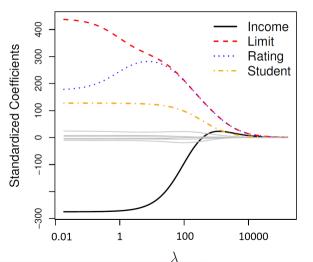


- The nature of ridge regression is similar to that of OLS: seek coefficient values that make the model fit the data well (by making RSS small).
- However, now we can no longer set values of coefficients to just significantly decrease RSS. This is because the second term $\lambda \sum_{j=1}^{p} \beta_{j}^{2}$, called a *shrinkage penalty*, will increase our loss function if values of $\beta_{0}, \beta_{1}, \ldots, \beta_{p}$ are far away from zero.

- The nature of ridge regression is similar to that of OLS: seek coefficient values that make the model fit the data well (by making RSS small).
- However, now we can no longer set values of coefficients to just significantly decrease RSS. This is because the second term $\lambda \sum_{j=1}^p \beta_j^2$, called a *shrinkage penalty*, will increase our loss function if values of $\beta_0, \beta_1, \ldots, \beta_p$ are far away from zero.
- Because loss function now has two terms to balance out, the extra second term has the effect of *shrinking* the estimates of β_i towards zero.

- The nature of ridge regression is similar to that of OLS: seek coefficient values that make the model fit the data well (by making RSS small).
- However, now we can no longer set values of coefficients to just significantly decrease RSS. This is because the second term $\lambda \sum_{j=1}^p \beta_j^2$, called a *shrinkage penalty*, will increase our loss function if values of $\beta_0, \beta_1, \ldots, \beta_p$ are far away from zero.
- Because loss function now has two terms to balance out, the extra second term has the effect of shrinking the estimates of β_j towards zero.
- The tuning parameter λ serves to control the relative impact of these two terms on the regression coefficient estimates. Selecting a good value for λ is critical, and is done via cross-validation.

Credit Card Data Example



- As can be seen from the picture, it is always possible to set λ to a value that will shrink all coefficients arbitrary close to zero.
- As such, we need to perform cross-validation testing to see which value of λ achieves minimal total value of ridge loss function.
- The process is usually done via a grid search algorithm (more on that later)

Feature Scaling and Standardization

• Standard least squares coefficient estimates are scale equivariant: multiplying X_j by a constant c simply leads to a scaling of the least squares coefficient estimate $\widehat{\beta}_j$ by a factor of 1/c. In other words, regardless of how the j-th predictor is scaled, $\widehat{\beta}_i X_j$ will always remain the same.

Feature Scaling and Standardization

- Standard least squares coefficient estimates are scale equivariant: multiplying X_j by a constant c simply leads to a scaling of the least squares coefficient estimate $\widehat{\beta}_j$ by a factor of 1/c. In other words, regardless of how the j-th predictor is scaled, $\widehat{\beta}_i X_j$ will always remain the same.
- In contrast, the ridge regression coefficient estimates can change substantially when multiplying a given predictor by a constant, due to the sum of squared coefficients term in the penalty part of the ridge regression loss function.
- Unlike the first term, which contains $\widehat{\beta}_j X_j$ parts, the shrinkage penalty contains values of only $\widehat{\beta}_j^2$, thus making is scale-dependent.

Feature Scaling and Standardization

- Standard least squares coefficient estimates are scale equivariant: multiplying X_j by a constant c simply leads to a scaling of the least squares coefficient estimate $\widehat{\beta}_j$ by a factor of 1/c. In other words, regardless of how the j-th predictor is scaled, $\widehat{\beta}_j X_j$ will always remain the same.
- In contrast, the ridge regression coefficient estimates can change substantially when multiplying a given predictor by a constant, due to the sum of squared coefficients term in the penalty part of the ridge regression loss function.
- Unlike the first term, which contains $\widehat{\beta}_j X_j$ parts, the shrinkage penalty contains values of only $\widehat{\beta}_j^2$, thus making is scale-dependent.
- Therefore, it is best to apply ridge regression after standardizing the predictors:

$$\widetilde{x}_{ij} = \frac{x_{ij}}{\sqrt{\frac{1}{n} \sum_{i=1}^{n} (x_{ij} - \overline{x}_j)^2}}$$



Why Does Ridge Regression Improve Over LS: Bias-variance trade-off

• Suppose our test data Te consists of a single data point (x_0, y_0) . Then

$$MSE = \mathbb{E}\left[\left(y_{0} - \hat{f}(x_{0})\right)^{2}\right] = \mathbb{E}\left[\left(f(x_{0}) - \hat{f}(x_{0})\right)^{2}\right] + \text{Var}(\epsilon_{0})$$

$$= \underbrace{\mathbb{E}\left[\left(\hat{f}(x_{0}) - \mathbb{E}\left[\hat{f}(x_{0})\right]\right)^{2}\right]}_{\text{Var}(\hat{f}(x_{0}))} + \underbrace{\mathbb{E}\left[\left(f(x_{0}) - \mathbb{E}\left[\hat{f}(x_{0})\right]\right)^{2}\right]}_{\mathbb{E}\left[\text{Bias}^{2}(\hat{f}(x_{0}))\right]} + \text{Var}(\epsilon_{0})$$

- ullet Variance refers to the amount by which \hat{f} would change if we estimated it using a different training data set
- Bias refers to the error that is introduced by approximating a real-life problem by a much simpler model

Why Does Ridge Regression Improve Over LS: Bias-variance trade-off

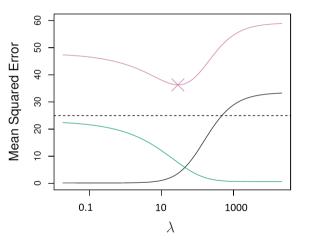
• Suppose our test data Te consists of a single data point (x_0, y_0) . Then

$$MSE = \mathbb{E}\left[\left(y_0 - \hat{f}(x_0)\right)^2\right] = \mathbb{E}\left[\left(f(x_0) - \hat{f}(x_0)\right)^2\right] + \text{Var}(\epsilon_0)$$

$$= \underbrace{\mathbb{E}\left[\left(\hat{f}(x_0) - \mathbb{E}\left[\hat{f}(x_0)\right]\right)^2\right]}_{\text{Var}(\hat{f}(x_0))} + \underbrace{\mathbb{E}\left[\left(f(x_0) - \mathbb{E}\left[\hat{f}(x_0)\right]\right)^2\right]}_{\mathbb{E}\left[\text{Bias}^2\left(\hat{f}(x_0)\right)\right]} + \text{Var}(\epsilon_0)$$

- ullet Variance refers to the amount by which \hat{f} would change if we estimated it using a different training data set
- Bias refers to the error that is introduced by approximating a real-life problem by a much simpler model
- Typically as the flexibility of \hat{f} increases, its variance increases, and its bias decreases. So choosing the flexibility based on MSE amounts to a bias-variance trade-off.

Why Does Ridge Regression Improve Over LS?



- Squared bias, variance and MSE.
- 45 predictors (p), 50 observations (n)
- Because OLS is free to choose any coefficient values, it tends to pick the ones that provide best fit, meaning less bias and more variance.
- Ridge regression, on the other hand, is penalized for choosing coefficients with high second moments, thus leading to less variance, slightly more bias, but lower MSE overall.

Lasso regression

• Unlike subset selection, which generally selects models that involve just a subset of all variables, ridge regression will include all p predictors in the final model. This makes ridge regression completely infeasible when p > n, as if often the case, for example, with Internet-related data.

Lasso regression

- Unlike subset selection, which generally selects models that involve just a subset of all variables, ridge regression will include all p predictors in the final model. This makes ridge regression completely infeasible when p > n, as if often the case, for example, with Internet-related data.
- The LASSO (Least Absolute Shrinkage and Selection Operator) is an alternative that overcomes this disadvantage. It achieves that by using a different type of shrinkage penalty:

$$\sum_{i=1}^{n} \left(y_i - \beta_0 - \sum_{j=1}^{p} \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^{p} |\beta_j| = RSS + \lambda \sum_{j=1}^{p} |\beta_j|$$

• In statistical lingo, this type of penalty is known as ℓ_1 -penalty, because it uses ℓ_1 norm of coefficient vector β given by $||\beta||_1 = \sum |\beta_j|$. Ridge regression, on the other hand, uses ℓ_2 norm as a penalty, given by $||\beta||_2 = \sum \beta_i^2$

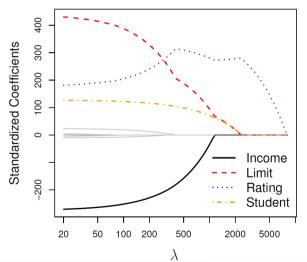
Lasso Variable Selection

- As with ridge regression, lasso shrinks all coefficient estimates towards zero.
- However, unlike ℓ_2 penalty in ridge regression, lasso penalty has the effect of forcing some of the coefficient estimates to be exactly equal to zero when the tuning parameter λ is sufficiently large.

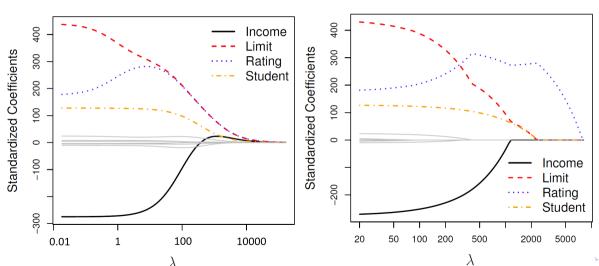
Lasso Variable Selection

- As with ridge regression, lasso shrinks all coefficient estimates towards zero.
- However, unlike ℓ_2 penalty in ridge regression, lasso penalty has the effect of forcing some of the coefficient estimates to be exactly equal to zero when the tuning parameter λ is sufficiently large.
- Hence, much like best subset selection, lasso performs variable selection, starting with the full set of
 p variables. We say that lasso yields sparse models that is, models that involve only a subset of
 the variables.
- ullet As in ridge regression, selecting a good value of λ for lasso is critical; cross-validation is again the method of choice.

Credit Card Data Example



- Similar to ridge regression, setting λ to sufficiently high value will shrink all coefficients to zero.
- Unlike ridge regression, lasso coefficients will get shrunk exactly to zero in a single jump, without smooth continuous decline.
- Additionally, while ridge regression shrinks all coefficients close to zero around the same values of λ, lasso sets some coefficients to zero much earlier than others.



• Why is it that in lasso regression we get some of the coefficients shrunk exactly to zero, but not in ridge regression?

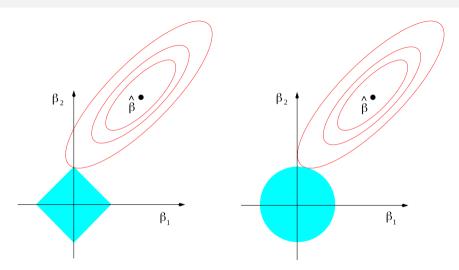
- Why is it that in lasso regression we get some of the coefficients shrunk exactly to zero, but not in ridge regression?
- One can show that lasso and ridge regression coefficient estimates solve the following problems:

$$\min_{\beta} \sum_{i=1}^{n} \left(y_i - \beta_0 - \sum_{j=1}^{p} \beta_j x_{ij} \right)^2 \quad \text{subject to } \sum_{j=1}^{p} |\beta_j| \le s$$

and

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

• These two problems have a useful geometric representation that shows exactly why lasso induces sparsity among coefficients.



Selecting the Tuning Parameter λ

• Both with ridge and with lasso we need to select the value for the tuning parameter λ or equivalently, the value of the constraint s in a way that will not lead to overfitting or other mistakes. Cross-validation provides a simple way to tackle this problem.

Selecting the Tuning Parameter λ

- Both with ridge and with lasso we need to select the value for the tuning parameter λ or equivalently, the value of the constraint s in a way that will not lead to overfitting or other mistakes. Cross-validation provides a simple way to tackle this problem.
- ullet We choose a grid of λ values and fit a separate model for every value from that grid using K-fold cross-validation.
- ullet We then compute the cross-validation error for each value of λ and select the one for which that error is smallest.

Selecting the Tuning Parameter λ

- Both with ridge and with lasso we need to select the value for the tuning parameter λ or equivalently, the value of the constraint s in a way that will not lead to overfitting or other mistakes. Cross-validation provides a simple way to tackle this problem.
- ullet We choose a grid of λ values and fit a separate model for every value from that grid using K-fold cross-validation.
- ullet We then compute the cross-validation error for each value of λ and select the one for which that error is smallest.
- Finally, the model is re-fit using all of the available observations and the selected value of the tuning parameter.

• In terms of overall fit neither ridge regression nor lasso will universally dominate the other.

- In terms of overall fit neither ridge regression nor lasso will universally dominate the other.
- In general, one might expect lasso to perform better when the response is a function of only a relatively small number of predictors. However, that is never known a priori with real-life data.

- In terms of overall fit neither ridge regression nor lasso will universally dominate the other.
- In general, one might expect lasso to perform better when the response is a function of only a relatively small number of predictors. However, that is never known a priori with real-life data.
- Ridge regression can perform better if p < n and there is no a priori reason for some of the variables to not be included in the model.

- In terms of overall fit neither ridge regression nor lasso will universally dominate the other.
- In general, one might expect lasso to perform better when the response is a function of only a relatively small number of predictors. However, that is never known a priori with real-life data.
- Ridge regression can perform better if p < n and there is no a priori reason for some of the variables to not be included in the model.
- Lasso can perform variable selection and model estimation with p > n, but has a known issue of ignoring groups of correlated variables (e.g. performance metrics of NBA players) and almost randomly selecting only one variable out of the group.

Lasso and Economics

- Despite its know flaws, over the past decade lasso has become very popular with both academic researchers and applied economists.
- From the theoretical perspective, multiple extensions and variations of lasso has been suggested, and today advanced versions of it can deal both with correlated regressors (group lasso, elastic net) and biased estimates (adaptive lasso, post-lasso).
- The main driving force behind is the ability to tackle datasets that previously were completely unusable due to number of variables p being close or even larger than sample size n.
- Additionally, lasso allows economists to utilize *sparse* structural models, e.g. consumer preferences across hundreds of product attributes with most of them having zero importance.