

## Lecture 10 — September 27

Lecturer: Purnamrita Sarkar

Scribe: Amin Anvari

**Disclaimer:** These scribe notes have been slightly proofread and may have typos etc.

**Note:** The latex template was borrowed from EECS, U.C. Berkeley.

## 10.1 Regularization

Regularization methods such as Ridge Regression and LASSO introduce a penalty term on model complexity to prevent overfitting.

### 10.1.1 Ridge Regression

For least squares regression, the model coefficients are selected by

$$\hat{\beta}_{LS} = \min_{\beta} \sum_{i=1}^n (y_i - \beta_0 - x_i^T \beta)^2 \quad (10.1)$$

For ridge regression, an additional term is added which penalizes all  $\beta_j$  for  $j > 0$

$$\hat{\beta}_{ridge} = \min_{\beta} \sum_{i=1}^n ((y_i - \beta_0 - x_i^T \beta)^2 + \lambda \beta^T \beta) \quad (10.2)$$

It is clear from this equation that if the variables are on different scales the Ridge Regression model will penalize them differently. Thus, centering must be performed to remove the means from all parameters. Additionally, the standard errors must be normalized to 1.

First let's forget about the normalization, and just consider the original problem. We will see that centering  $\mathbf{X}$  and  $\mathbf{y}$  ( $X_c(i, j) = X_{ij} - \bar{x}_j$  and  $y_c(i) = y_i - \bar{y}$ ) will not change the  $\beta_i$ 's other than the intercept term.

$$\begin{aligned} & \sum_i (y_i - \beta_0 - \sum_j \beta_j x_{ij})^2 + \lambda \beta_j^2 \\ &= \sum_i (y_i - (\beta_0 + \sum_j \beta_j \bar{x}_j) + \sum_j \beta_j (x_{ij} - \bar{x}_j))^2 + \lambda \beta_j^2 \end{aligned}$$

$$\beta'_0 = \beta_0 + \sum_{j=1}^k x_{ij} \beta_j \quad (10.3)$$

$$\beta'_j = \beta_j, \forall j > 0 \quad (10.4)$$

As it turns out, deriving the above expression w.r.t  $\beta'_0$  and setting to zero gives:

$$\hat{\beta}'_0 = \bar{y}$$

Thus the reparameterized equation for  $\hat{\beta}_{ridge}$  is

$$OBJ_{ridge} = \min_{\beta} (\mathbf{y}_c - \mathbf{X}_c \beta)^\top (\mathbf{y}_c - \mathbf{X}_c \beta) + \lambda \beta^\top \beta, \quad (10.5)$$

where  $\mathbf{y}_c$  and  $\mathbf{X}_c$  are centered versions of  $\mathbf{y}$  and  $\mathbf{X}$ . Setting the derivative of this expression to zero gives the solution

$$\hat{\beta}_{ridge} = (\mathbf{X}^\top \mathbf{X} + \lambda I)^{-1} \mathbf{X}^\top \mathbf{y} \quad (10.6)$$

Note that this ridge regularization is the same as adding  $\lambda$  to all eigenvalues of  $\mathbf{X}^\top \mathbf{X}$ . Furthermore, exactly corresponding to the hat matrix, in this case if we write  $\hat{y} = S y$  we get:

$$S = \mathbf{X}(\mathbf{X}^\top \mathbf{X} + \lambda I)^{-1} \mathbf{X}^\top \quad (10.7)$$

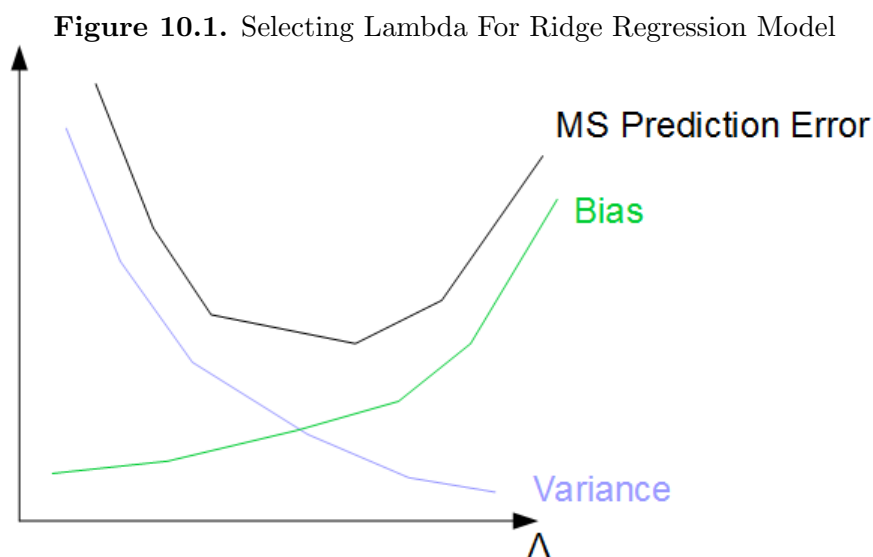
and so the effective degrees of freedom is:

$$\text{trace}(S) = \sum_{i=1}^p \frac{\sigma_i^2}{\sigma_i^2 + \lambda} \quad (10.8)$$

In which  $\sigma_i$ 's are the singular values of matrix  $\mathbf{X}$ .

**How to Choose  $\lambda$ :** Selection of  $\lambda$  is a tradeoff between bias, variance and mean square error (Fig. 9.3).

1. Perform a grid search over  $\lambda$  or  $\log(\lambda)$ .
2. For each  $\lambda$  perform a cross-validation and calculate the mean and standard error on the estimated model prediction error
3. Identify the  $\lambda$  with the lowest mean predictive error
4. Apply One Standard Error Rule to select most parsimonious model whose mean lies within one standard error (in this case more parsimonious means smaller effective degrees of freedom, so a larger  $\lambda$ )



**Equivalent Bayesian Interpretation** Assume that  $y$  is normally distributed, and apply a Gaussian prior to  $\beta$

$$y \sim N(\mathbf{X}\beta, \sigma^2 I) \quad (10.9)$$

$$\beta \sim N(0, \tau^2 I) \quad (10.10)$$

Then the posterior distribution of  $\beta$  given  $y$  can be calculated as

$$f(\beta|y) = \exp\left(\frac{-(\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta)}{2\sigma^2}\right) \exp\left(\frac{-\beta^\top \beta}{\tau^2}\right) \quad (10.11)$$

From this form it can be observed that the relationship between  $\sigma$  and  $\tau$  defines  $\lambda$ . Ridge regression generally has smaller LSE than linear regression, and Ridge Regression will make coefficients small for variables which are not highly correlated to the output. However, Ridge Regression will not drive coefficients to zero (unless  $\lambda = \infty$ , in which case it drives all coefficients to zero), so it cannot be used for variable selection.

### 10.1.2 LASSO

$$OBJ_{LASSO} = \min_{\beta} (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) + \lambda \|\beta\|_1 \quad (10.12)$$

The LSE of LASSO is comparable to Ridge Regression, and LASSO will drive some coefficients to zero with a large  $\lambda$ . However, there is no closed form solution. If the problem is reformulated as

$$\hat{\beta}_{LASSO} = \min_{\beta} (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) \quad (10.13)$$

$$\text{s.t. } \sum_{j=1}^k |\beta_j| \leq t \quad (10.14)$$

This is a convex optimization problem, which will have a solution even without a closed-form.

Note that we can also devise an objective function using the zero-norm instead of one-norm (All Subsets objective function):

$$OBJ_{\text{All subsets}} = \min_{\beta} (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) \text{ s.t. } \|\beta\|_0 \leq k \quad (10.15)$$

The problem with this objective function is that it does not yield a convex problem anymore. In fact, one would need to search through all possible subsets of size smaller than  $k + 1$  to evaluate the best possible subset.

Consider the easier problem where I want the best subset of size  $k$ . We will denote this by  $OBJ_{BS}$  from now on. Typically this is a NP complete problem, which is computationally intractable. However, as it turns out in the case where the design matrix has orthogonal columns, i.e. the features are orthonormal to each other, i.e.  $\mathbf{X}^T \mathbf{X} = \mathbf{I}$  one can actually write down effect of best subset of size  $k$ ,  $L_1$  and  $L_2$  norm penalizations easily.

## 10.2 Orthogonal design matrix

Consider an orthogonal design matrix such that  $\mathbf{X}^T \mathbf{X} = \mathbf{I}$ . Now we will look at Best subset of size  $k$ , Lasso and Ridge to get an intuition about how the different regularizations work. The main trick is to see that in this case the loss function becomes separable w.r.t the coefficients of  $\beta$ . Also note that in this case  $\beta^{ols} = \mathbf{X}^T \mathbf{y}$ . So

$$(\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) = \mathbf{y}^T \mathbf{y} + \beta^T \beta - 2\mathbf{y}^T \mathbf{X}\beta = (\beta^{ols} - \beta)^T (\beta^{ols} - \beta) + \mathbf{y}^T (\mathbf{I} - \mathbf{X}\mathbf{X}^T) \mathbf{y}$$

Since the last term is independent of  $\beta$ , we only need to consider the first term.

### 10.2.1 Best subset of size $k$

$$\min_{\beta} (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) \quad (10.16)$$

$$\text{s.t. } \|\beta\|_0 = k \quad (10.17)$$

which is equivalent to:

$$\min_{\beta} (\beta - \beta^{ols})^T (\beta - \beta^{ols}) = \min_{\beta} \sum_{i=1}^p (\beta_i - \beta_i^{ols})^2 \quad (10.18)$$

$$\text{s.t. } \|\beta\|_0 = k \quad (10.19)$$

Now think for a second. I can only put nonzero values in  $k$  positions, and at the same time minimize the above sum of squares. Can I put them on the smallest  $k$   $\beta_i$  values? That

will incur a lot of error from the large coefficients. So, we should use:

$$\beta_i^{\hat{BS}} = \begin{cases} \hat{\beta}_i^{OLS} & \text{If } |\beta_i^{ols}| \geq |\beta_{(k)}^{ols}| \\ 0 & \text{Otherwise} \end{cases}$$

Here  $x^{(k)}$  means the  $k^{th}$  largest element in absolute value.

### 10.2.2 Lasso

Now lets think about Lasso. Now the objective function is

$$\min_{\beta} \sum_i (\beta_i - \beta_i^{ols})^2 + \lambda \sum_i |\beta_i| \quad (10.20)$$

$$(10.21)$$

Since the function is now separable w.r.t the elements, we can reason about each coordinate separately. Lets take coordinate  $i$ . If  $\beta_i^{ols} > 0$ , it makes sense to set  $\beta_i > 0$  in order to minimize the above function. So that gives us  $\min_{\beta_i} (\beta_i^{ols} - \beta_i)^2 + \lambda \beta_i$ . Derive w.r.t  $\beta_i$  and set to zero to solve.  $\beta_i^{LASSO} = \beta_i^{ols} - \lambda/2$ , but in this case  $\beta_i^{LASSO} > 0$  as well. So, use:  $\beta_i^{LASSO} = \max(\beta_i^{ols} - \lambda/2, 0)$ . A similar argument for the negative OLS coefficients give the following:

$$\text{for } \beta_i^{ols} > 0 : \quad \hat{\beta}_i = \max \left\{ \beta_i^{ols} - \frac{\lambda}{2}, 0 \right\} \quad (10.22)$$

$$\text{for } \beta_i^{ols} \leq 0 : \quad \hat{\beta}_i = \min \left\{ \beta_i^{ols} + \frac{\lambda}{2}, 0 \right\} \quad (10.23)$$

### 10.2.3 Ridge

Finally, lets figure out ridge regression. This is actually quite straight forward. We want:

$$\min_{\beta} \sum_i (\beta_i - \beta_i^{ols})^2 + \lambda \sum_i \beta_i^2 \quad (10.24)$$

Derive both sides w.r.t  $\beta_i$  and set to zero. This gives:

$$\hat{\beta}_i = \frac{\beta_i^{ols}}{1 + \lambda}$$