

# Homework 3

Daniel Hartig

March 7, 2017

**e. 2.**

For the ridge regression function

$$|\mathbf{y}_{\text{train}} - \Phi_{\text{train}} w|_2^2 + \lambda |w|_2^2,$$

then we can solve for a minimum by taking the zeroes of the derivative.

$$\begin{aligned} \frac{d}{dw} \left( |\mathbf{y}_{\text{train}} - \Phi_{\text{train}} w|_2^2 + \lambda |w|_2^2 \right) &= -2\Phi_{\text{train}}^\top (\mathbf{y}_{\text{train}} - \Phi w) + 2\lambda w = 0 \\ -2\Phi_{\text{train}}^\top \Phi_{\text{train}} w - 2\lambda w &= -2\Phi_{\text{train}}^\top \mathbf{y}_{\text{train}} \\ w &= (\Phi_{\text{train}}^\top \Phi_{\text{train}} + \lambda \mathbf{I}_D)^{-1} \Phi_{\text{train}}^\top \mathbf{y}_{\text{train}} \end{aligned}$$

Using the singular value decomposition of  $\Phi = \mathbf{U}\mathbf{S}\mathbf{V}^\top$ , and the spectral decomposition

$$\begin{aligned} \Phi \Phi^\top &= (\mathbf{U}\mathbf{S}\mathbf{V}^\top) (\mathbf{U}\mathbf{S}\mathbf{V}^\top)^\top \\ &= \mathbf{V}\mathbf{S}^\top \mathbf{U}^\top \mathbf{U} \mathbf{S} \mathbf{V}^\top \\ &= \mathbf{V}\mathbf{S}^2 \mathbf{V}^\top \end{aligned}$$

Plugging this into the formula for  $w$  gives, and utilizing the property that the transpose of a square diagonal matrix is itself ( $\mathbf{S}^\top = \mathbf{S}$ ),

$$\begin{aligned} w &= (\mathbf{V}\mathbf{S}^2 \mathbf{V}^\top + \lambda \mathbf{I}_D) (\mathbf{U}\mathbf{S}\mathbf{V}^\top)^\top \mathbf{y}_{\text{train}} \\ &= \mathbf{V} (\mathbf{S}^2 + \lambda \mathbf{I}_D) \mathbf{V}^\top \mathbf{V} \mathbf{S}^\top \mathbf{U}^\top \mathbf{y}_{\text{train}} \\ &= \mathbf{V} (\mathbf{S}^2 + \lambda \mathbf{I}_D) \mathbf{S} \mathbf{U}^\top \mathbf{y}_{\text{train}} \end{aligned}$$

$(\mathbf{S}^2 + \lambda \mathbf{I}_D) \mathbf{S}$  resolves to a diagonal matrix whose trace is

$$\text{tr}(\mathbf{S}_\lambda) = \sum_{j=1}^D \frac{d_j}{d_j^2 + \lambda}$$

yielding the final form

$$\mathbf{V} \mathbf{S}_\lambda \mathbf{U}^\top \mathbf{y}_{\text{train}}.$$

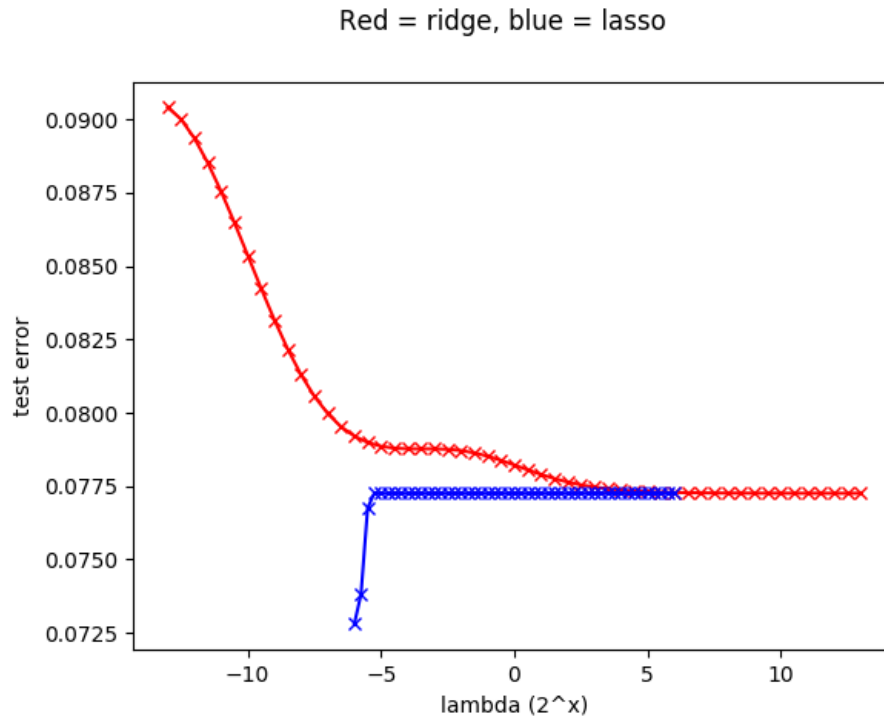
**e. 3.**

The computational complexity of multiplying an  $n \times m$  matrix by an  $m \times p$  matrix is  $O(nmp)$ . Multiplying  $\mathbf{S}_\lambda$  (a  $D \times D$  matrix) by  $\mathbf{U}^\top$  ( $D \times n$ ) has a complexity of  $O(nD^2)$  flops.

However, if we pre-compute  $\mathbf{U}^\top \mathbf{y}_{\text{train}}$  only one time, then each matrix multiplication operation in the loop is  $(D \times D) \times (D \times 1)$ ; by multiplying  $\mathbf{U}^\top \mathbf{y}_{\text{train}}$  first by  $\mathbf{S}_\lambda$  then by  $\mathbf{V}$ . These operations have a complexity of  $O(D^2)$  flops.

i.)

The test errors for both methods combined are:



We can see from the graphs that both methods converge to a similar error rate as  $\lambda$  increases, but that Lasso regression holds that low error rate over a wider range of lambda values. Lasso regression also has a lower error rate at low values of lambda. Lasso regression holds the lowest error rate of 0.072 when  $\lambda = \sqrt{\frac{\log D}{n_{\text{train}}}} 2^{-6}$