

1. (Augmented dataset) Ridge regression is equivalent to applying OLS on an expanded dataset that has additional examples. Describe these additional examples in detail. Intuitively, what effect do these additional examples have?

The ridge regression estimate is defined as  $\beta_{\text{RR}} = (XX^T + \lambda I)^{-1}Xy$ , where  $X \in \mathbb{R}^{p \times n}$ ,  $y \in \mathbb{R}^n$ ,  $\beta \in \mathbb{R}^p$ , which can be reformulated as a modified least-square problem

$$\beta_{\text{RR}} = \arg \min_{\beta} \left\| \begin{bmatrix} y \\ 0 \end{bmatrix} - \begin{bmatrix} X^T \\ \sqrt{\lambda} I_{n \times p} \end{bmatrix} \beta \right\|_2^2$$

It seems like we have added  $p$  vectors to the original  $n$  datapoints  $[x_1, \dots, x_n]$ ,  $x_i \in \mathbb{R}^p$ ,  $i = 1, \dots, n$ . By doing so, Ridge regression embedded the original problem in  $\mathbb{R}^n$  into a larger space  $\mathbb{R}^{n+p}$  by moving into  $p$  different directions with a small amount  $\sqrt{\lambda}$  which could decrease any collinearity present in the original data points  $X$ . As we have seen in the notes of linear regression, when the number of training data is small,  $\lambda$  neutralizes the large variance of the errors between the true  $\beta$  coefficients and the ridge regression coefficients due to small singular values of the sample covariance matrix.

2. (Correlated features) Consider a regression problem where the response only depends on one feature, but we don't know it, so we incorporate an additional feature into the model that happens to be very correlated with the first feature. More specifically, let  $y \in \mathbb{R}^n$  be defined by

$$y := \beta_{\text{true}} w_1 + z, \quad (1)$$

where  $\beta_{\text{true}} \in \mathbb{R}$  is the true coefficient,  $w_1 \in \mathbb{R}^n$  is the first feature vector, and  $z \in \mathbb{R}^n$  is additive noise. The second feature vector is given by  $w_2 \in \mathbb{R}^n$  and can be decomposed into

$$w_2 = \alpha w_1 + \sqrt{1 - \alpha^2} w_{\perp}, \quad (2)$$

where  $w_{\perp}$  is orthogonal to  $w_1$ . The vectors  $w_1$ ,  $w_2$ ,  $w_{\perp}$  and  $z$  all have unit  $\ell_2$  norm. In addition, we assume

$$w_1^T z = 0.1, \quad (3)$$

$$w_{\perp}^T z = 0.1. \quad (4)$$

We fit a linear regression model to  $y$  using the feature matrix

$$X = \begin{bmatrix} w_1^T \\ w_2^T \end{bmatrix}. \quad (5)$$

- (a) What does the OLS estimator of the coefficients  $\beta_{\text{OLS}}$  equal to when  $\alpha \rightarrow 1$ ? Explain what is happening.

*Hint:* Use the fact that for any  $a$ ,  $b$ ,  $c$ , and  $d$  such that  $ad \neq bc$

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix}^{-1} = \frac{1}{ad - bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}. \quad (6)$$

The OLS estimator of the coefficients of  $\beta_{\text{OLS}}$  is given by  $\beta_{\text{OLS}} = (XX^T)^{-1}Xy$ . Expand-

ing each term, we have:

$$\begin{aligned}
(XX^T) &= \begin{bmatrix} w_1^T \\ w_2^T \end{bmatrix} [w_1 \ w_2] \\
&= \begin{bmatrix} w_1^T w_1 & w_1^T w_2 \\ w_2^T w_1 & w_2^T w_2 \end{bmatrix} \\
&= \begin{bmatrix} \|w_1\|_2^2 & w_1^T w_2 \\ w_2^T w_1 & \|w_2\|_2^2 \end{bmatrix} \\
w_1^T w_2 &= w_1^T (\alpha w_1 + \sqrt{1 - \alpha^2} w_\perp) \\
&= \alpha \|w_1\|_2^2 + \sqrt{1 - \alpha^2} w_1^T w_\perp \\
&= \alpha \quad \text{since } \|w_1\|_2 = 1, w_1 \perp w_\perp \\
w_2^T w_2 &= (\alpha w_1 + \sqrt{1 - \alpha^2} w_\perp)^T (\alpha w_1 + \sqrt{1 - \alpha^2} w_\perp) \\
&= \alpha^2 w_1^T w_1 + 2\alpha \sqrt{1 - \alpha^2} w_1^T w_\perp + (1 - \alpha^2) w_\perp^T w_\perp \\
&= \alpha^2 \|w_1\|_2^2 + (1 - \alpha^2) \|w_\perp\|_2^2 \quad \text{knowing that } w_\perp \perp w_1 \\
&= \alpha^2 + (1 - \alpha^2) = 1 \quad \text{by assumptions } \|w_1\|_2 = \|w_\perp\|_2 = 1 \\
\Rightarrow (XX^T) &= \begin{bmatrix} 1 & \alpha \\ \alpha & 1 \end{bmatrix} \\
\Rightarrow (XX^T)^{-1} &= \frac{1}{1 - \alpha^2} \begin{bmatrix} 1 & -\alpha \\ -\alpha & 1 \end{bmatrix} \\
Xy &= \begin{bmatrix} w_1^T \\ w_2^T \end{bmatrix} (\beta_{\text{true}} w_1 + z) \\
&= \begin{bmatrix} \beta_{\text{true}} w_1^T w_1 + w_1^T z \\ \beta_{\text{true}} w_2^T w_1 + w_2^T z \end{bmatrix} \\
&= \begin{bmatrix} \beta_{\text{true}} + 0.1 \\ \alpha \beta_{\text{true}} + 0.1 \end{bmatrix}
\end{aligned}$$

Substituting each of the previous terms back into the expression of  $\beta_{\text{OLS}}$ , we find that:

$$\begin{aligned}
\beta_{\text{OLS}} &= \frac{1}{1 - \alpha^2} \begin{bmatrix} 1 & -\alpha \\ -\alpha & 1 \end{bmatrix} \begin{bmatrix} \beta_{\text{true}} + 0.1 \\ \alpha \beta_{\text{true}} + 0.1 \end{bmatrix} \\
&= \frac{1}{1 - \alpha^2} \begin{bmatrix} \beta_{\text{true}} + 0.1 - \alpha^2 \beta_{\text{true}} - 0.1\alpha \\ -\alpha \beta_{\text{true}} - 0.1\alpha + \alpha \beta_{\text{true}} + 0.1 \end{bmatrix} \\
&= \frac{1}{1 - \alpha^2} \begin{bmatrix} (1 - \alpha^2) \beta_{\text{true}} + 0.1(1 - \alpha) \\ 0.1(1 - \alpha) \end{bmatrix} \\
&= \begin{bmatrix} \beta_{\text{true}} + \frac{0.1}{1 + \alpha} \\ \frac{0.1}{1 + \alpha} \end{bmatrix}
\end{aligned}$$

When  $\alpha \rightarrow 1$ ,  $\beta_{\text{OLS}} \rightarrow \begin{bmatrix} \beta_{\text{true}} + 0.05 \\ 0.05 \end{bmatrix}$ . The OLS estimator ignores the correlated feature  $w_2$  and adds a fixed bias 0.05 to the true  $\beta_{\text{true}}$  coefficient which could be significant compared to  $\beta_{\text{true}}$ . Notice that in this case  $XX^T$  is rank 1, and the algorithm used to find the OLS estimator may be unable to find a solution.

- (b) What does the corresponding estimate of the response  $y_{\text{OLS}} := X^T \beta_{\text{OLS}}$  equal to when  $\alpha \rightarrow 1$ ? Is it collinear with the true feature  $w_1$  when  $\alpha \rightarrow 1$ ? Explain what is happening.

Taking  $\alpha \rightarrow 1$ ,  $w_2 \rightarrow w_1$  and we have

$$\begin{aligned} y_{\text{OLS}} &:= X^T \beta_{\text{OLS}} \rightarrow [w_1 \ w_2] \begin{bmatrix} \beta_{\text{true}} + 0.05 \\ 0.05 \end{bmatrix} \\ &\rightarrow [w_1 \ w_1] \begin{bmatrix} \beta_{\text{true}} + 0.05 \\ 0.05 \end{bmatrix} \\ &\rightarrow (\beta_{\text{true}} + 0.1)w_1 \end{aligned}$$

When  $\alpha \rightarrow 1$ , the correlated feature  $w_2 \rightarrow w_1$ , and the response variable is collinear with the true feature  $w_1$ , the OLS linear model ignores  $w_2$  and estimates  $y_{\text{OLS}}$  as a linear scaling of  $w_1$  up to a factor of 0.1 which could lead to an important error depending the magnitude of the linear coefficient  $\beta_{\text{true}}$  compared to 0.1.

- (c) What does the ridge regression estimator of the coefficients  $\beta_{\text{RR}}$  equal to when  $\alpha \rightarrow 1$  and the regularization parameter  $\lambda > 0$  is fixed? Describe the difference with the OLS estimate.

By definition the ridge regression estimator of the coefficients  $\beta_{\text{RR}}$  is:

$$\begin{aligned} \beta_{\text{RR}} &= (XX^T + \lambda I)^{-1} Xy \quad \lambda > 0 \\ XX^T + \lambda I &= \begin{bmatrix} 1 + \lambda & \alpha \\ \alpha & 1 + \lambda \end{bmatrix} \\ (XX^T + \lambda I)^{-1} &= \frac{1}{(1 + \lambda)^2 - \alpha^2} \begin{bmatrix} 1 + \lambda & -\alpha \\ -\alpha & 1 + \lambda \end{bmatrix} \\ \Rightarrow \beta_{\text{RR}} &= \frac{1}{(1 + \lambda)^2 - \alpha^2} \begin{bmatrix} 1 + \lambda & -\alpha \\ -\alpha & 1 + \lambda \end{bmatrix} \begin{bmatrix} \beta_{\text{true}} + 0.1 \\ \alpha \beta_{\text{true}} + 0.1 \end{bmatrix} \\ &= \frac{1}{(1 + \lambda)^2 - \alpha^2} \begin{bmatrix} (1 + \lambda)(\beta_{\text{true}} + 0.1) - \alpha^2 \beta_{\text{true}} - 0.1\alpha \\ -\alpha \beta_{\text{true}} - 0.1\alpha + \alpha(1 + \lambda)\beta_{\text{true}} + 0.1(1 + \lambda) \end{bmatrix} \\ &= \frac{1}{(1 + \lambda)^2 - \alpha^2} \begin{bmatrix} (1 + \lambda - \alpha^2)\beta_{\text{true}} + 0.1(1 + \lambda - \alpha) \\ \alpha\lambda\beta_{\text{true}} + 0.1(1 + \lambda - \alpha) \end{bmatrix} \\ &= \begin{bmatrix} \frac{(1 + \lambda - \alpha^2)\beta_{\text{true}}}{(1 + \lambda)^2 - \alpha^2} + \frac{0.1}{1 + \lambda + \alpha} \\ \frac{\lambda\alpha\beta_{\text{true}}}{(1 + \lambda)^2 - \alpha^2} + \frac{0.1}{1 + \lambda + \alpha} \end{bmatrix} \\ \alpha \rightarrow 1 &\rightarrow \frac{\beta_{\text{true}} + 0.1}{2 + \lambda} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \end{aligned}$$

If we compare to  $\beta_{\text{OLS}}$  estimator, the ridge regression estimator has coefficients on both  $w_1$  and  $w_2$  and this coefficient can be regularized using  $\lambda$ . Even if  $XX^T$  is singular (rank 1), the  $l_2$  regularization using  $\lambda$  makes the matrix  $XX^T$  non singular and we are still able to find a solution which includes both features. In addition we have a new set of coefficients each time we tune  $\lambda$ : when  $\lambda \rightarrow 0$ ,  $\beta_{\text{RR}} = \frac{\beta_{\text{OLS}}}{2}$  and when  $\lambda \rightarrow \infty$ ,  $\beta_{\text{RR}} \rightarrow 0$ . So the ridge

regression estimator with the regularization parameter,  $\lambda$ , governs the trade-off between a model that has good fit on the training set (low bias) and a model which reduces the test error (low variance).

- (d) What does the corresponding estimate of the response  $y_{\text{RR}} := X^T \beta_{\text{RR}}$  equal to when  $\alpha \rightarrow 1$ ? Is it collinear with the true feature  $w_1$ ?

When  $\alpha \rightarrow 1$ ,  $w_2 \rightarrow w_1$ , which leads to

$$\begin{aligned} y_{\text{RR}} := X^T \beta_{\text{RR}} &\rightarrow \frac{\beta_{\text{true}} + 0.1}{2 + \lambda} [w_1 w_2] \begin{bmatrix} 1 \\ 1 \end{bmatrix} \\ &\rightarrow \frac{\beta_{\text{true}} + 0.1}{2 + \lambda} [w_1 w_1] \begin{bmatrix} 1 \\ 1 \end{bmatrix} \\ &\rightarrow 2 \frac{\beta_{\text{true}} + 0.1}{2 + \lambda} w_1 \end{aligned}$$

$y_{\text{RR}}$  is collinear with the true feature  $w_1$ , compared to  $y_{\text{OLS}}$ , the parameter  $\lambda$  allows to control the amount of linearity between the response and control variable (when  $\lambda = 0$ ,  $y_{\text{RR}} = y_{\text{OLS}}$ ), which might be desirable if the test data points, not known in advance, are not totally dependent on the feature  $w_1$ .

3. (Prior knowledge) Consider a linear regression problem where we have prior information indicating that the coefficients should be close to a certain value  $\beta_{\text{prior}}$ .

- (a) How can you incorporate this prior knowledge if you are using ridge regression? Write the corresponding optimization problem.

If we want to include that the coefficients should be close to a  $\beta_{\text{prior}}$ , the ridge-regression estimator is the minimizer of the optimization problem:

$$\beta_{\text{RR}} := \arg \min_{\beta} \|y - X^T \beta\|_2^2 + \lambda \|\beta - \beta_{\text{prior}}\|_2^2$$

where  $\lambda > 0$  is a fixed regularization parameter.

- (b) Assume that the data are generated according to a linear model  $\tilde{y} := X^T \beta_{\text{true}} + \tilde{z}$ , where  $\beta_{\text{true}} \in \mathbb{R}^p$  and  $X \in \mathbb{R}^{p \times n}$  are fixed and  $\tilde{z}$  is an iid Gaussian random vector with zero mean and variance  $\sigma^2$ . Does the modification change the mean or the covariance matrix of the estimator? If so, report the new value.

$$\begin{aligned} \|y - X^T \beta\|_2^2 &= (y - X^T \beta)^T (y - X^T \beta) + \lambda (\beta - \beta_{\text{prior}})^T (\beta - \beta_{\text{prior}}) \\ &= (y^T - \beta^T X)(y - X^T \beta) + \lambda (\beta^T - \beta_{\text{prior}}^T)(\beta - \beta_{\text{prior}}) \\ &= \beta^T X X^T \beta - 2(Xy)^T \beta + y^T y + \lambda (\beta^T \beta - 2\beta_{\text{prior}}^T \beta + \beta_{\text{prior}}^T \beta_{\text{prior}}) \\ &= \beta^T X X^T \beta - 2(Xy)^T \beta + \lambda (\beta^T \beta - 2\beta_{\text{prior}}^T \beta) + y^T y + \lambda \beta_{\text{prior}}^T \beta_{\text{prior}} := f(\beta) \end{aligned}$$

$f$  is a quadratic form in  $\beta$  and its gradient and Hessian equal:

$$\begin{aligned} \nabla_{\beta} f(\beta) &= 2(X X^T \beta - Xy + \lambda(\beta - \beta_{\text{prior}}))I \\ \nabla_{\beta}^2 f(\beta) &= 2(X X^T + \lambda I) \end{aligned}$$

$v^T (X X^T + \lambda I) v = v^T X X^T v + \lambda v^T v = \|X^T v\|_2^2 + \lambda \|v\|_2^2 \geq 0$  and it is zero only when  $v$  is the zero vector (by property of the norm operator), thus the matrix  $X X^T + \lambda I$  is positive definite and invertible. The unique minimum can be found by setting the gradient to zero, which gives that the ridge regression estimator is:

$$\beta_{\text{RR}} = (X X^T + \lambda I)^{-1} (Xy + \lambda \beta_{\text{prior}} I)$$

with mean:

$$\begin{aligned} \mathbb{E}[\beta_{\text{RR}}] &= \mathbb{E}[(X X^T + \lambda I)^{-1} (X X^T \beta_{\text{true}} + X \tilde{z} + \lambda \beta_{\text{prior}} I)] \\ &= (X X^T + \lambda I)^{-1} (X X^T \beta_{\text{true}} + \lambda \beta_{\text{prior}} I) \end{aligned}$$

by linearity of the expectation and  $\tilde{z}$  has zero mean

Let  $X = U S V^T$  the SVD of  $X$ , we can then expand the previous expression:

$$\begin{aligned} \mathbb{E}[\beta_{\text{RR}}] &= (U S^2 U^T + \lambda U U^T)^{-1} (U S^2 U^T \beta_{\text{true}} + \lambda \beta_{\text{prior}}) \\ &= U (S^2 + \lambda I)^{-1} S^2 U^T \beta_{\text{true}} + \lambda \beta_{\text{prior}} U (S^2 + \lambda I)^{-1} U^T \end{aligned}$$

and variance:

$$\begin{aligned}
\text{Var}(\beta_{\text{RR}}) &= \text{Var}((XX^T + \lambda I)^{-1}(Xy + \lambda\beta_{\text{prior}}I)) \\
&= \text{Var}((XX^T + \lambda I)^{-1}Xy + \lambda\beta_{\text{prior}}(XX^T + \lambda I)^{-1}) \\
&= \text{Var}((XX^T + \lambda I)^{-1}Xy) \\
&= (XX^T + \lambda I)^{-1}X\text{Var}(y)((XX^T + \lambda I)^{-1}X)^T \\
&= \sigma^2(XX^T + \lambda I)^{-1}XX^T(XX^T + \lambda I)^{-1} \\
&= \sigma^2U(S^2 + \lambda I)^{-1}S^2(S^2 + \lambda I)^{-1}U^T
\end{aligned}$$

- (c) How can you incorporate this prior knowledge if you are using gradient descent with early stopping? Write the corresponding update equation as a function of  $\beta_{\text{prior}}$ .
- (d) Assume that the data are generated according to the linear model described above. Does the modification change the mean or the covariance matrix of the estimator? If so, report the new value.
4. The code you will implement in this question is located in the `regress.py` file in the time folder of `hw5.zip`. Define a sequence of random variables as follows:

$$\begin{aligned}
\vec{x}[0] &= 1 \\
\vec{x}[1] &= \vec{x}[0] + \vec{z}[1] \\
\vec{x}[2] &= \vec{x}[1] + \vec{z}[2] \\
\vec{x}[3] &= \vec{x}[2] + \vec{z}[3] \\
\vec{x}[4] &= \vec{x}[3] + \vec{z}[4],
\end{aligned}$$

where  $\vec{z}[1], \vec{z}[2], \vec{z}[3], \vec{z}[4]$  are independent,  $\vec{z}[1] \sim \mathcal{N}(0, 1)$  and  $\vec{z}[2], \vec{z}[3], \vec{z}[4] \sim \mathcal{N}(0, 0.01^2)$ . There is a function  $f : \mathbb{R}^5 \rightarrow \mathbb{R}$  of the form  $f(x) = \vec{\beta}^T x$  where  $\vec{\beta}$  is unknown. We are given a training sample of independent draws

$$(\vec{x}_1, f(\vec{x}_1) + \tilde{w}_1), \dots, (\vec{x}_n, f(\vec{x}_n) + \tilde{w}_n) \in \mathbb{R}^5 \times \mathbb{R},$$

where  $\tilde{w}_i$  are iid standard normal random variables corrupting our measurements of  $f$ . Using this training data, we will estimate  $\vec{\beta}$  and test our performance on a validation set drawn from the same distribution. Below we refer to the square loss function  $L : \mathbb{R}^5 \times \mathbb{R}^{n \times 5} \times \mathbb{R}^n \rightarrow \mathbb{R}$  defined by

$$L(\hat{\beta}, X, y) = \sum_{i=1}^n (X[i, :] \hat{\beta} - y[i])^2$$

where  $X \in \mathbb{R}^{n \times 5}$  denotes a matrix of data (training or validation; each row is a data point), and  $y$  is the corresponding vector of  $f$ -values.

- (a) Using least squares (i.e., minimizing the square loss on the training set) compute an estimate for  $\vec{\beta}$ . Include your estimate for  $\vec{\beta}$ , your square loss on the training set, and your square loss on the validation set in your submission. [Hint: If computed correctly your training loss should be larger than 30 and your validation loss should be larger than 10.]

- (b) Compute the singular values of the training data matrix  $X \in \mathbb{R}^{n \times 5}$ .
- (c) The true value of  $\vec{\beta}$  can be found at the top of `regress.py`. Give an explanation as to why the least squares estimates aren't close to the true  $\vec{\beta}$ -values.
- (d) Use ridge regression to produce a new estimate of  $\vec{\beta}$  and report the resulting estimate of  $\vec{\beta}$ , and your square loss on the training and validation sets. Here  $\hat{\beta}$  should solve

$$\text{minimize}_{\vec{\eta}} \quad \|X\vec{\eta} - \vec{y}\|_2^2 + 0.5\|\vec{\eta}\|_2^2.$$

You're not required to include your code in your submission, but you are free to do so.