

Exercise 7

Fundamentals of Machine Learning 8 ECTS
Heidelberg University WiSe 20/21

Catherine Knobloch, Elias Olofsson, Julia Siegl

February 9, 2021

2 Proof - Ridge Regression - Primal vs. Dual (10 pts)

In the primal formulation, the ridge regression problem takes the following form

$$\hat{\beta} = \operatorname{argmin}_{\beta} \|\mathbf{y} - X\beta\|_2^2 + \tau\|\beta\|_2^2, \quad (1)$$

where X is an $N \times D$ matrix, β is a D -dimensional vector and \mathbf{y} is an N -dimensional vector. As you saw in the lecture, the optimal $\hat{\beta}$ is given by

$$\hat{\beta} = (X^T X + \tau \mathbb{1}_D)^{-1} X^T \mathbf{y}. \quad (2)$$

Here $\mathbb{1}_D$ is the D -dimensional unit matrix. You also know that the dual formulation of the problem is given by

$$\hat{\alpha} = \operatorname{argmax}_{\alpha} -\alpha^T (X X^T + \tau \mathbb{1}_N) \alpha + 2\alpha^T \mathbf{y}, \quad (3)$$

with the solution for $\hat{\alpha}$

$$\hat{\alpha} = (X X^T + \tau \mathbb{1}_N)^{-1} \mathbf{y}. \quad (4)$$

For each feasible α , a corresponding feasible β is given by

$$\beta = X^T \alpha. \quad (5)$$

Prove that the optimal $\hat{\beta}$ corresponds to the optimal $\hat{\alpha}$

$$\hat{\beta} = X^T \hat{\alpha}. \quad (6)$$

Hint: Prove the following lemma (e.g. using the SVD of X), which will be useful in your derivation

$$(X^T X + \tau \mathbb{1}_D)^{-1} X^T = X^T (X X^T + \tau \mathbb{1}_N)^{-1}. \quad (7)$$

To prove the relation in Eq.(6), we use the individually derived expressions for the optimal $\hat{\beta}$ and $\hat{\alpha}$ in Eq.(2) and Eq.(4) and substitute these into Eq.(6), obtaining

$$(X^T X + \tau \mathbb{1}_D)^{-1} X^T \mathbf{y} = X^T (X X^T + \tau \mathbb{1}_N)^{-1} \mathbf{y}, \quad (8)$$

where we can simplify the equation by removing the response vector \mathbf{y} from each side, yielding

$$(X^T X + \tau \mathbb{1}_D)^{-1} X^T = X^T (X X^T + \tau \mathbb{1}_N)^{-1}, \quad (9)$$

which is the same equation as the one given in the hint, Eq.(7). To prove that this equality holds, we left-multiply Eq.(9) with $(X^T X + \tau \mathbb{1}_D)$ and right-multiply with $(X X^T + \tau \mathbb{1}_N)$, which then gives

$$X^T (X X^T + \tau \mathbb{1}_N) = (X^T X + \tau \mathbb{1}_D) X^T. \quad (10)$$

By finally multiplying the X^T :s into each parenthesis, we can trivially see that

$$\text{LHS} = X^T X X^T + \tau X^T = X^T X X^T + \tau X^T = \text{RHS}, \quad (11)$$

for any training set X and any regularization parameter τ . We have thus shown that the relation in Eq.(6) holds true, given the expressions Eq.(2) and Eq.(4) for the optimal $\hat{\beta}$ and $\hat{\alpha}$ of the ridge regression problem. \square