

1 Bias and variance of ridge regression (8 points)

Ridge regression solves the regularized least squares problem

$$\hat{\beta}_\tau = \arg\min_{\beta} (y - X\beta)^T (y - X\beta) + \tau \beta^T \beta$$

with regularization parameter $\tau \geq 0$. Regularization introduces some bias into the solution in order to achieve a potentially large gain in variance. Assume that the true model is $y = X\beta^* + \epsilon$ with zero-mean Gaussian noise $\epsilon \sim \mathcal{N}(0, \sigma^2)$ and centered features $\frac{1}{N} \sum_i X_i = 0$ (note that these assumptions imply that y is also centered in expectation). Prove (e.g. using the SVD of X) that expectation and covariance matrix of the regularized solution (both taken over all possible training sets of size N) are then given by

$$\mathbb{E}[\hat{\beta}_\tau] = S_\tau^{-1} S \beta^* \quad \text{Cov}[\hat{\beta}_\tau] = S_\tau^{-1} S S_\tau^{-1} \sigma^2$$

where S and S_τ are the ordinary and regularized scatter matrices:

$$S = X^T X \quad S_\tau = X^T X + \tau I_D$$

Notice that expectation and covariance reduce to the corresponding expressions of ordinary least squares (as derived in the lecture) when $\tau = 0$:

$$\mathbb{E}[\hat{\beta}_{\tau=0}] = \beta^* \quad \text{Cov}[\hat{\beta}_{\tau=0}] = S^{-1} \sigma^2$$

Since S_τ is greater than S (in any norm), regularization has a shrinking effect on both expectation and covariance.

$$\begin{aligned} \mathbb{E}[\hat{\beta}_\tau] &= (X^T X + \tau I_D)^{-1} X^T y \\ &= (X^T X + \tau I_D)^{-1} X^T (X \beta^* + \epsilon) \\ &= \underbrace{(X^T X + \tau I_D)^{-1} X^T X}_{S_\tau^{-1} S} \underbrace{\beta^*}_{\mathbb{E}[\beta^*]} + \underbrace{(X^T X + \tau I_D)^{-1} X^T \epsilon}_{\mathbb{E}[\epsilon] = 0} \\ &\Rightarrow \mathbb{E}[\hat{\beta}_\tau] = S_\tau^{-1} S \beta^* \end{aligned}$$

$$\text{bias} = \hat{\beta}_\tau - \beta^* = (X^T X + \tau I_D)^{-1} (X^T X \beta^* + X^T \epsilon) - \beta^* \Rightarrow \text{In the lecture, this step was introduced and therefore I used it again. But I should have just used the above proven term instead.}$$

$$\text{Cov}(\hat{\beta}_\tau) = \mathbb{E}[(\hat{\beta}_\tau - \mathbb{E}[\hat{\beta}_\tau])(\hat{\beta}_\tau - \mathbb{E}[\hat{\beta}_\tau])^T] = \mathbb{E}[(X^T X + \tau I_D)^{-1} (X^T X \beta^* + X^T \epsilon) - \beta^*] (X^T X + \tau I_D)^{-1} X^T X (X^T X + \tau I_D)^{-1} \sigma^2$$

Most probably, there is something wrong here which I can't seem to find/solve

Number 2:

$$\frac{\partial}{\partial \beta} \sum_{i=1}^N (y_i^* - X_i \cdot \beta)^2 \stackrel{!}{=} 0 \quad (1)$$

$$\Sigma \cdot \beta + \frac{1}{4} (\mu_1 - \mu_{-1})^T \cdot (\mu_1 - \mu_{-1}) \cdot \beta = \frac{1}{2} (\mu_1 - \mu_{-1})^T \quad (2)$$

$$\frac{\partial}{\partial \beta} \sum_{i=1}^N (y_i^* - X_i \cdot \beta)^2 = 0$$

$$\Rightarrow \sum_{i=1}^N (-2 y_i^* X_i + 2 X_i^2 \beta) = 0$$

$$\Rightarrow 2 \cdot \sum_{i=1}^N (X_i^2 \beta - y_i^* X_i) = 0$$

$$\Rightarrow \sum_{i=1}^N (X_i^2 \beta) = \sum_{i=1}^N y_i^* X_i$$

$$\Rightarrow \beta \cdot \sum_{i=1}^N X_i^2 = \sum_{i=1}^N y_i^* X_i$$

$$\Rightarrow \sum_{i=1}^N \beta X_i^2 = \sum_{i=1}^N y_i^* X_i \quad ?$$

up to this step quite correct, but instead of X_i^2 $X_i X_i^T$ and instead of $y_i^* X_i$, $y_i^* X_i^T$ would have been correct. The above error leaves me at this point as I am unable to rewrite the solution as explained in the sample-solution

I have no clue how to go on from this or if it is even remotely correct so far