

IML - 5 תרגיל  
 ממשל: נניח

1. (a) נניח  $\sqrt{\frac{1}{2m} \ln\left(\frac{2}{\delta}\right)} = \epsilon$  ונניח  $\epsilon$  קטן מספיק

$$P\left(|L_{\text{Sall}}(h) - L_D(h)| \leq \sqrt{\frac{\ln\left(\frac{2}{\delta}\right)}{2m}}\right) \geq 1 - \delta$$

לכן עבור  $h \in H_n$  וכל  $\delta$  קטן מספיק

$$P\left(|L_{\text{Sall}}(h_i) - L_D(h_i)| \geq \sqrt{\frac{\ln\left(\frac{2|H_n|}{\delta}\right)}{2m}}\right) \leq \frac{\delta}{|H_n|}$$

נניח  $\delta = \frac{\epsilon^2}{|H_n|}$  ונניח  $\epsilon$  קטן מספיק

$$\delta \leq \frac{|H_n| \epsilon^2}{|H_n|} = \epsilon^2$$

$$|L_{\text{Sall}}(h) - L_D(h)| \leq \sqrt{\frac{\ln\left(\frac{2|H_n|}{\delta}\right)}{2m}}$$

$$\Rightarrow P\left(L_D(h^*) \leq \min_{h \in H_n} L_D(h) + \sqrt{\frac{2 \ln\left(\frac{2|H_n|}{\delta}\right)}{m}}\right) \geq 1 - \delta$$

1. (b) נניח  $\epsilon$  קטן מספיק

$$P\left(L_D(h^*) \leq \min_{h \in H_n} L_D(h) + \sqrt{\frac{2 \ln\left(\frac{4|H_n|}{\delta}\right)}{2m}}\right) \geq 1 - \frac{\delta}{2}$$

נניח  $\delta = \frac{\epsilon^2}{|H_n|}$  ונניח  $\epsilon$  קטן מספיק

$$P\left(L_D(h) \leq \min_{h \in H_n} L_D(h) + \sqrt{\frac{2 \ln\left(\frac{4|H_n|}{\delta}\right)}{(1-\alpha)m}}\right) \geq 1 - \frac{\delta}{2}$$

נניח  $\delta = \frac{\epsilon^2}{|H_n|}$  ונניח  $\epsilon$  קטן מספיק

$$L_D(h^*) \leq L_D(h_j) + \sqrt{\frac{2 \ln\left(\frac{4|H_n|}{\delta}\right)}{2m}} \leq \min_{h \in H_n} L_D(h) + \sqrt{\frac{2 \ln\left(\frac{4|H_n|}{\delta}\right)}{(1-\alpha)m}} + \sqrt{\frac{2 \ln\left(\frac{4|H_n|}{\delta}\right)}{2m}}$$



$$= \min_{h \in \mathcal{H}_n} L_b(h) + \sqrt{\frac{2 \ln(4|\mathcal{H}_n|/r)}{(1-\alpha)m}} + \sqrt{\frac{2 \ln(4n/r)}{2m}}$$

הנה מתקן הנוכחי  $|\mathcal{H}_n| = 2^{2^n}$ . נראה שהערך  $L_b(h)$  הוא קבוע (constant).

$$L_b(h^*) \leq \min_{h \in \mathcal{H}_n} L_b(h) + \sqrt{\frac{2 \ln(2^{2^n}/r)}{m}}$$

אם  $B > 1/m$  אז

~~$$L_b(h^*) \leq \min_{h \in \mathcal{H}_n} L_b(h) + \sqrt{\frac{2 \ln(2^{2^n}/r)}{m}}$$~~

$$L_b(h^*) \leq \min_{h \in \mathcal{H}_n} L_b(h) + \sqrt{\frac{2 \ln(2^{2^n}/r)}{(1-\alpha)m}} + \sqrt{\frac{2 \ln(4n/r)}{2m}}$$

אם  $B > 1/m$

אם  $B > 1/m$  אז  $V \rightarrow B$  ונראה שהערך  $L_b(h)$  הוא קבוע.

$$\lim_{n \rightarrow \infty} \frac{B}{V} = \infty$$

העקף של  $B$  הוא קבוע  $B$  ונראה שהערך  $L_b(h)$  הוא קבוע. לכן, הערך  $L_b(h)$  הוא קבוע.



2. (א) ridge regression

$$\hat{w}_\lambda^{\text{ridge}} = (X^T X + \lambda I)^{-1} X^T y$$

(ב) orthogonal design

$$\hat{w}_\lambda^{\text{ridge}} = X y (I + \lambda I)^{-1} = X y ((1 + \lambda) I)^{-1} = \frac{\hat{w}^{\text{LS}}}{1 + \lambda}$$

(ב) ridge regression is orthogonal design  $\hat{w}^{\text{LS}} = X^T y$

we have a loss function  $\|y - Xv\|^2$  over  $v \in \mathbb{R}^d$

$$\|y - Xv\|^2 = \sum_{i=1}^n (\hat{w}_i^{\text{LS}} - v_i)^2$$

we want to minimize the loss

$$\argmin_{v \in \mathbb{R}^d} \sum_{i=1}^n (\hat{w}_i^{\text{LS}} - v_i)^2 + \lambda \|v\|_0$$

(c) ridge regression

$$v_i = 0 \quad \text{if} \quad |\hat{w}_i^{\text{LS}}| \leq \sqrt{\lambda} \quad \text{and} \quad (\hat{w}_i^{\text{LS}})^2 \leq \lambda$$

$$v_i = \hat{w}_i^{\text{LS}} \quad \text{if} \quad (\hat{w}_i^{\text{LS}})^2 > \lambda$$

$$\hat{w}_\lambda^{\text{subset}} = \frac{1}{\sqrt{\lambda}} \eta_{\sqrt{\lambda}}(\hat{w}^{\text{LS}})$$

$$A_\lambda \hat{w} = (X^T X + \lambda I)^{-1} (X^T X) (X^T X)^{-1} X^T y$$

$$= (X^T X + \lambda I)^{-1} X^T y = \hat{w}(\lambda)$$

$$E[\hat{w}(\lambda)] = A_\lambda E[\hat{w}] = (X^T X + \lambda I)^{-1} (X^T X) w$$

$$E[\hat{w}(\lambda)] \neq w \quad \text{for} \quad \lambda > 0$$

$$\text{Var}(\hat{w}(\lambda)) = A_\lambda \text{Var}(\hat{w}) A_\lambda^T = \sigma^2 A_\lambda (X^T X)^{-1} A_\lambda^T \quad (c)$$

א) ממונה על המבחן וההערכה

- $bias^2(\lambda) = \|(A_\lambda - I)w\|^2 = w^T (A_\lambda - I)^T (A_\lambda - I) w$
- $Var(\lambda) = \sigma^2 Tr(A_\lambda (X^T X)^{-1} A_\lambda^T)$

ז"ל  
הי  
הערכה (הערכה)  
הערכה

$$MSE(\lambda) = bias^2(\lambda) + Var(\lambda)$$

הערכה על המבחן וההערכה

$$\frac{d}{d\lambda} bias^2(\lambda) \Big|_{\lambda=0} = \frac{d}{d\lambda} \left( \sum_i \left( \sum_j (A_\lambda - I)_{ij} w_j \right)^2 \right) = 0$$

הערכה

$$\frac{d}{d\lambda} Var(\lambda) \Big|_{\lambda=0} = \sum_{ij} \frac{\partial Var(\lambda)}{\partial (X^T X + \lambda I)_{ij}} \cdot \frac{d(X^T X + \lambda I)_{ij}}{d\lambda} =$$

$$Tr \left( \frac{d Var(\lambda)}{d(X^T X + \lambda I)} \cdot \frac{d(X^T X + \lambda I)}{d\lambda} \right) =$$

$$= -2\sigma^2 Tr((X^T X)^{-1} (X^T X)^{-1}) < 0$$

הערכה על המבחן וההערכה

$$\frac{d Var(\lambda)}{d\lambda} = -2\sigma^2 ((X^T X)^{-1} (X^T X)^{-1})$$

$$\frac{d}{d\lambda} MSE(\lambda) \Big|_{\lambda=0} < 0 \quad \text{על המבחן וההערכה}$$

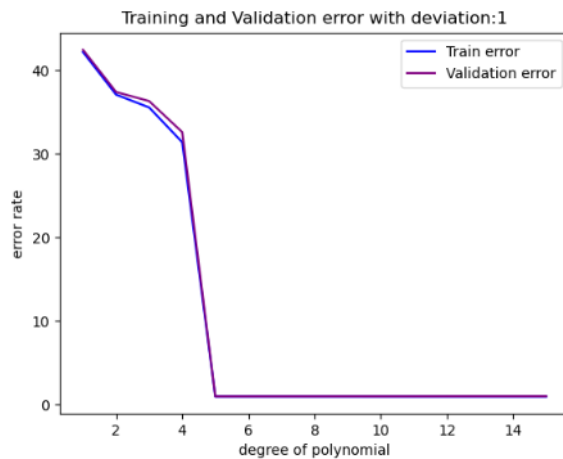
הערכה על המבחן וההערכה

הערכה על המבחן וההערכה

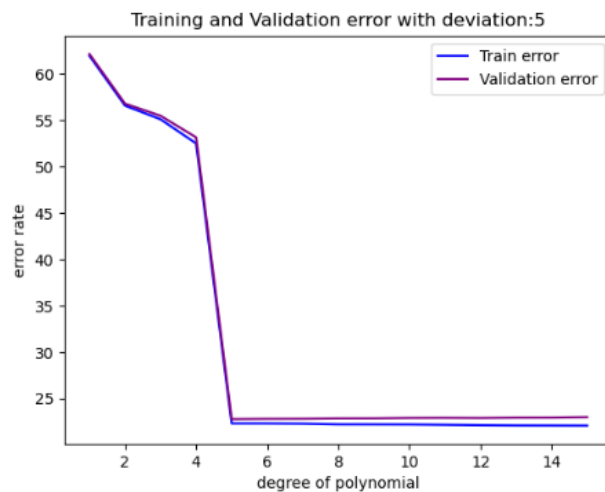
4.

e)

The degree of the polynomial with the lowest validation error for deviation 1 is 5



The degree of the polynomial with the lowest validation error for deviation 5 is 5





g)

Both polynomials ( $h^*$ ) that performed the best by the ERM principal were of degree 5, and the following is their test error:

test error for deviation 1: 0.5244495414916898

test error for deviation 5: 13.764926054734158

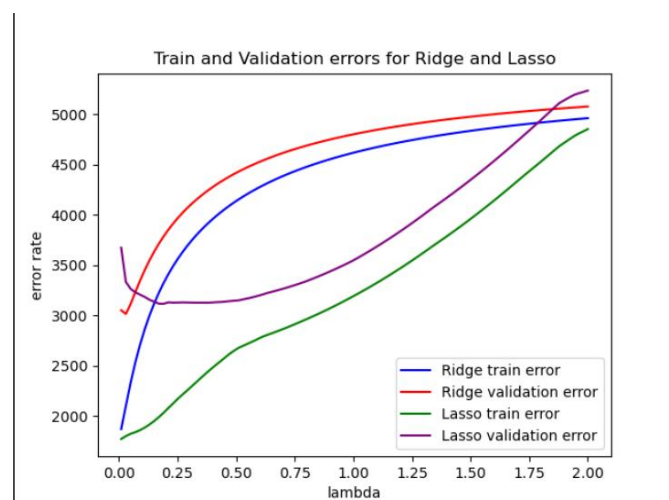
The validation error in both cases was similar to the test error

5.

cii)

The values of lambda that I chose to check were between 0.01(close to zero) and 2. I started with such a low number as to see if there is a difference between using a very small regularization term that is close to zero, which is very close to regular Linear regression as we have learned. By checking enough examples of lambda between these two numbers, we can clearly tell by the graph that there is a ceil to how much we can improve our classification if our regularization term(lambda) is higher.

d)



e)

Best lambda for Ridge is 0.0301010101010101

Best lambda for Lasso is 0.19090909090909092

f)

Test error for best Ridge hypothesis 3216.0093244988275

Test error for best Lasso hypothesis  
3392.0050650611543

Test error for best Linear Regression 3612.249688324898

g)

We can clearly tell from section f) that the regularization has indeed helped in getting better results, compared to regular Linear Regression.