

DAC

2.1

ר. ו. ג. 3.2

(b) \leq (a)

ר. ו. ג. 3.2

1

$$\lim_{n \rightarrow \infty} \mathbb{P}_{S_n \sim \Omega^n} [L_0(A(S)) \leq \varepsilon] \geq 1 - \delta$$

$$\lim_{n \rightarrow \infty} \mathbb{P}_{S_n \sim \Omega^n} [L_0(A(S)) < \varepsilon] \geq 1 - \delta$$

$$\forall \varepsilon > 0 \exists m_0 \in \mathbb{N} \text{ s.t. } \forall n \geq m_0, \mathbb{E}_{S_n \sim \Omega^n} [L_0(A(S))] < \varepsilon$$

$$\mathbb{P}_0 = P(L_0(A(S)) \leq \varepsilon) \geq 1 - \delta$$

$$m \geq m_0 \Rightarrow \mathbb{P}_{S_m \sim \Omega^m} [L_0(A(S)) \leq \varepsilon]$$

$$\mathbb{E}_{S_m \sim \Omega^m} [L_0(A(S))] = \sum_{\varepsilon=0}^{\frac{\varepsilon}{2}} \mathbb{P}_{S_m \sim \Omega^m} [L_0(A(S)) = \varepsilon] \cdot \varepsilon$$

$$\leq \sum_{\varepsilon=0}^{\frac{\varepsilon}{2}} \mathbb{P}_{S_m \sim \Omega^m} [L_0(A(S)) = \varepsilon] \cdot \varepsilon + \sum_{\varepsilon=\frac{\varepsilon}{2}}^1 \mathbb{P}_{S_m \sim \Omega^m} [L_0(A(S)) = \varepsilon] \cdot \varepsilon$$

$$= \frac{\varepsilon}{2} \cdot \mathbb{P}_{S_m \sim \Omega^m} [L_0(A(S)) \leq \varepsilon/2] + \mathbb{P}_{S_m \sim \Omega^m} [L_0(A(S)) > \varepsilon/2] \cdot \frac{\varepsilon}{2} = \delta.$$

$$\mathbb{P}_{S_m \sim \Omega^m} [L_0(A(S)) \leq \varepsilon/2] \geq 1 - \delta$$

$$\mathbb{P}_{S_m \sim \Omega^m} [L_0(A(S)) > \varepsilon] \leq \frac{\mathbb{E}_{S_m \sim \Omega^m} [L_0(A(S))]}{\varepsilon} = \frac{\mathbb{E}_{S_m \sim \Omega^m} [L_0(A(S))]}{\varepsilon} = \delta$$

אך נשים $\delta \rightarrow 0$ נסובב

לעבב א ERM מילון 100% 2

$$\text{לעבב } S = \{(x_i, y_i)\}_{i=1}^m \subseteq \mathbb{R}^2 \times \{0,1\}^m \quad |N| \approx 60$$

לעבב $x_i \in \mathbb{R}^2$ מילון נסוב בדרכו (אך לא כפולה)

לעבב $y_i \in \{0,1\}^m$ (אך לא כפולה)

לעבב (x_i, y_i) מילון נסוב בדרכו (אך לא כפולה)

לעבב $(x_i, y_i) \in \mathbb{R}^2 \times \{0,1\}^m$ מילון נסוב בדרכו (אך לא כפולה)

לעבב $x_i \in \mathbb{R}^2$ מילון נסוב בדרכו (אך לא כפולה)

$$\Pr_{r \sim \mathcal{R}}(\{y_i : r \in \mathcal{R} \mid y_i = 1\}) = \varepsilon$$

לעבב $E = \{x \in \mathbb{R}^2 \mid r \leq \|x\| \leq R\}$ מילון נסוב בדרכו (אך לא כפולה)

לעבב $E = \mathbb{R}^2$ מילון נסוב בדרכו (אך לא כפולה)

לעבב $e^{-\epsilon m} \geq (1-\varepsilon)^m$ מילון נסוב בדרכו (אך לא כפולה)

לעבב $e^{-\epsilon m} \leq \frac{1}{2}$ מילון נסוב בדרכו (אך לא כפולה)

לעבב $\sum_{i=1}^m \delta(x_i) = 100$ מילון נסוב בדרכו (אך לא כפולה)

לעבב $\sum_{i=1}^m \delta(x_i) = 100$ מילון נסוב בדרכו (אך לא כפולה)

$$z = \frac{23}{3}$$

m_{VC}^{upper}

פוקט λ נייד δ , S מילוי ϵ, δ (ככל ש)

אנו מודים $m_{VC}(\epsilon/2, \delta) = m$

אנו מודים $m_{VC}(\epsilon/2, \delta) = m$

אנו מודים $m_{VC}(\epsilon/2, \delta) = m$

$$L_0(h_s) \leq L_0(h) + \frac{\epsilon}{2}$$

$$\sum_i l_i(h_s) \leq \sum_i l_i(h) + \frac{\epsilon}{2} \quad \forall i \in H \quad \text{כפי}$$

$$\leq L_0(h) + \frac{\epsilon}{2} + \frac{\epsilon}{2} = L_0(h) + \epsilon$$

לפיכך $L_0(h_s) \leq L_0(h) + \epsilon$

הוכחה סהרה

$m_{VC}(\epsilon/2, \delta) = m_{VC}(\epsilon, \delta)$

$L_0(h_s) \leq L_0(h) + \epsilon$

$H \subset \mathcal{C}$ (ולכן)

$$L_0(h) \leq \min_{h \in H} l_i(h) + \epsilon$$

$m_{VC}(\epsilon/2, \delta) \leq m_{VC}(\epsilon, \delta)$

$$m_{VC}(\epsilon, \delta) \leq m_{VC}(\epsilon/2, \delta)$$

VC-Dimension

2.2

1) $B \rightarrow D^3 \cup T^2 \cup \{B\}$ for $\{e_1, \dots, e_n\}$

$h_1(e_i) = y_i$ for $y_i \in \{0, 1\}$ for $i \in I = \{i \mid y_i = 1\}$

$|H| = 2^n - 2$ for H

$\nabla Cdim(H) \leq \log_2 |H| / 10^{10}$

• $Vcbim(H) = n - 2$

2) $H_1 \cup H_2$ for $H_1, H_2 \subseteq \{1, \dots, m\}$

H_1 for H_2 for $C = \{C_1, \dots, C_m\}$

$H_2 \subseteq P_2$ for $P_2 \subseteq \{1, \dots, m\}$

$\nabla Cdim(H_1) \leq \nabla Cdim(H_2)$

REGULARIZATION

2.3

$$\begin{aligned} A_2 \hat{w} &= (X^T X + \lambda I_d)^{-1} (X^T X)^{-1} X^T y \quad (1) \\ &= (X^T X + \lambda I_d)^{-1} X^T y = \hat{w}_2 \end{aligned}$$

$$\begin{aligned} E[\hat{w}_2] &= E[A_2 \hat{w}] = A_2 E[\hat{w}] \quad (2) \\ &= (X^T X + \lambda I_d)^{-1} (X^T X) w \\ \bullet E[\hat{w}_2] &\neq w \text{ (not } \Rightarrow \text{)} \quad |21| \end{aligned}$$

$$\begin{aligned} \text{Var}(\hat{w}_2) &= \text{Var}(A_2 \hat{w}) \quad (3) \\ &= A_2 \text{Var}(w) A_2^T = G^2 A_2 (X^T X)^{-1} A_2^T \end{aligned}$$

$$\begin{aligned} \text{MSE}(\bar{y}) &= E[|\hat{y} - y|^2] = E[(\hat{y} - \bar{y})^2] + |\bar{y} - y|^2 \quad (4) \\ &= \text{Var}(\hat{y}) + \text{bias}^2(\bar{y}) \end{aligned}$$

$$y = E[\hat{w}], \bar{y} = \hat{w}_0, \bar{y} = E[\hat{w}_0]: \text{Pf 3.1C}$$

$$\begin{aligned} \text{bias}^2(\bar{y}) &= |E[\hat{w}_0] - E[w]|^2 \quad |21| \\ &= |(A_2 - I)w|^2 \end{aligned}$$

$$\text{Var}(\bar{y}) = \text{Tr}(\text{Var}(\hat{w}_2)) = \bar{\lambda} \text{Tr}(A_2 (X^T X)^{-1} A_2^T)$$

$$\text{MSE}(\hat{w}_2) = \bar{\lambda} \text{Tr}(A_2 (X^T X)^{-1} A_2^T) + |(A_2 - I)w|^2$$

וריאנט $\hat{W}_{\lambda=0}$ מינימום MSE $\lambda=0$ ופה יזכיר
 $\hat{W}_{\lambda=0}$ מינימום MSE $\lambda \neq 0$
 $\text{MSE}(\hat{W}_{\lambda=0}) < \text{MSE}(\hat{W}_{\lambda>0}) = \text{MSE}(\hat{W})$
 $\hat{W} = X(X^T X)^{-1} X^T Y$

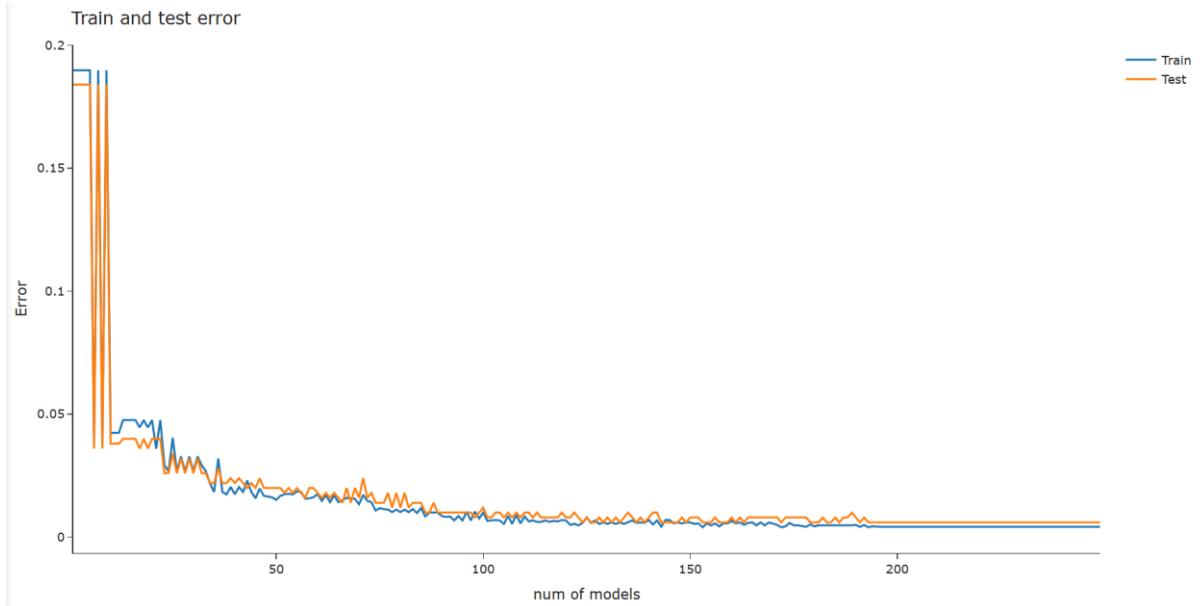
$\textcircled{1}$ $(X^T X)^{-1} X^T A = (X^{-1} X^T)^{-1} X^T A = (X^T)^{-1} A =$
 $X^T X$ מינימום MSE $\hat{W} = X^T X$

$\textcircled{2}$ $(X^T X)^{-1} X^T X V = (X^T)^{-1} X V =$
 $X^T X$ מינימום MSE $\hat{W} = X^T X$

$\textcircled{3}$ $f_{\hat{W}}(X) = f_{\hat{W}}(X^T X) + f_{\hat{W}}(X^T X V) =$
 $f_{\hat{W}}(X^T X) + f_{\hat{W}}(X^T X V)$

PRACTICAL PART-3.1.1

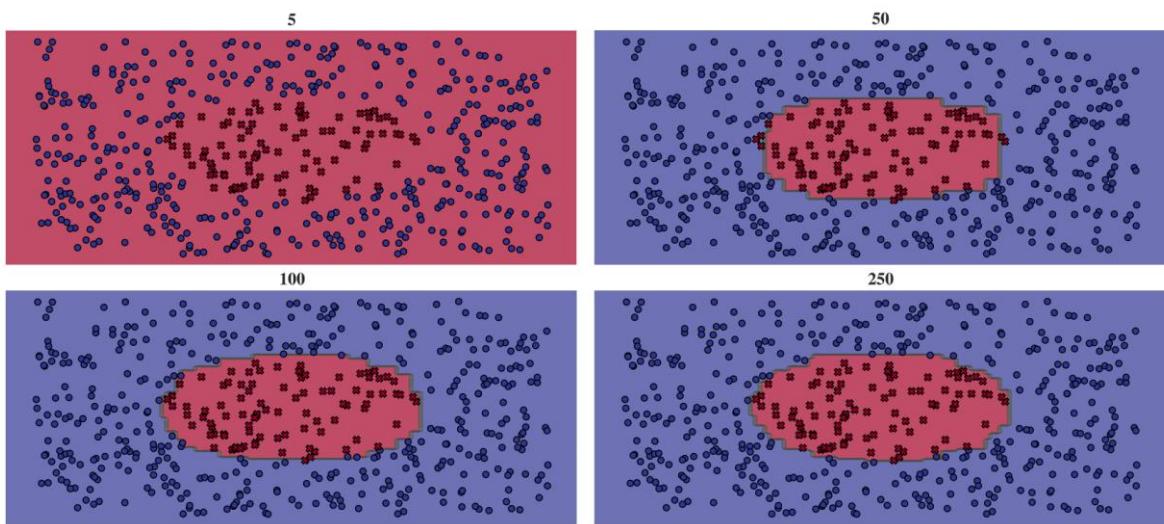
According to the graph, it can be seen that as the number of iterations increases, the error decreases. This demonstrates what we have learned, that boosting reduces bias more rapidly than it increases variance. Another interesting point observed in the graph is that overfitting does not occur. Although the test error is slightly higher than the training error, at no point do we see a significant increase in test error due to overfitting. (graph 1)



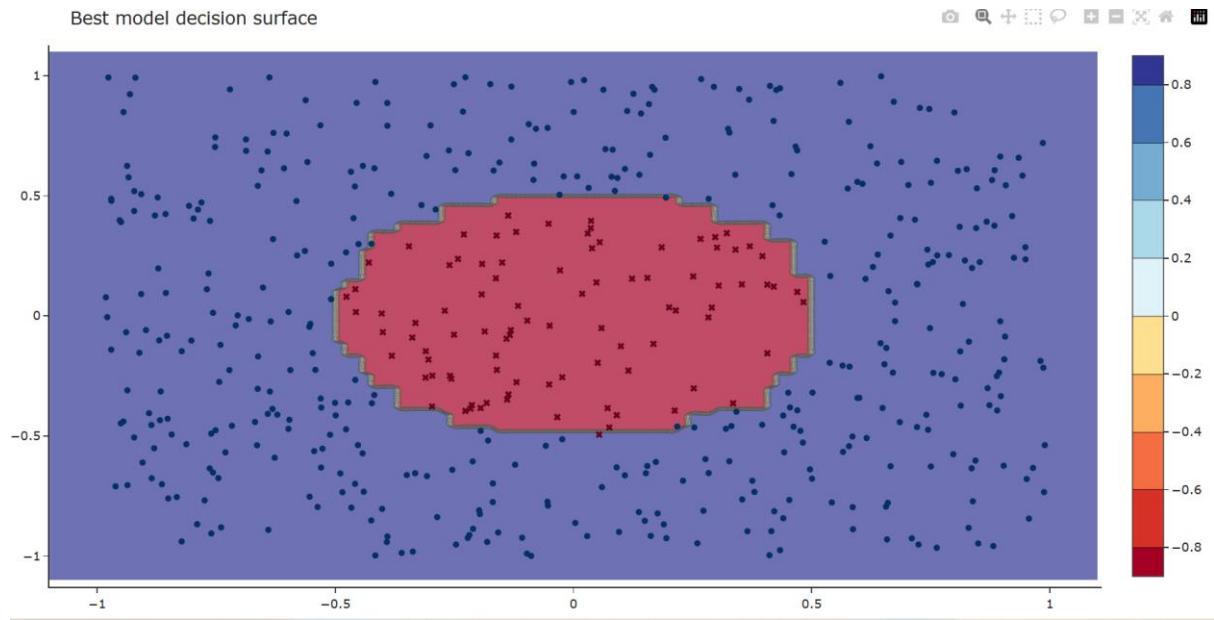
3.1.2

It can be seen that for a small number of iterations (5), the model fails to classify between the red and blue points (as seen in the graph from section 1 where the error is very high). However, when the number of iterations is slightly increased (to 50), it can be observed that the model performs very well in separating the points, as shown in the previous graph. Further increasing the number of iterations only slightly improves the result.

Decision surface for different iterations

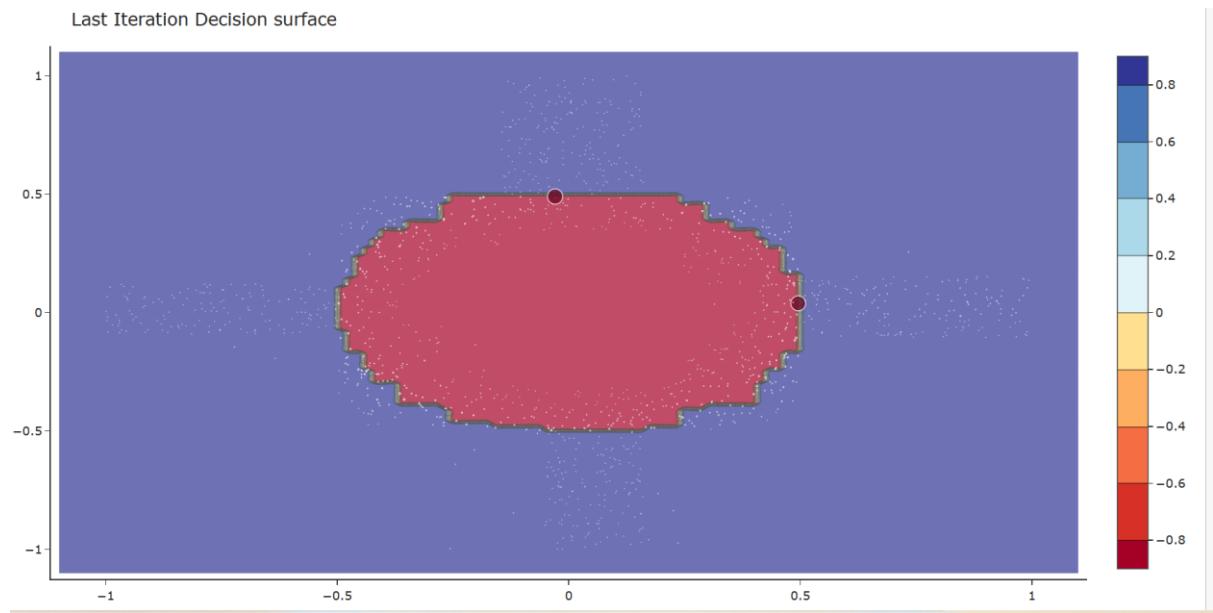


3.1.3



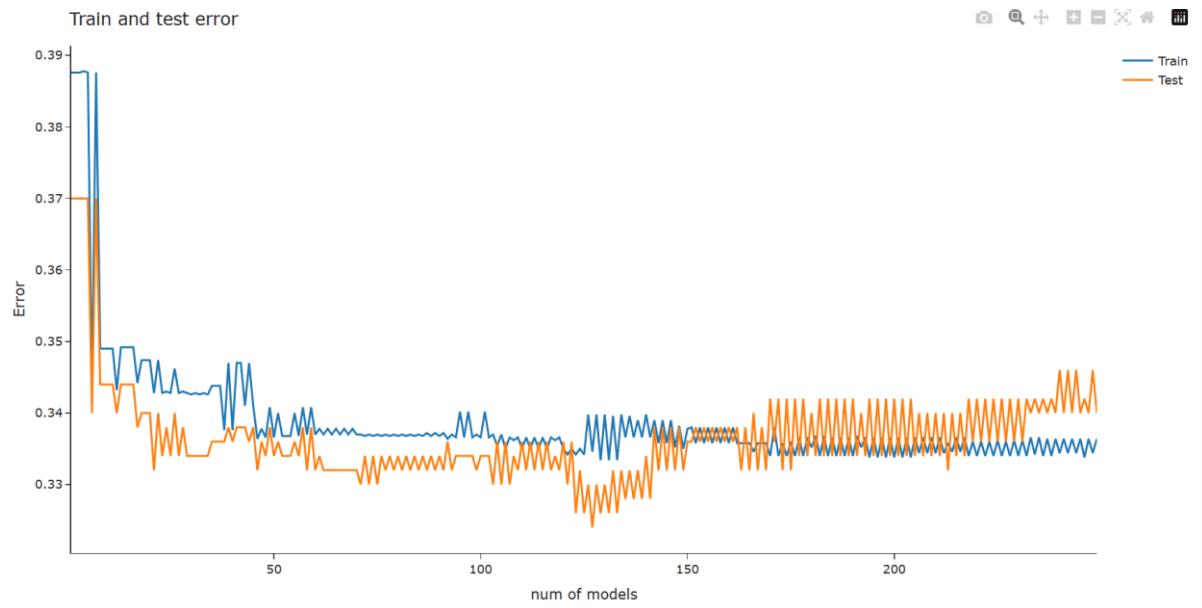
3.1.4

It can be observed that the points on the boundary receive higher weights (the points that are not close to the boundary are hardly visible). In other words, it is challenging for the model to classify them.



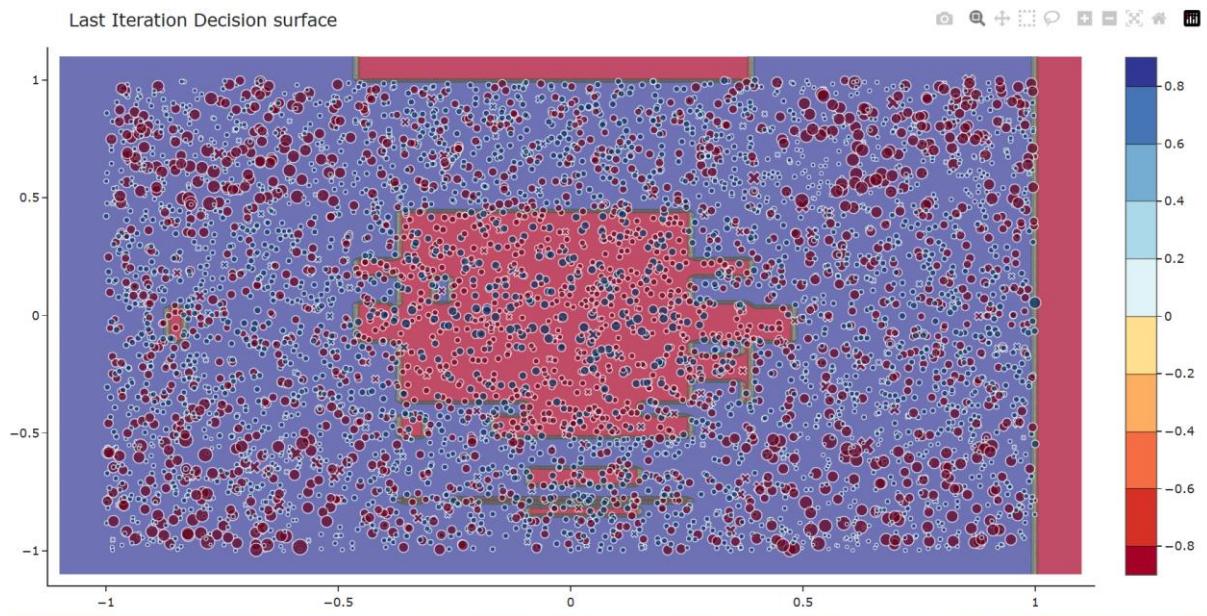
3.1.5

Due to the addition of a large amount of noise, we expect each sample to be significantly different from the other. However, despite this, in the graph, we can see that although the error has increased for both the training set and the test set, it is reasonable because, indeed, due to the high noise, many of the red points are now outside the red circle. Nevertheless, it can be observed that the error for the test set has not increased much compared to the error for the training set. This means that the generalization error has not increased significantly. This is further evidence that AdaBoost reduces bias significantly more than it increases variance



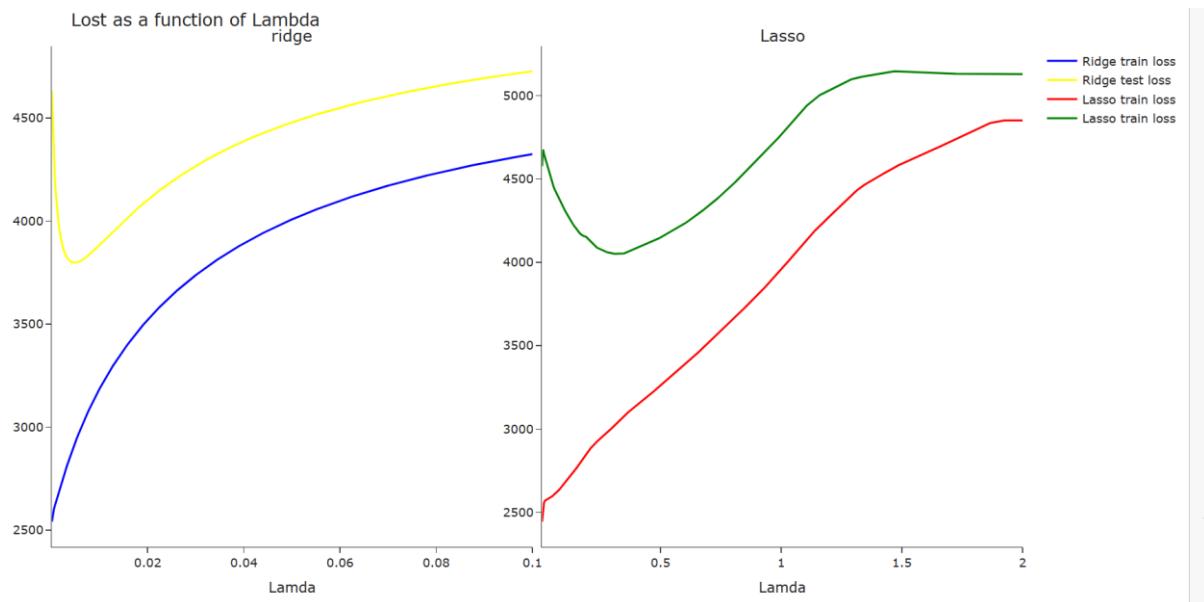
In this graph, it can be observed that despite the noise, the model manages to provide responses that are close to the truth, even though many of the red points are located outside the red circle and .many of the blue points are inside it

It can also be seen that the points in the corners receive the highest weight, which is because they are outside the range of the circle in both coordinates. Therefore, it is logical that most of the weak learners made mistakes regarding these points



3.2

It can be observed that in both graphs, the train_loss is minimal when λ (lambda) is minimal. This is because when $\lambda=0$, we obtain the solution of least squares, which minimizes the empirical error. When λ is increased, we get a simpler model, but its empirical error and generalization error also increase.



3.2.3

```
Best lambda for ridge is: 0.004904809619238477
Best lambda for lasso is: 0.3170741482965932
Best ridge loss is: 4330.988920435104
Best lasso loss is: 4372.179666085873
Linear Regression loss is: 8146.267790957025
```

```
>>>
```

