

1. Theoretical Part

1.1 Regularization

1)

a) $\hat{w}_\lambda = (X^T X + \lambda I_d)^{-1} (X^T y)$ as shown in class

$$\hat{w} = (X^T X)^{-1} X^T y \Leftrightarrow X^T y = X^T X \hat{w}$$

$$A_\lambda = (X^T X + \lambda I_d)^{-1} X^T X$$

$$\hat{w}_\lambda = (X^T X + \lambda I_d)^{-1} (X^T X) \hat{w} = A_\lambda \hat{w} \quad \blacksquare$$

d) $E[\hat{w}_\lambda] \stackrel{a}{=} E[(X^T X + \lambda I_d)^{-1} (X^T X) \hat{w}] = A_\lambda E[\hat{w}] = A_\lambda \cdot w + w$

c) $\hat{w}_\lambda = A_\lambda \hat{w} \Leftrightarrow \text{Var}(\hat{w}_\lambda) = \text{Var}(A_\lambda \hat{w})$

From the hint: where B is A_λ and z : \hat{w}

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

given that $\text{Var}(\hat{w}) = \sigma^2 (X^T X)^{-1}$

d) $\text{Bias}(\hat{w}_\lambda) = E[\hat{w}_\lambda] - w = A_\lambda w - w = (A_\lambda - I) w$

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

$$\text{Squarred Bias} = \| (A_\lambda - I) w \|^2 = w^T (A_\lambda - I)^T (A_\lambda - I) w$$

$$E[\| \hat{w}_\lambda - w \|^2] = E[(\hat{w}_\lambda - w)^T (\hat{w}_\lambda - w)] = E\left[\sum_{i=1}^d (\hat{w}_{\lambda i} - w_i)^2\right]$$

$$= \sum_{i=1}^d E[(\hat{w}_{\lambda i} - w_i)^2] = \sum_{i=1}^d \left(E[(w_{\lambda i} - E(w_i))^2] + (E(w_i) - w_i)^2 \right)$$

$$= E[\| w_\lambda - w \|^2] + \sum_{i=1}^d (E(w_i) - w_i)^2$$

$$= E[\| w_\lambda - w \|^2] + \| E(w) - w \|^2$$

↓

$$\text{Bias}(\hat{w}_\lambda) = E(\hat{w}_\lambda) - w = (A_\lambda - I) w$$

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

$$\Rightarrow \text{Bias}^2(\lambda_i) = (A_\lambda - I) w \Big|_i$$

$$\Rightarrow \text{Var}(\lambda_i) = \text{Var}(\hat{w}_\lambda) \Big|_i$$

b) Because when $\lambda=0$ then it's the normal ls so we search the best solution only relying on the data adding λ increases variance and bias but causes its negative we add more bias but remove a lot more variance

2)

$$a) Z = W_L^{Tik}$$

$$LS(w) = \|x_2 - y\|^2 + \|L_2\|^2 = (x_2 - y)^T(x_2 - y) + (L_2)^T(L_2)$$

\Downarrow $\overbrace{\hspace{10em}}$

$$= (x_2)^T \cdot x_2 - 2y^T x_2 + y^T y + z^T L^T L z$$

$$= z^T (X^T X + L^T L) z + 2y^T X z + y^T y$$

$$\nabla_{\mathbf{z}} L(\mathbf{w}) = 2(X^T X + L^T L)\mathbf{z} - 2y^T X = 0$$

$$w_L^{Tik} = (X^T X + L^T L)^{-1} X^T y \quad \text{回}$$

b) No it's not necessary because $X^T X$ is always PSD

so by adding LTL we can make $X^T X + L^T L$ invertible

Thus the new condition would be $X^T X + L^T L$ invertible.

$$c) i. \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}^T \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} = \begin{bmatrix} 1 & 3 \\ 2 & 4 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} = \begin{bmatrix} 10 & 14 \\ 14 & 20 \end{bmatrix} = X^T X$$

$$(X^T X)^{-1} = \frac{1}{10 \cdot 20 - 14^2} \begin{bmatrix} 20 & -14 \\ -14 & 10 \end{bmatrix} \quad X^T y = \begin{bmatrix} 23 \\ 34 \end{bmatrix}$$

$$L_S(w) = (X^T X)^{-1} X^T y = \begin{bmatrix} -4 \\ 9/2 \end{bmatrix} = w$$

$$\text{ii. } X^T X + I = \begin{bmatrix} 10 & 14 \\ 14 & 20 \end{bmatrix} + \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 11 & 14 \\ 14 & 21 \end{bmatrix}$$

$$LS(W) = \frac{1}{11 \cdot 21 - 14^2} \begin{bmatrix} 21 & -4 \\ -14 & 10 \end{bmatrix} \cdot \begin{bmatrix} 23 \\ 34 \end{bmatrix} = \begin{bmatrix} 1/5 \\ \frac{52}{35} \end{bmatrix} = \hat{w}_{\text{Ridge}}$$

$$\text{iii. } L^T L = \begin{bmatrix} 0 & 0 \\ 0 & 2 \end{bmatrix} \begin{bmatrix} 0 & 0 \\ 0 & 2 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 0 & 4 \end{bmatrix}$$

$$X^T X + L^T L = \begin{bmatrix} 10 & 14 \\ 14 & 24 \end{bmatrix}$$

$$LS(w) = \frac{1}{10 \cdot 29 - 19^2} \begin{bmatrix} 29 & -19 \\ -19 & 10 \end{bmatrix}, \begin{bmatrix} 2 \\ 3 \end{bmatrix} = \begin{bmatrix} 19/11 \\ 9/22 \end{bmatrix} = \hat{w}_L^{T+1}$$

$$i. \|\hat{w}\| = \sqrt{9 + 9/2} \left(\frac{-4}{9/2} \right) = \sqrt{\frac{9.16}{4} + \frac{81}{4}} = \frac{\sqrt{145}}{2} \approx 6$$

$$\|\hat{w}_{\text{Ridge}}\| = \left(\frac{1}{5}, \frac{32}{35} \right) \left(\frac{115}{32/35} \right) = \frac{\sqrt{2753}}{35} \approx 1.5$$

$$\|\hat{w}_L^{\text{Tik}}\| = \left(\frac{19}{11}, \frac{9}{22} \right) \left(\frac{19/11}{9/22} \right) = \frac{\sqrt{1525}}{22} \approx 1.8$$

The largest is \hat{w} because there is no penalization for large weights we just trying to fit the data in Ridge its small because we penalize both weights and w_L^{Tik} is in the middle because we punish just for one weight

$$ii. \hat{w}_2 = 9/2 = 4.5$$

$$\hat{w}_{\text{Ridge}} = 52/35 \approx 1.5$$

$$\hat{w}_L^{\text{Tik}} = 9/22 \approx 0.5$$

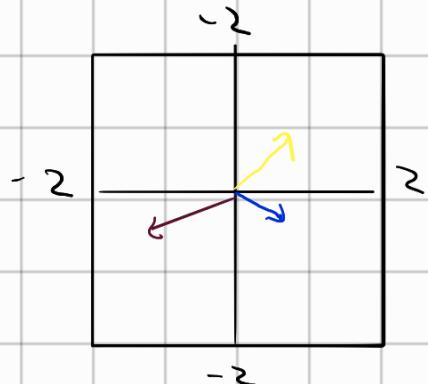
This happens for the same reasons above and because in Tik we punish more than Ridge
The second coordinate is smaller

1.2 PCA

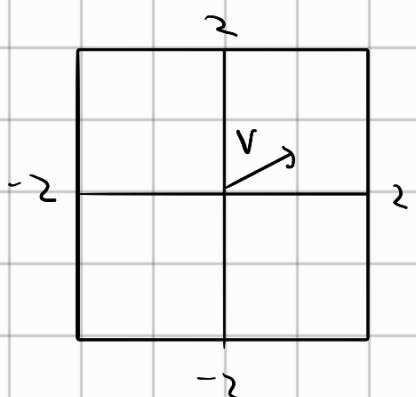
1.

a) $\bar{X} = \frac{1}{3} \sum_{i=1}^3 X_i = \begin{bmatrix} \frac{1}{3} \\ \frac{4}{3} \\ \frac{4}{3} \end{bmatrix}$

$$z_1 = \begin{bmatrix} \frac{2}{3} \\ \frac{1}{3} \\ \frac{2}{3} \end{bmatrix} \quad z_2 = \begin{bmatrix} -\frac{4}{3} \\ \frac{1}{3} \\ -\frac{1}{3} \end{bmatrix} \quad z_3 = \begin{bmatrix} \frac{2}{3} \\ \frac{1}{3} \\ -\frac{1}{3} \end{bmatrix}$$



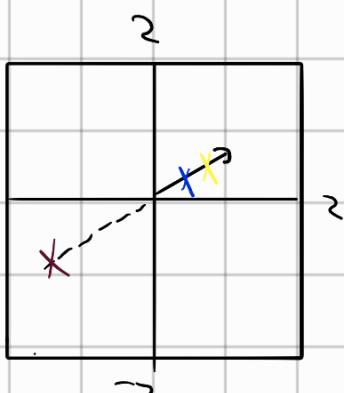
b) $\Sigma = \begin{bmatrix} 9/9 & 2/9 \\ 2/9 & 2/9 \end{bmatrix} \quad V \propto \begin{pmatrix} 0.95 \\ 0.29 \end{pmatrix}$



c) $\sqrt{\lambda_1} = 0.92$ X

$\sqrt{\lambda_2} = -1.37$ X

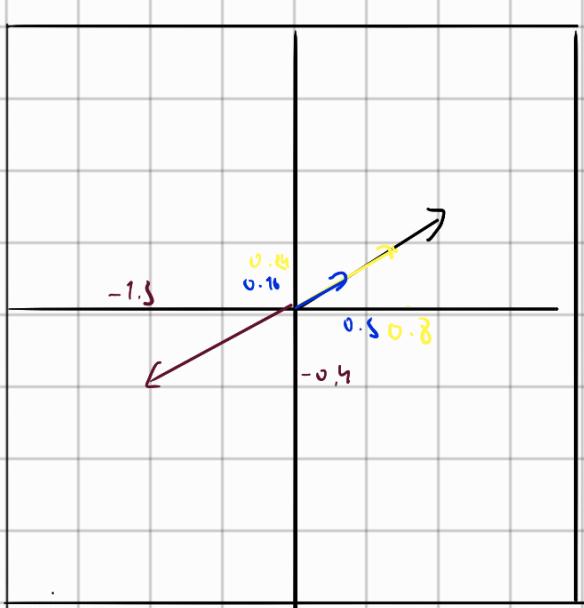
$\sqrt{\lambda_3} = 0.55$ -



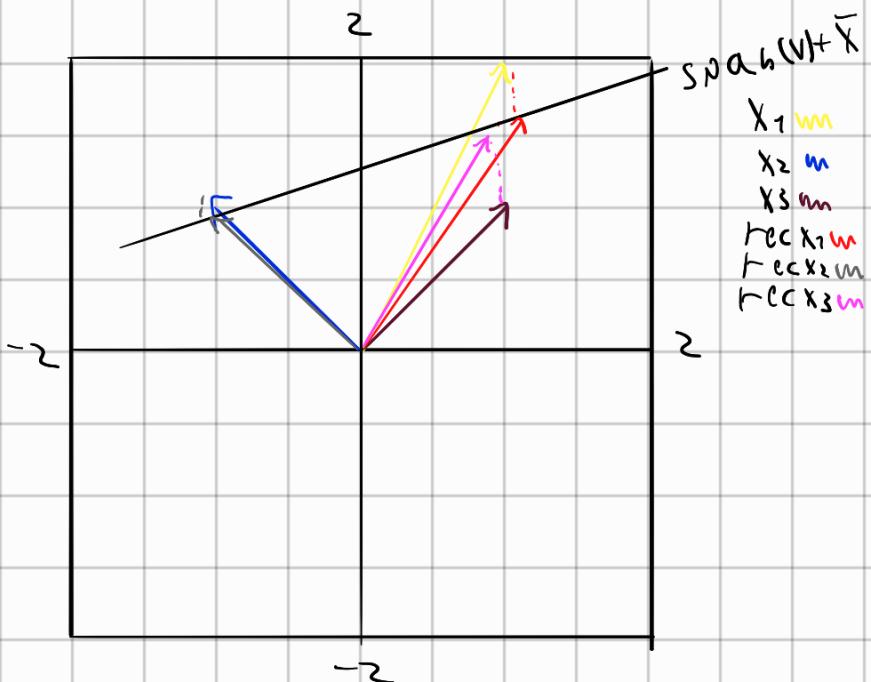
d) $\text{proj}_V(z_1) = \begin{bmatrix} 0.79 \\ 0.24 \end{bmatrix}$ z1

$\text{proj}_V(z_2) = \begin{bmatrix} -1.3 \\ -0.4 \end{bmatrix}$ z2

$\text{proj}_V(z_3) = \begin{bmatrix} 0.5 \\ 0.16 \end{bmatrix}$ z3



e) $\text{rec } x_1 = \begin{bmatrix} 1.1 \\ 1.6 \end{bmatrix} \quad \text{rec } x_2 = \begin{bmatrix} -0.98 \\ 0.936 \end{bmatrix} \quad \text{rec } x_3 = \begin{bmatrix} 0.852 \\ 1.49 \end{bmatrix}$



$$2) \text{ a)} \sum_{i=1}^m \|x_i - \bar{x} - w^T w(x_i - \bar{x})\|^2 = \sum_{i=1}^m \|z_i - w^T w z_i\|^2$$

$$\begin{aligned} &= \sum_{i=1}^m z_i^T z_i - 2 z_i^T (w^T w z_i) + (w^T w z_i)^T (w^T w z_i) \\ &= \sum_{i=1}^m \|z_i\|^2 - 2 z_i^T w^T w z_i + z_i^T w^T w z_i \\ &= \sum_{i=1}^m \|z_i\|^2 - z_i^T w^T w z_i \end{aligned}$$

$$* = z_i^T w^T \underbrace{w w^T}_{I_K} w z_i = z_i^T w^T w z_i$$

□

because w rows is orthonormal vectors

$$\begin{aligned} b) \sum_{i=1}^m \underbrace{z_i^T w^T w z_i}_{\text{scalar}} &= \sum_{i=1}^m \text{Tr}(w z_i z_i^T w^T) = \text{Tr}(w \left(\sum_{i=1}^m z_i z_i^T \right) w) \\ &= \text{Tr}(w (m \Sigma) w^T) = m \cdot \text{Tr}(w \Sigma w^T) \end{aligned}$$

$$c) \text{Tr}(w \Sigma w^T) = \text{Tr}\left(\left(\sum_{i=1}^k v_i v_i^T\right) \sum\left(\sum_{i=1}^k v_i\right)\right) = \sum_{i=1}^k \text{Tr}(v_i \Sigma v_i^T)$$

d) From a) we know

$$\sum_{i=1}^m \|x_i - \bar{x} - w^T w(x_i - \bar{x})\|^2 = \sum_{i=1}^m \|z_i\|^2 - z_i^T w^T w z_i$$

$$\min\left(\sum_{i=1}^m \|z_i\|^2 - z_i^T w^T w z_i\right) \stackrel{a,b}{=} \min\left(\sum_{i=1}^m \|z_i\|^2 - m \cdot \text{Tr}(w \Sigma w^T)\right)$$

$$\max \sum_{j=1}^k v_j^T \Sigma v_j \stackrel{c}{=} \max(\text{Tr}(w \Sigma w^T))$$

To minimize the LS error we need to maximize the orthogonal projection

3. We need to prove: $\forall V \in \mathbb{R}^d$ with $\|V\|_2 = 1 \Rightarrow \text{Var}(V^T X) \leq \text{Var}(V_1^T X)$

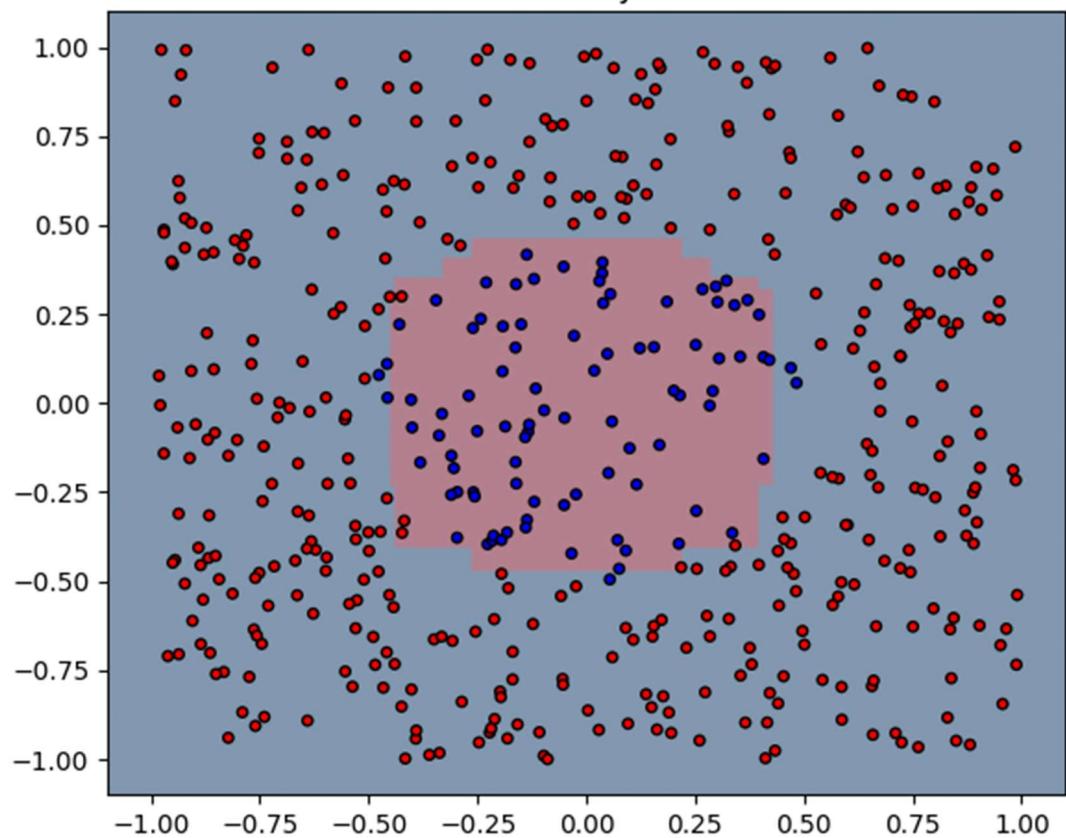
which is equivalent to $V^T \Sigma V \leq V_1^T \Sigma V_1$

PCA returns $\sum_{i=1}^k V_i^T \Sigma V_i$ in descending order that is
the Variance of $\langle V_1, X \rangle$ is the largest

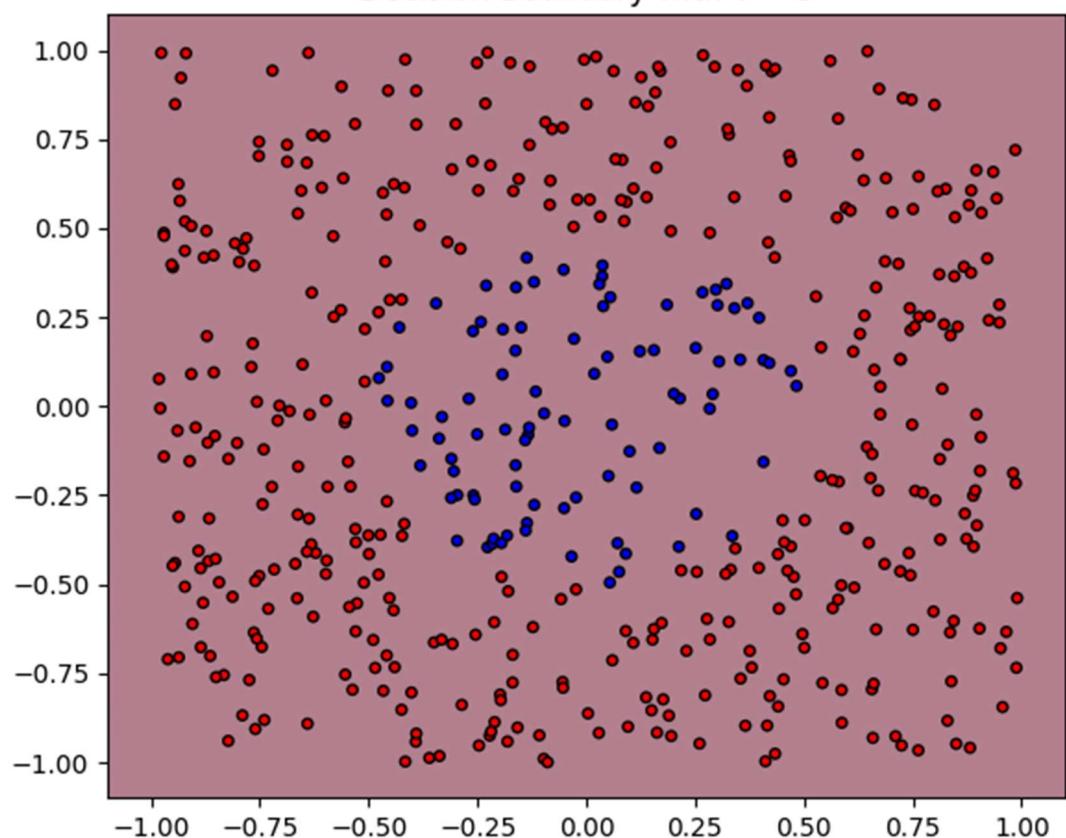
1.3 Clustering

1. a) True if $m_{min}=1$ Then for cluster of size one can be pushed into final C
- b) False its relying on k-means which is not optimal and can give different output
- c) True because we making assumption for noise
- d) True because got different starting point and has different output

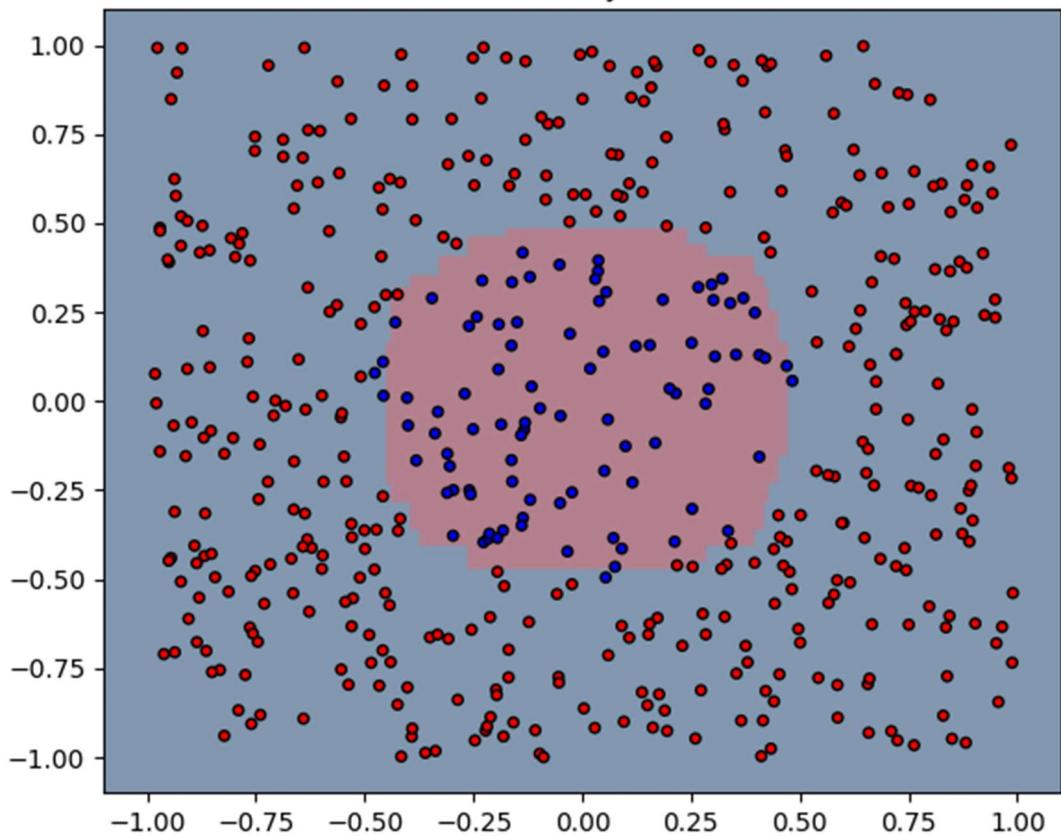
Decision Boundary with $T = 50$



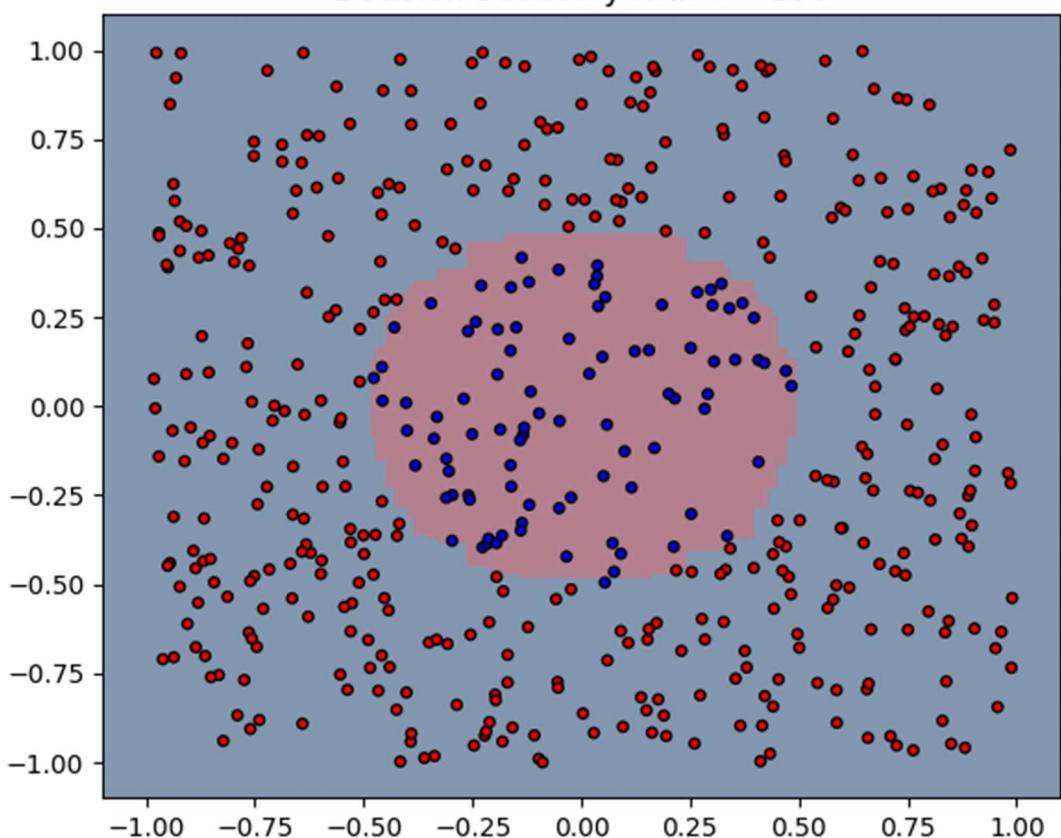
Decision Boundary with $T = 5$



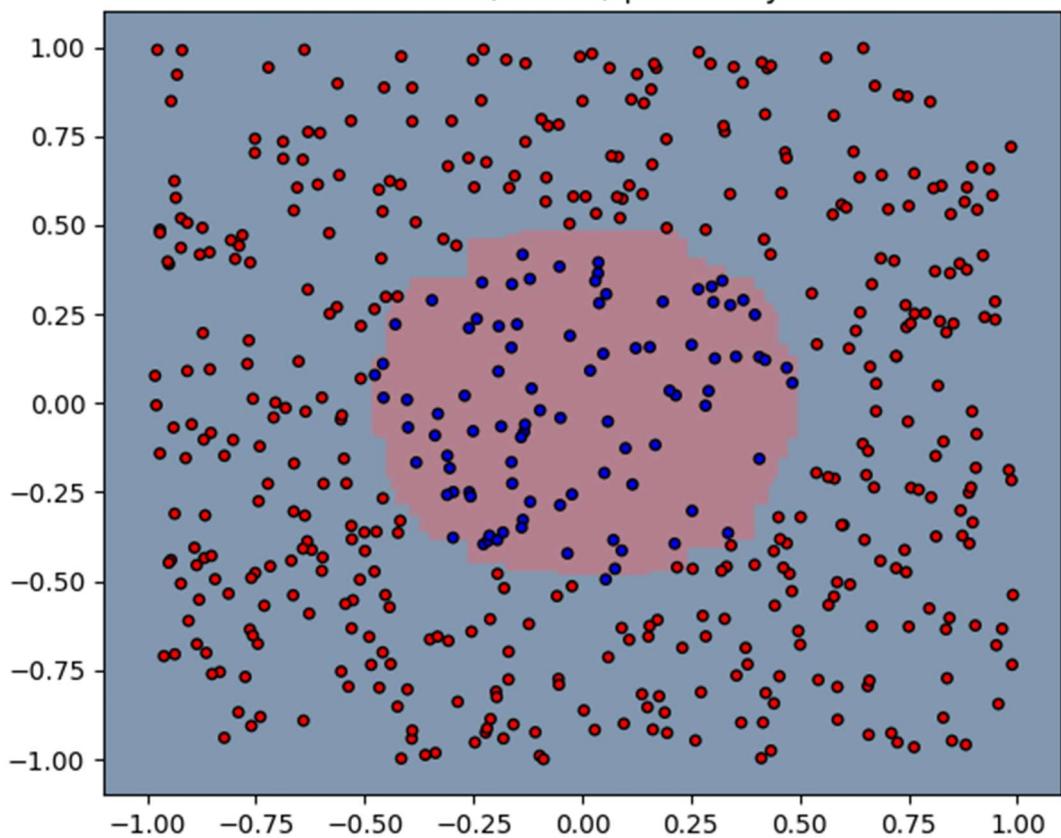
Decision Boundary with $T = 100$



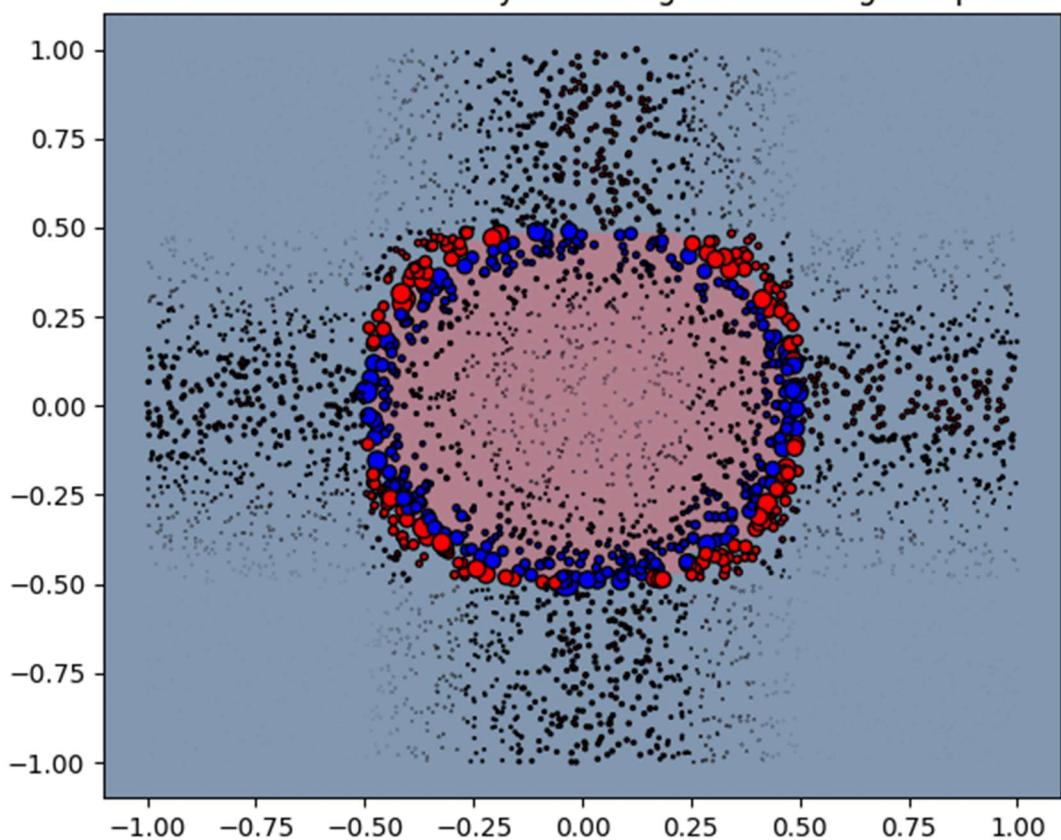
Decision Boundary with $T = 250$

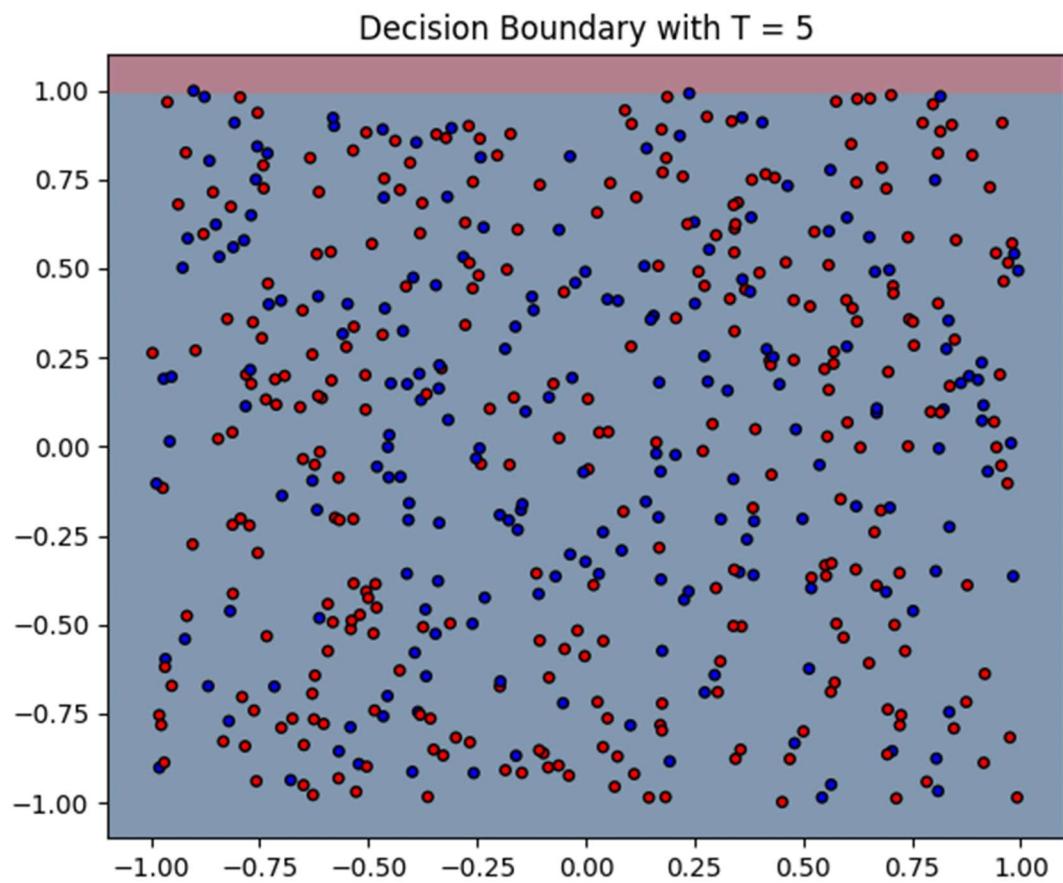
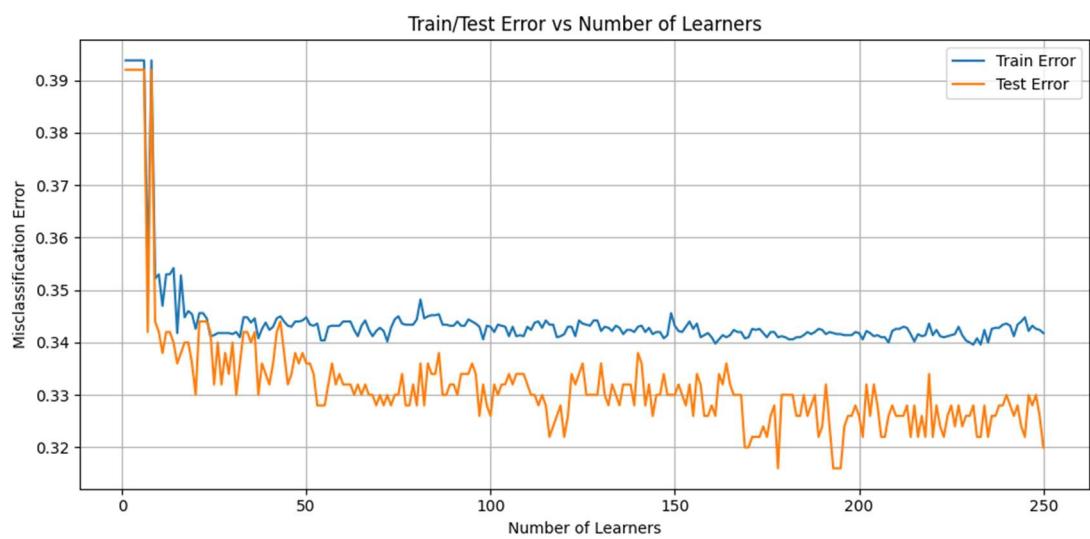


Best Ensemble (T=172) | Accuracy = 0.996

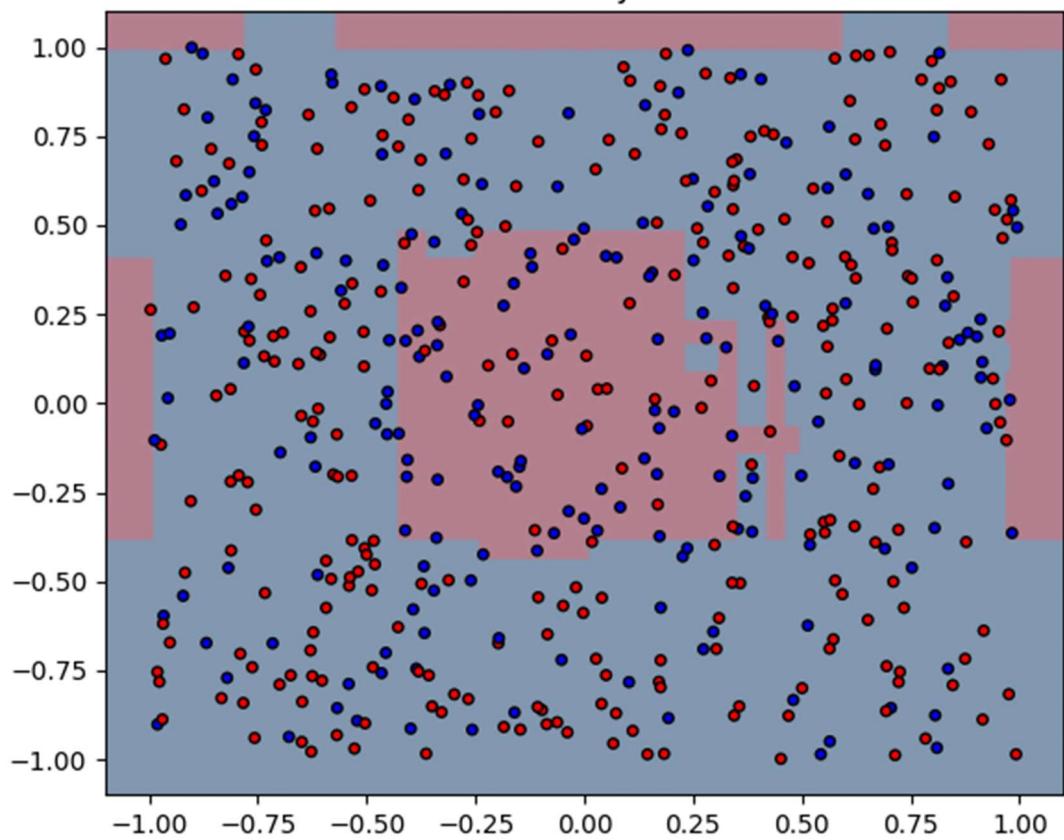


Final Decision Boundary with Weighted Training Samples

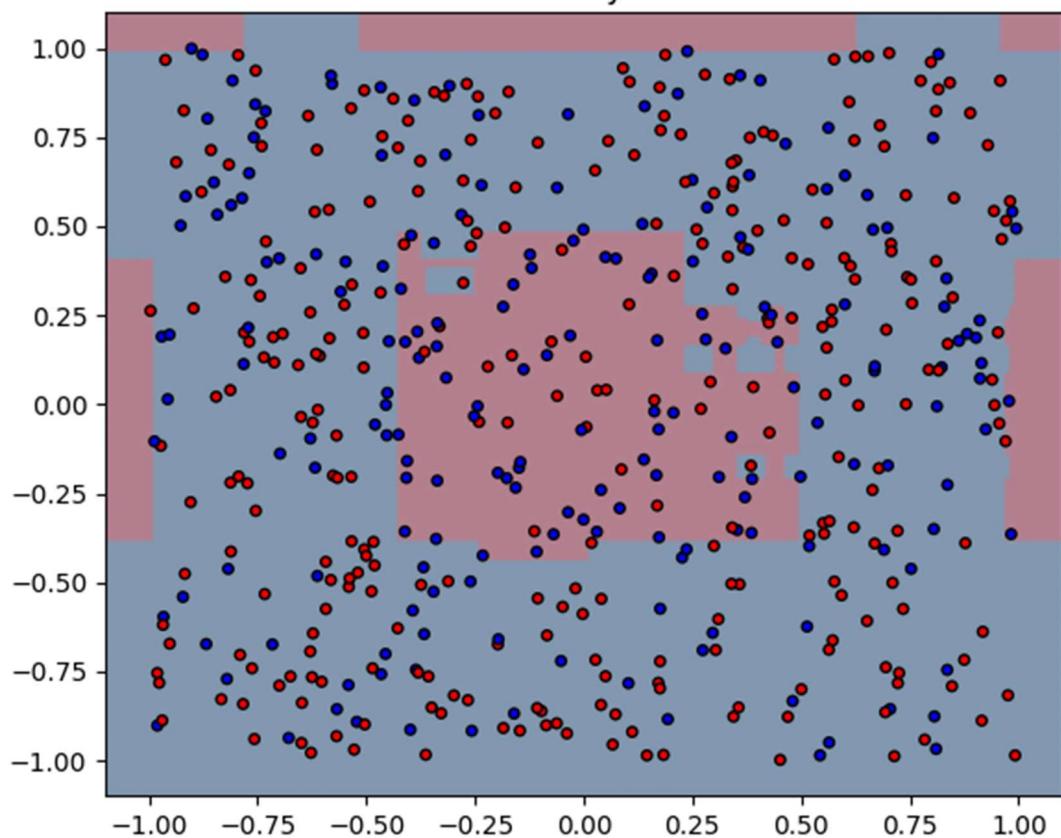




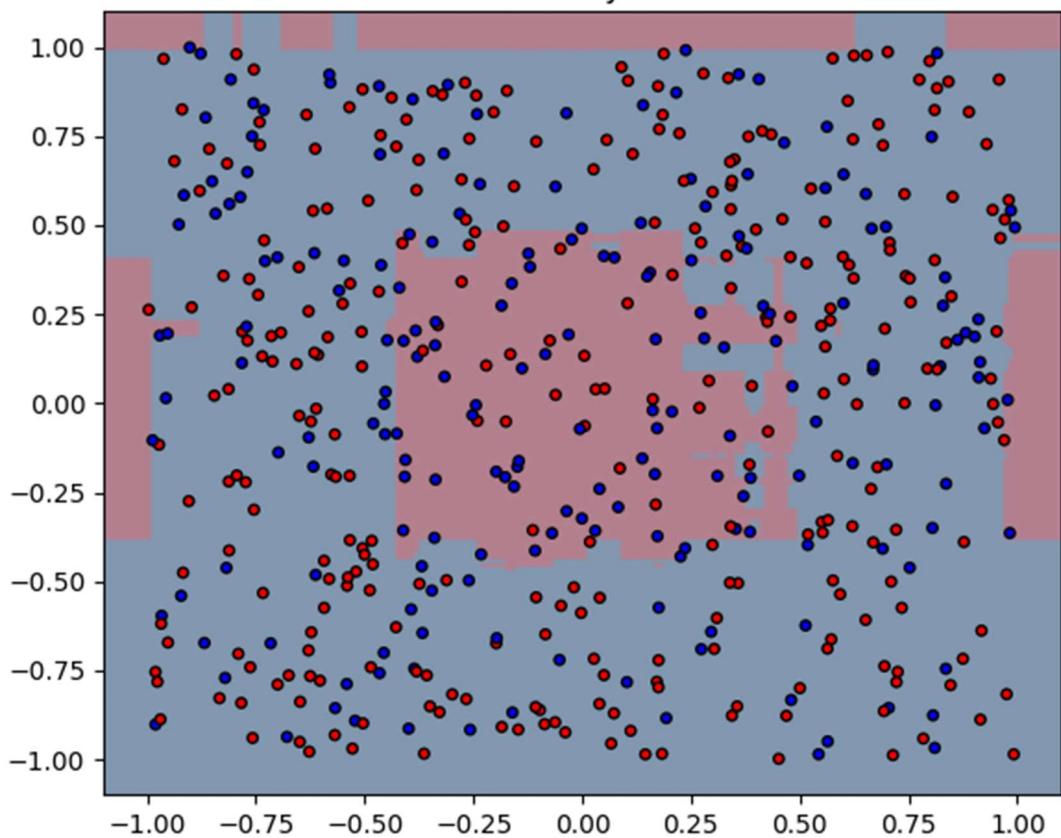
Decision Boundary with $T = 50$



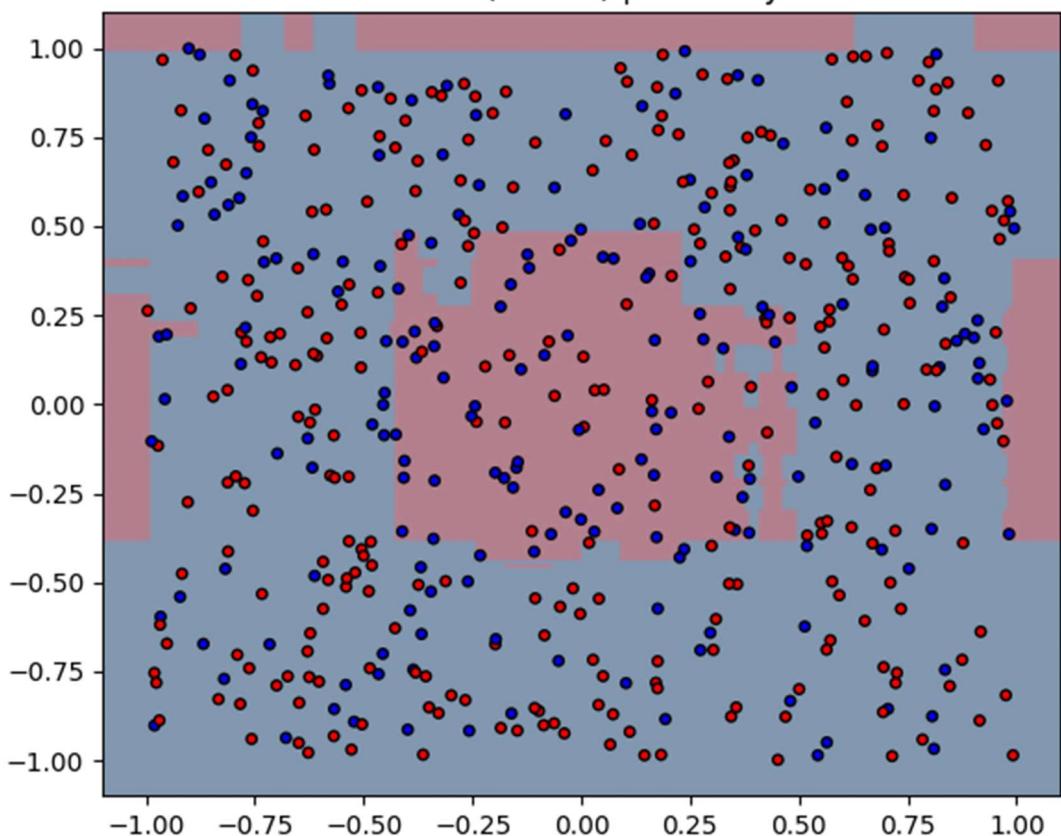
Decision Boundary with $T = 100$



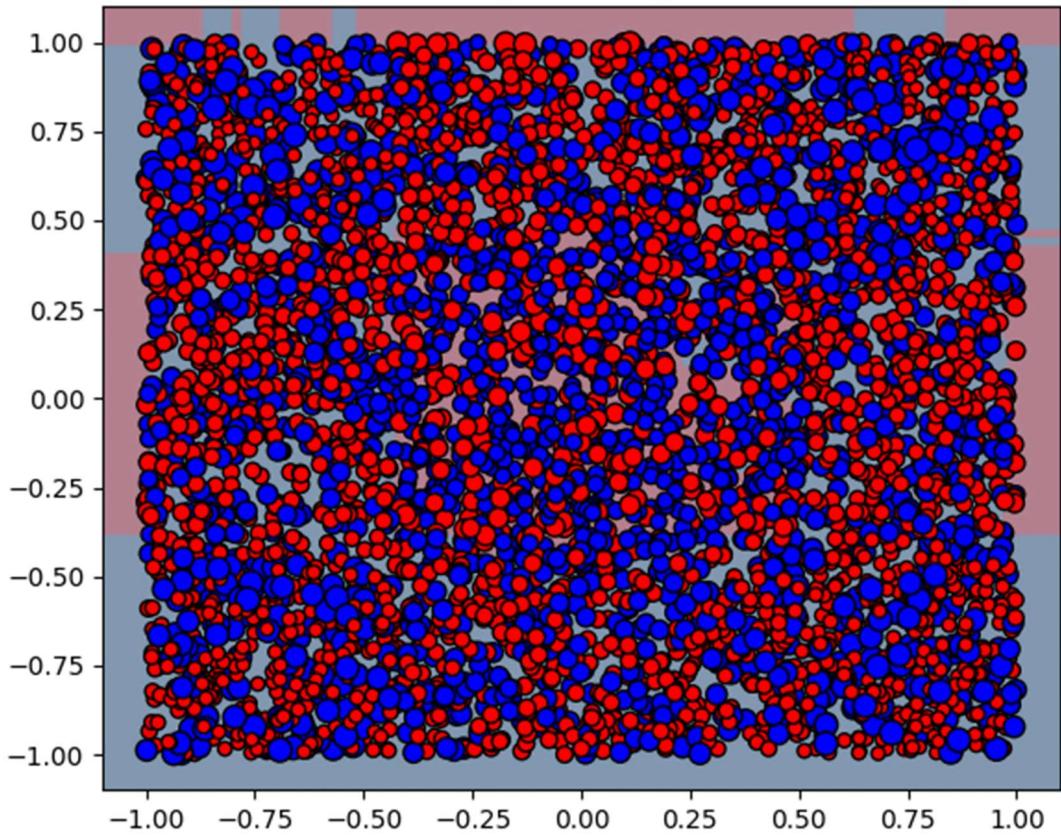
Decision Boundary with $T = 250$



Best Ensemble ($T=178$) | Accuracy = 0.684



Final Decision Boundary with Weighted Training Samples



2 Practical Part

2.1 Boosting

1. מכיוון שאין כלל רעש המודל לומד את הנתונים לגמרי עד שאין לו שום טעות ולכן השכיהה מגיעה לכמעט 0 מכיוון שהנתונים נלמדים לגמרי

4. נראה שהמודל מתקשה להצליח עבור נקודות שנמצאות על התפר בין הנתונים מכיוון שפעם הגבול לא ברור ולכן זה קשה עבורו

לא ממש הצלחתי להמשיך מפני לצערי לא היה לי את הזמן) :

