

Ex 2

1. Theoretical Part

1.1 Hard & Soft SVM

1. Let's define: $V = \begin{pmatrix} b \\ w_1 \\ \vdots \\ w_d \end{pmatrix} \in \mathbb{R}^{d+1}$

From the hint we know: $\|w\|^2 = w^T I_d w$ and now we can use V instead of w , let's define $Q = \begin{bmatrix} I_d & 0 \\ 0 & 0 \end{bmatrix} \in \mathbb{R}^{(d+1) \times (d+1)}$

$$[w_1 \dots w_d] \begin{bmatrix} 1 & 0 \\ \vdots & \vdots \\ 0 & 1 \end{bmatrix} \begin{pmatrix} w_1 \\ \vdots \\ w_d \end{pmatrix} = [w_1 \dots w_d, b] \begin{bmatrix} I_d & 0 \\ 0 & 0 \end{bmatrix} \begin{pmatrix} w_1 \\ \vdots \\ w_d \\ b \end{pmatrix} = V^T Q V$$

for the condition: $y_i(\langle w, x_i \rangle + b) \geq 1 \Leftrightarrow y_i(V^T[1, x_i]) \geq 1 \Leftrightarrow$

$$-y_i(V^T[1, x_i]) \leq -1$$

so let's define $d = \begin{bmatrix} -1 \\ \vdots \\ 1 \end{bmatrix} \in \mathbb{R}^m$ and $A \in \mathbb{R}^{m \times d}$ s.t. $A_i = -y_i \cdot (V^T[1, x_i])$

and $a=0$ occurs that isn't non-linear term in $\|w\|^2$

how we get that $\arg \min_{(w,b)} \|w\|^2 = \arg \min_{(w,b)} \frac{1}{2} V^T Q V$

we can multiply by $\frac{1}{2}$ because this will give the same solution for the optimization problem

1.2 PAC Learnability

$$1. (A) \quad \left(1 + \frac{x}{n}\right)^n \leq e^{nx}$$

$$x \in \mathbb{R}, x > -1 \Rightarrow \ln(1+x) \leq x \Rightarrow 1+nx \leq e^{nx} \quad x_1, \dots, x_m \in S_m$$

$$P_D[S_m \cap X' = Q] = P_D[\bigcap_{i=1}^m x_i \notin X'] = \prod_{i=1}^m P_D[x_i \notin X'] = (1 - P_D[x \in X'])^m$$

*

We know that $P_D[x_i \in X'] > \varepsilon \Leftrightarrow P_D[x_i \notin X'] \leq \varepsilon$

$$1 - P_D[x \in X'] \leq 1 - \varepsilon \Rightarrow (1 - P_D[x \in X'])^m \leq (1 - \varepsilon)^m$$

know from Bernoulli's inequality $(1 - \varepsilon)^m \leq e^{-m\varepsilon}$

□

(b) Let's define $f = \text{The real function}$

$$Z_{f_s} = \{x \in X \text{ s.t } \hat{f}_s(x) \neq f(x)\}$$

$$H_{\text{Bad}} := \{\hat{f} \in H \text{ s.t } L_{D,f}(\hat{f}) > \varepsilon\}$$

We can see that by definition that mistake can

happen if $S_m \cap Z_{f_s} = \emptyset$. Then we can get $L_{S,f}(\hat{f}_{S_m}) = 0$ and $\hat{f} \neq \hat{f}_s$

The probability to get a set like that

is the one from (a) meaning: $P[S_m \cap Z_{f_s} = \emptyset] \leq e^{-m\varepsilon}$

$$\text{So } \{L_{D,f}(\hat{f}_s) > \varepsilon\} \subseteq \{\forall h \in H_{\text{Bad}}, L_{S_m}(\hat{f}_s) = 0\}$$

Because by definition this is all the functions that

$L_{D,f}(\hat{f}) > \varepsilon$ mistake will be made $L_{S_m}(\hat{f}_s) = 0$

$$\Rightarrow P_D(\{L_{D,f}(\hat{f}_s) > \varepsilon\}) \leq P_D(\{\forall h \in H_{\text{Bad}}, L_{S_m}(\hat{f}_s) = 0\})$$

$$\leq \sum_{\hat{f}_s \in H_{\text{Bad}}} P_D(S_m \cap Z_{f_s} = \emptyset) \leq |H| \cdot e^{-m\varepsilon}$$

(c) From (b) we know that $P_D[L_{DF}(A(S_m)) > \varepsilon] \leq |H| \cdot e^{-m\varepsilon}$

We need m that $|H| \cdot e^{-m\varepsilon} \leq \delta$

↑

$$-m\varepsilon = \log(e^{-m\varepsilon}) \leq \log\left(\frac{\delta}{|H|}\right)$$

↑

$$m\varepsilon \geq -\log\left(\frac{\delta}{|H|}\right) = \log|H| - \log\delta =$$

↑

$$\log 1 - \log\delta = \log\left(\frac{1}{\delta}\right)$$

Hands for every $m > m$

$$m \geq \frac{\log|H| + \log\frac{1}{\delta}}{\varepsilon}$$

We still hold the conditions

$$P_D[L_{DF}(A(S_m)) > \varepsilon] \leq |H| \cdot e^{-m\varepsilon} \leq \delta$$

↑

$$P_D[L_{DF}(A(S_m)) > \varepsilon] \geq 1 - \delta$$

◻

$$2. \quad H := \{h_r : r \in \mathbb{R}^+\} \quad h(x) = \prod_{\{||x_i||_2 \leq r\}} 1$$

Let A be an algorithm that for every S_m sample over X choose $x_0 \in S_m$ s.t. $(\forall x \in S_m \quad ||x_0||_2 \geq ||x|| \wedge g_0 = 1)$

We want the max-r that we know that his label is 1

the algorithm will label sample $h(x) = 1$ if $X \leq \text{max}_r$

Otherwise how we know that for other sample labeled 1 the algorithm is correct but we still have the margin between $\text{max}_r \leq ||x|| \leq \text{real}_r$

that mistakenly labeled 0

lets define $Z_n = \{x : \text{max}_r \leq ||x|| \leq \text{real}_r\}$ the "bad samples"

the algorithm hold to all the conditions in Q.1

From 1.(a) we know

$$\Pr(L_{D, \text{real}_r}(h_x) > \varepsilon) \leq 1 - \Pr(L_{D, \text{real}_r}(h_x) \leq \varepsilon) \leq (1 - \varepsilon)^m \leq e^{-m\varepsilon}$$

and all we need now is to show that holds for this

$$\text{condition } e^{-m\varepsilon} \leq S \Leftrightarrow m(S, \varepsilon) \geq \frac{\log(\frac{1}{S})}{\varepsilon}$$

OK

3. Given that \mathcal{H} has VCP that is $m_{\mathcal{H}}^{VC} : (0, 1)^2 \rightarrow \mathbb{N}$
s.t. $\forall \epsilon, \delta \in (0, 1) \ \forall D$ over $X \times Y \ P_D^m(\{S \subseteq (X \times Y)^m \mid S \rightarrow \text{ECP}\}) \geq 1 - \delta$

lets define A with learning principle ERM

$$\text{thus } A(S_n) = \hat{h}_S = \underset{h \in \mathcal{H}}{\operatorname{arg\,min}} L_S(h)$$

S is ϵ -representative, from the assumption \mathcal{H} -VCP

thus we can define $m_{\mathcal{H}}^{CV}(\epsilon, \delta)$

$$S \text{ is ECP} \Rightarrow \forall h \in \mathcal{H} \ |L_S(h) - L_D(h)| < \epsilon \text{ with } P(\cdot) \geq 1 - \delta$$

let $\hat{h} = \underset{h \in \mathcal{H}}{\operatorname{arg\,min}} (L_D(h))$ from VCP then $\forall h \in \mathcal{H} \ |L_S(h) - L_D(h)| < \epsilon$

thus for \hat{h} and for \hat{h}_S

$$|L_S(\hat{h}) - L_D(\hat{h})| < \epsilon \Leftrightarrow -\epsilon < L_D(\hat{h}) - L_S(\hat{h}) < \epsilon \Rightarrow L_S(\hat{h}) < L_D(\hat{h}) + \epsilon$$

$$\text{for } \hat{h} \text{ also } L_D(\hat{h}) < L_S(\hat{h}) + \epsilon$$



From Erm \hat{h}_S has the min loss over $S \quad L_S(\hat{h}_S) \leq L_S(\hat{h})$

$$\text{thus } L_D(\hat{h}_S) \leq L_S(\hat{h}_S) + \epsilon \leq L_S(\hat{h}) + \epsilon \leq L_D(\hat{h}) + \epsilon + \epsilon$$

$$\text{so for } \epsilon_0 = \epsilon/2 \text{ we get } L_D(\hat{h}_S) \leq \min_{h \in \mathcal{H}} L_D(h) + \epsilon$$

and that's happening in probability $> 1 - \delta$

which satisfies the requirement of APAC \square

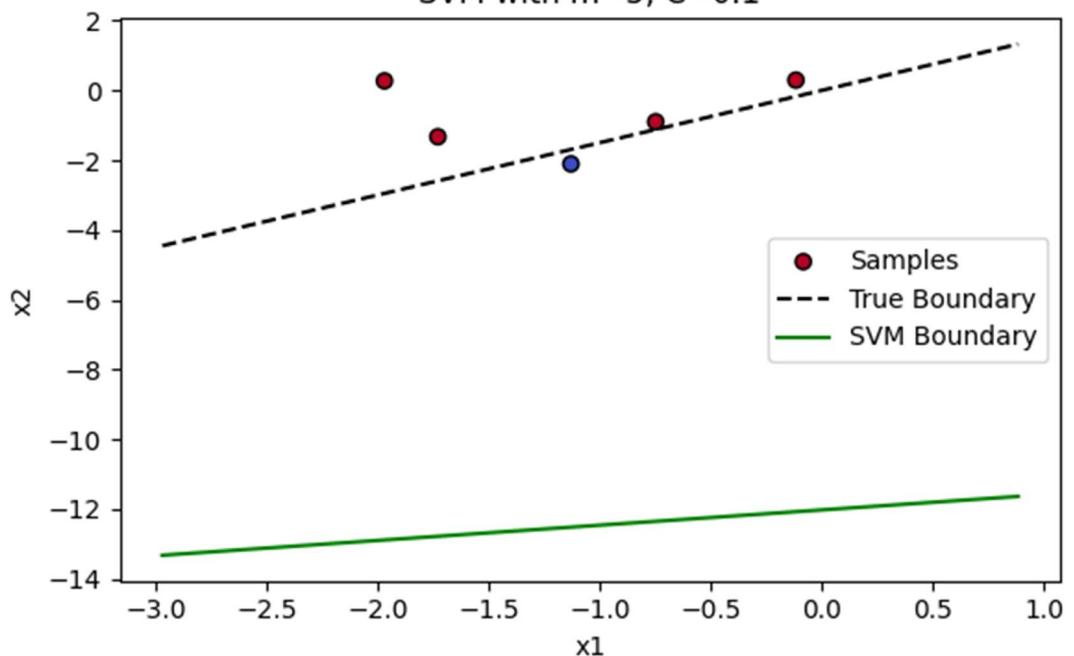
1.3 Vc Dimension

1. For sample of size n we need to show that some set can be shattered so we choose $\{c_1, c_2, \dots, c_n\}$ and define $I := \{i | y_i = 1\}$ then if $y_{i=1} \in I$ which means $h_I(c_i) = y_i = 1$ so it always odd. The size of H is all the vector $[0, 1]^n \subset \mathbb{R}^n$ we know that the size of subsets $I \subseteq [n]$ is 2^n and in the reaction we saw that $Vc\text{-dim}$ can't be bigger than $\log(|H|) \Rightarrow Vc\text{-dim}(H) \leq \log(2^n) = n$

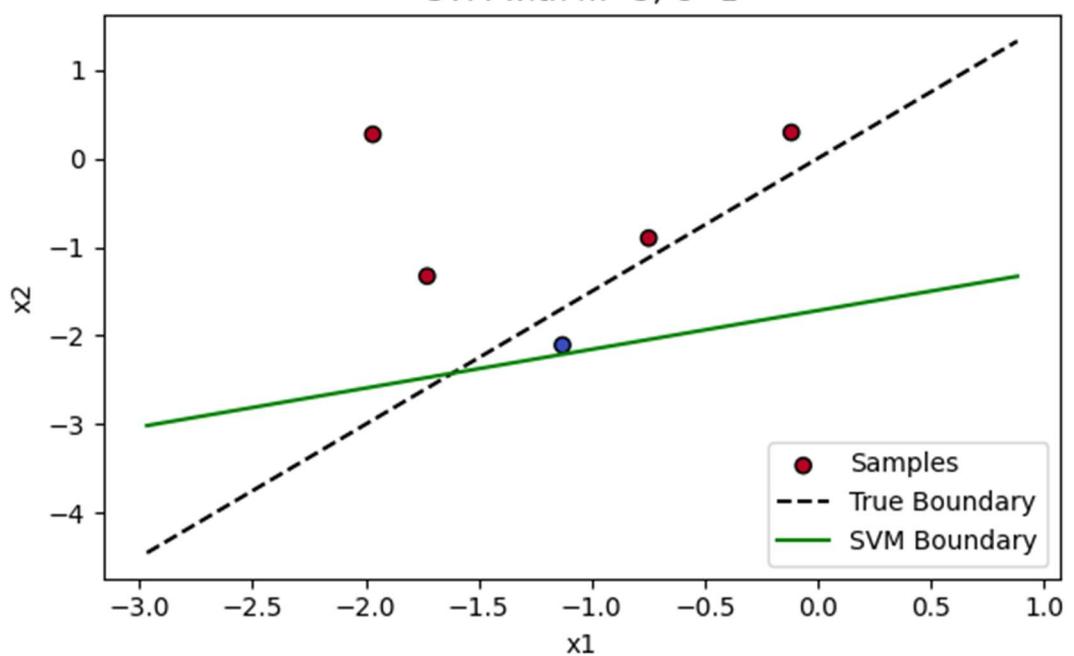
$$\begin{aligned} * \quad & Vc\text{-dim}(H) \geq n \Rightarrow Vc\text{-dim}(H) = n \\ * \quad & Vc\text{-dim}(H) \leq n \end{aligned}$$

2. $H_1 \subseteq H_2$ thus every $x \subseteq X$ s.t. $x \in H_1$ if x shattered by H_1 , $H_1 \subseteq H_2$ thus x shattered by $H_2 \Rightarrow Vc\text{-dim}(H_1) \leq Vc\text{-dim}(H_2)$

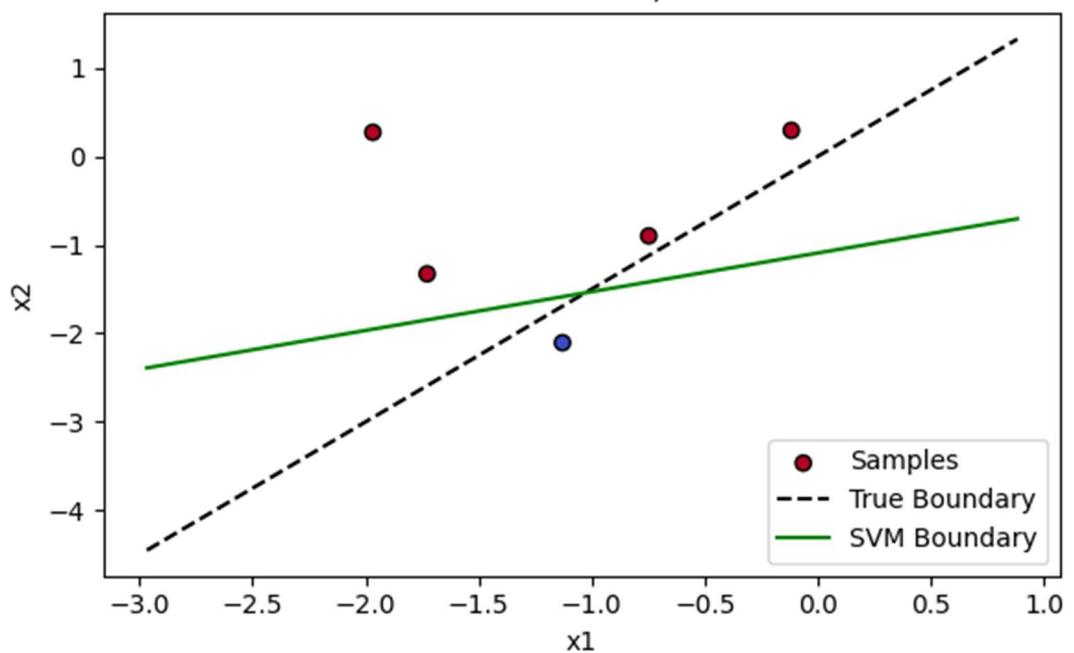
SVM with $m=5$, $C=0.1$



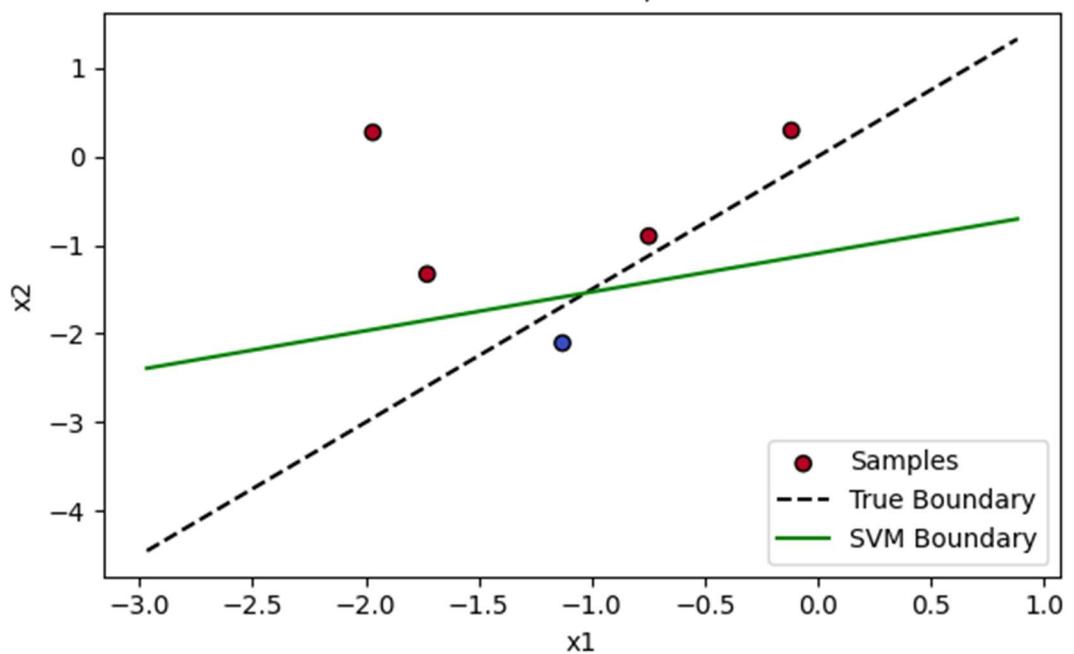
SVM with $m=5$, $C=1$



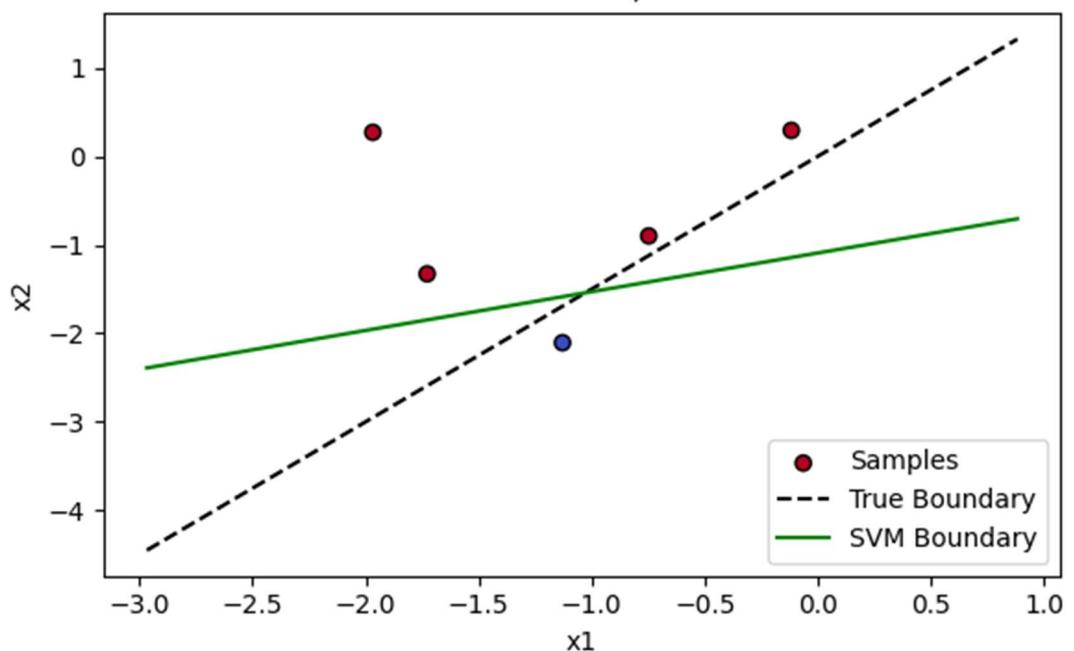
SVM with $m=5$, $C=5$



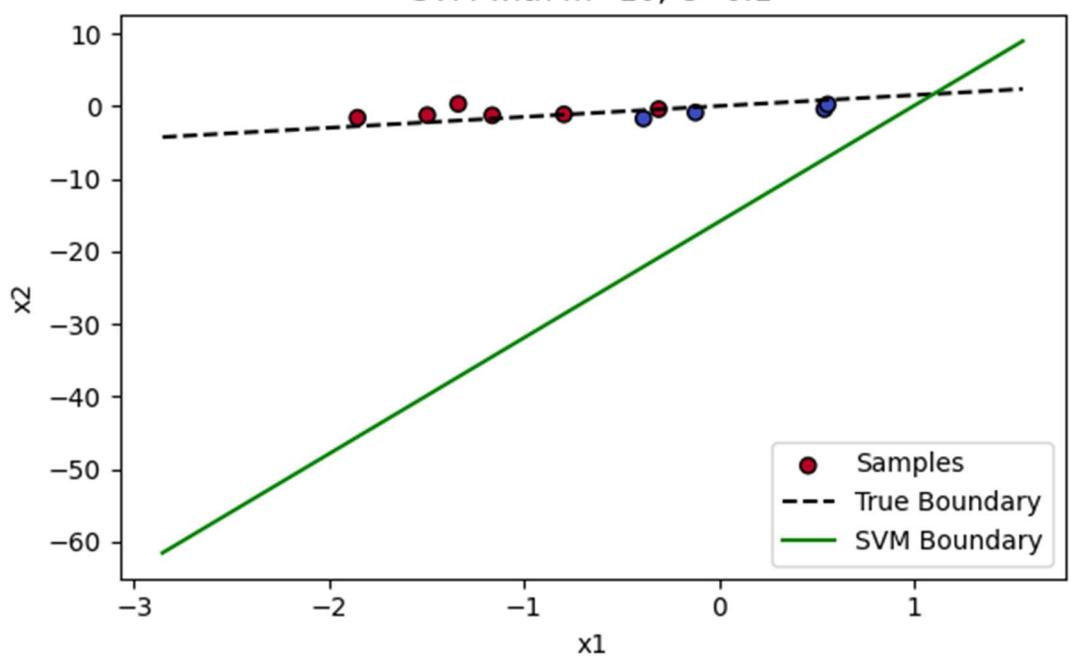
SVM with $m=5$, $C=10$



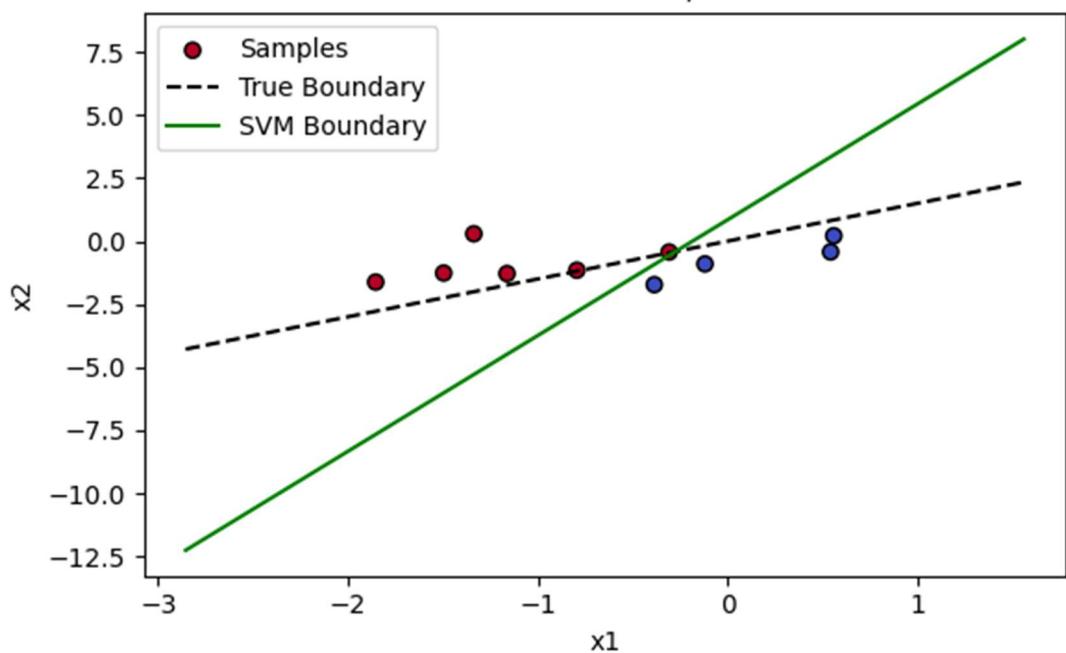
SVM with $m=5$, $C=100$



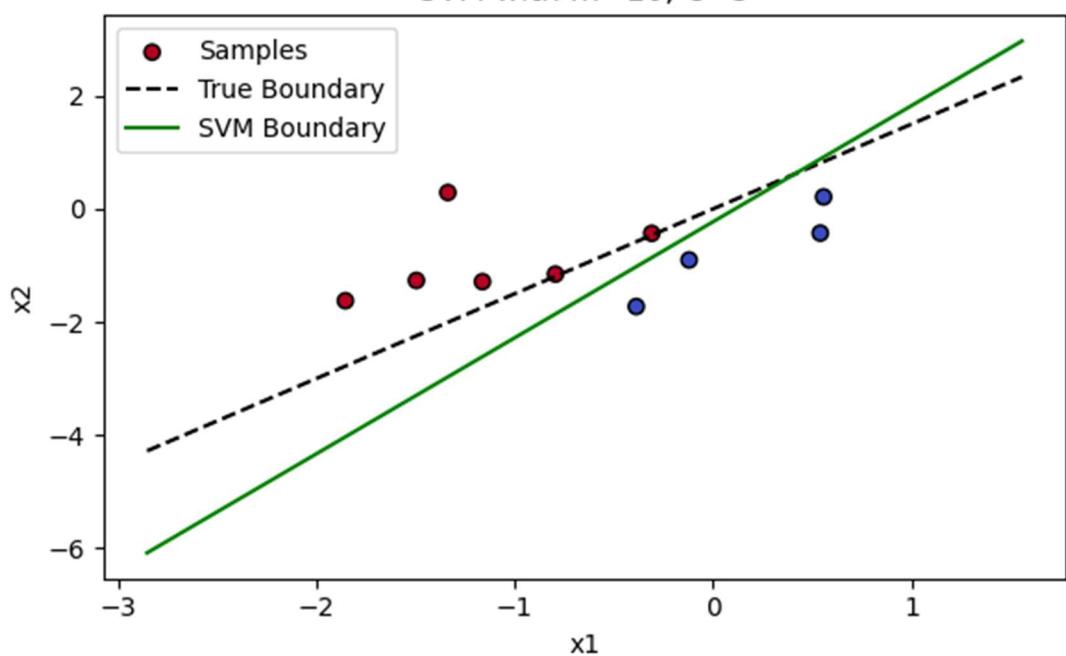
SVM with $m=10$, $C=0.1$



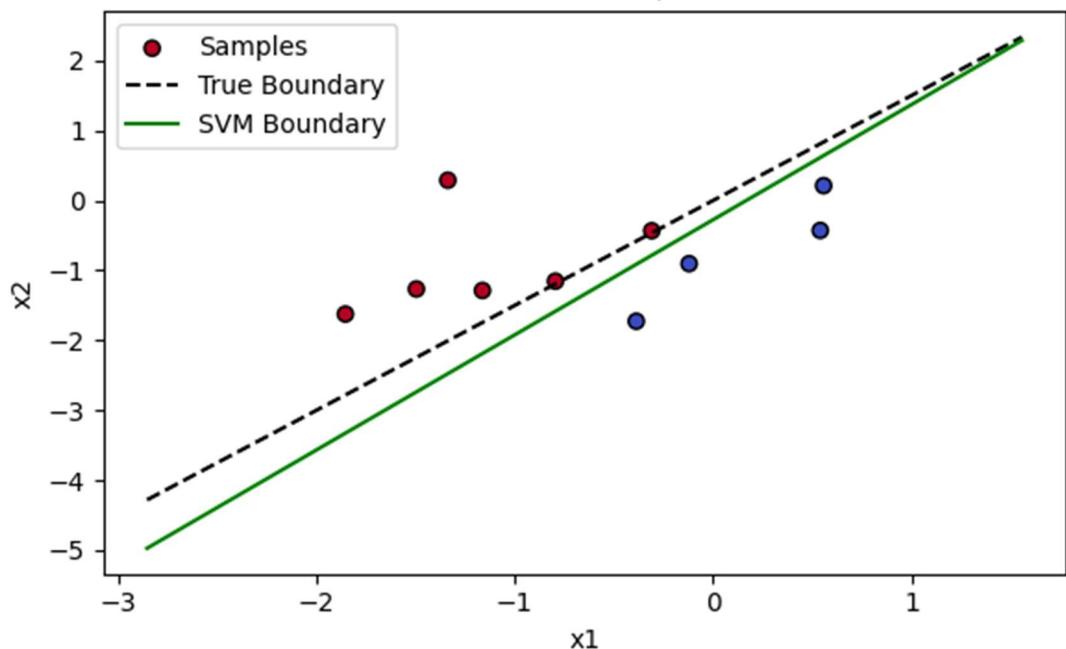
SVM with $m=10, C=1$



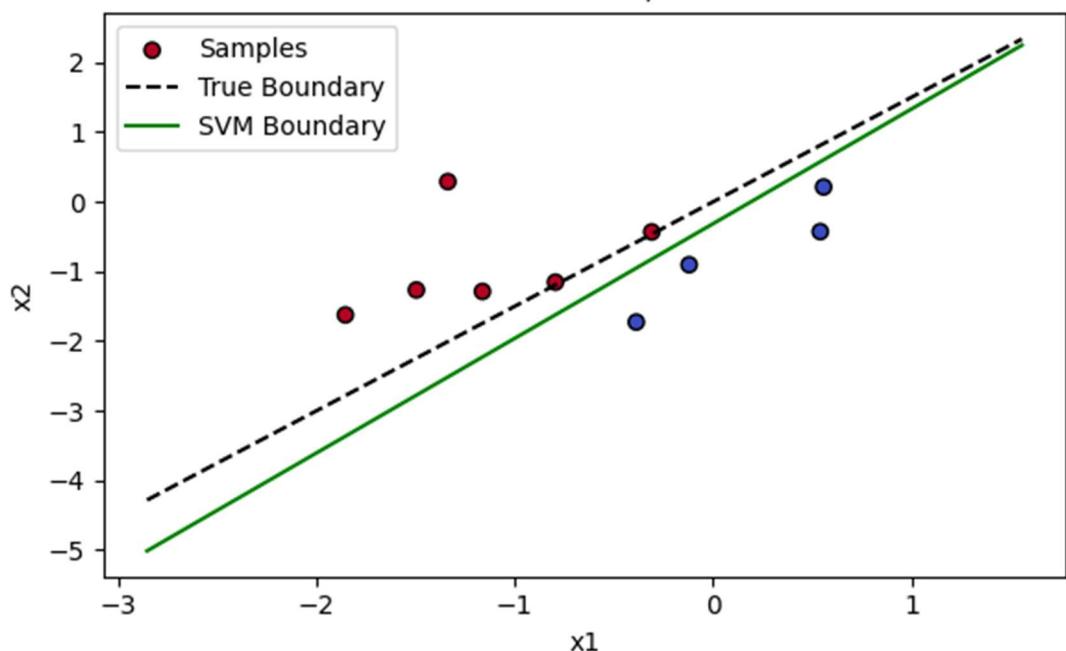
SVM with $m=10, C=5$



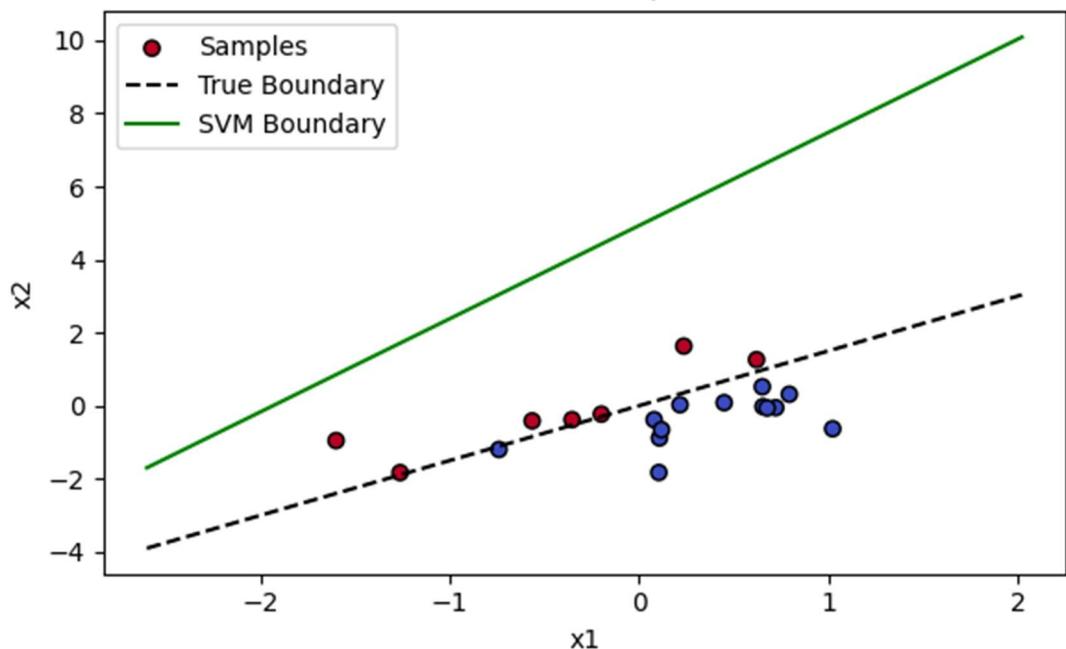
SVM with $m=10$, $C=10$



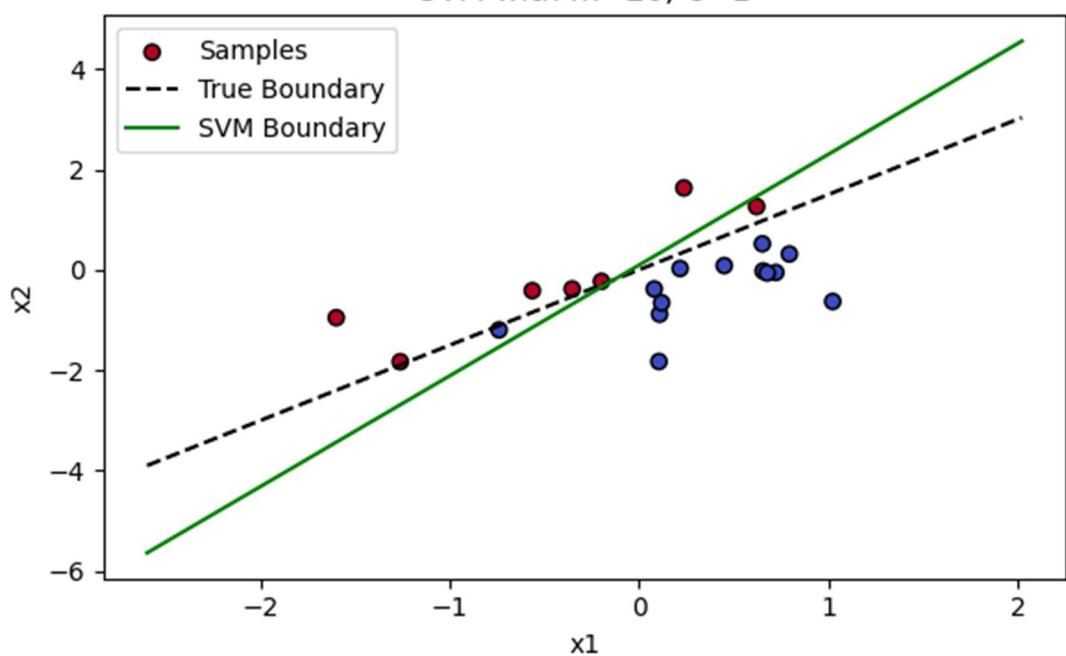
SVM with $m=10$, $C=100$



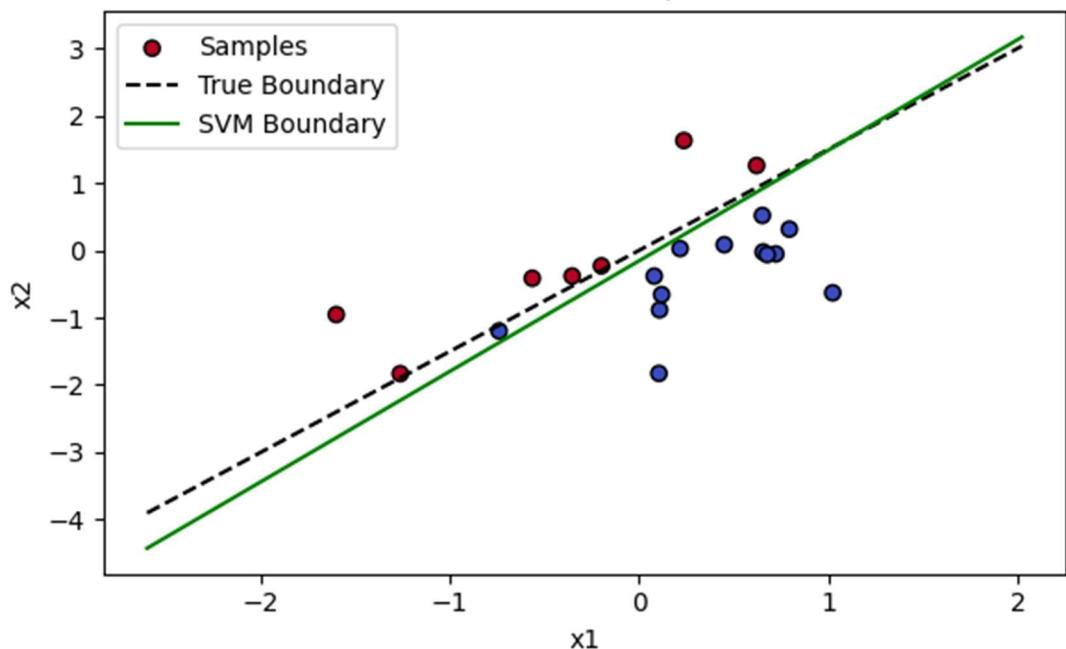
SVM with $m=20$, $C=0.1$



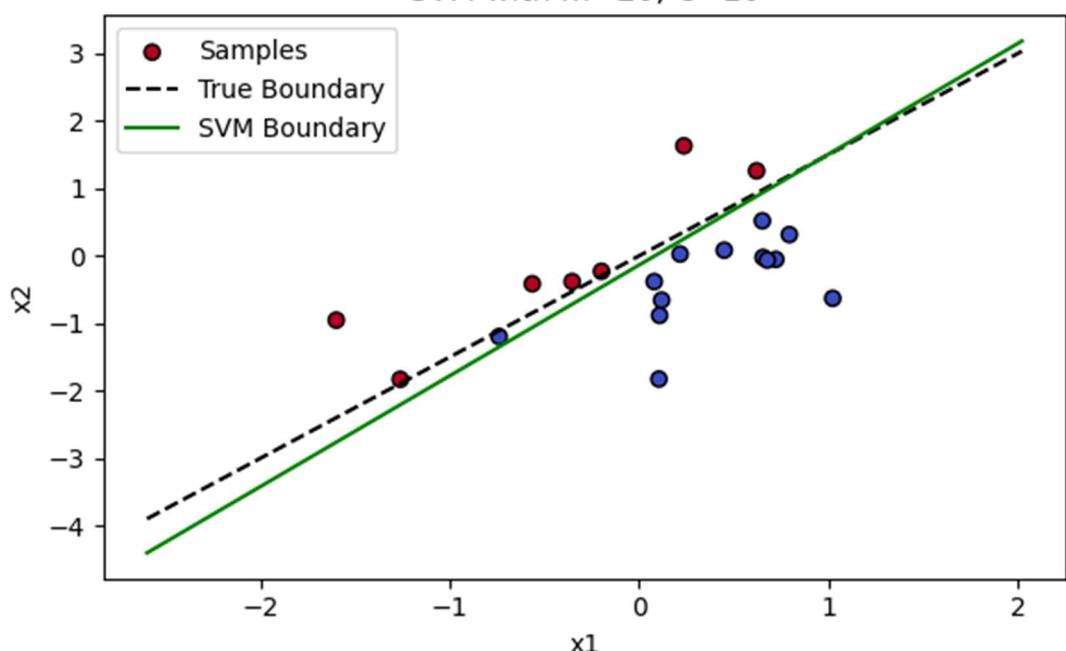
SVM with $m=20$, $C=1$



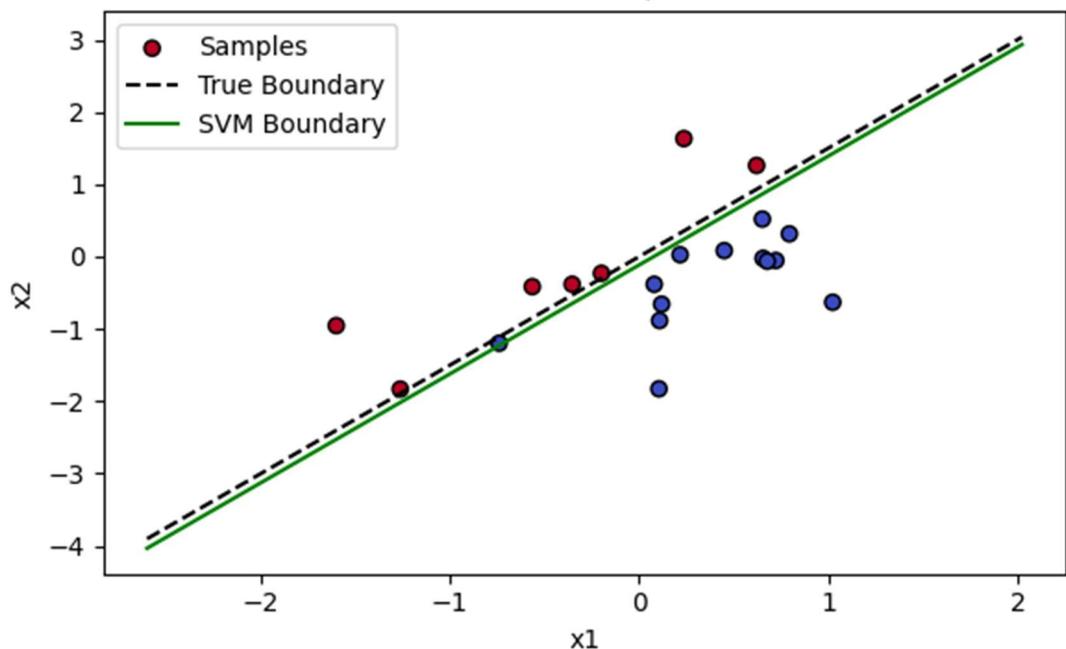
SVM with $m=20$, $C=5$



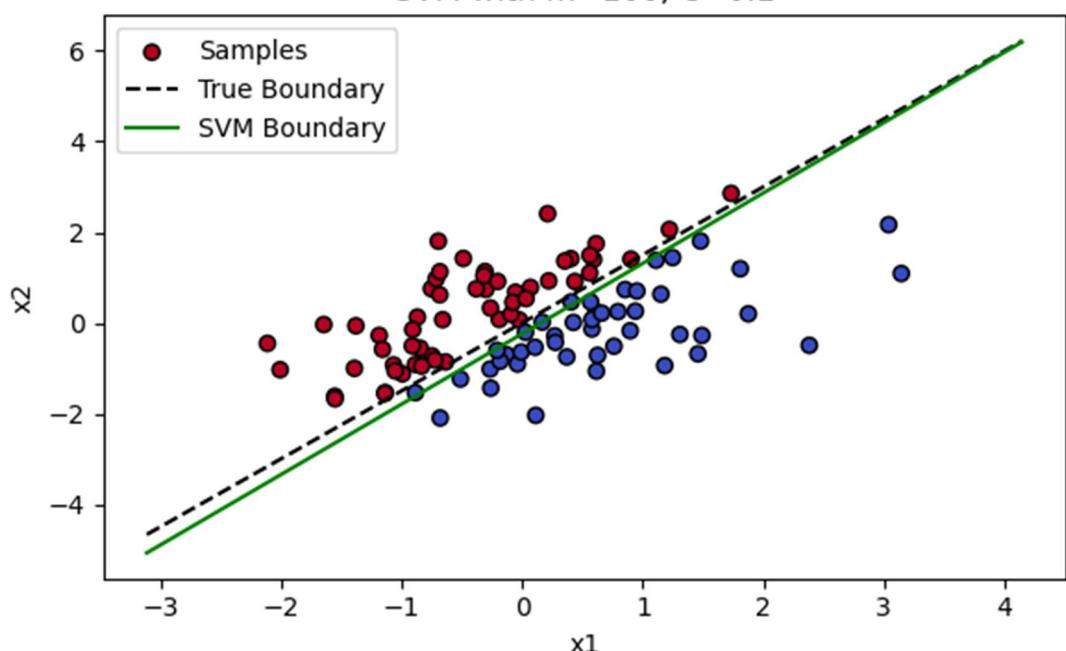
SVM with $m=20$, $C=10$



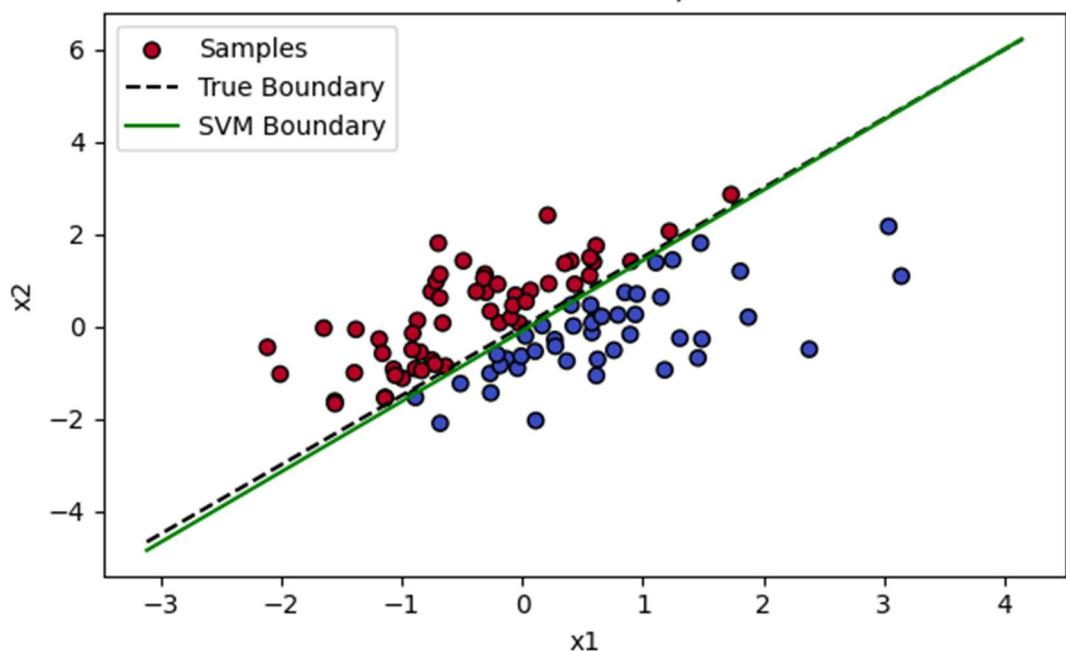
SVM with $m=20$, $C=100$



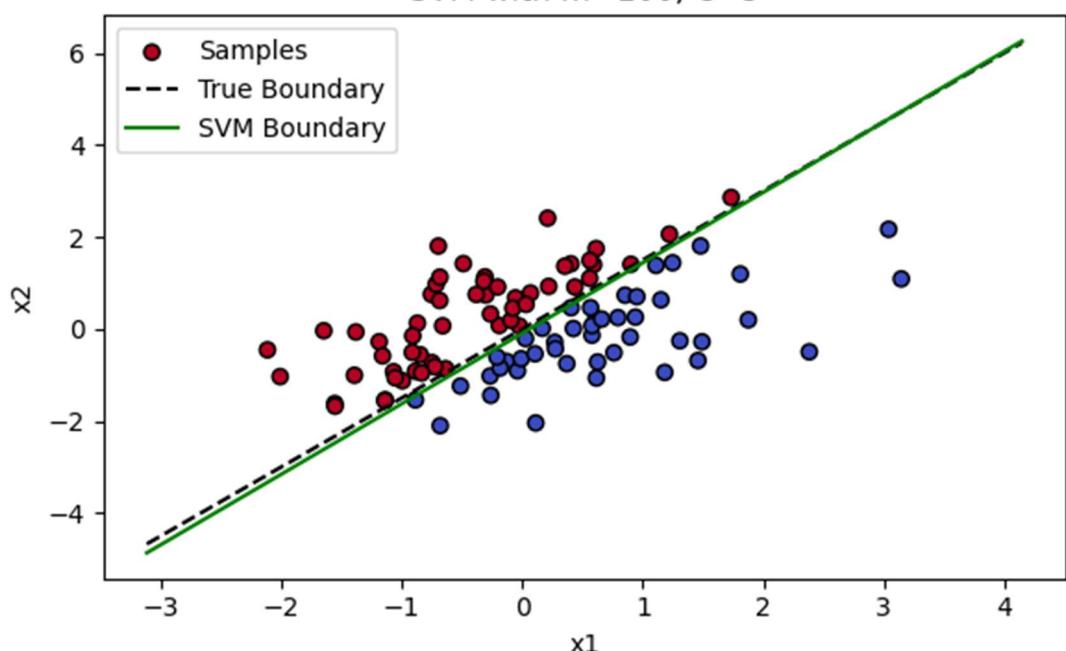
SVM with $m=100$, $C=0.1$



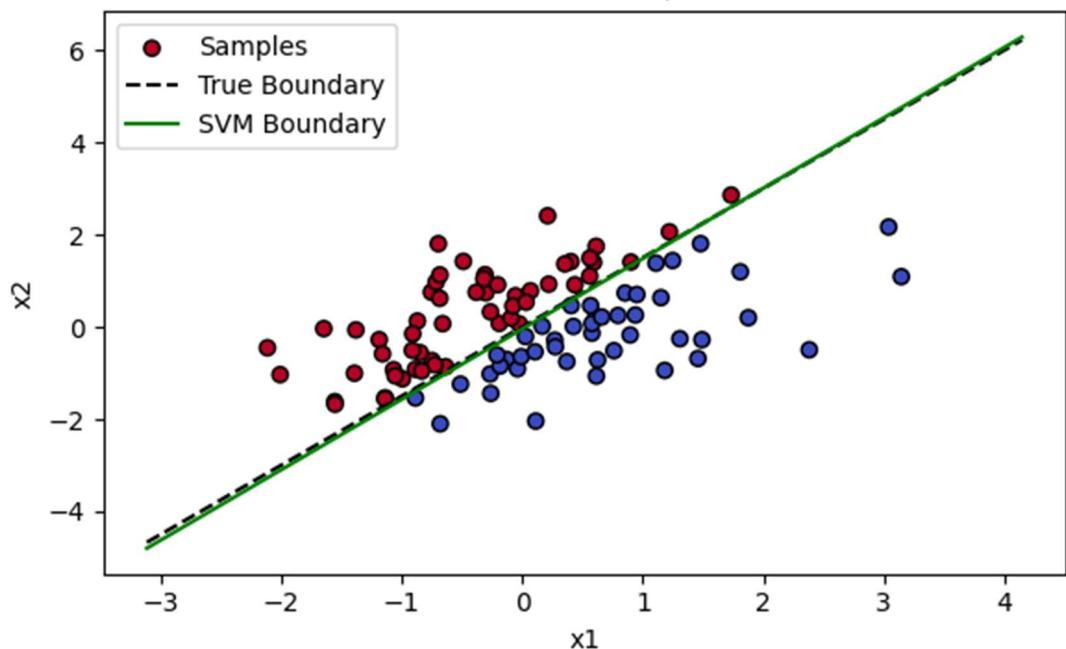
SVM with $m=100$, $C=1$



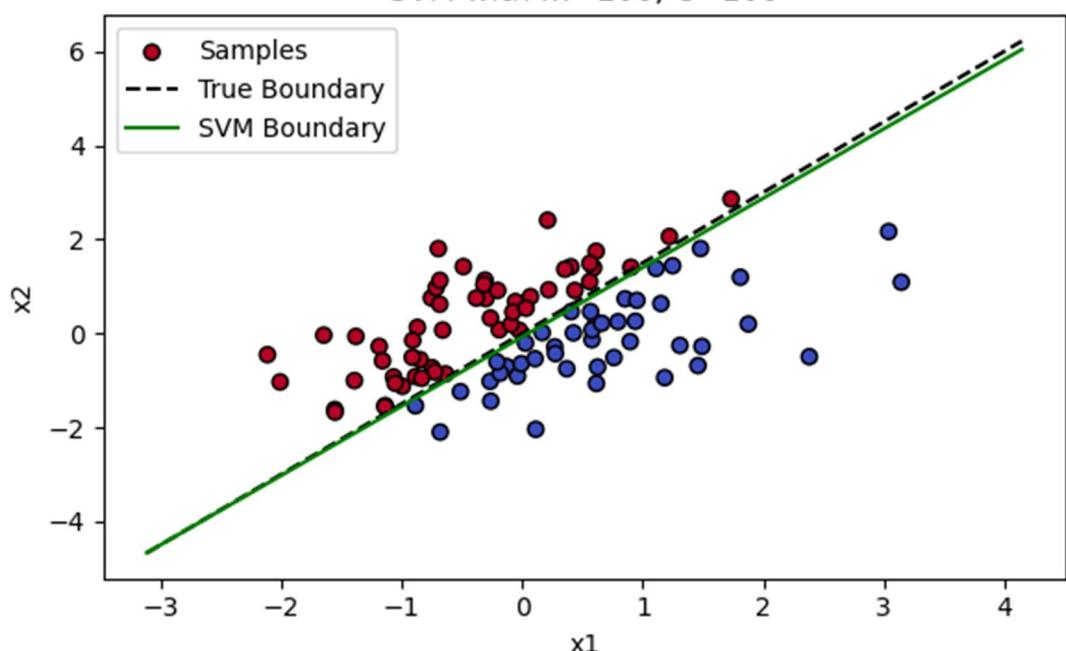
SVM with $m=100$, $C=5$



SVM with $m=100$, $C=10$



SVM with $m=100$, $C=100$



2

2.1 SVM

1.above

2.

We can see that as **C** gets bigger, the SVM line is more accurate.

For **C = 0.1**, we almost don't care about the mistakes — we just want the **w** with minimum norm,

so the learner will choose ψ (slack variables) that are very big without getting punished.

On the other hand, if **C = 100**,

then for every ψ we pay more, so the learner will search for the minimum ψ possible, which could lead to overfitting.

3.

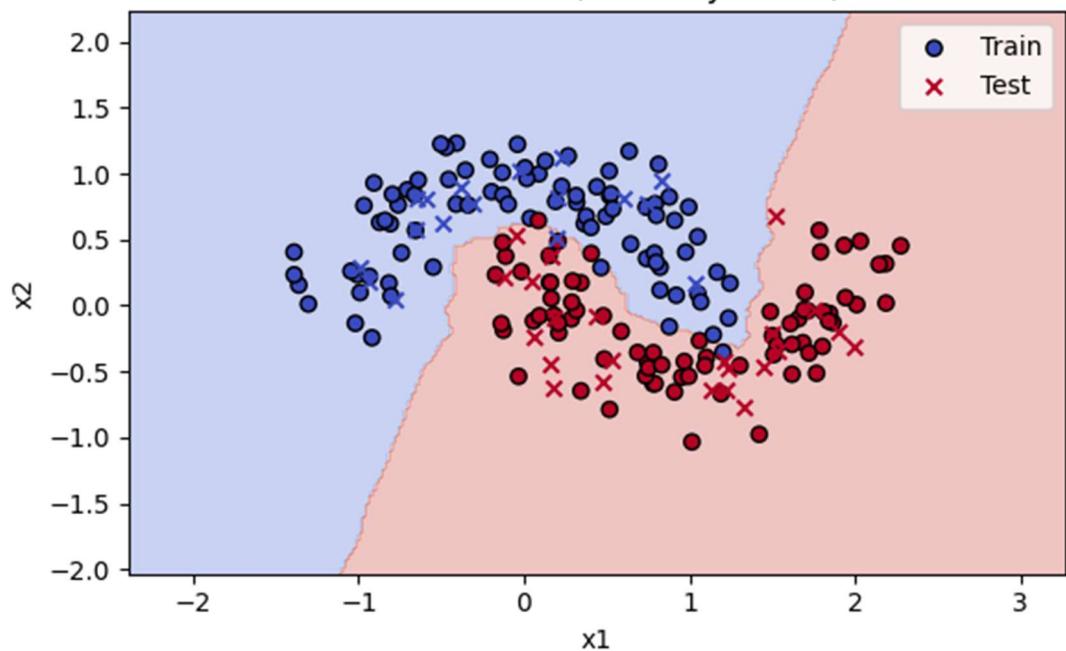
For **m**, as the sample size gets bigger, the SVM becomes more accurate for the same **C**.

That's because in PAC learning there is a minimum **m** required for **H** to be PAC-learnable.

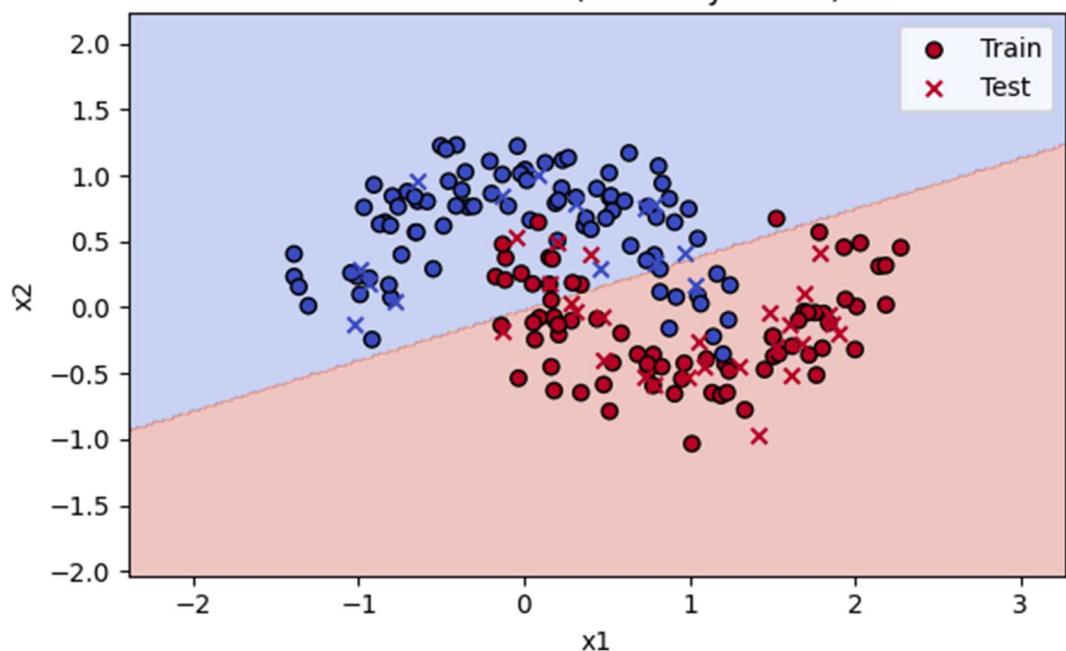
So as **m** increases, we approach that boundary such that:

$$m \geq \frac{(\log(|H|) * \log(1/\delta))}{\epsilon}$$

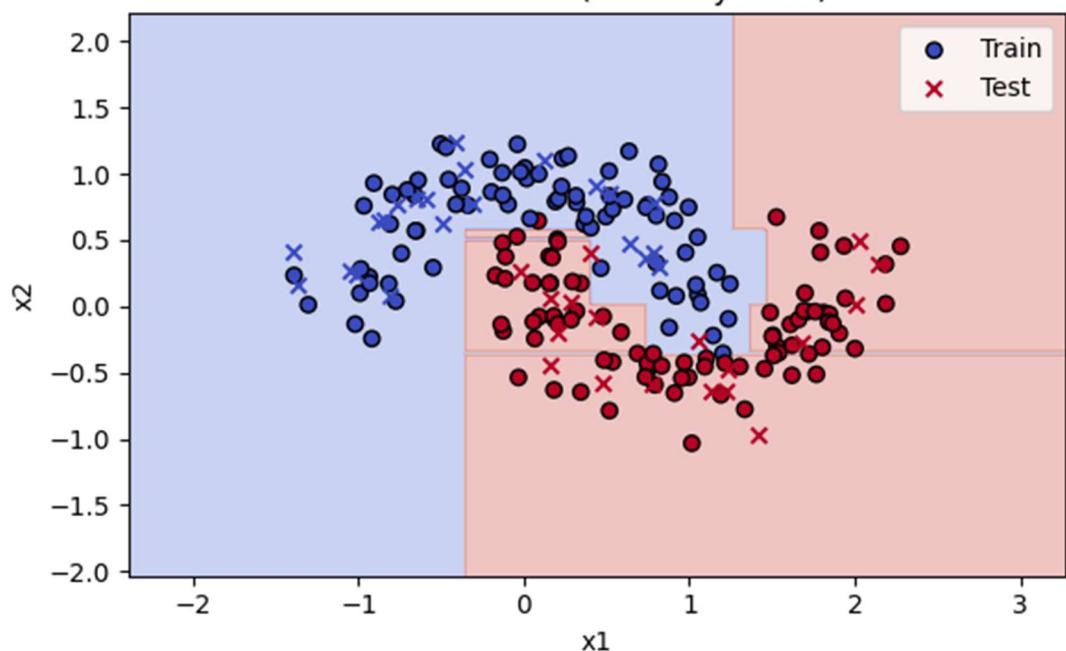
KNN on moon (Accuracy: 0.975)



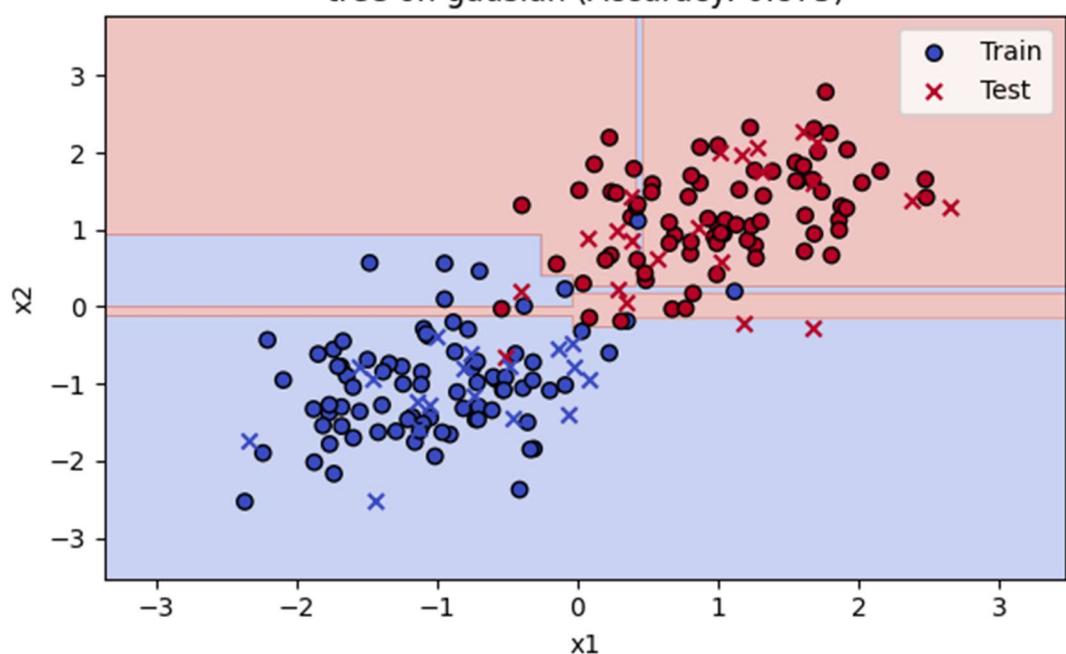
SVM on moon (Accuracy: 0.875)



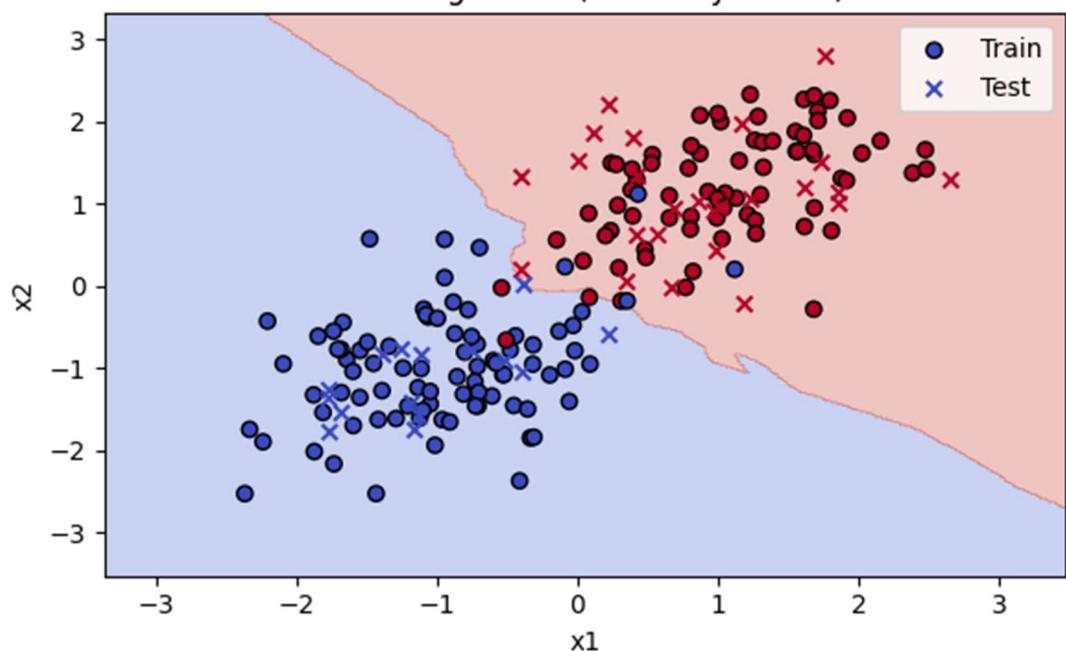
tree on moon (Accuracy: 0.95)



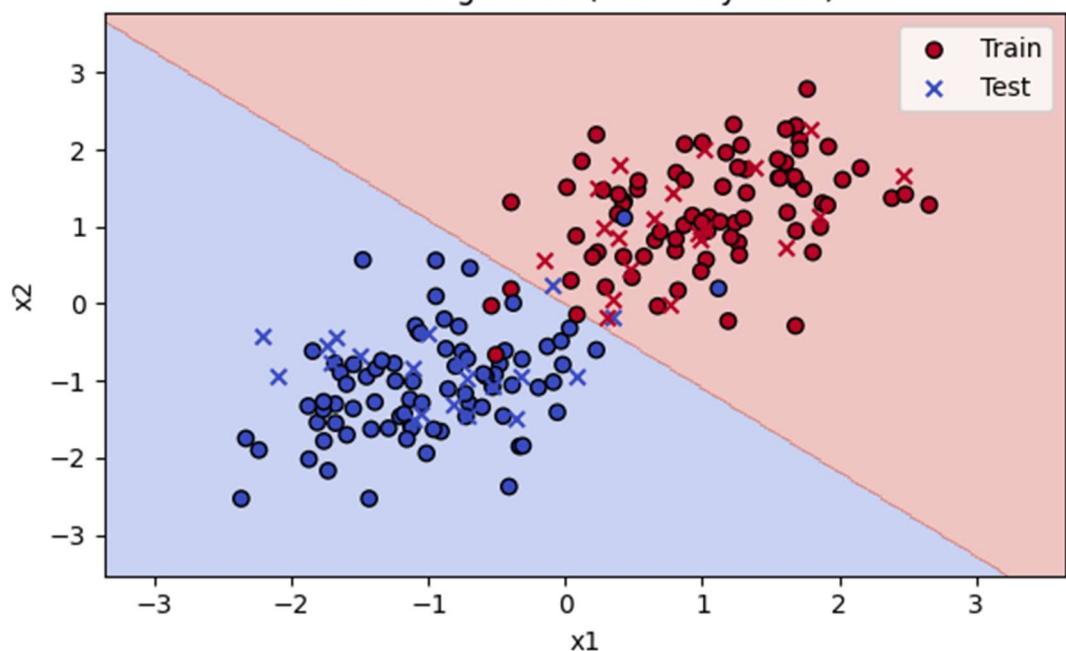
tree on gausian (Accuracy: 0.875)



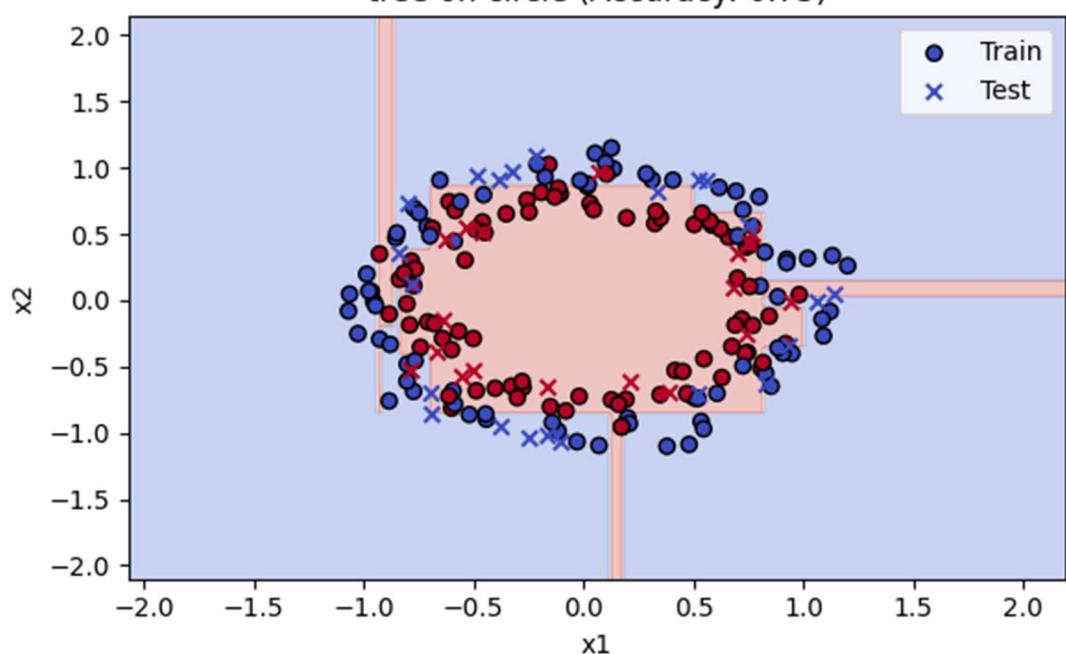
KNN on gaussian (Accuracy: 0.975)



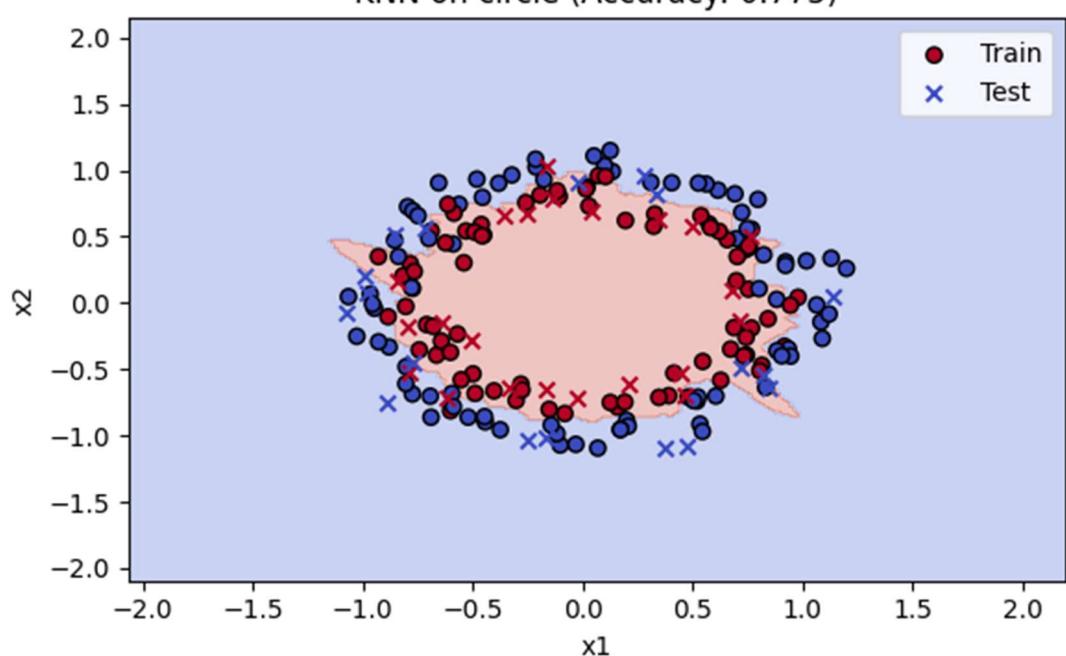
SVM on gaussian (Accuracy: 0.95)



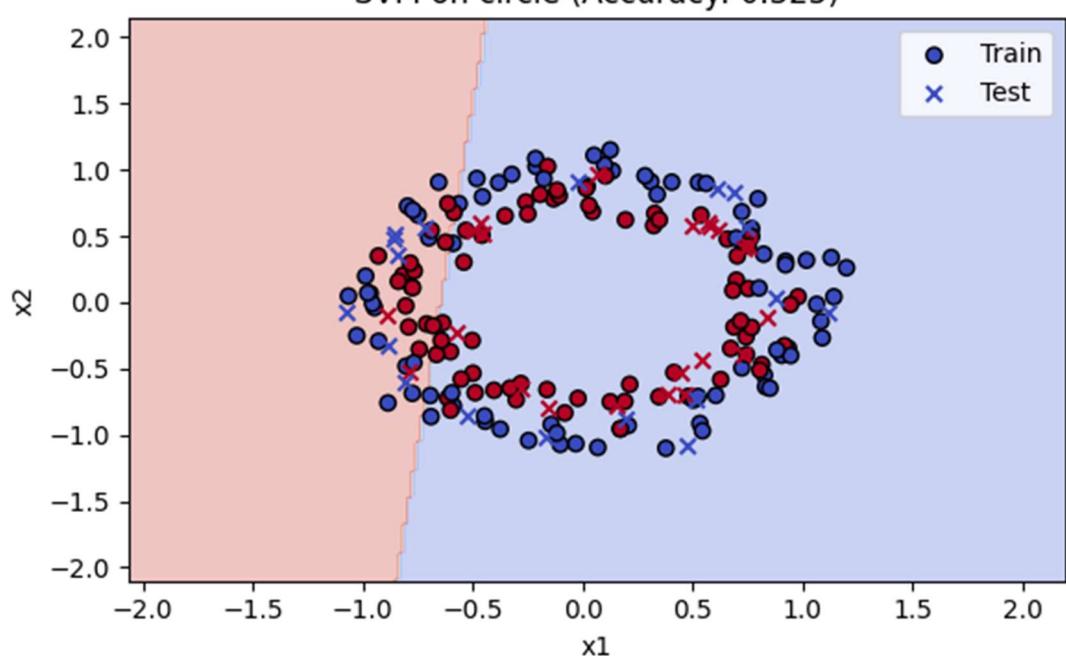
tree on circle (Accuracy: 0.75)



KNN on circle (Accuracy: 0.775)



SVM on circle (Accuracy: 0.325)



2.2 Classification Boundaries and Model Comparison

1.abov

2.

Moon – KNN (acc = 0.975), Tree (acc = 0.95), SVM (acc = 0.875)

Because the data is not linear, SVM can't completely separate them.

Tree performs better because we can separate better by dividing into groups, and KNN gives the best result because the data isn't mixed, which means the closest points to a point are the right choice.

Circle – KNN (acc = 0.775), Tree (acc = 0.75), SVM (acc = 0.325)

Here KNN performs best but still less accurately than KNN on moon because it's much less clear where one class starts and the other ends.

Tree also performs okay for the same reason.

But SVM performs poorly because for every w we choose, the distance from the data is the same.

Gaussian – KNN (acc = 0.975), Tree (acc = 0.875), SVM (acc = 0.95)

Tree performs worst because it only creates rectangles, which causes it to misclassify whole regions

unlike KNN that works best because the different classes are grouped in the same region.

That also allows SVM to stretch a line between them.

3.

SVM – It maximizes the margin between the two classes and is sensitive to noise and outliers.

We can use **delta** or **C** to decide how much we care about noise.

Tree – Stretches lines parallel to the axes with a pre-decided depth and splits them recursively. High k will cause overfitting while small one will misclassify

KNN – Based on the k nearest neighbors, we decide how to classify a point.

If we choose small k , it can cause overfitting, but for large k , we can misclassify a lot.

I used chat gpt for typos here part of the plotting such as how to draw background and some documentation

k

