

Constrained Optimization in Machine Learning



Dr. Pritam Anand.
Assistant Professor,
DA-IICT, Gandhinagar.

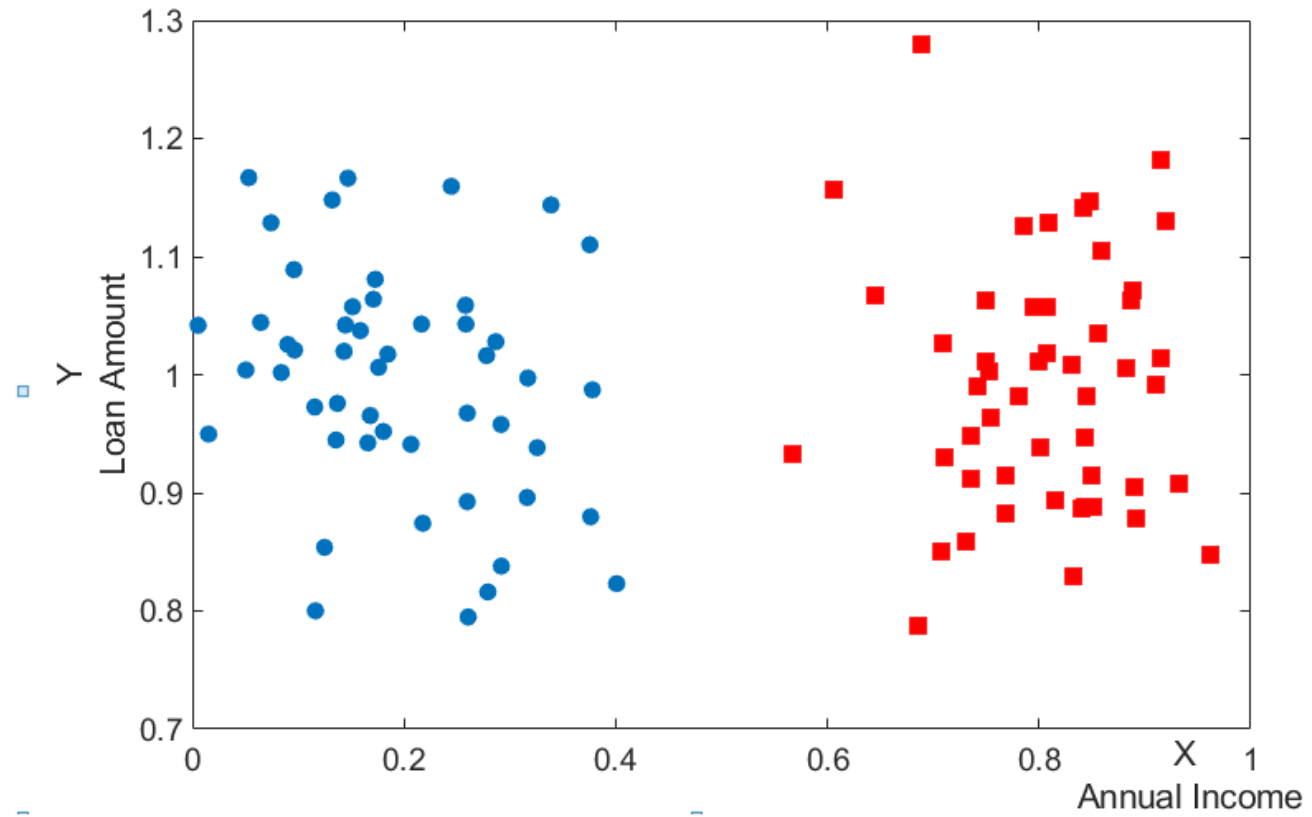
Takeaways

- *Develop the understanding of the Convex Programming Problems.*
- *Aware with the constrained optimization methods used in the solution of Support Vector Machine.*
- *Develop the detailed understanding of the Support Vector Machine.*
- *Understand the dual and primal relationship in context of Support Vector Machine problem.*
- *Ability to use the Support Vector Machine in different domain of applications.*

Load Defaulter Dataset

Index	Employed	Bank Balance (in thousands rupees)	Annual Salary (in million rupees)	Loan Amount (in thousands rupees)	Defaulted
1	1	0.4721	0.1358	0.9448	-1
2	0	0.8412	0.3169	0.9972	-1
3	0	0.3687	0.1249	0.8539	-1
4	0	0.2547	0.8416	1.1406	1
5	1	0.3111	0.7502	1.014	1

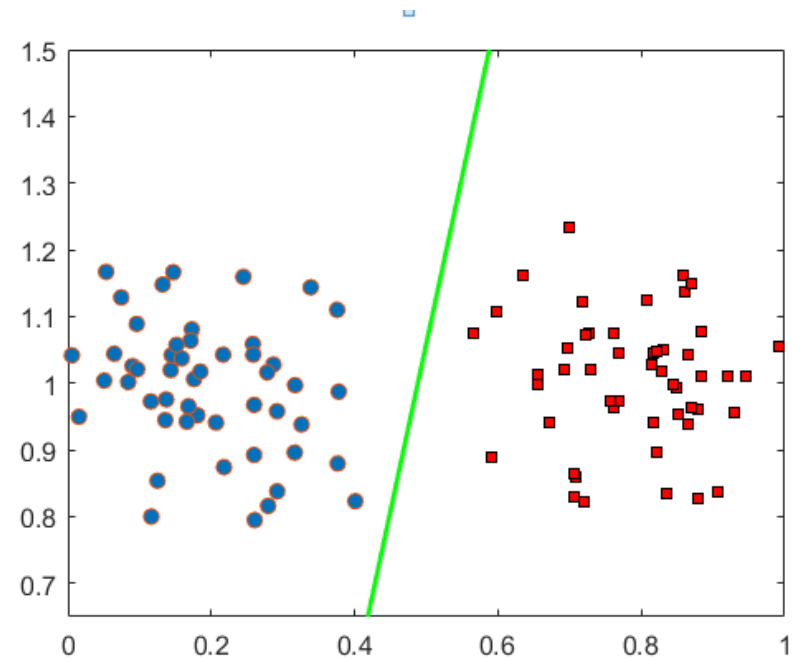
A random classifier



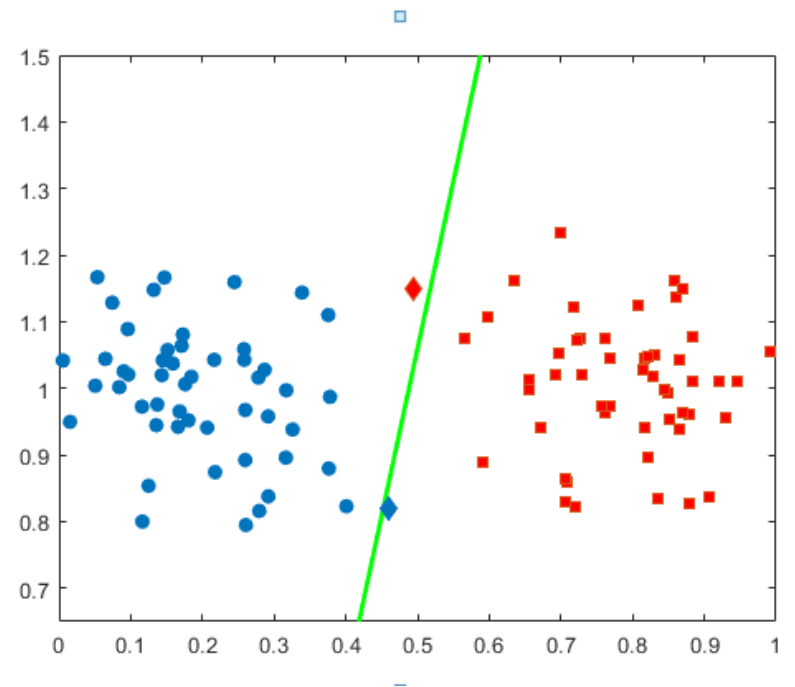
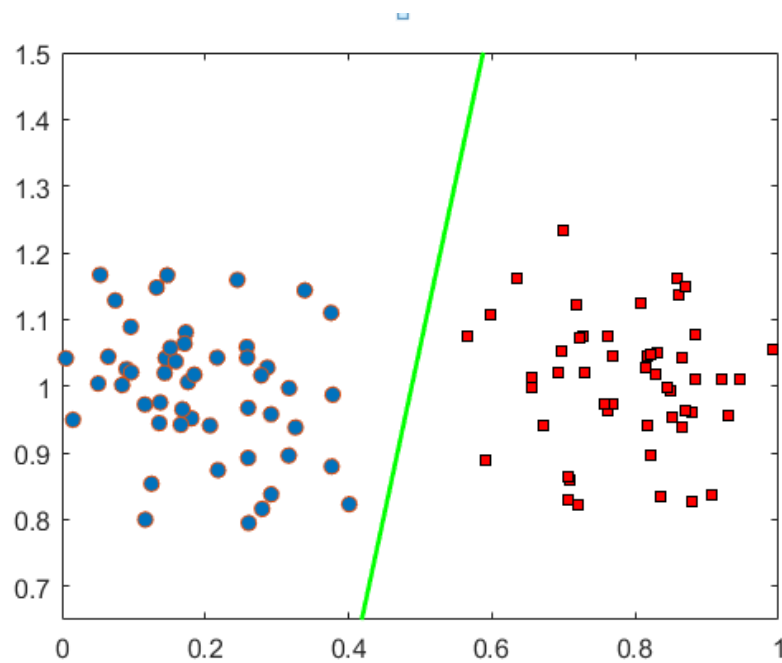
A random Classifier

$$\beta_1 x_1 + \beta_2 x_2 + \beta_0 = 0$$

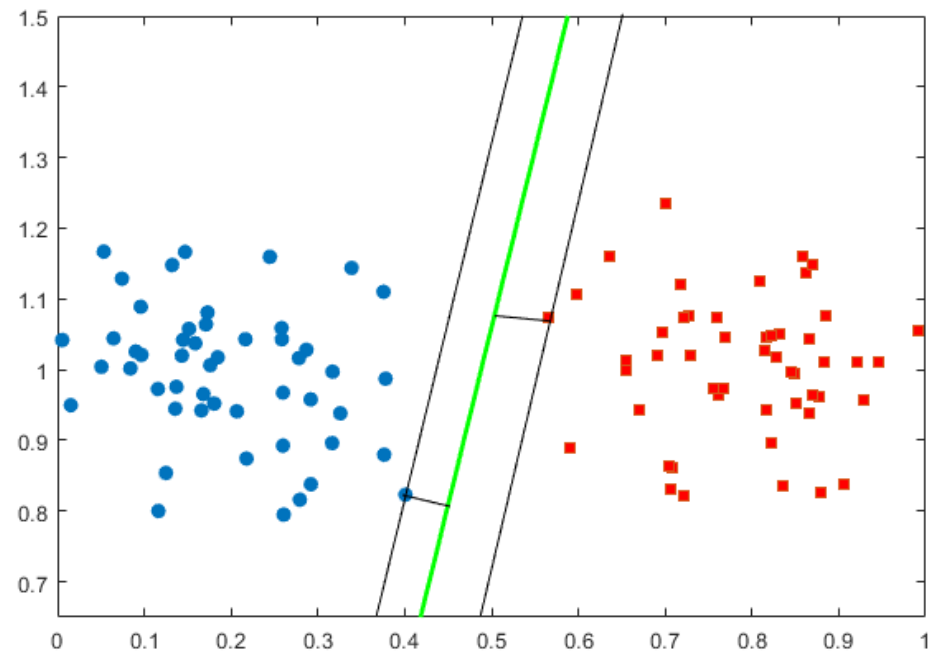
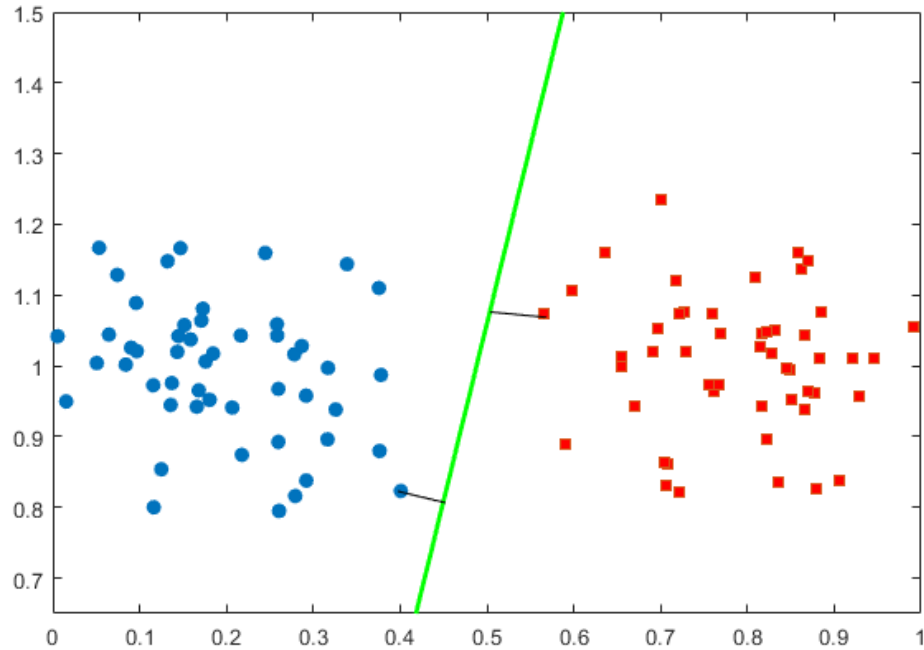
What is Problem ??

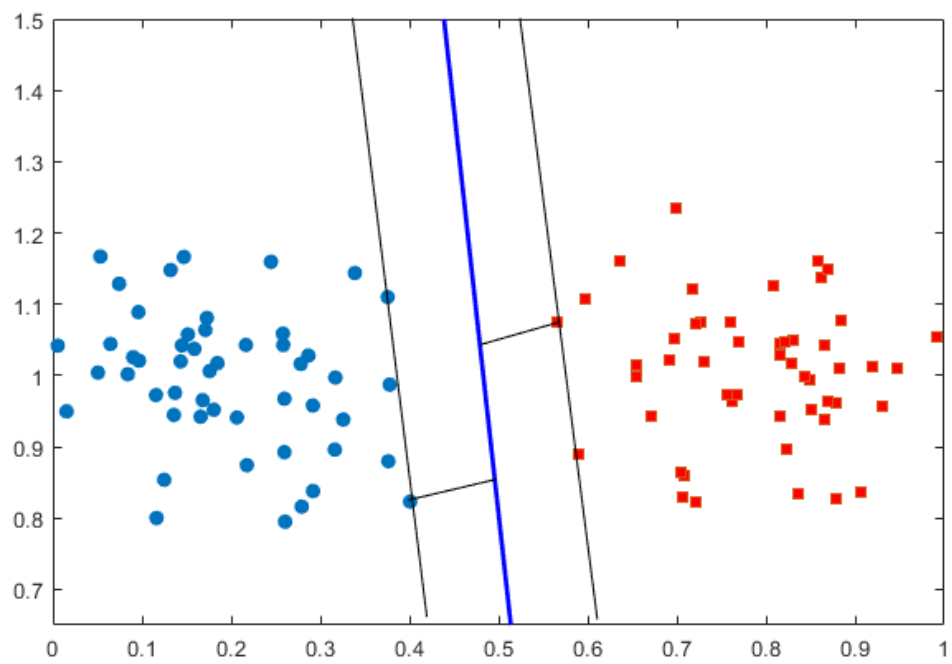


More Susceptible to misclassification

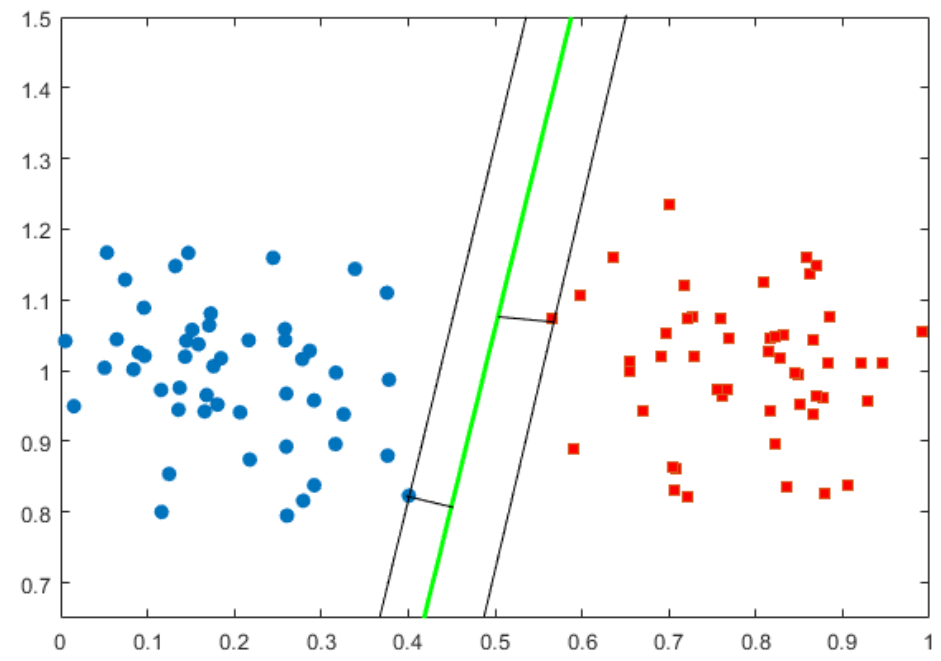
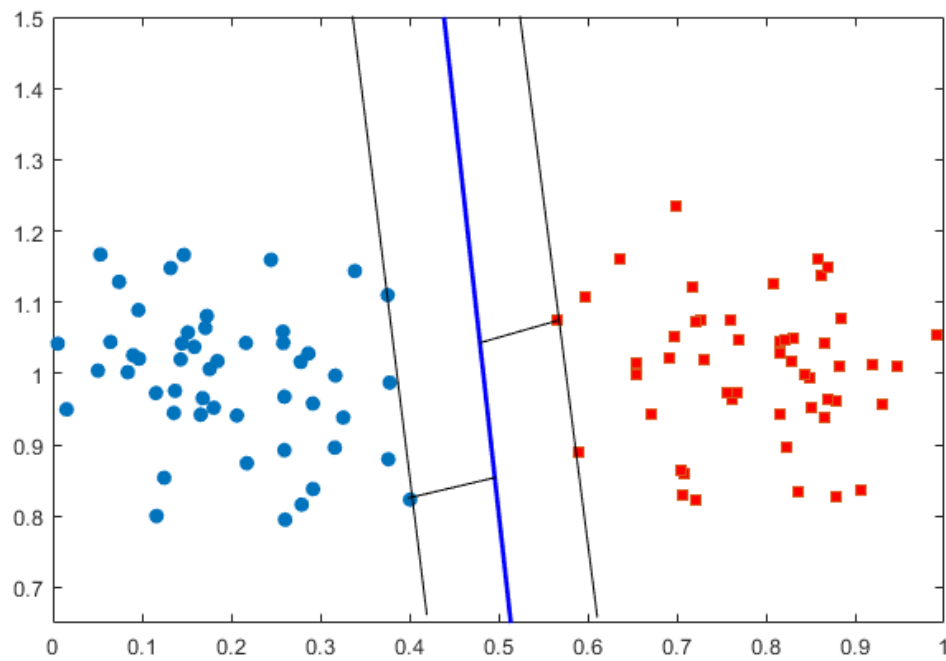


Margin of random classifier

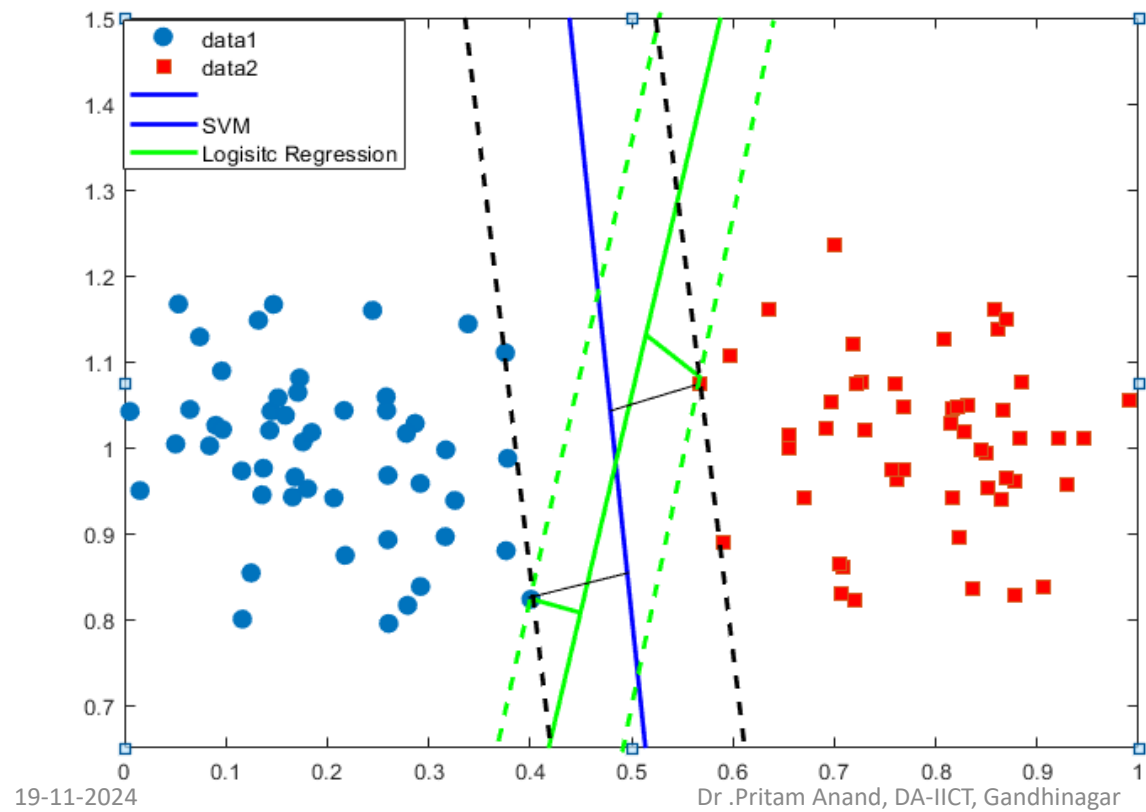




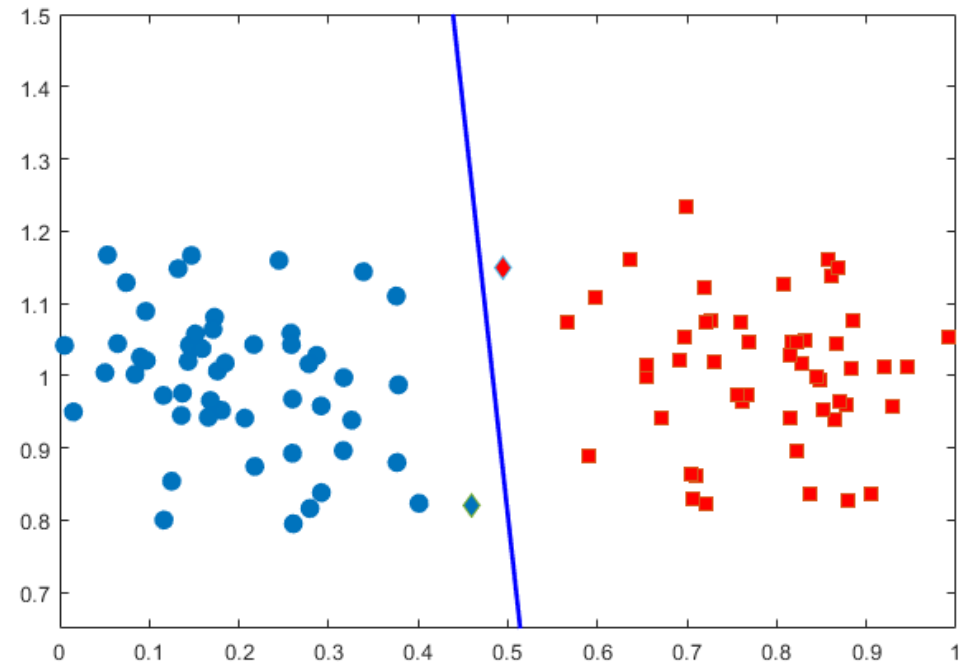
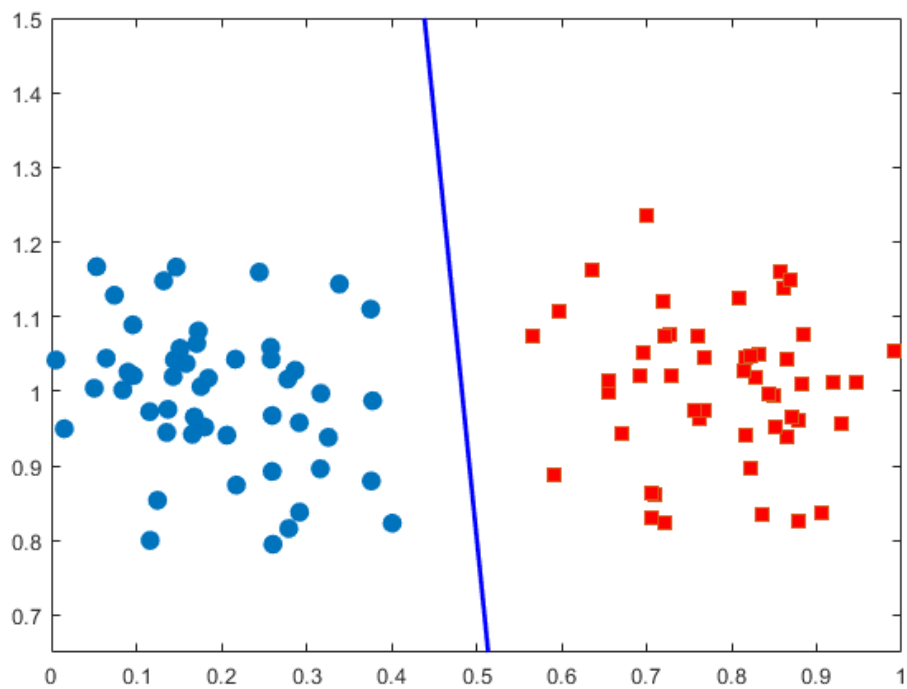
Maximal Margin Classifier



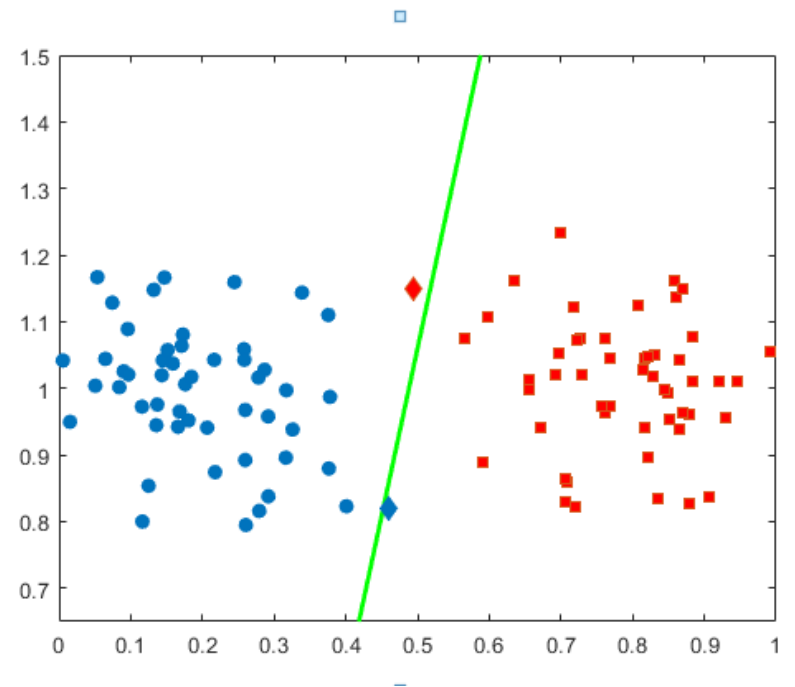
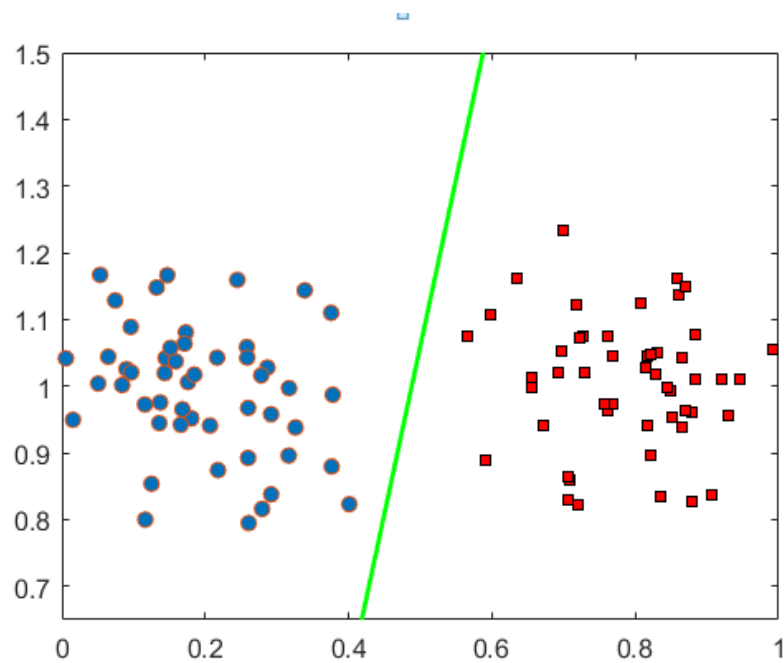
Width of margin



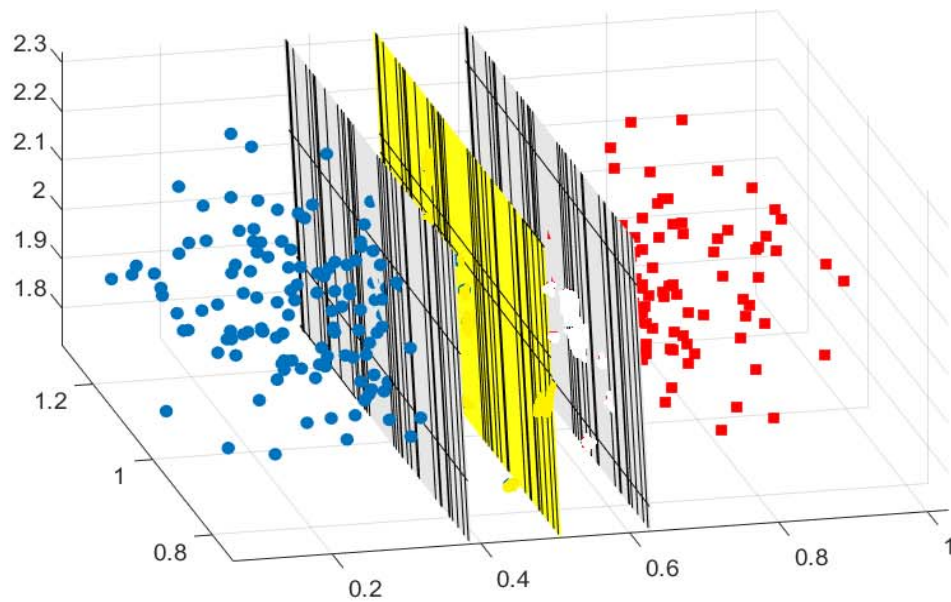
Checking the SVM



More Susceptible to misclassification



SVM margin



Hyperplane

$$\hat{w}^T x = c$$

$$\hat{w}^T x - c = 0$$

$$\frac{w}{\|w\|}^T x - c = 0$$

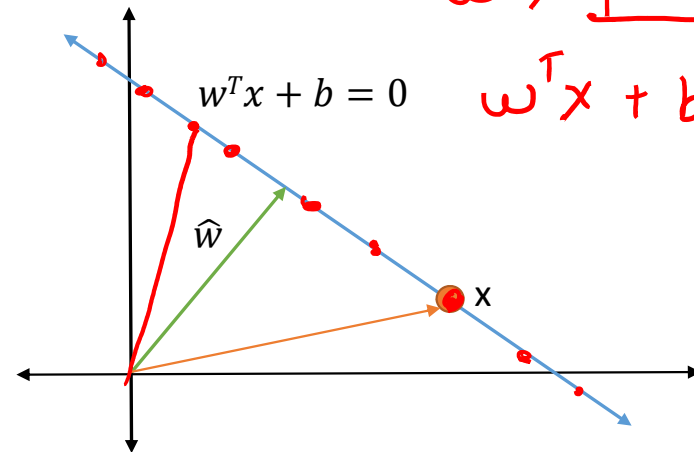
• Hyperplane:- A set of points in \mathbb{R}^n satisfying $w^T x + b = 0, w \in \mathbb{R}^n, b \in \mathbb{R}$.

• For $n = 2$, it is a line \mathbb{R}^2 .

• For $n = 3$, it is a plane in \mathbb{R}^3 .

$$w^T x \boxed{- c \|w\|} = 0$$

$$w^T x + b = 0$$



$$\hat{w}^T x = \frac{-b}{\|w\|}, \text{ where } \hat{w} = \frac{w}{\|w\|}$$

$$\hat{w}^T x = c$$

$$\hat{w}^T x = -b$$

$$\hat{w}^T x$$

$$\hat{w} = \frac{w}{\|w\|}$$

Hyperplane , Projection and Distances

$$\hat{w}^T x + b = 0$$

$$\frac{\hat{w}^T}{\|\hat{w}\|} \|w\| x + b \|w\| = 0$$

$$\Rightarrow w^T x + b \|w\| = 0$$

$$w^T x + b$$

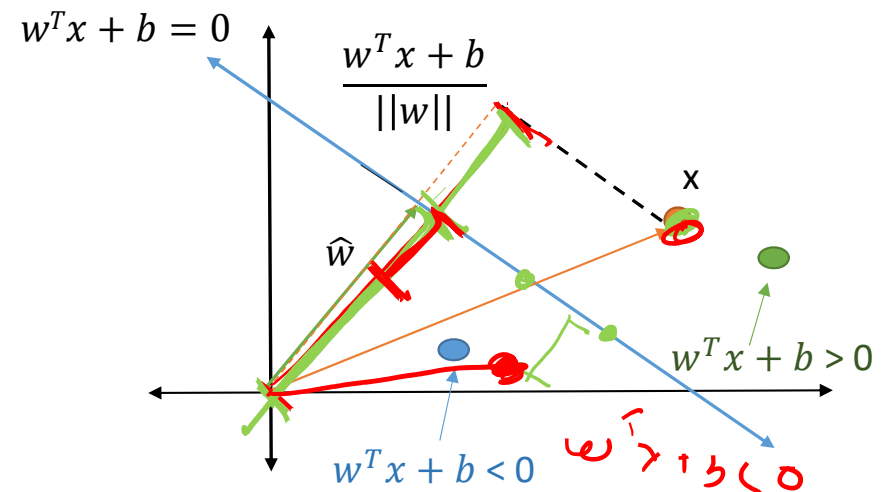
$$\hat{w}^T x - C$$

$$\frac{w^T}{\|w\|} x - C = \frac{w^T x - C \|w\|}{\|w\|} = \frac{w^T x + b}{\|w\|}$$

$$\hat{w}^T x$$

$$\frac{w^T x + b}{\|w\|} > 0$$

$$w^T x + b > 0$$



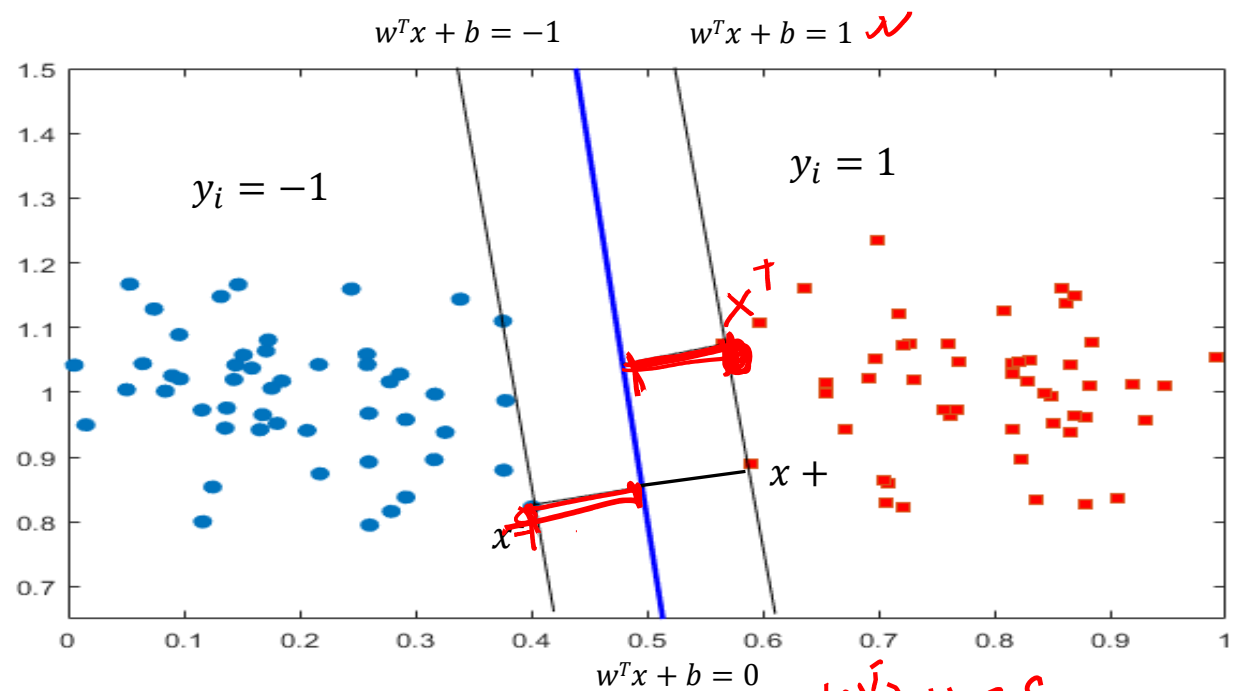
Margin in SVM

$$\frac{(w^T x^+ + b) - (w^T x^- + b)}{\|w\|} = \frac{1 - (-1)}{\|w\|} = \frac{2}{\|w\|}$$

Width of the margin:-

$$\frac{(w^T x^+ + b) - (w^T x^- + b)}{\|w\|}$$

$$= \frac{2}{\|w\|}$$



$$w^T x + b = \delta$$

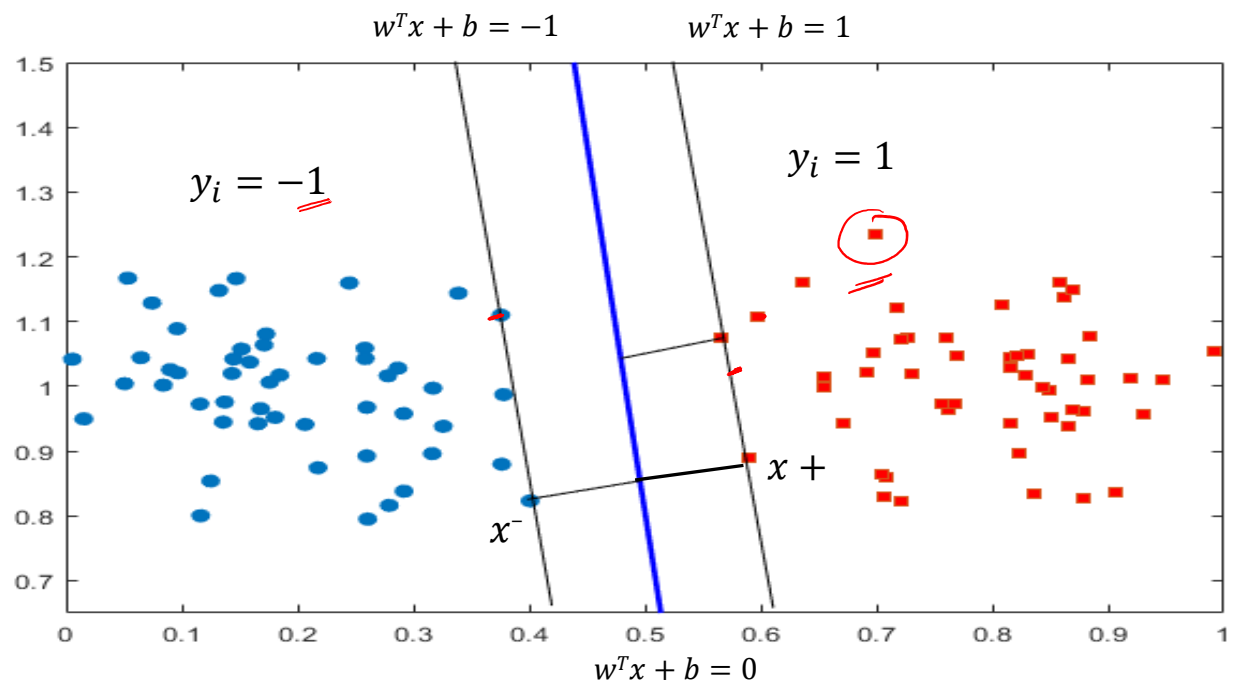
$$w^T x + b = -\delta$$

$$\left(\frac{w}{\delta}\right)^T x + \frac{b}{\delta} = -1$$

Optimization in SVM

Need to maximize the width of the margin but, subject to certain constraints

$$\begin{aligned} & \max_{(w,b)} \frac{2}{\|w\|} \\ & \text{subject to,} \\ & (w^T x_i + b) \geq 1, \text{ if } y_i = 1. \\ & (w^T x_i + b) \leq -1, \text{ if } y_i = -1. \end{aligned}$$

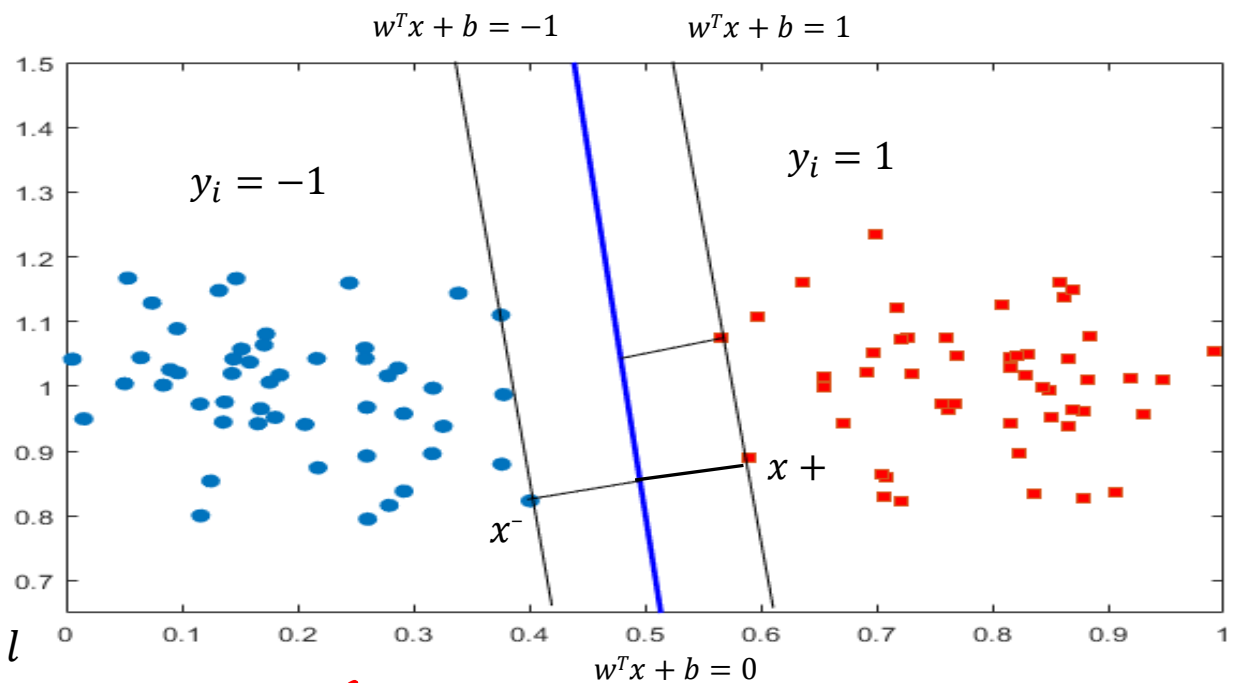


Optimization in SVM

$$\begin{aligned} & \max_{(w,b)} \frac{2}{\|w\|} \\ & \text{subject to,} \\ & (w^T x_i + b) \geq 1, \text{ if } y_i = 1. \\ & (w^T x_i + b) \leq -1, \text{ if } y_i = -1. \end{aligned}$$

is equivalent to

$$\begin{aligned} & \max_{(w,b)} \frac{2}{\|w\|} \\ & \text{subject to,} \\ & y_i (w^T x_i + b) \geq 1, i = 1, 2, \dots, l \end{aligned}$$



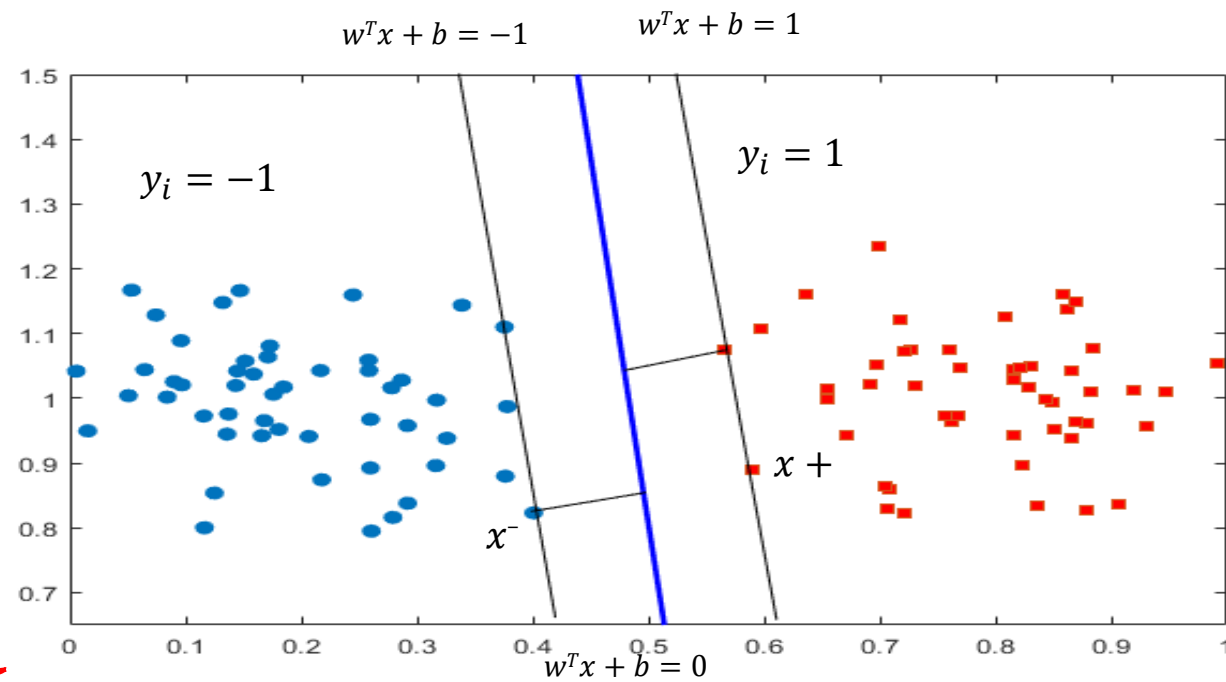
$$-w^T x_i + b \leq -1$$

$$\begin{aligned} & (w^T x_i + b) \geq 1 \text{ if } y_i = 1 \\ & \leftarrow -(w^T x_i + b) \geq -1 \text{ if } y_i = -1 \end{aligned}$$

Optimization in SVM

$$\max \frac{2}{\|w\|} \quad \min \frac{\|w\|^2}{2}$$

- $\max_{(w,b)} \frac{2}{\|w\|}$
subject to,
 $y_i (w^T x_i + b) \geq 1, i = 1, 2, \dots, l.$
- $\min_{(w,b)} \frac{\|w\|^2}{2}$
subject to,
 $y_i (w^T x_i + b) \geq 1, i = 1, 2, \dots, l.$



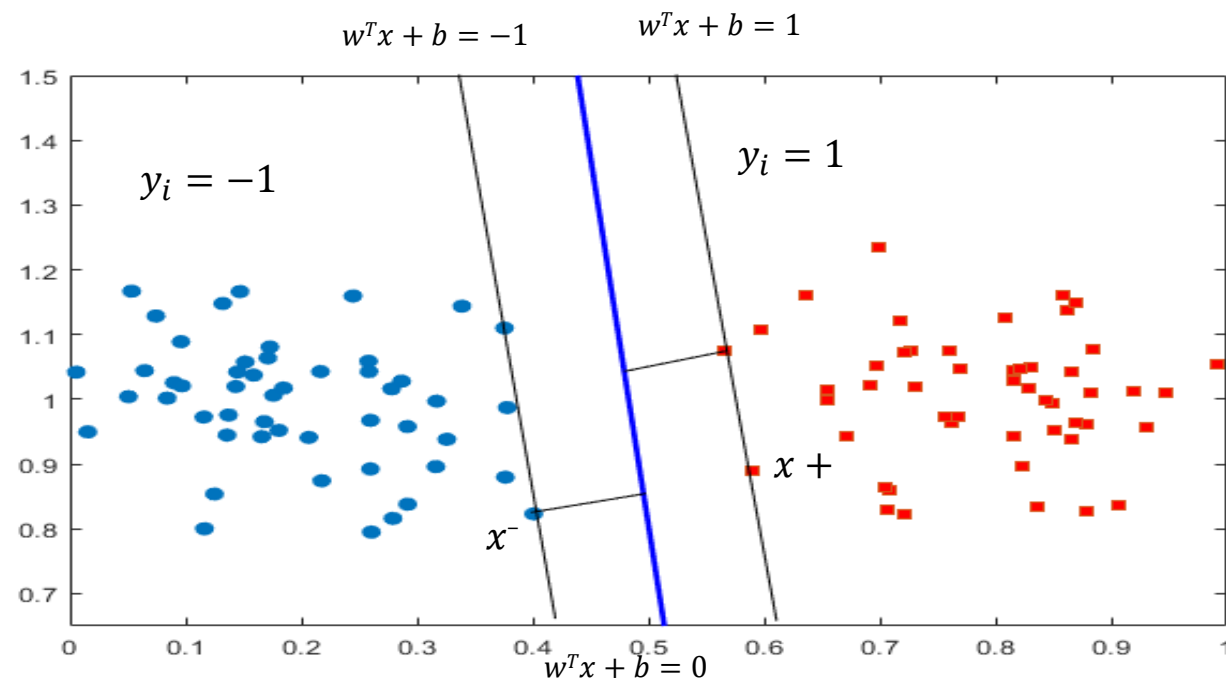
$$\sqrt{w_1^2 + w_2^2 + \dots + w_n^2}$$

$$W^T W$$

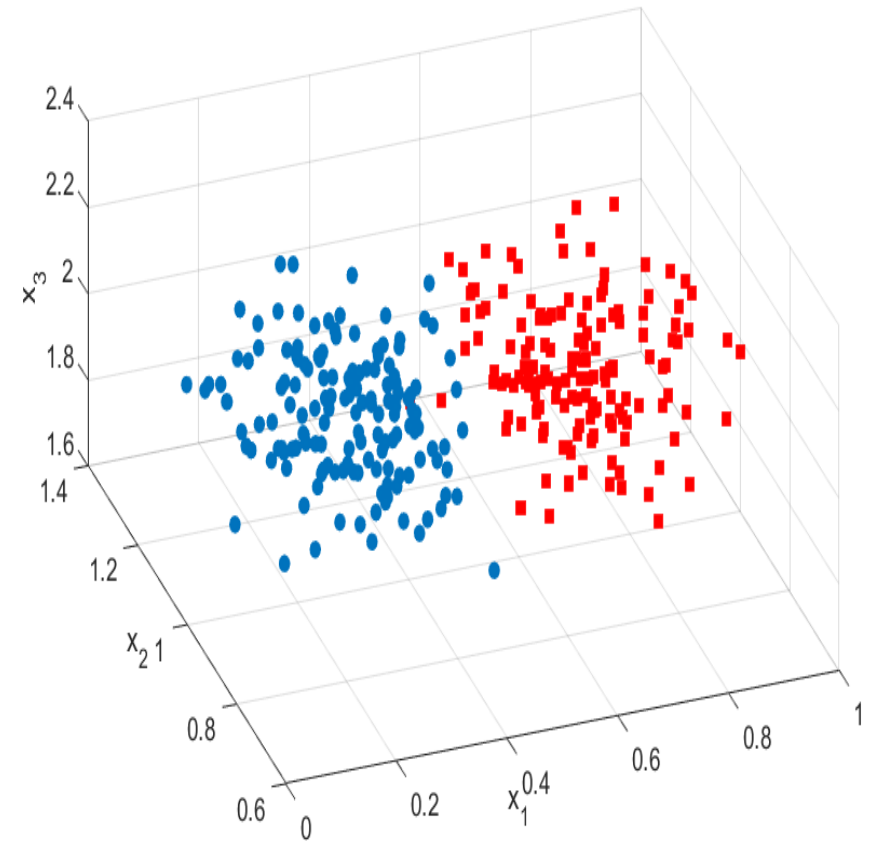
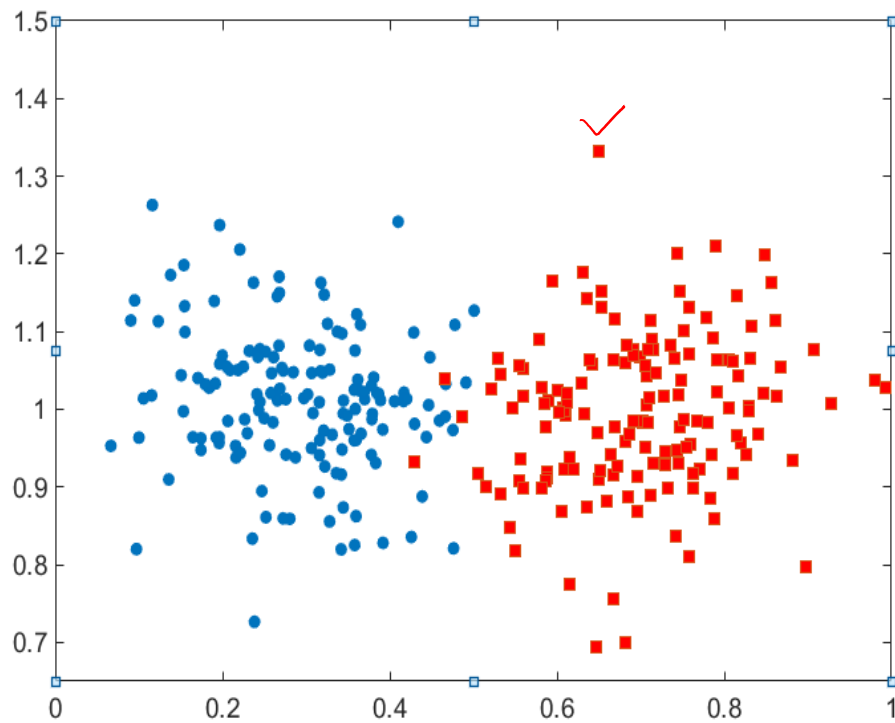
Optimization in SVM

- $\max_{(w,b)} \frac{2}{||w||}$
subject to,
 $y_i (w^T x_i + b) \geq 1, i = 1, 2, \dots, l.$
- $\min_{(w,b)} \frac{||w||^2}{2}$
subject to,
 $y_i (w^T x_i + b) \geq 1, i = 1, 2, \dots, l.$

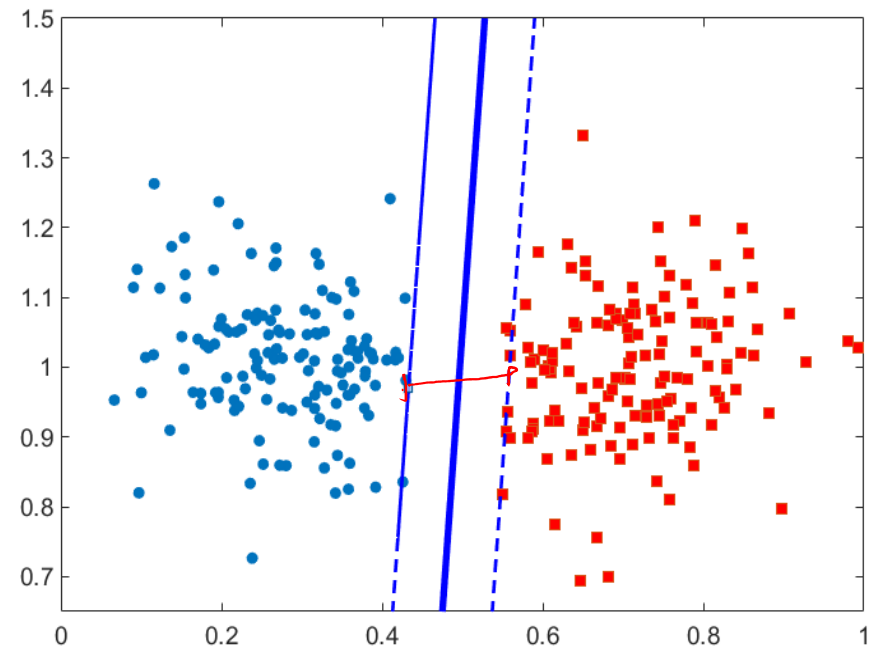
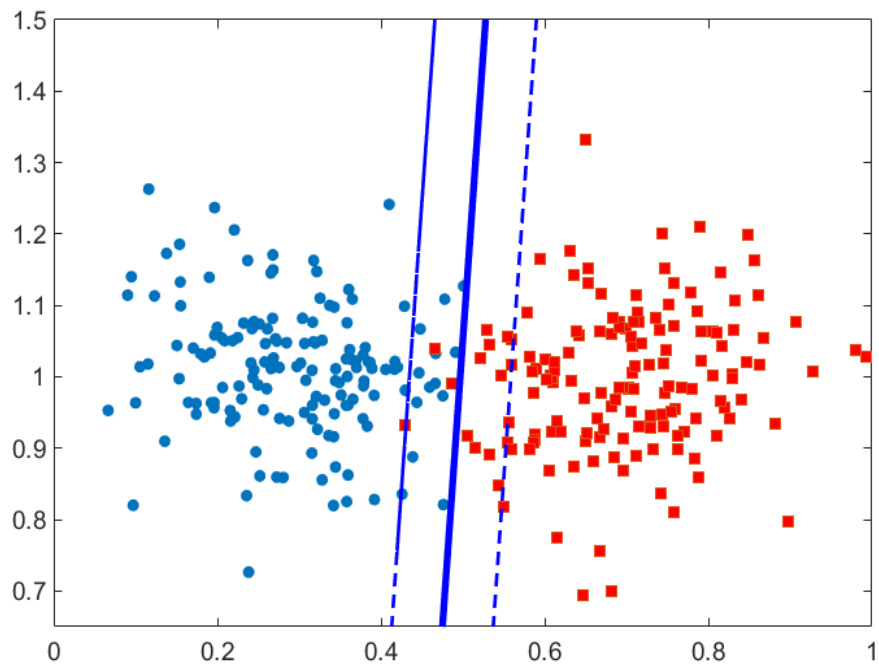
Hard-margin SVM problem



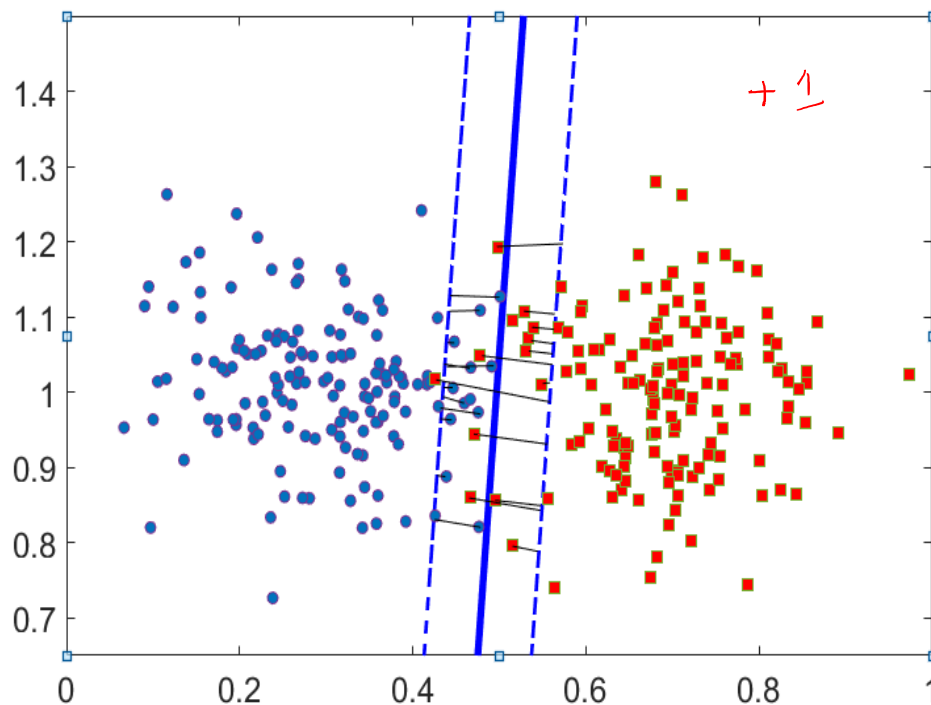
What to do for this



Soft-margin SVM



Soft-margin SVM



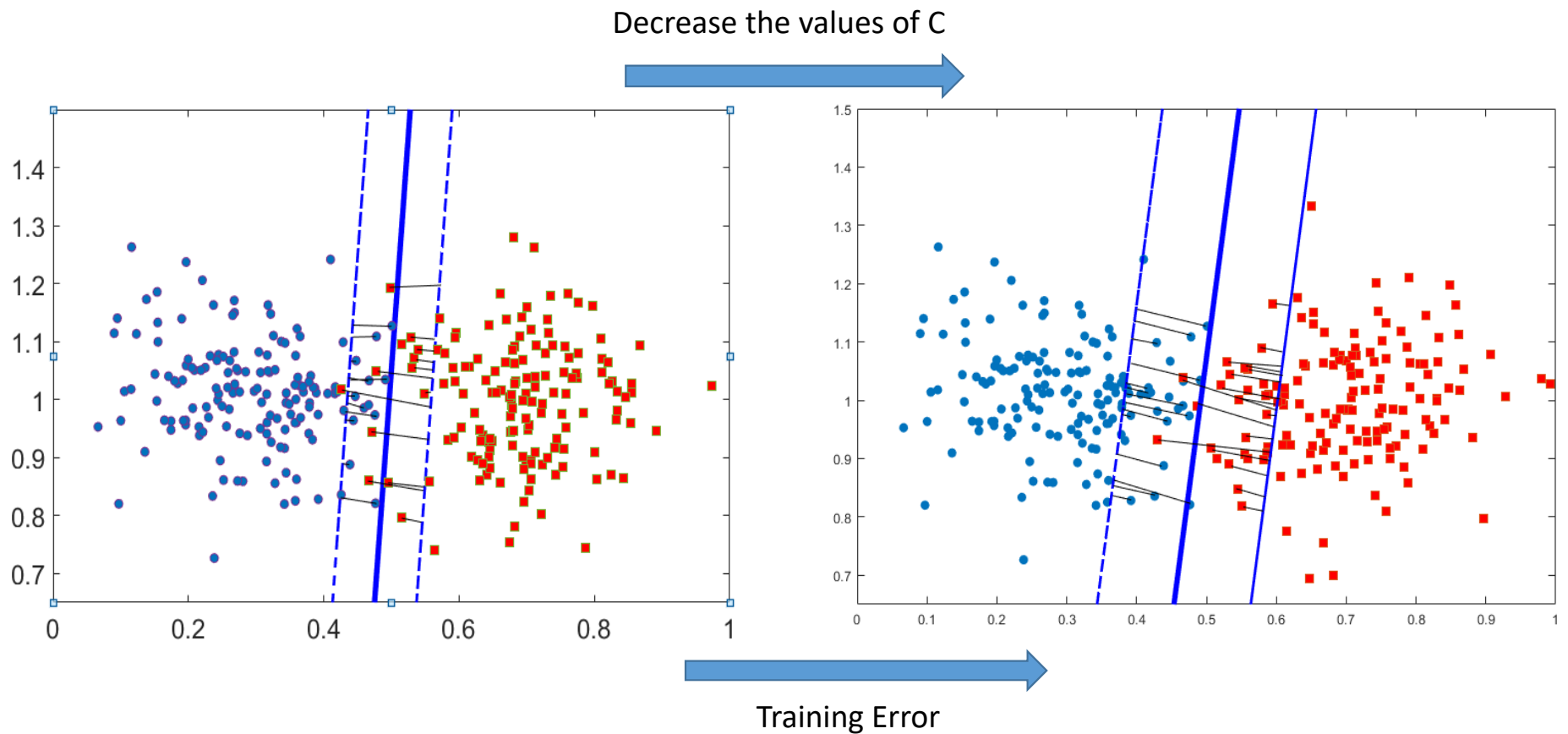
Relax the constraints by introducing slack variable ξ .

$$\begin{aligned} & \min_{(w,b)} \frac{1}{2} w^T w + C \sum_{i=1}^l \xi_i \\ & \text{subject to,} \\ & y_i (w^T x_i + b) \geq 1 - \xi_i, \\ & \xi_i \geq 0, \quad i = 1, 2, \dots, l. \end{aligned}$$

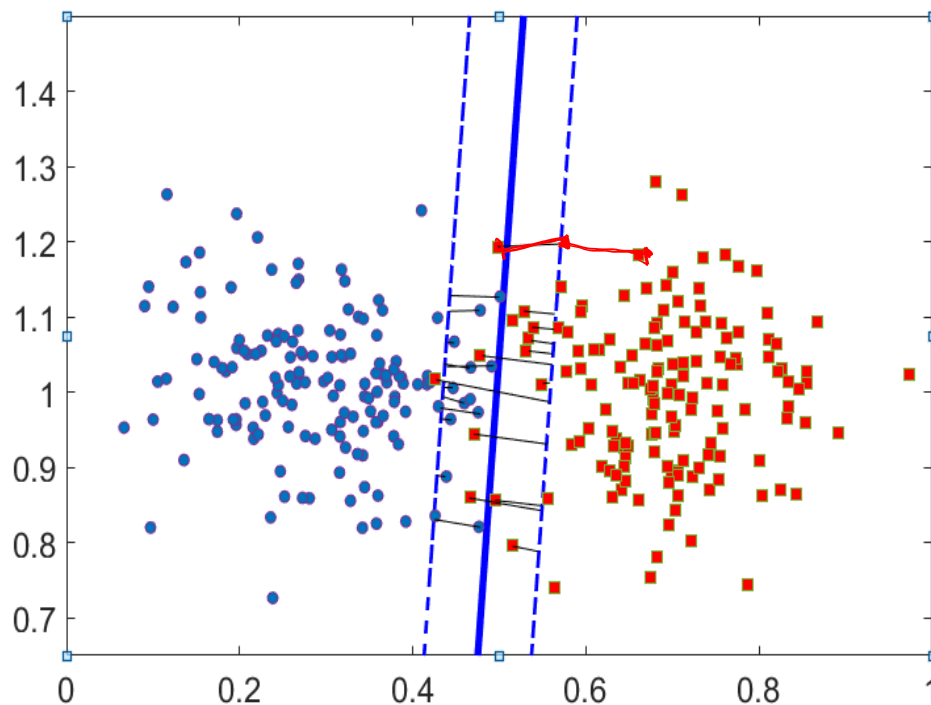
where $C > 0$ is the user defined parameter

Soft -margin SVM problem

Trading-off the training error and width of margin



Soft-margin SVM optimization problem



$$\begin{aligned} \min_{(w,b)} \quad & \frac{1}{2} w^T w + C \sum_{i=1}^l \xi_i \\ \text{subject to,} \quad & y_i (w^T x_i + b) \geq 1 - \xi_i, \\ & \xi_i \geq 0, \quad i = 1, 2, \dots, l. \end{aligned}$$

where $C > 0$ is the user defined parameter.

- Limitation :- The number of constraint becomes huge for large-scale datasets

$$\min_{w, b} \frac{1}{2} w^T w + C \sum_{i=1}^L \xi_i$$

$$\min (\max(a, b))$$

$$\min c$$

$$c > a$$

$$c > b$$

①

$$y_i(w^T x_i + b) \geq 1 - \xi_i, \quad i=1, 2, \dots, L$$

$$\xi_i \geq 0, \quad i=1, 2, \dots, L$$

$$\min_{w, b} \frac{1}{2} w^T w + C \sum_{i=1}^L \max(0, 1 - y_i(w^T x_i + b))$$

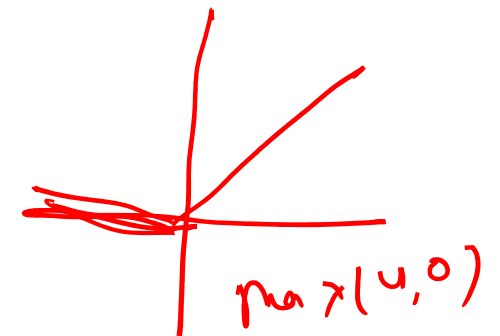
②

$$\xi_i =$$

$$\xi_i = \max(0, 1 - y_i(w^T x_i + b))$$

$$\min_{w, b} \frac{1}{2} w^T w + C \sum_{i=1}^L \xi_i$$

$$\xi_i \geq 0, \quad (\xi_i \geq 1 - y_i(w^T x_i + b))$$



$\min_{\substack{w, b \\ w \in \mathbb{R}^n \\ b \in \mathbb{R}}} J(w, b) = \frac{1}{2} w^T w + C \sum_{i=1}^L \max(\underbrace{1 - y_i(w^T x_i + b)}, 0)$

$y_i(w^T x_i + b) < 1$
 \swarrow
 $\underbrace{1 - y_i(w^T x_i + b)}_0$ if $y_i(w^T x_i + b) < 1$
 otherwise $\underline{0}$

$$\nabla_w J(w, b) = w + C \sum_{i=1}^L \begin{matrix} -x_i y_i \\ 0 \end{matrix} \quad \begin{matrix} \text{if } y_i(w^T x_i + b) < 1 \\ \text{otherwise} \end{matrix}$$

$$\nabla_b J(w, b) = C \sum_{i=1}^L \begin{matrix} -y_i \\ 0 \end{matrix} \quad \begin{matrix} \text{if } y_i(w^T x_i + b) < 1 \\ \text{otherwise} \end{matrix}$$

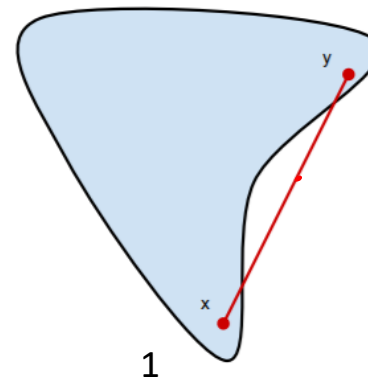
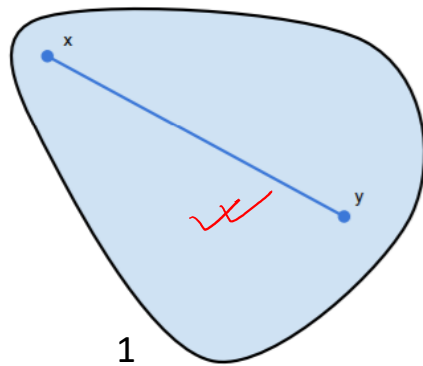
Convex Programming Problem

Convex Sets

- $C \subseteq \mathbb{R}^n$ is convex

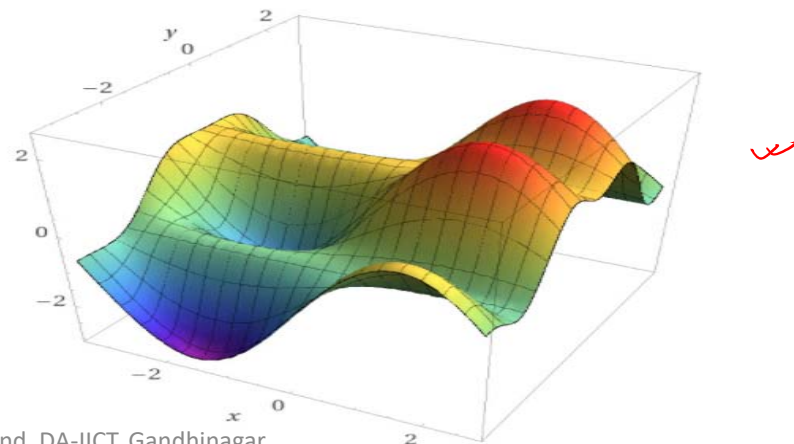
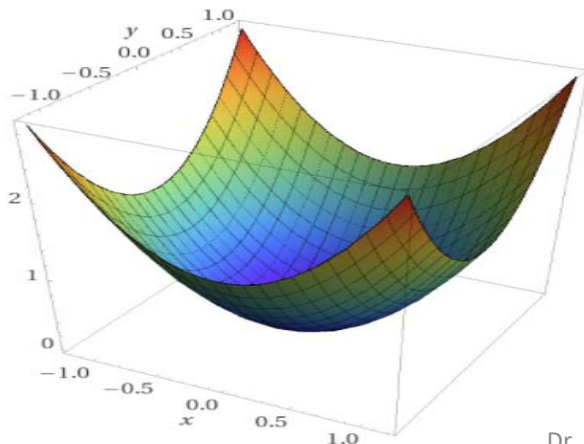
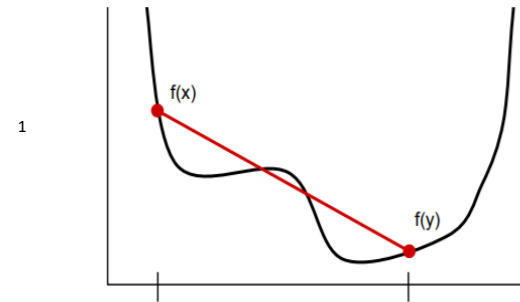
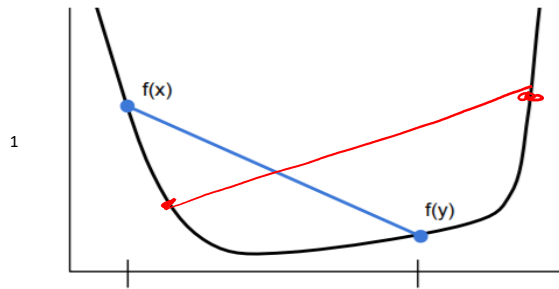
if $\lambda x + (1 - \lambda)y \in C$ for any $x, y \in C$ and $0 \leq \lambda \leq 1$.

that is, a set is convex if the line connecting any two points in the set is entirely inside the set.



Convex Functions

- $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex
 - if $\text{dom}(f)$ (the domain of f) is a convex set,
 - and if $f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y)$ for any $x, y \in \text{dom}(f)$ and $0 \leq \lambda \leq 1$.
 - that is, the line connecting any two points on the graph of the function stays above the graph.



19-11-2024

Dr .Pritam Anand, DA-IICT, Gandhinagar

30

Convex Programming Problem

A convex programming problem is an optimization problem in the form

$$\left\{ \begin{array}{l} \min_x \quad \underline{f_0(x)} \rightarrow \text{Convex func of } x \\ \text{subject to,} \\ \leftarrow \underline{f_i(x)} \leq 0, i = 1, 2, \dots, m, \\ \underline{h_i(x)} = 0, i = 1, 2, \dots, p, \end{array} \right.$$

Convex func

$$\begin{array}{l} g_i(x) \geq 0 \\ \text{---} g_i(x) \leq 0 \end{array} \quad f(x)$$

where $f_0(x)$ and $f_i(x), i = 1, \dots, m$ are continuous convex functions of \mathbb{R}^n , and $h_i(x), i = 1, \dots, p$ are linear functions.

Convex Programming Problem

- A Convex Programming Problem (CPP) is an optimization problem in the form

$$\min_x f_0(x)$$

subject to ,

$$f_i(x) \leq 0, i = 1, 2, \dots, m,$$

$$h_i(x) = 0, i = 1, 2, \dots, p,$$

where $f_0(x)$ and $f_i(x), i = 1, \dots, m$ are continuous convex functions of \mathbb{R}^n , and $h_i(x), i = 1, \dots, p$ are linear functions.

Convex Programming Problem

- A Convex Programming Problem (CPP) is an optimization problem in the form

$$\begin{aligned} \min_x \quad & f_0(x) \\ \text{subject to,} \quad & \\ & f_i(x) \leq 0, i = 1, 2, \dots, m, \\ & h_i(x) = 0, i = 1, 2, \dots, p, \end{aligned} \tag{1}$$

where $f_0(x)$ and $f_i(x), i = 1, \dots, m$ are continuous convex functions of \mathbb{R}^n , and $h_i(x), i = 1, \dots, p$ are linear functions.

- If x^* is its local solution of CPP (1), then x^* is also its global solution.

Duality Theory

$$D = \left\{ x : f_i(x) \leq 0, i = 1, 2, \dots, m \text{ and } h_j(x) = 0, j = 1, 2, \dots, p \right\}$$

- A Convex Programming Problem (CPP) is an optimization problem in the form

$$\begin{cases} \min_x f_0(x) \\ \text{subject to,} \\ f_i(x) \leq 0, i = 1, 2, \dots, m, \\ h_i(x) = 0, i = 1, 2, \dots, p, \end{cases} \quad \begin{matrix} \text{Min } x^2 + y^2 \\ \text{Subject} \\ f_1(x^*) \leq 0 \rightarrow x + y \leq 1 \\ h_1(x^*) = 0 \end{matrix} \quad (1)$$

where $f_0(x)$ and $f_i(x), i = 1, \dots, m$ are continuous convex functions of \mathbb{R}^n , and $h_i(x), i = 1, \dots, p$ are linear functions.

- Let p^* is the optimal solution of CPP (1), i.e., $p^* = \inf\{f_0(x) \mid x \in D\}$, where, $D = \{x \mid f_i(x) \leq 0, i = 1, 2, \dots, m, h_i(x) = 0, i = 1, 2, \dots, p\}$

$$\max_{\lambda, \nu} \text{Lagrangian Function} \quad \inf_{x \in \mathbb{R}^n} L(x, \lambda, \nu) \leq \inf_{x \in D} L(x, \lambda, \nu) \leq \inf_{x \in D} f_0(x)$$

The Lagrangian function is given as

$$L(x, \lambda, \nu) = f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^p \nu_i h_i(x),$$

where $\lambda = (\lambda_1, \lambda_2, \dots, \lambda_m)$ and $\nu = (\nu_1, \nu_2, \dots, \nu_p)$ are the Lagrangian multipliers.

$[\lambda, \nu]$

$f(x) + \lambda x$

Lagrangian Function

The Lagrangian function is given as

$$L(x, \lambda, \nu) = f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^p \nu_i h_i(x),$$

where $\lambda = (\lambda_1, \lambda_2, \dots, \lambda_m)$ and $\nu = (\nu_1, \nu_2, \dots, \nu_p)$ are the Lagrangian multipliers.

Obviously, when $x \in D, \lambda \geq 0$, we have

$$L(x, \lambda, \nu) \leq f_0(x).$$

Thus, $\inf_{x \in \mathbb{R}^n} L(x, \lambda, \nu) \leq \inf_{x \in D} L(x, \lambda, \nu) \leq \inf_{x \in D} f_0(x) = p^*.$

Lagrangian dual function

- $\inf_{x \in R^n} L(x, \lambda, \nu) \leq \inf_{x \in D} L(x, \lambda, \nu) \leq \inf_{x \in D} f_0(x) = p^* .$
Consider the dual function

$$g(\lambda, \nu) = \inf_{x \in R^n} L(x, \lambda, \nu) , \text{ then}$$
$$g(\lambda, \nu) \leq p^*$$

- The above inequality indicates that, for any $\lambda \geq 0$, $g(\lambda, \nu)$ is a lower bound of p^* .
- Among these lower bounds, finding the best one leads to the optimization problem

$$\begin{aligned} \max g(\lambda, \nu) &= \inf_{x \in R^n} L(x, \lambda, \nu), \\ \text{subject to, } \lambda &\geq 0 \end{aligned}$$

(2)

Weak Duality

- Let d^* is solution of the dual problem

$$\begin{aligned} \max (g(\lambda, \nu) &= \inf_{x \in R^n} L(x, \lambda, \nu)), \\ \text{subject to, } \lambda &\geq 0 \end{aligned}$$



(2)

- Let p^* is the optimal solution of the CPP(1)

$$\begin{aligned} \min_x \quad & f_0(x) \\ \text{subject to,} \quad & \\ & f_i(x) \leq 0, i = 1, 2, \dots, m, \\ & h_i(x) = 0, i = 1, 2, \dots, p, \end{aligned}$$

then $d^* \leq p^*$

Slater's condition

The CPP(1)

$$\begin{aligned} \min_x \quad & f_0(x) \\ \text{subject to,} \quad & \\ & f_i(x) \leq 0, i = 1, 2, \dots, m, \\ & h_i(x) = 0, i = 1, 2, \dots, p, \end{aligned}$$

is said to satisfy the Slater's condition if there exists a feasible point x such that

$$\begin{aligned} f_i(x) &< 0, i = 1, 2, \dots, m \text{ and} \\ h_i(x) &= 0, i = 1, 2, \dots, p. \end{aligned}$$

Strong Duality Theorem

- Let p^* is the optimal value of the CPP (1) satisfying the Slater's condition and d^* is the solution of the dual problem

$$\begin{aligned} \max (g(\lambda, v) &= \inf_{x \in R^n} L(x, \lambda, v)), \\ \text{subject to, } \lambda &\geq 0 \end{aligned}$$

then $p^* = d^*$.

- Furthermore, if p^* is attained, i.e. there exists a solution x^* to the primal CPP(1), then d^* is also attained, i.e. there exists a global solution (λ^*, v^*) to the dual problem such that $p^* = f_0(x^*) = g(\lambda^*, v^*) = d^* < \infty$.

KKT conditions

For the given CPP (1), the Point x^* is said to satisfy the Karush-Kuhn-Tucker(KKT) conditions, if there exist the multipliers $\lambda^* = (\lambda_1^*, \lambda_2^* \dots, \lambda_m^*)$ and $v = (v_1^*, v_2^*, \dots, v_p^*)$ corresponding to constraints

$$\begin{aligned} f_i(x) &\leq 0, i = 1, 2, \dots, m, \\ h_i(x) &= 0, i = 1, 2, \dots, p, \end{aligned}$$

respectively, such that the Lagrangian function

$$L(x, \lambda, v) = f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^p v_i h_i(x)$$

satisfies

$$\max_{\lambda, v} \left(\inf_x L(x, \lambda, v) \right)$$

min $f(x)$
 Subject to
 $f_i(x) \leq 0, i=1, 2, \dots, m$
 $h_j(x) = 0, j=1, 2, \dots, p$

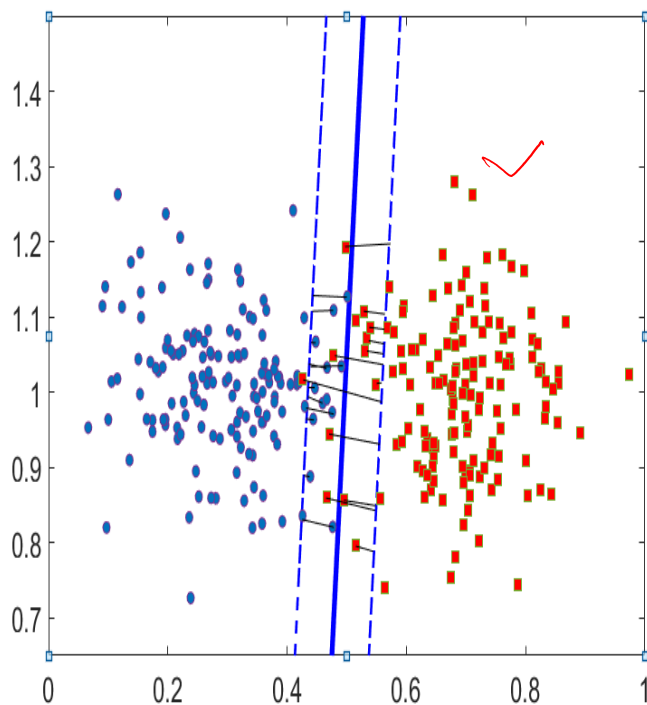
KKT conditions

$$\left\{ \begin{array}{l} f_i(x^*) \leq 0, i = 1, 2, \dots, m, \\ h_i(x^*) = 0, i = 1, 2, \dots, p, \\ \lambda_i^* \geq 0, i = 1, 2, \dots, m, \\ \lambda_i^* f_i(x^*) = 0 \\ \nabla_x L(x, \lambda, \nu) = \nabla f_0(x) + \sum_{i=1}^m \lambda_i \nabla f_i(x) + \sum_{i=1}^p \nu_i \nabla h_i(x) = 0 \end{array} \right.$$

- If x^* is the solution of the CPP(1) satisfying the Slater condition, then it must satisfy the KKT condition.
- If x^* satisfies the KKT condition, then it is the solution of the CPP(1).

Back to SVM Optimization Problem

SVM Optimization problem



SVM primal problem

$$\begin{aligned} \min_{(w,b,\xi)} & \frac{1}{2} w^T w + C \sum_{i=1}^l \xi_i \\ \text{subject to,} & y_i (w^T x_i + b) \geq 1 - \xi_i, \\ & \xi_i \geq 0, \quad i = 1, 2, \dots, l. \end{aligned}$$

where $C > 0$ is the user defined parameter

- Convex Programming Problem, hence global optimal solution.
- Quadratic Programming Problem, hence can be solved with QPP solver.
- But, need to handle too many constraints, which makes the solution computationally expensive.

$$d_i (1 - y_i (w^T x_i + b) - \xi_i) \leq 0$$

$$-d_i (y_i (w^T x_i + b) - 1 + \xi_i)$$

Lagrangian Function for SVM problem

$$\begin{aligned} \xi_i &\geq 0 \\ -\xi_i &\leq 0 \\ \beta_i \xi_i \end{aligned}$$

The Lagrangian function is given as

$$L(w, b, \alpha, \beta, \xi) = \frac{1}{2} w^T w + C \sum_{i=1}^l \xi_i - \sum_{i=1}^l \alpha_i (y_i (w^T x_i + b) - 1 + \xi_i) - \sum_{i=1}^l \beta_i \xi_i$$

where, $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_l)$ and $\beta = (\beta_1, \beta_2, \dots, \beta_l)$ are the positive Lagrangian multipliers.

Dual Problem:-

$$\begin{aligned} &\max_{\alpha \geq 0, \beta \geq 0} (\inf_{(w, b, \xi)} L(w, b, \alpha, \beta, \xi)) \\ &\frac{1}{2} \left(\sum_{i=1}^l \alpha_i y_i x_i \right)^T \left(\sum_{j=1}^l (\alpha_j y_j x_j) \right) - \sum_{i=1}^l \alpha_i y_i \left(\sum_{j=1}^l \alpha_j y_j x_j \right)^T x_i \\ &\quad + \sum_{i=1}^l \alpha_i \\ &\text{Subject to, } C - \alpha_i - \beta_i = 0 \quad \sum_{i=1}^l \alpha_i y_i = 0 \end{aligned}$$

max
 $\alpha \geq 0$
 $\beta \geq 0$

KKT conditions for optimality

$$\alpha_i = C \quad \beta_i = 0$$

$$y_i(w^T x_i + b) - 1 + \xi_i = 0$$

$$y_i(w^T x_i + b) - 1 = -\xi_i$$

$$1 - y_i(w^T x_i + b) = \xi_i$$

$$1 - y_i(w^T x_i + b) > 0$$

$$y_i(w^T x_i + b) < 1$$

$$\nabla_w L(w, b, \alpha, \beta, \xi) = 0 \Rightarrow w = \sum_{i=1}^l \alpha_i y_i x_i \quad (1)$$

$$\nabla_b L(w, b, \alpha, \beta, \xi) = 0 \Rightarrow \sum_{i=1}^l \alpha_i y_i = 0 \quad (2)$$

$$\nabla_\xi L(w, b, \alpha, \beta, \xi) = 0 \Rightarrow C - \alpha_i - \beta_i = 0 \quad (3)$$

$$y_i (w^T x_i + b) \geq 1 - \xi_i, \quad i = 1, 2, \dots, l. \quad (4)$$

$$\xi_i \geq 0, \quad i = 1, 2, \dots, l. \quad (5)$$

$$\alpha_i (y_i (w^T x_i + b) - 1 + \xi_i) = 0, \quad i = 1, 2, \dots, l. \quad (6)$$

$$\beta_i \xi_i = 0, \quad i = 1, 2, \dots, l. \quad (7)$$

$$\alpha_i \geq 0, \quad i = 1, 2, \dots, l. \quad (8)$$

$$\beta_i \geq 0, \quad i = 1, 2, \dots, l. \quad (9)$$

Dual SVM Problem

The Lagrangian function is given as

$$L(w, b, \alpha, \beta, \xi) = \frac{1}{2} w^T w + C \sum_{i=1}^l \xi_i - \sum_{i=1}^l \alpha_i (y_i (w^T x_i + b) - 1 + \xi_i) - \sum_{i=1}^l \beta_i \xi_i,$$

where, $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_l)$ and $\beta = (\beta_1, \beta_2, \dots, \beta_l)$ are the positive Lagrangian multipliers.

Dual Problem:-

$$\begin{aligned} & \max_{\alpha \geq 0, \beta \geq 0} \left(\inf_{(w, b, \xi)} L(w, b, \alpha, \beta, \xi) \right) \\ & \frac{1}{2} \left(\sum_{i=1}^l \alpha_i y_i x_i \right)^T \sum_{j=1}^l \left(\alpha_j y_j x_j \right) + \sum_{i=1}^l \alpha_i y_i - \sum_{i=1}^l \alpha_i y_i \left(\sum_{j=1}^l \alpha_j y_j x_j^T x_i \right) \\ & \text{Subject to,} \quad \sum_{i=1}^l \alpha_i y_i = 0, \\ & \quad \quad \quad C - \alpha_i - \beta_i = 0, \quad i = 1, 2, \dots, l \end{aligned}$$

$$\max_{\alpha, \beta} \quad -\frac{1}{2} \sum_{i=1}^2 \sum_{j=1}^2 \alpha_i \alpha_j y_i y_j (x_i^T x_j) + \sum_{i=1}^2 \alpha_i y_i$$

Subject to, $\sum_{i=1}^2 \alpha_i y_i = 0$

$$C - \alpha_i - \beta_i = 0 \Rightarrow 0 \leq \alpha_i \leq C$$

$$\alpha_i, \beta_i \geq 0$$

Dual SVM Problem

Dual Problem:- $\max_{\alpha \geq 0, \beta \geq 0} (\inf_{(w, b, \xi)} L(w, b, \alpha, \beta, \xi))$

where, $L(w, b, \alpha, \beta, \xi) = \frac{1}{2} w^T w + C \sum_{i=1}^l \xi_i - \sum_{i=1}^l \alpha_i (y_i (w^T x_i + b) - 1 + \xi_i) - \sum_{i=1}^l \beta_i \xi_i$,

$$\max_{\alpha \geq 0, \beta \geq 0} \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j x_i^T x_j - \sum_{i=1}^l \alpha_i (y_i (\sum_{j=1}^l \alpha_j y_j x_j)^T x_i) - 1)$$

Subject to, $\sum_{i=1}^l \alpha_i y_i = 0$,

$$C - \alpha_i - \beta_i = 0,$$

$$\alpha_i \geq 0, \beta_i \geq 0, i = 1, 2, \dots, l$$

Dual SVM Problem

$$k(x_i, x_j) =$$

$$\leftarrow \frac{\phi(x_i)^T \phi(x_j)}{}$$

$$\max_{\alpha} \frac{-1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j x_i^T x_j + \sum_{i=1}^l \alpha_i$$

Subject to, $\sum_{i=1}^l \alpha_i y_i = 0$,
 $0 \leq \alpha_i \leq C, i = 1, 2, \dots, l.$

$$\min_{\alpha} \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j x_i^T x_j - \sum_{i=1}^l \alpha_i$$

Subject to, $\sum_{i=1}^l \alpha_i y_i = 0$,
 $0 \leq \alpha_i \leq C, i = 1, 2, \dots, l.$

Dual SVM Problem

$$\min_{\alpha} \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j x_i^T x_j - \sum_{i=1}^l \alpha_i$$

Subject to, $\sum_{i=1}^l \alpha_i y_i = 0,$
 $0 \leq \alpha_i \leq C, i = 1, 2, \dots, l.$

$$\min_{(w,b)} \frac{1}{2} w^T w + C \sum_{i=1}^l \xi_i$$

subject to,

$$y_i (w^T x_i + b) \geq 1 - \xi_i,$$

$$\xi_i \geq 0, \quad i = 1, 2, \dots, l.$$

- The Dual problem is a Quadratic Programming Problem (QPP) but, less number of constraints.
- Suppose that $\alpha^* = (\alpha_1^*, \dots, \alpha_l^*)$ is a solution to the dual problem, Then the training point (x_i, y_i) , is said to be a support vector if the corresponding component α_i^* is non-zero and otherwise it is a non-support vector.

$$\begin{aligned} (x_1, y_1) &\rightarrow \alpha_1^* \\ (x_2, y_2) &\rightarrow \alpha_2^* \\ &\vdots \\ (x_l, y_l) &\rightarrow \alpha_l^* \end{aligned}$$

SVM Solution

$$\min_{(w,b)} \quad \frac{1}{2} w^T w + C \sum_{i=1}^l \xi_i$$

subject to,

$$y_i (w^T x_i + b) \geq 1 - \xi_i, \\ \xi_i \geq 0, \quad i = 1, 2, \dots, l.$$

$$\min_{\alpha} \quad \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j x_i^T x_j - \sum_{i=1}^l \alpha_i$$

$$\text{Subject to, } \sum_{i=1}^l \alpha_i y_i = 0, \\ 0 \leq \alpha_i \leq C, i = 1, 2, \dots, l.$$

- After obtaining the solution of dual problem $\alpha^* = (\alpha_1^*, \dots, \alpha_l^*)$, the w^* can be obtained by using the KKT condition (1) as

$$w^* = \sum_{i=1}^l \alpha_i^* y_i x_i$$

$$\frac{1}{y_0} = y_j$$

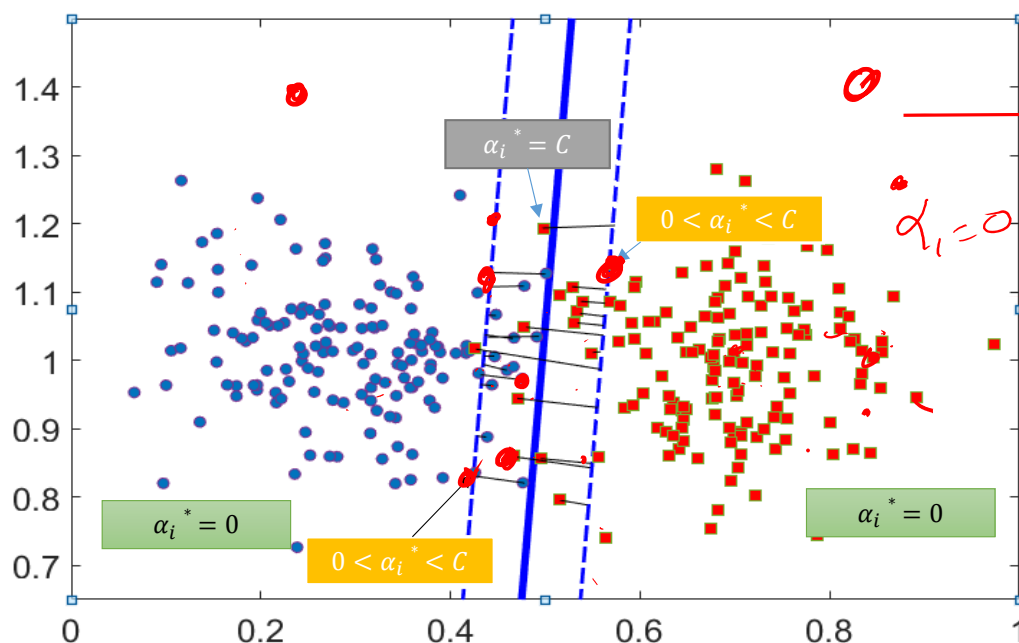
- Further, for a given $0 < \alpha_j^* < C$, we have $y_j (w^{*T} x_j + b^*) = 1$ from KKT condition (6) and (7) which gives,

$$b^* = y_j - \sum_{i=1}^l \alpha_i^* y_i x_i^T x_j$$

$$x_j$$

- In practice, we consider all α_j^* , satisfying $0 < \alpha_j^* < C$ and compute the value of b^* . The final value of b^* is considered as the mean of all computed values.

SVM Solution



$$d_i = C$$

$$y_i(w^T x_i + b) \geq 1 - \xi_i$$

$$w = \sum_{i=1}^L \alpha_i y_i x_i$$

$$d_i = 0$$

$$\rightarrow \xi_i > 0$$

$$0 < d_i < C$$

- For data point (x_i, y_i) , satisfying $y_i(w^{*T} x_i + b^{*}) < 1$, the corresponding $\alpha_i^* = C$.
- For data point (x_i, y_i) , satisfying $y_i(w^{*T} x_i + b^{*}) = 1$, the corresponding $0 < \alpha_i^* < C$.
- For data point (x_i, y_i) , satisfying $y_i(w^{*T} x_i + b^{*}) > 1$, the corresponding $\alpha_i^* = 0$.

SVM is a sparse classification model

Proof using the KKT conditions

$$\nabla_w L(w, b, \alpha, \beta, \xi) = 0 \implies w = \sum_{i=1}^l \alpha_i y_i x_i \quad (1)$$

$$\nabla_b L(w, b, \alpha, \beta, \xi) = 0 \implies \sum_{i=1}^l \alpha_i y_i x_i = 0 \quad (2)$$

$$\nabla_\xi L(w, b, \alpha, \beta, \xi) = 0 \implies C - \alpha_i - \beta_i = 0 \quad (3)$$

$$y_i (w^T x_i + b) \geq 1 - \xi_i, \quad i = 1, 2, \dots, l. \quad (4)$$

$$\xi_i \geq 0, \quad i = 1, 2, \dots, l. \quad (5)$$

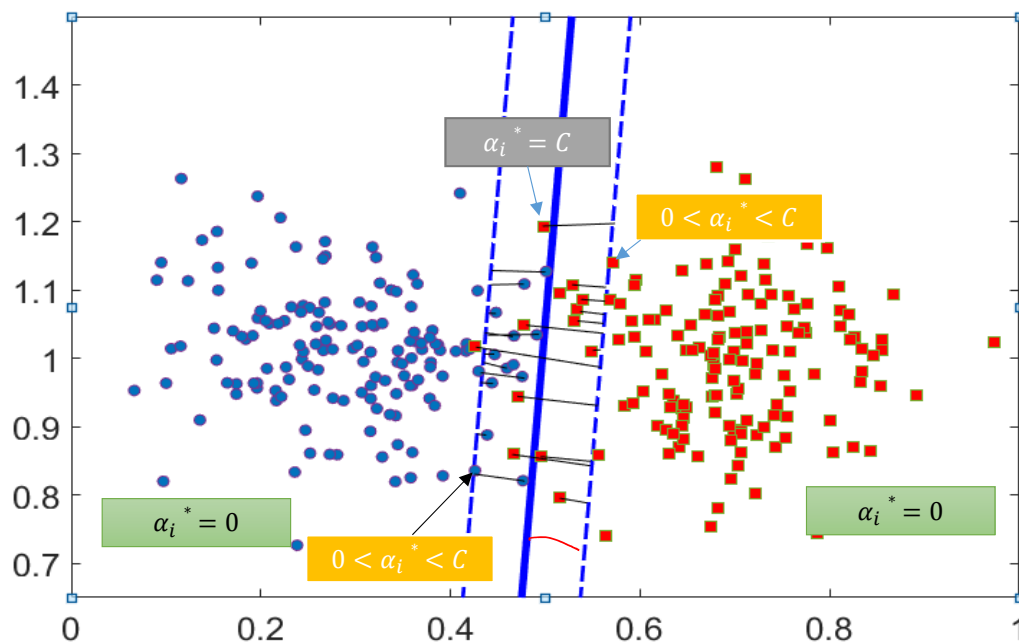
$$\alpha_i (y_i (w^T x_i + b) - 1 + \xi_i) = 0, \quad i = 1, 2, \dots, l. \quad (6)$$

$$\beta_i \xi_i = 0, \quad i = 1, 2, \dots, l. \quad (7)$$

$$\alpha_i \geq 0, \quad i = 1, 2, \dots, l. \quad (8)$$

$$\beta_i \geq 0, \quad i = 1, 2, \dots, l. \quad (9)$$

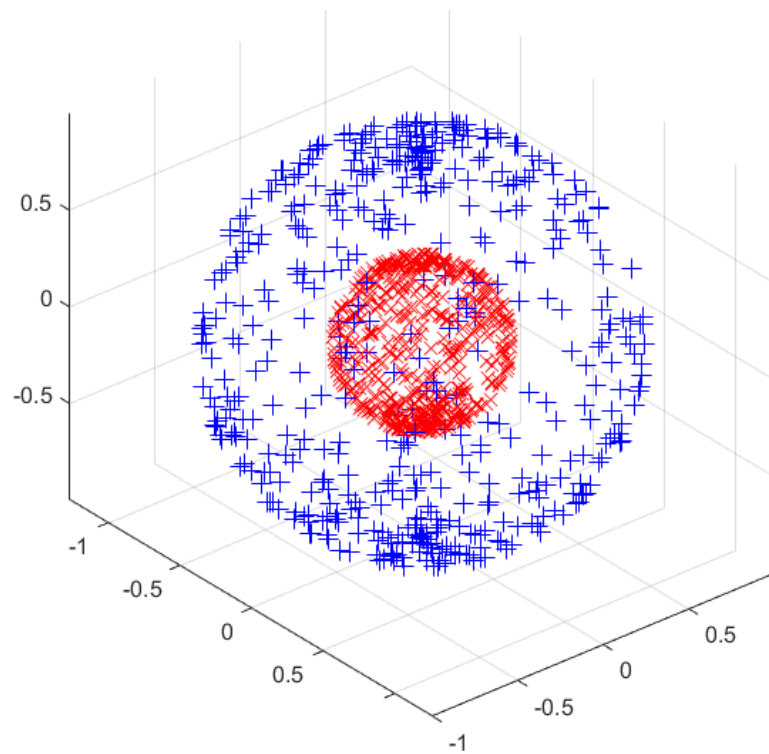
SVM Solution



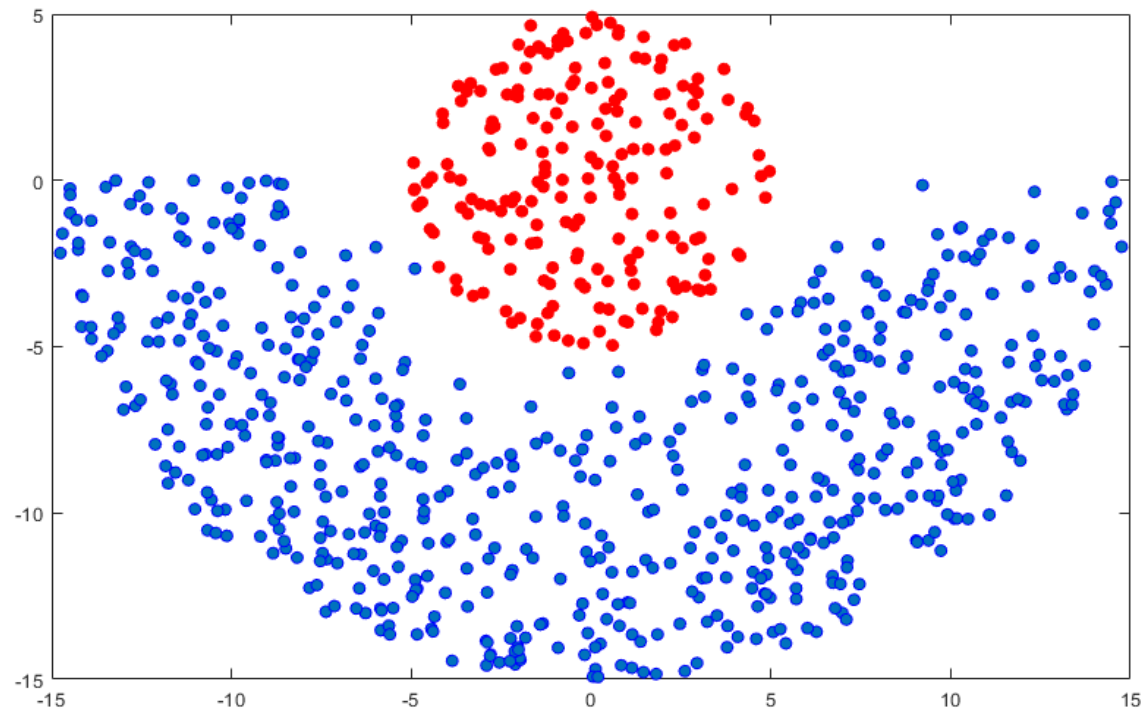
- For data point (x_i, y_i) , satisfying $y_i(w^{*T}x_i + b^*) < 1$, the corresponding $\alpha_i^* = C$.
- For data point (x_i, y_i) , satisfying $y_i(w^{*T}x_i + b^*) = 1$, the corresponding $0 < \alpha_i^* < C$.
- For data point (x_i, y_i) , satisfying $y_i(w^{*T}x_i + b^*) > 1$, the corresponding $\alpha_i^* = 0$.

For a test point say \tilde{x} , the SVM prediction is $\text{sign}(w^{*T}\tilde{x} + b^*)$

Non-linear separable data

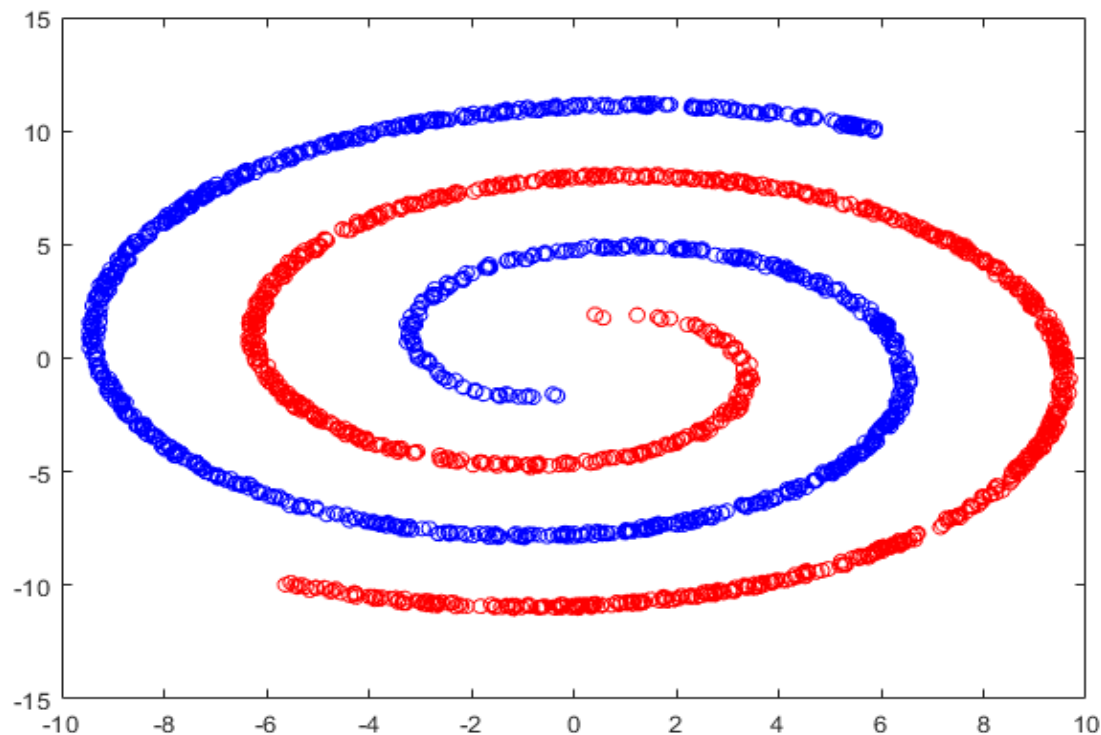


Non-linear separable data



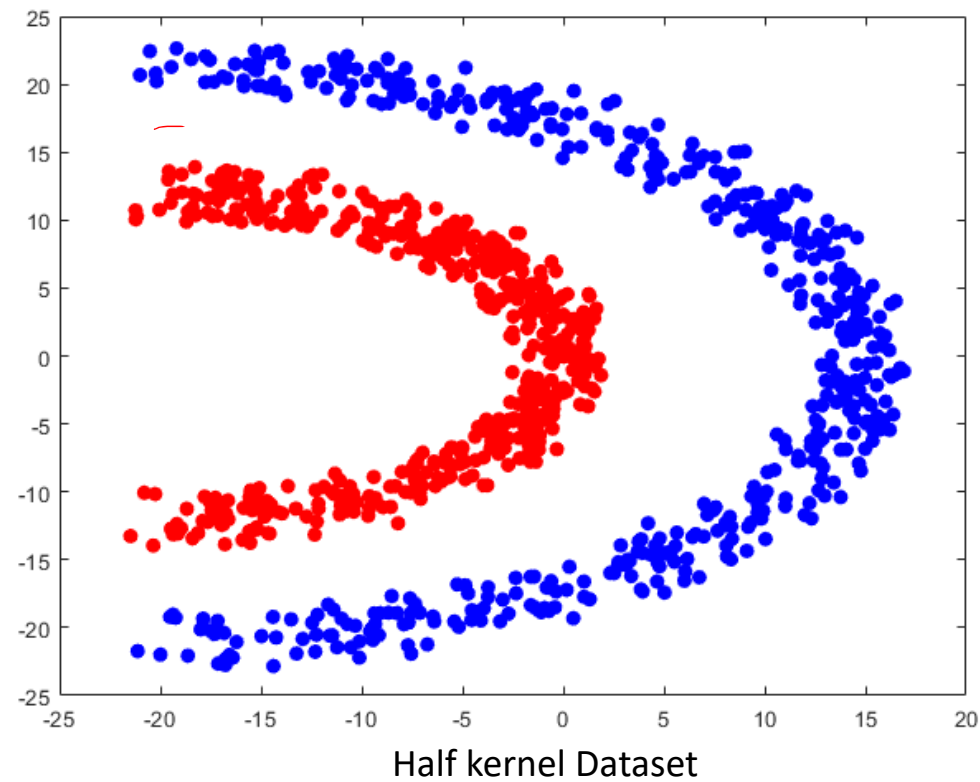
Crescent full moon dataset

Non-linear separable data

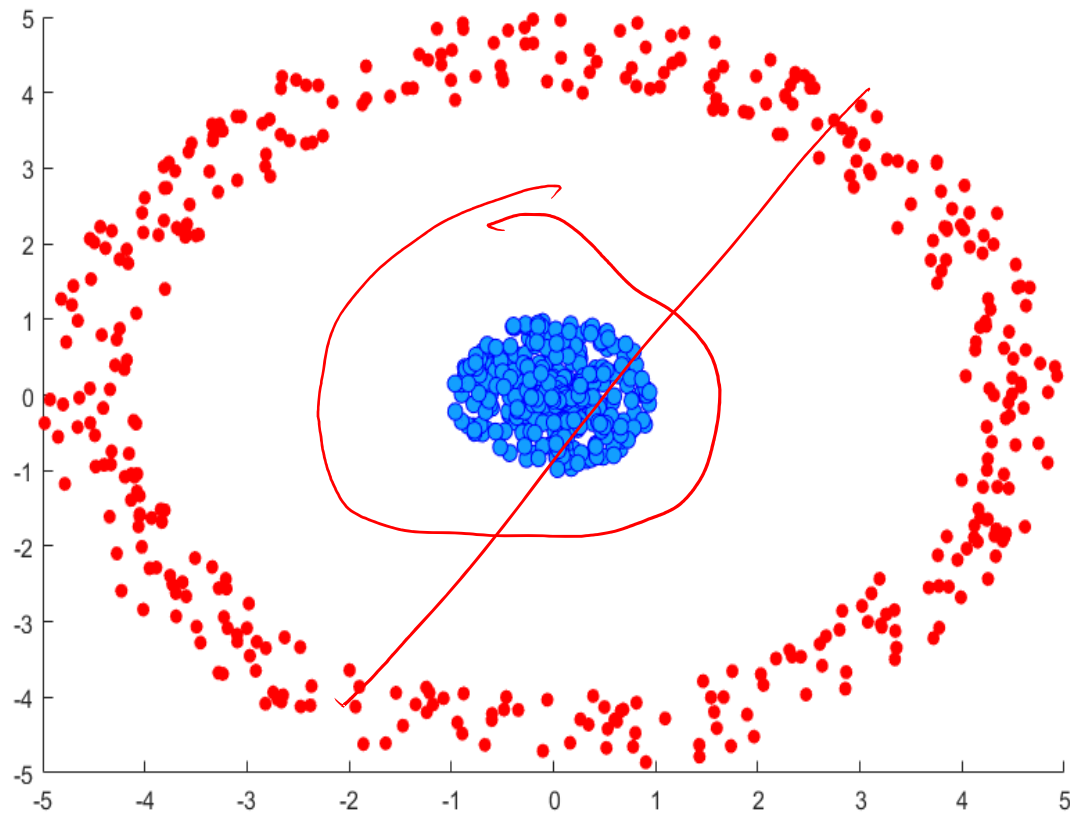


Spiral Dataset

Non-linear separable data



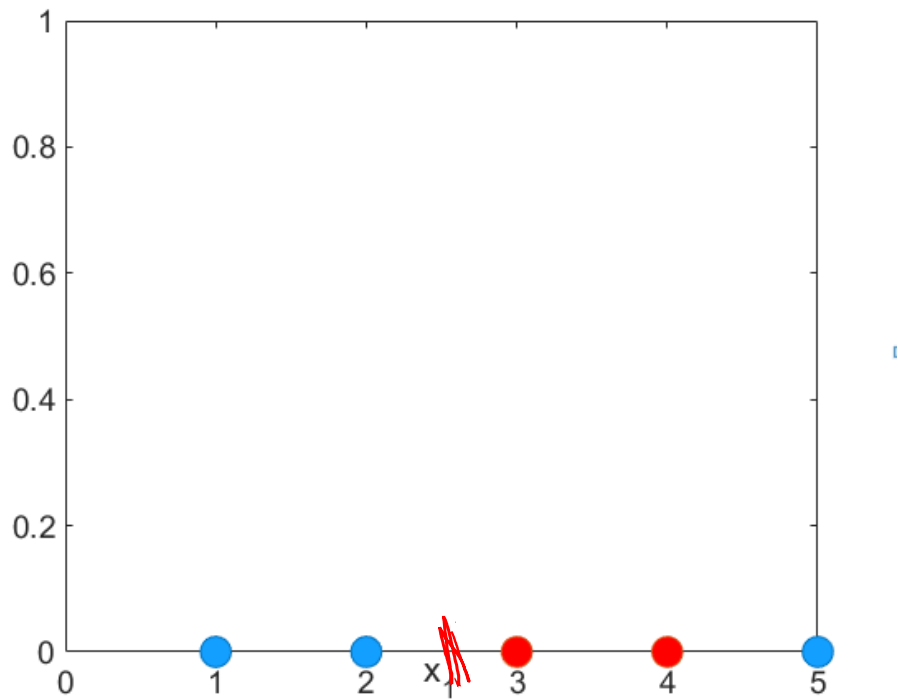
Non-linear separable data



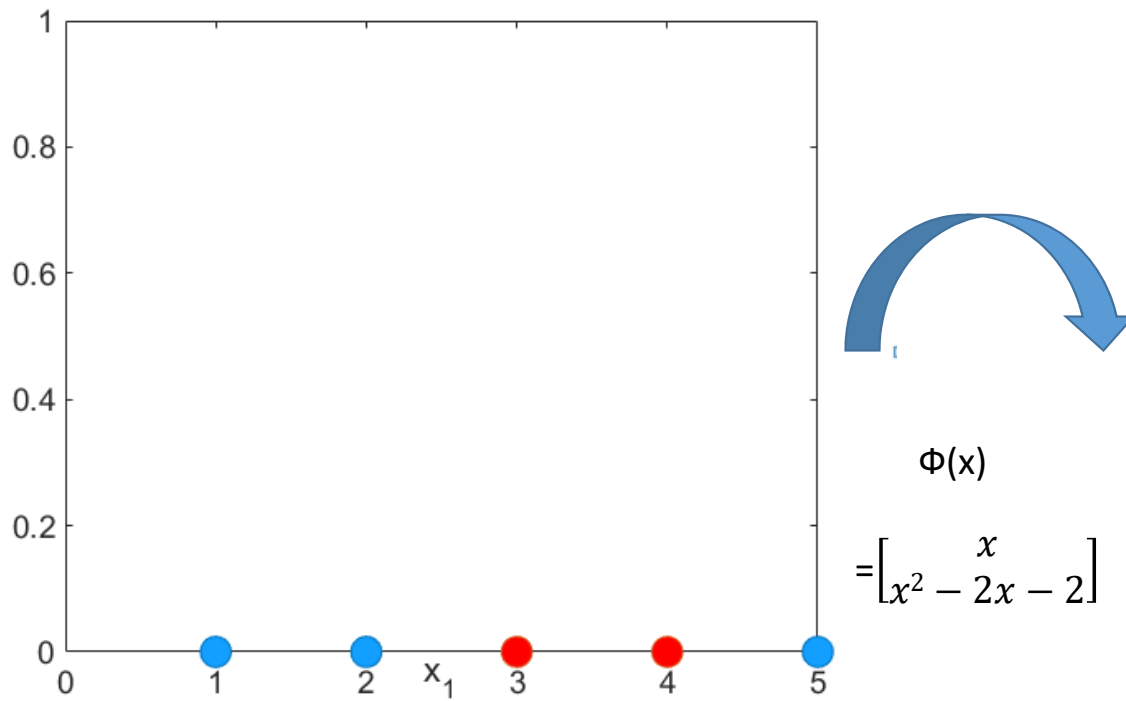
Co-centric Circle Dataset

Dr. Pritam Anand, DA-IICT, Gandhinagar

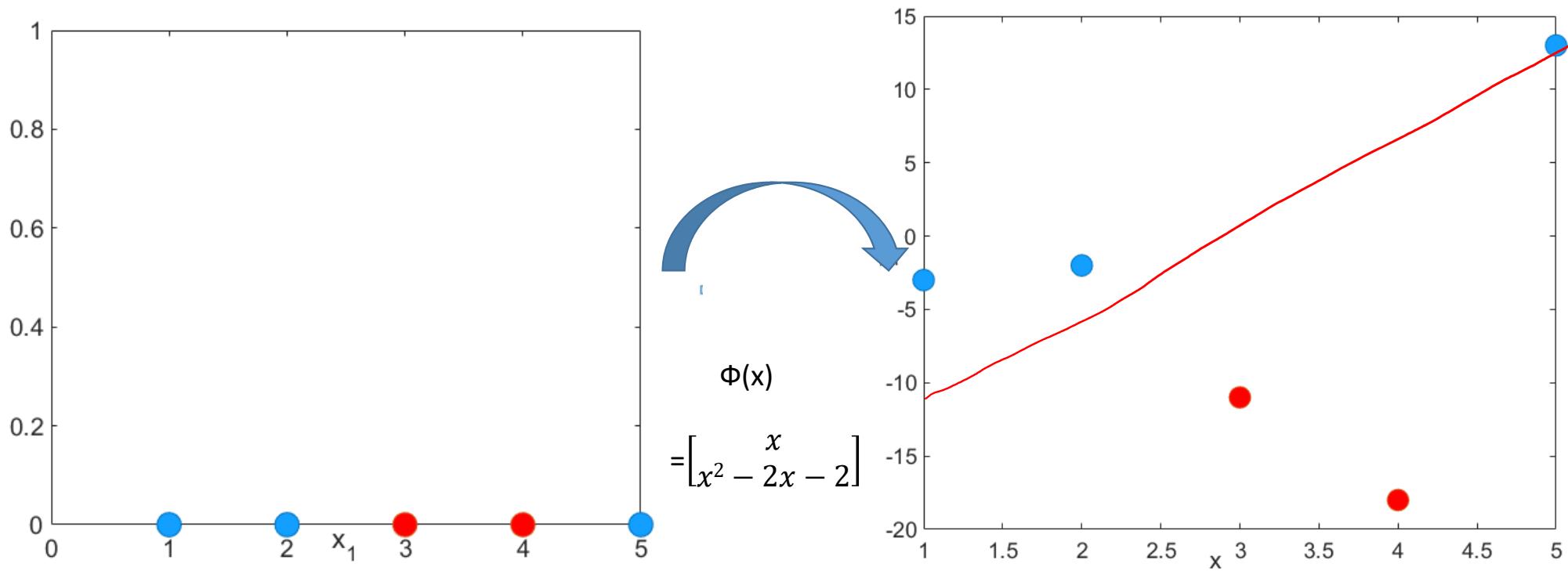
Kernel Trick



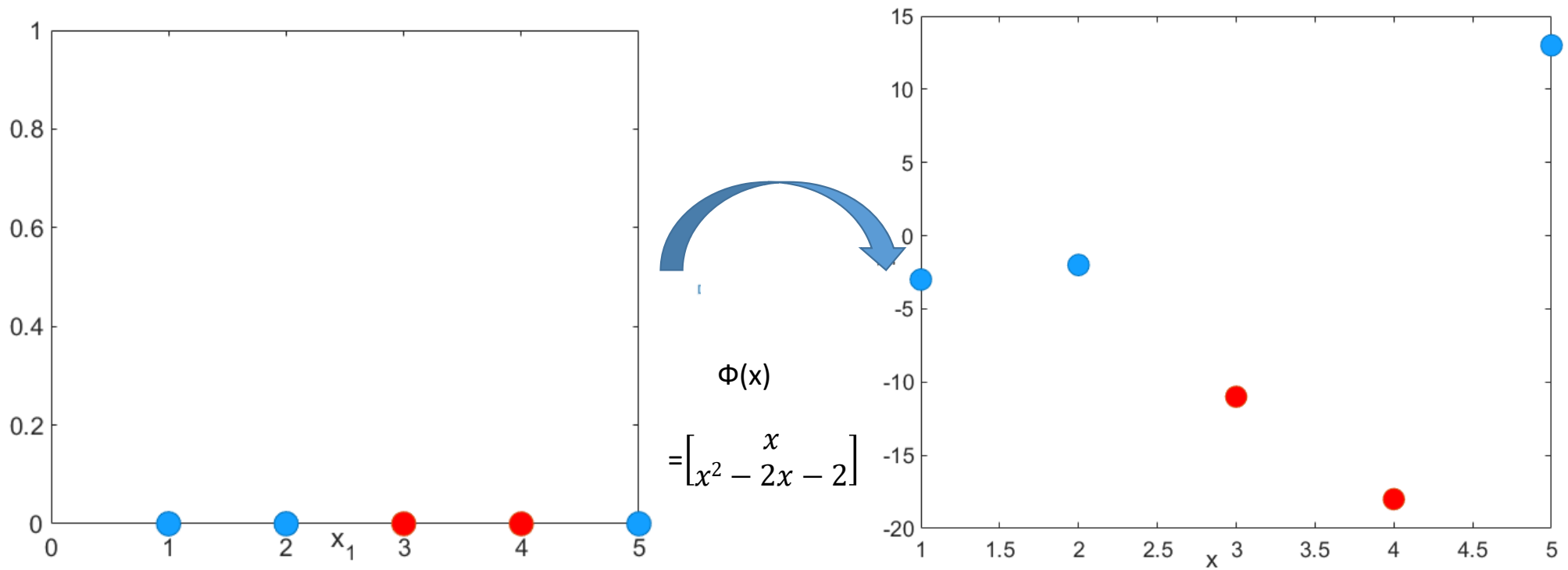
Kernel Trick



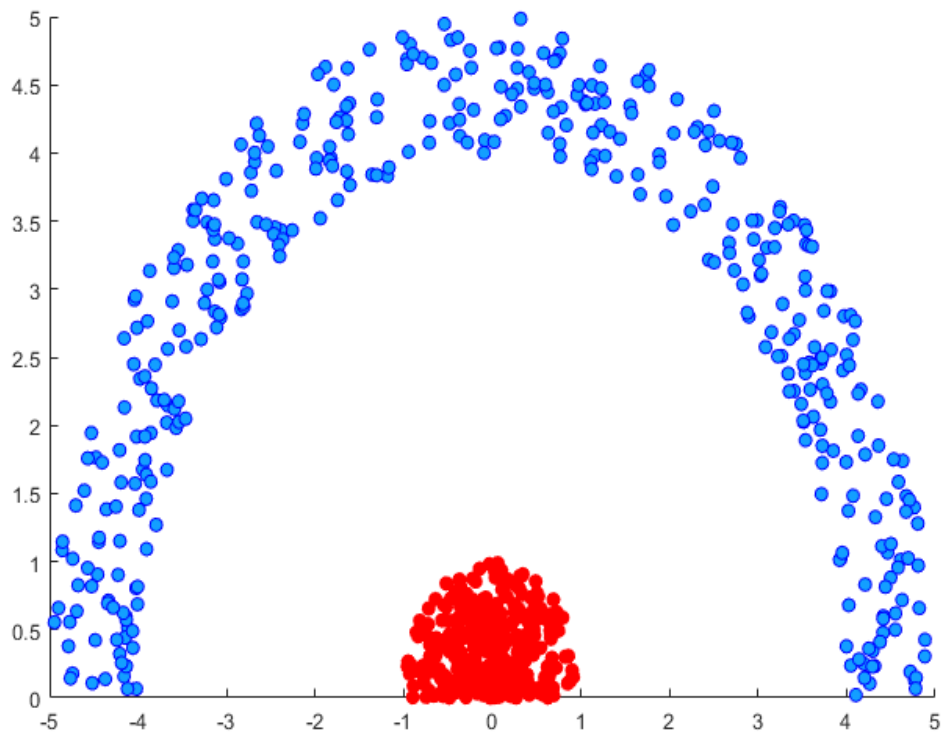
Kernel Trick



Kernel Trick

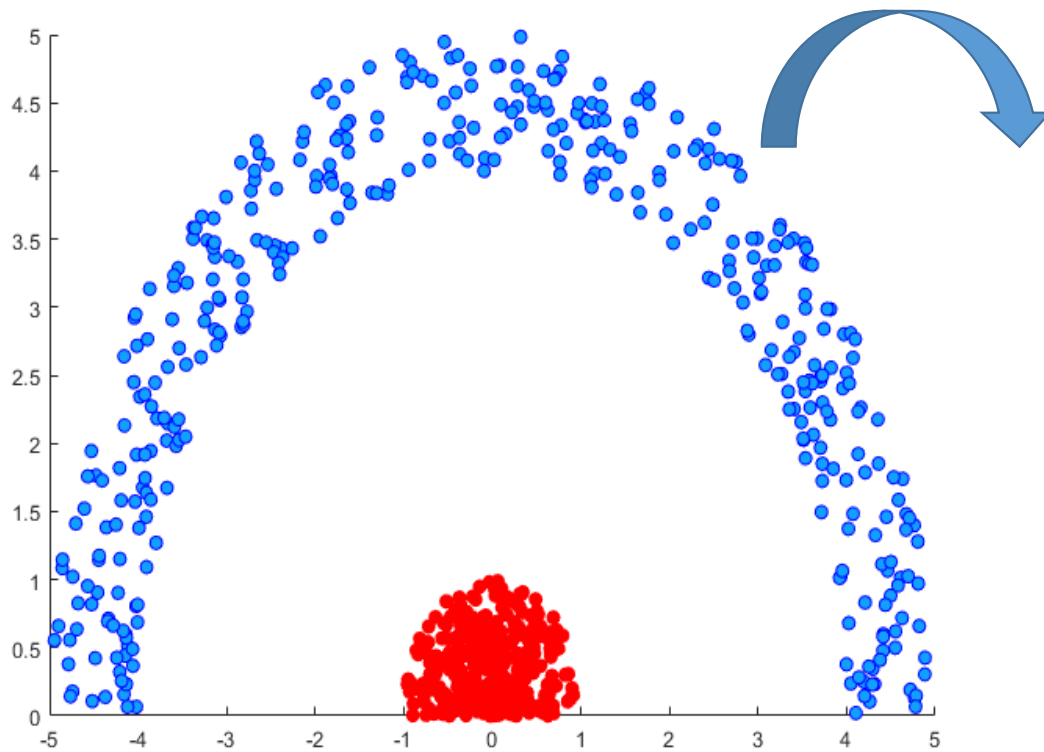


Kernel Trick



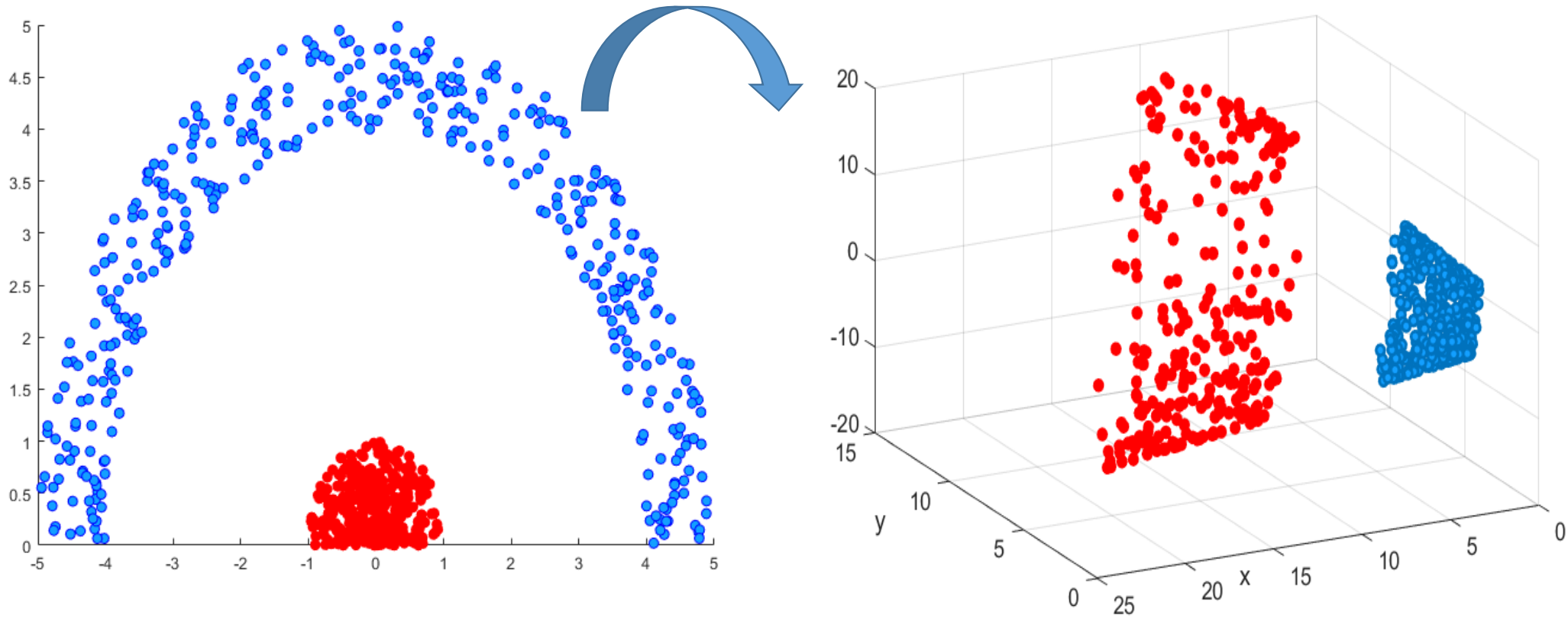
Kernel Trick

$$\Phi(x) = \Phi \left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right) = \begin{bmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2}x_2x_1 \end{bmatrix}$$



Kernel Trick

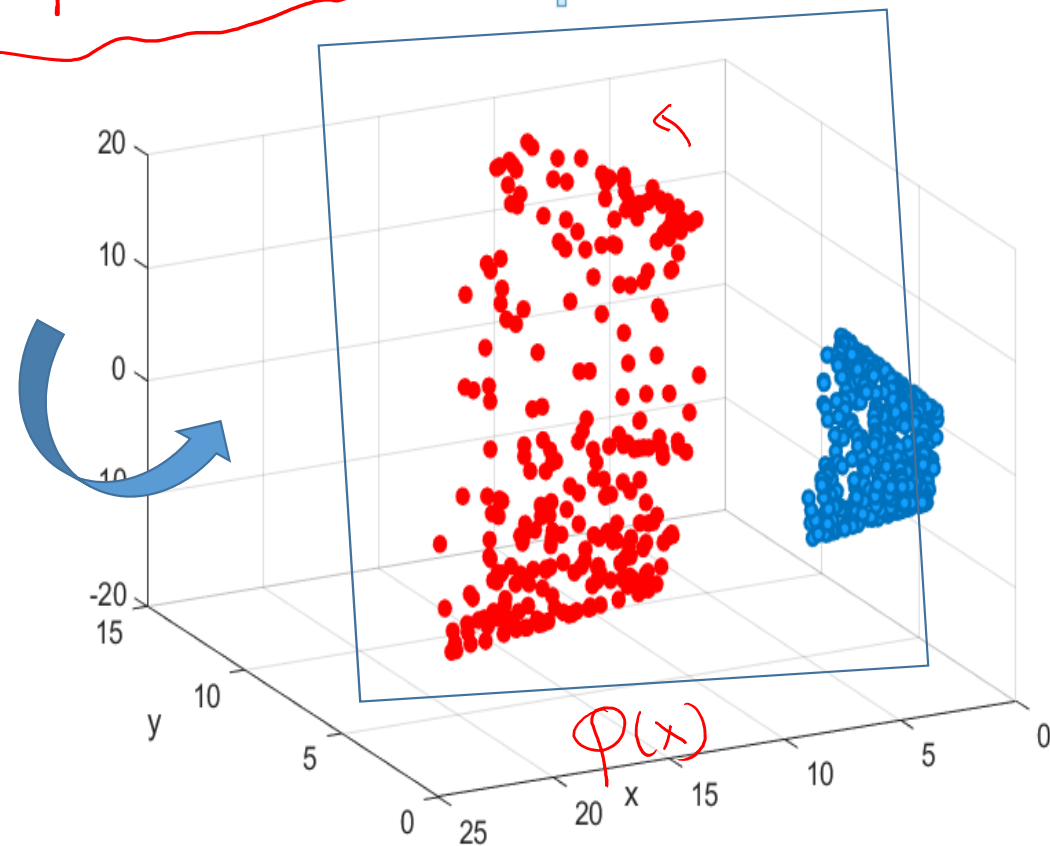
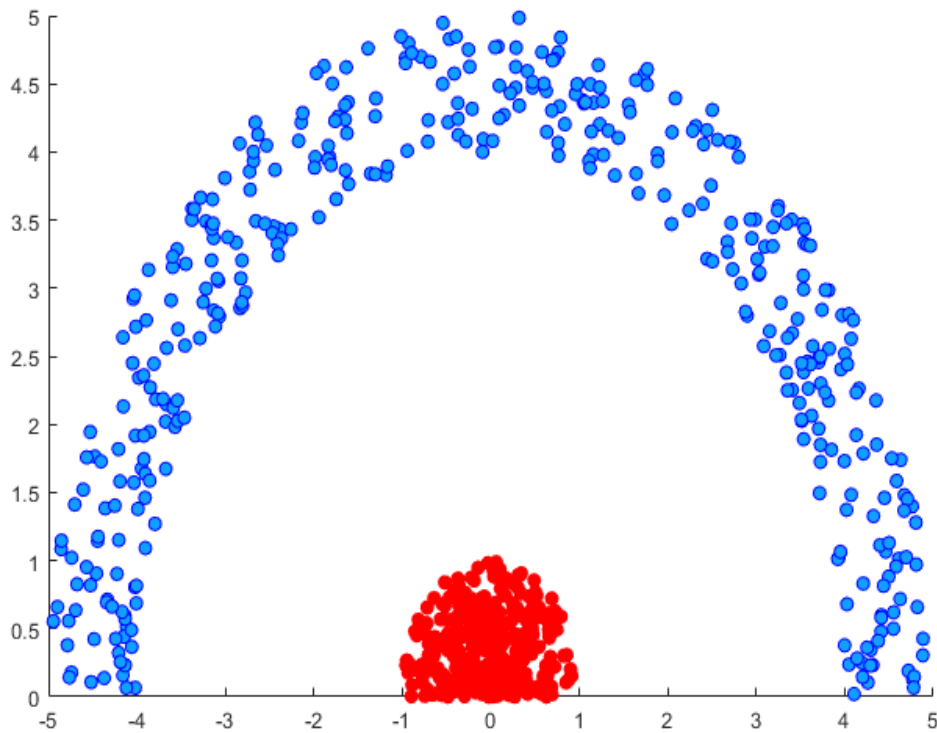
$$\Phi(x) = \Phi \left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right) = \begin{bmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2}x_2x_1 \end{bmatrix}$$



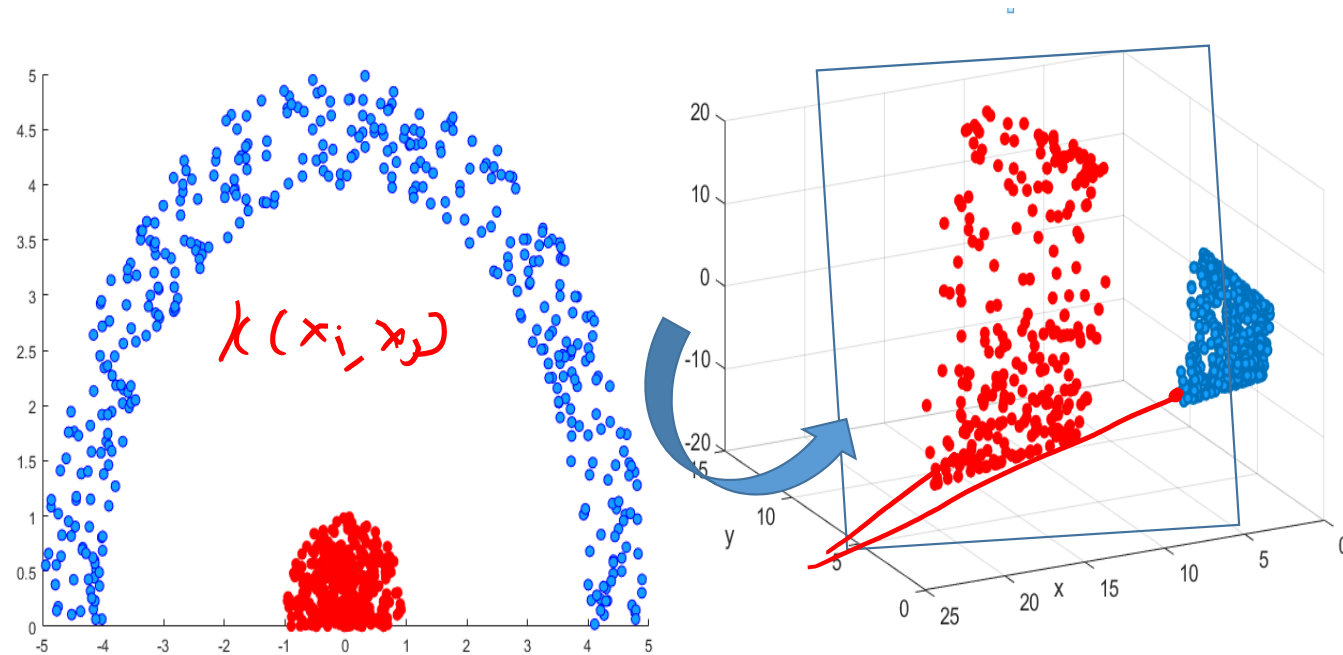
Kernel Trick

$$\Phi(x) = \Phi \left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right) = \begin{bmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2}x_2x_1 \end{bmatrix}$$

Handwritten red text: $w^T \Phi(x) + b = 0$



Non-linear SVM



Primal Problem

$$\min_{(w,b)} \frac{1}{2} w^T w + C \sum_{i=1}^l \xi_i$$

subject to,

$$y_i (w^T \phi(x_i) + b) \geq 1 - \xi_i,$$

$$\xi_i \geq 0, \quad i = 1, 2, \dots, l.$$

19-11-2024

Dual Problem

$$\min_{\alpha} \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j \phi(x_i)^T \phi(x_j) - \sum_{i=1}^l \alpha_i$$

Handwritten red annotations: $\phi(x_i)^T \phi(x_j)$ is circled, and $\alpha_j y_i y_j$ is underlined.

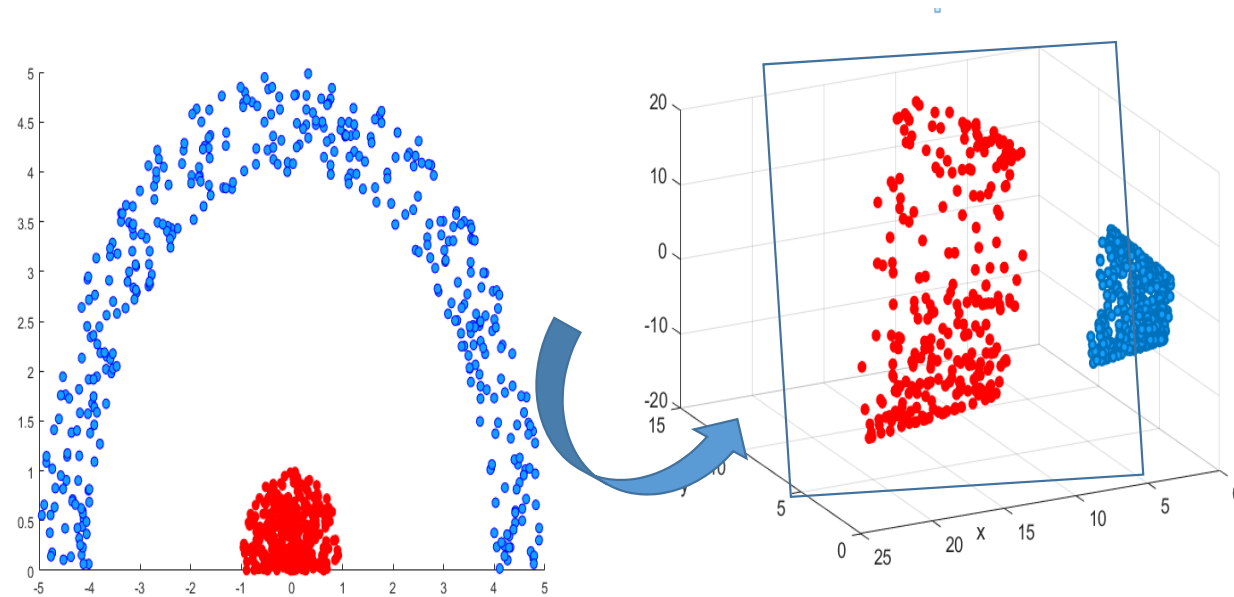
Subject to, $\sum_{i=1}^l \alpha_i y_i = 0,$

$$0 \leq \alpha_i \leq C, i = 1, 2, \dots, l.$$

Dr .Pritam Anand, DA-IICT, Gandhinagar

69

Non-linear SVM



- The dual problem requires the knowledge of only $\phi(x_i)^T \phi(x_j)$.
- We can use a kernel function such that $k(x_i, x_j) = \phi(x_i)^T \phi(x_j)$.
- Dual Problem can be solved without the explicit knowledge of mapping ϕ .

Dual Problem

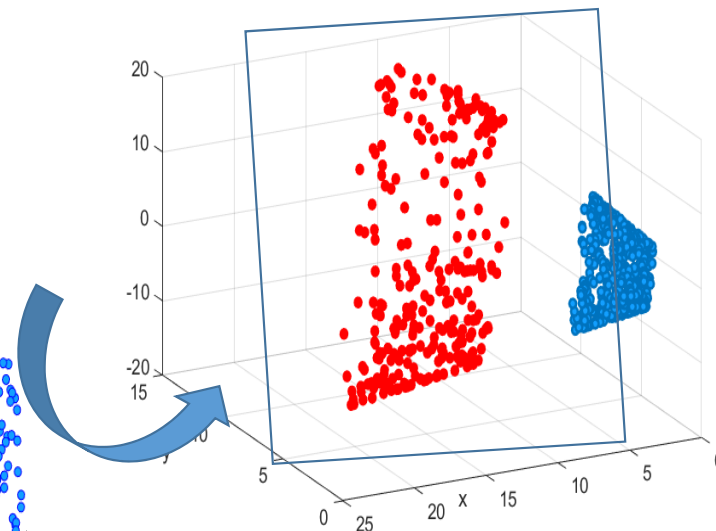
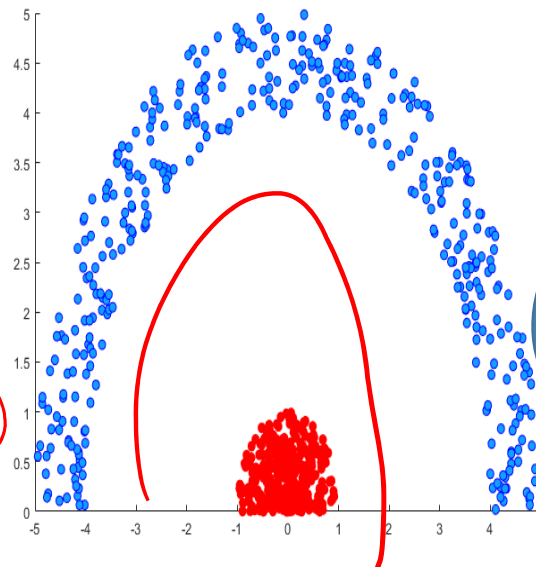
$$\min_{\alpha} \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j \overbrace{\phi(x_i)^T \phi(x_j)}^{K(x_i, x_j)} - \sum_{i=1}^l \alpha_i$$

$$\text{Subject to, } \sum_{i=1}^l \alpha_i y_i = 0, \\ 0 \leq \alpha_i \leq C, i = 1, 2, \dots, l.$$

Non-linear SVM solution

$$w = \left(\sum_{i=1}^l \alpha_i y_i \phi(x_i) \right)^T \phi(x)$$

$$\sum_{i=1}^l y_i \alpha_i \underbrace{\phi(x_i)^T \phi(x)}_{k(x_i, x)}$$



Dual Problem

$$\min_{\alpha} \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j k(x_i, x_j) - \sum_{i=1}^l \alpha_i$$

Subject to, $\sum_{i=1}^l \alpha_i y_i = 0,$
 $0 \leq \alpha_i \leq C, i = 1, 2, \dots, l.$

After obtaining the optimal solution of the dual problem $\alpha^* = (\alpha_1^*, \dots, \alpha_l^*)$, the decision function can be obtained as

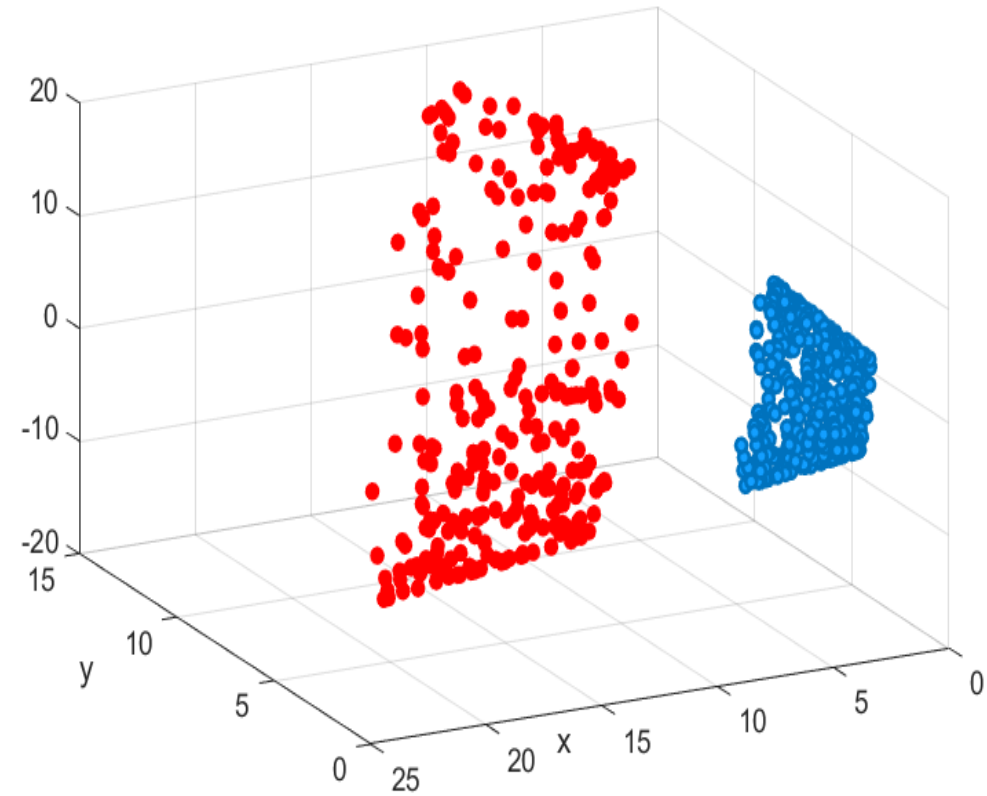
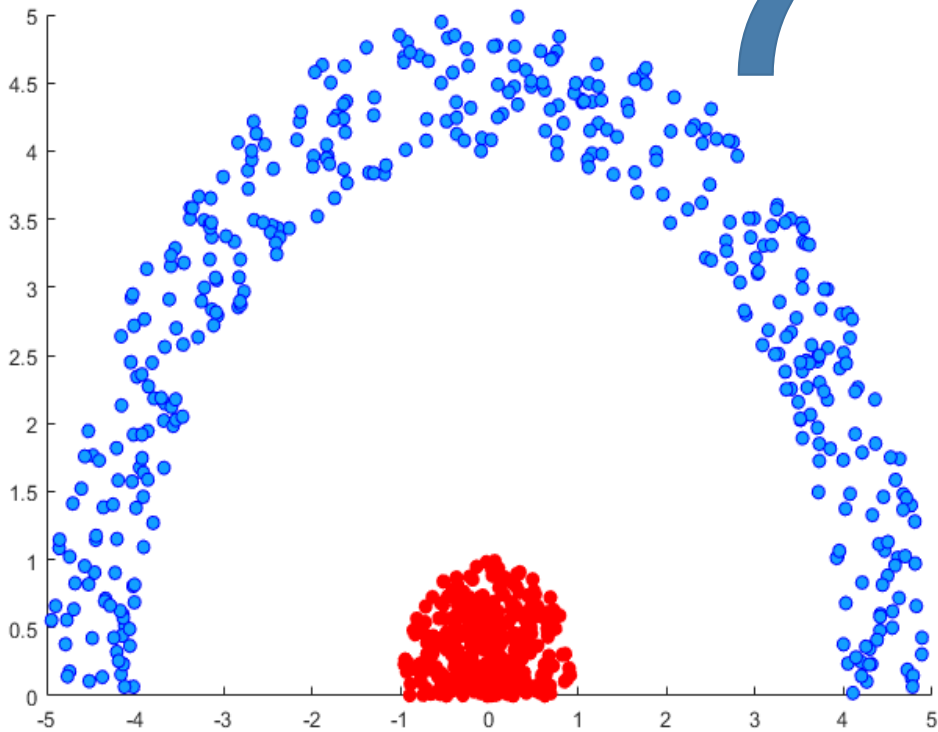
$$f(x) = \text{sign}(w^{*T} \phi(x) + b^*)$$

$$= \text{sign}(\sum_{i=1}^l \alpha_i^* y_i \phi(x_i)^T \phi(x) + b^*)$$

$$= \text{sign}(\sum_{i=1}^l \alpha_i^* y_i k(x_i, x) + b^*)$$

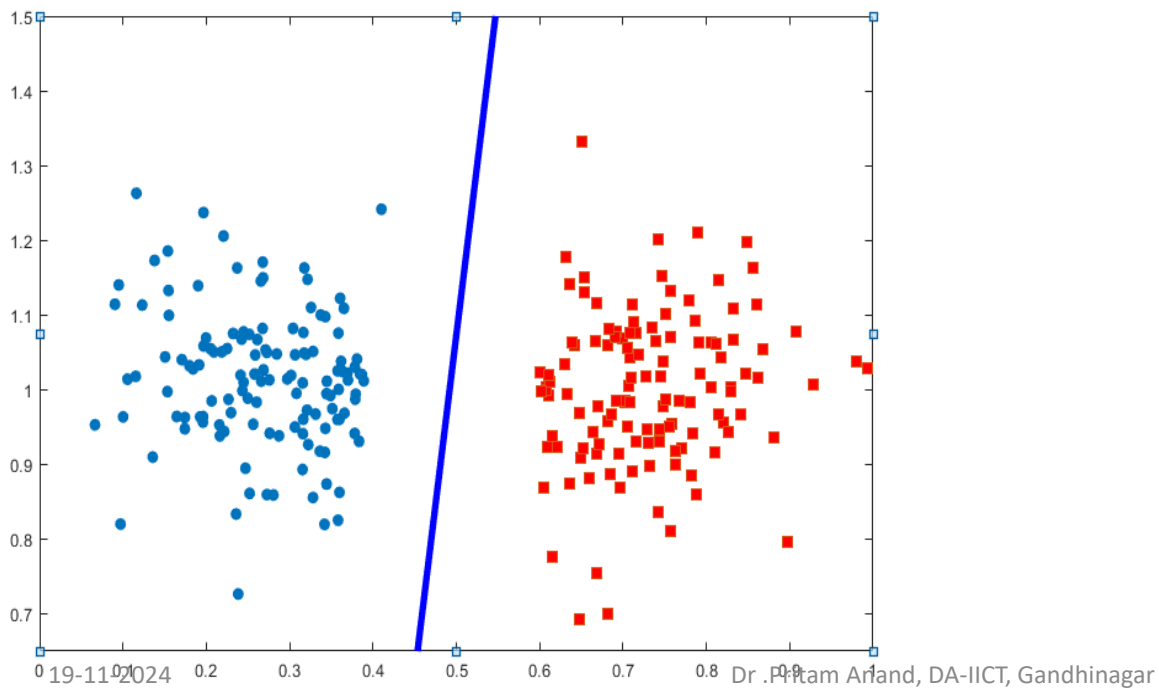
Kernel Trick

$$k(x,y) = \Phi(x)^T \Phi(y) = \underbrace{(x^T y)}_{\text{dot product}}^2$$



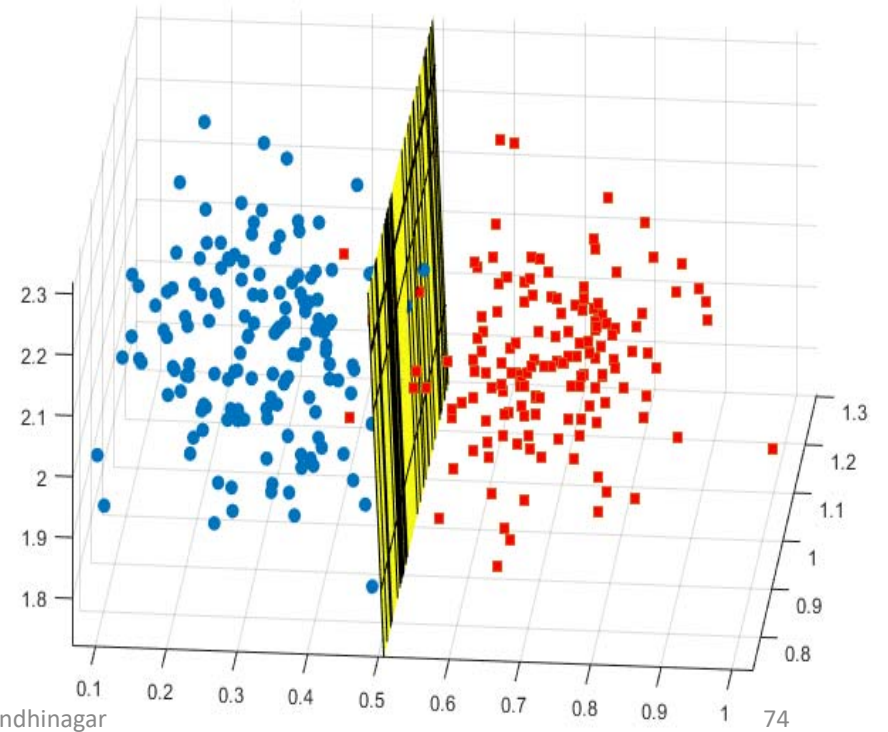
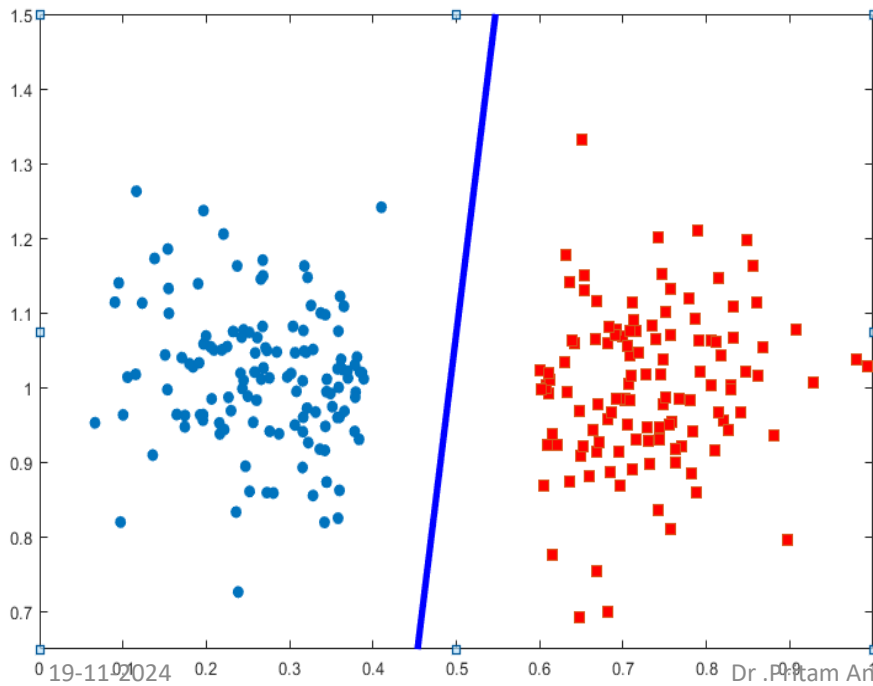
Kernel Types.

- Linear Kernel: $k(x,y) = \Phi(x)^T \Phi(y) = (x^T y)$



Kernel Types.

- Linear Kernel: $k(x,y) = \Phi(x)^T \Phi(y) = (x^T y)$

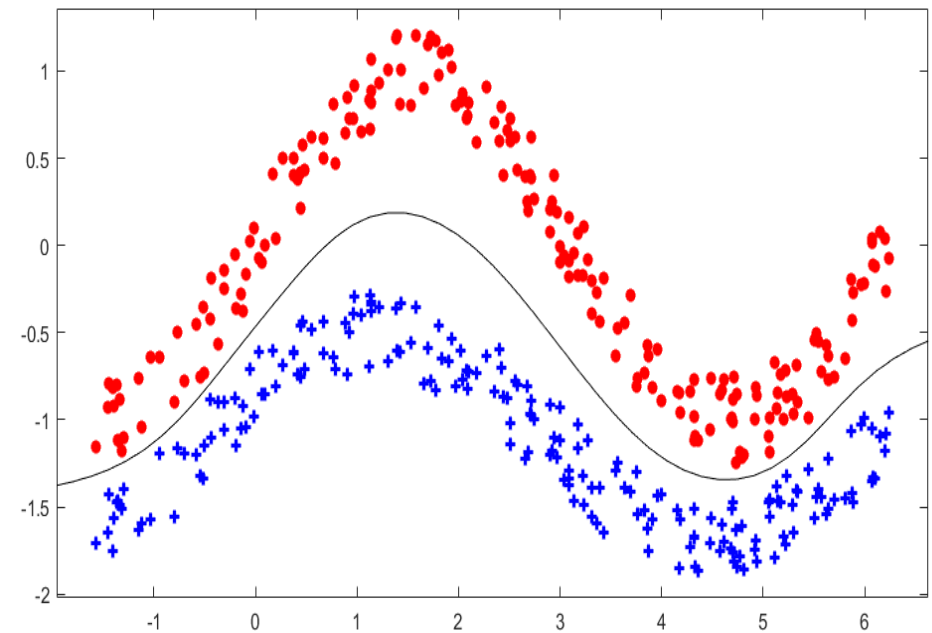


Kernel Types

- Quadratic Kernel:

$$k(x,y) = \Phi(x)^T \Phi(y) = (x^T y + c)^2$$

, where c is the user-defined parameter.

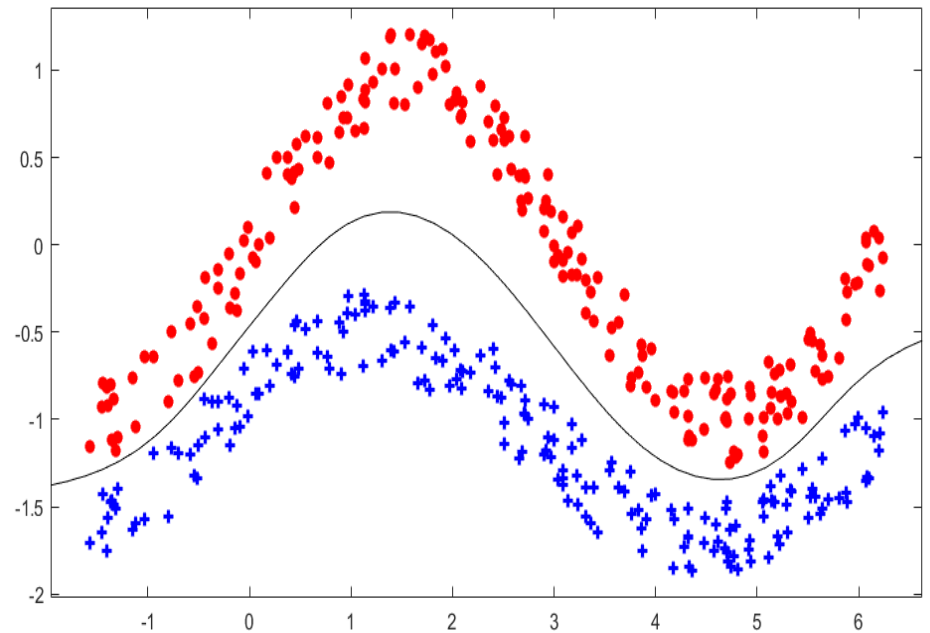


Kernel Types

- Polynomial Kernel:

$$k(x,y) = \Phi(x)^T \Phi(y) = (x^T y + c)^p$$

, where c is the user-defined parameter.



- It can generate any type of polynomial surfaces.

Kernel Types.

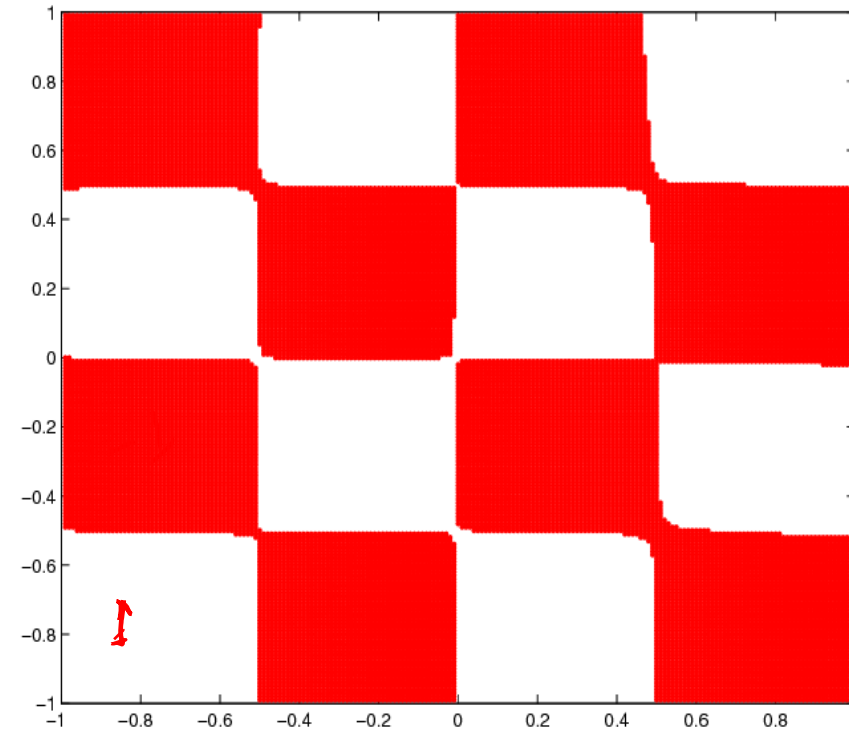
- RBF kernel:

$$k(x,y) = \Phi(x)^T \Phi(y)$$

$$= e^{-q||x-y||^2},$$

where q is the user-defined parameter.

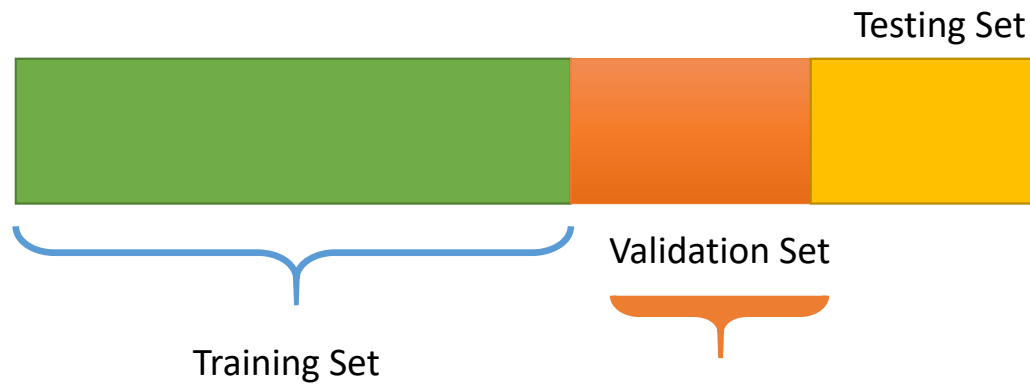
- It can generate any type of continuous surfaces.



TOY SVM

<https://greitemann.dev/svm-demo>

Model selection in SVM



Model selection in SVM

2^{-5}	2^{-4}	2^{-3}	2^{-2}	2^{-1}	2^0	2^1	2^2	2^3	2^4	2^5	c
2^{-5}	2^{-4}	2^{-3}	2^{-2}	2^{-1}	2^0	2^1	2^2	2^3	2^4	2^5	q

Implementation

- A number of libraries in MATLAB and Python is available.
- The *fitcsvm* in MATLAB provides an efficient implementation.
- The LIBSVM library was once popular for SVM . Link-<https://www.csie.ntu.edu.tw/~cjlin/libsvm/index.html>
- You can also code the SVM from scratch. You need to solve the dual problem ,which is a QPP. You can use the *quadprog* function for solving the QPP. For coding the SVM from scratch, you can refer the Appendices of the <https://svms.org/tutorials/Gunn1998.pdf>

Popular Tutorial

- Gunn, Steve R. "Support vector machines for classification and regression." *ISIS technical report* 14.1 (1998): 5-16.
- Burges, Christopher JC. "A tutorial on support vector machines for pattern recognition." *Data mining and knowledge discovery* 2.2 (1998): 121-167.

Popular Books

- Deng, Naiyang, Yingjie Tian, and Chunhua Zhang. *Support vector machines: optimization based theory, algorithms, and extensions*. CRC press, 2012. (Optimization Prespective)
- Scholkopf, Bernhard, and Alexander J. Smola. *Learning with kernels: support vector machines, regularization, optimization, and beyond*. MIT press, 2018.

Questions ??

Thanks