

Mathematical Programming Applications in Machine Learning



E V N Sai Bharadwaja
Machine Learning and Computing
IIST

Q. What is the goal of any Machine Learning algorithm?

- Training Data
- Testing Data

Three basic basic problems that are often Studied in Machine Learning are :

- Regression
- Classification
- Clustering

There are so many approaches to solve above problem

- But **Mathematical Programming** has now become very popular in Machine Learning Community and Mathematical Programming Community

Q. What is the Advantage of Studying Mathematical Programming?

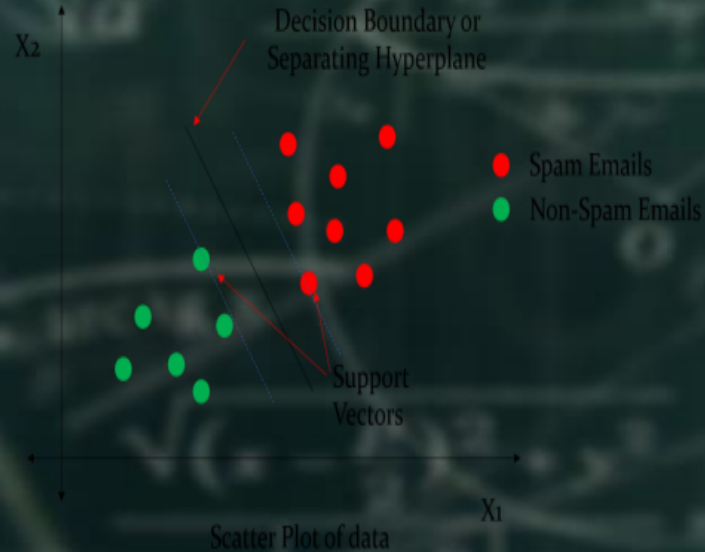
- Most of the Machine Learning Models result into Linear Programming , Quadratic Programming , Convex Programming.
- But Mathematical Programming has
 - Theoretical results(KKT conditions and Duality Theory)
 - Efficient algorithms

Machine Learning Techniques

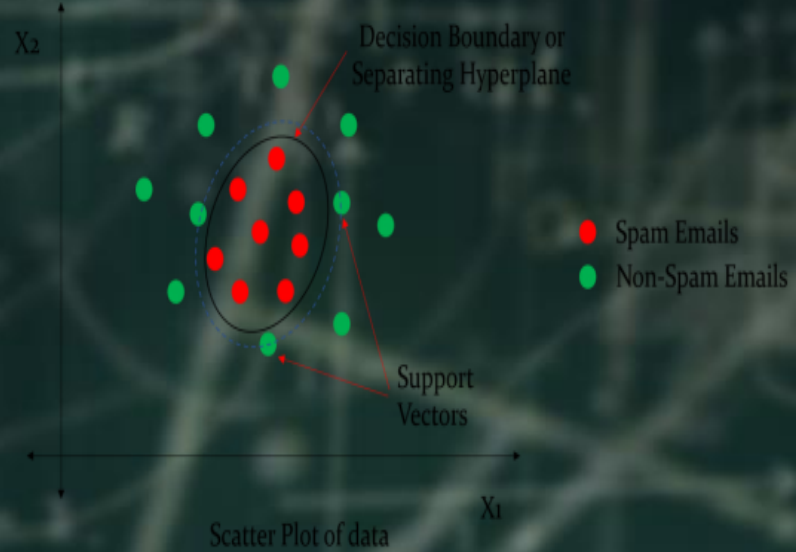
- Supervised Learning (Ex : Classification and Regression)
- Unsupervised Learning (Ex : Clustering)

Binary (Two-Class) Pattern Classification Problems

- $X \in \mathbb{R}^n$
- m - patterns \in Class +1
- k - patterns \in Class -1



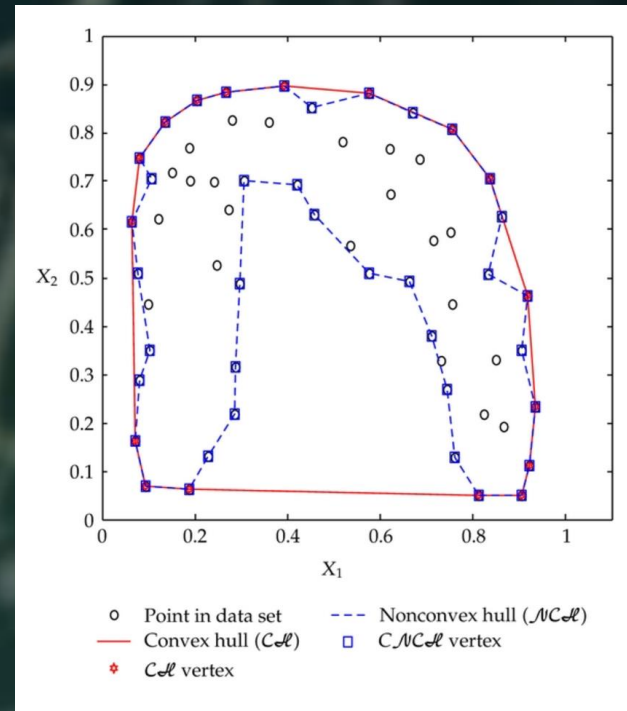
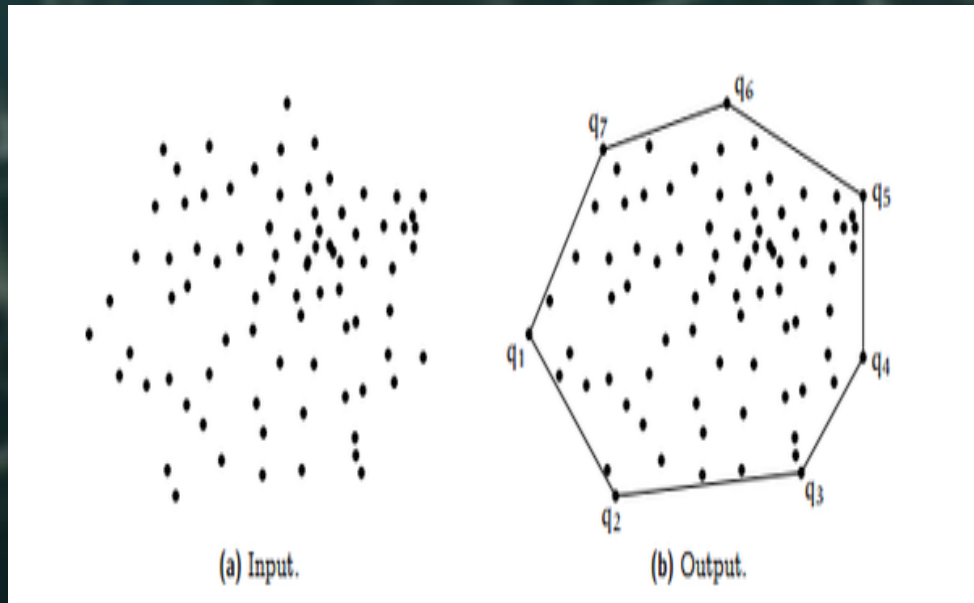
Linear Decision Boundary



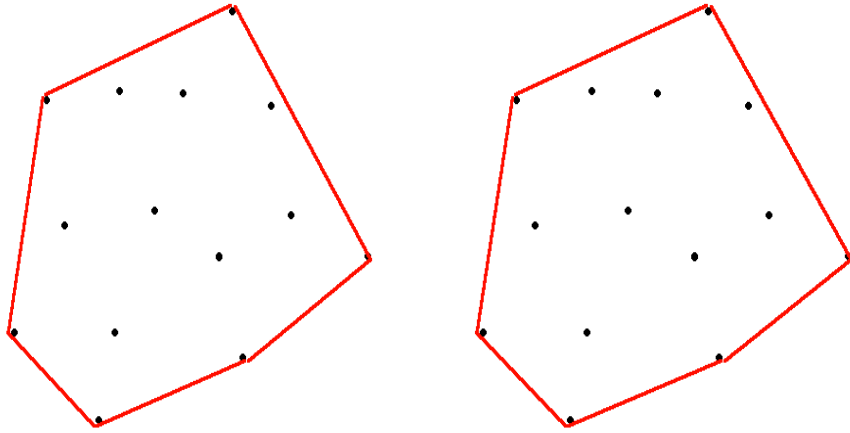
Non- Linear Decision Boundary

Convex Hull :

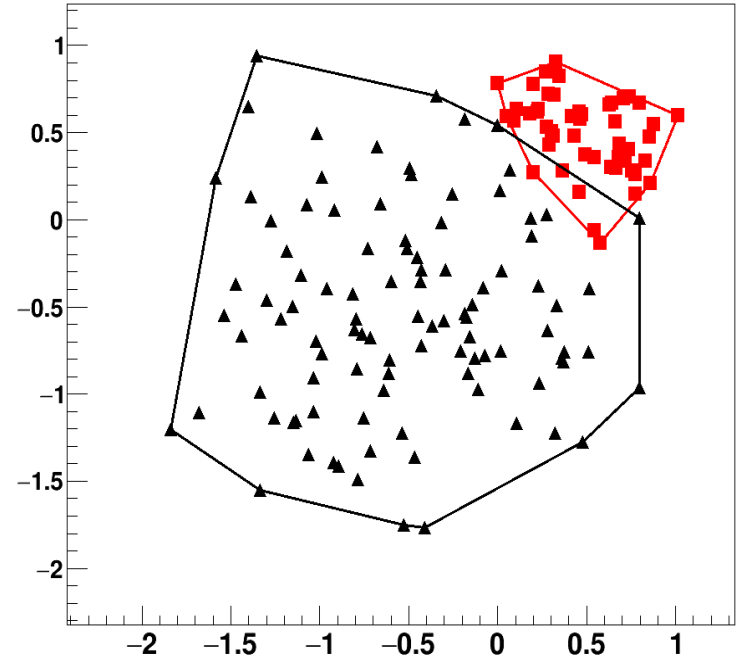
The Convex hull of a set of points X in euclidean space is the smallest convex set containing X .



Disjoint Convex hulls
Or Linearly Separable



Joint Convex hulls
Or Linearly non Separable



Q. $A = \{ (-1,0), (0,1), (1,0) \}$,

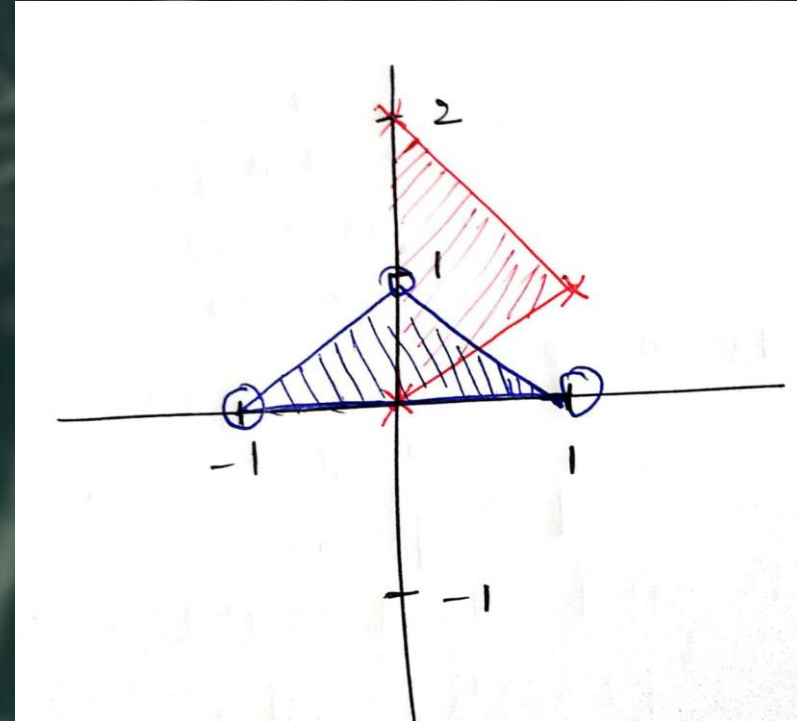
$B = \{ (0,0), (1,1), (0,2) \}$

Are A and B linearly Separable ?

Ans :

- Convex Hulls are NOT Disjoint Hence given Problem is NOT Linearly Separable.
- But Finding convex hulls is not an easy task

So we make use of **Linear Programming** to decide whether the problem is Linearly Separable or not



Solving the above problem using Linear Programming

$$X = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m1} & x_{m2} & \dots & x_{mn} \end{bmatrix}_{N \times n}$$

m patterns Belongs to class 1.

K patterns Belongs to class -1

→ From the given Dataset X.
points those belong to class 1

$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix}_{m \times n}$$

$$\sum_{j=1}^n a_{ij} w_j \geq b \quad \forall i = 1, 2, \dots, m$$

$$\sum_{j=1}^n a_{ij} w_j \geq b + \alpha \quad \forall i = 1, 2, \dots, m$$

$$\sum_{j=1}^n a_{ij} w_j \geq b + \alpha^*$$

where $\alpha^* = \min[\alpha_i]$

$$Aw \geq eb + e$$

→ From the given Data set X.
points belongs to class -1

$$B = \begin{bmatrix} b_{11} & b_{12} & \dots & b_{1n} \\ b_{21} & b_{22} & \dots & b_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ b_{k1} & b_{k2} & \dots & b_{kn} \end{bmatrix}$$

$$\sum_{j=1}^n b_{rj} w_j \leq b \quad \forall r = 1, 2, \dots, k$$

$$\sum_{j=1}^n b_{rj} w_j \leq b - \beta \quad \forall r = 1, 2, \dots, k$$

$$\sum_{j=1}^n b_{rj} w_j \leq b - \beta^*$$

where $\beta^* = \min[\beta_r]$

$$Bw \leq eb - e$$

Error Minimizing Linear Programming Problem

→ For a Real number $a_+ = \max(a, 0)$.

→ For $x \in \mathbb{R}^n$ $x_+ = [(x_+)_1, (x_+)_2, \dots, (x_+)_n]^T \in \mathbb{R}^n$

→ L-1 norm $\|x\|_1 = \sum_{i=1}^n |x_i|$

So the optimization problem.

$$\begin{aligned} \min_{(w,b)} \quad & \frac{1}{n} \|(-Aw + eb + e)_+\|_1 \\ & + \frac{1}{k} \|(Bw - eb + e)_+\|_1 \end{aligned}$$

If A and B are **Linearly Separable** then Error is zero and **Optimal value** for the Optimization problem is **zero** and the converse is also True.

Converse Statement :

If **Optimal Value** for the optimization problem is **zero** then A and B are **Linearly Separable** otherwise **NOT Linearly Separable**

Lemma : Transforms into LPP

Lemma: Consider the problems

$$\begin{aligned} \text{i)} \quad & \min_{x \in S} \|g(x)_+\|_1 + \|h(x)_+\|_1 \\ \text{ii)} \quad & \min_{x \in S} \left\{ e^T y + e^T z : \begin{aligned} y &\geq g(x), y \geq 0 \\ z &\geq h(x), z \geq 0 \end{aligned} \right\} \end{aligned}$$

Where $S \subset \mathbb{R}^n$, $g: S \rightarrow \mathbb{R}^m$, $h: S \rightarrow \mathbb{R}^k$,
 $y \in \mathbb{R}^m$ and $z \in \mathbb{R}^k$.

→ Then both ~~pro~~ problem [i] and [ii] have identical solution sets.

Applying Lemma to our Optimization Problem

Our optimization problem is

$$\begin{aligned} \min_{(w,b)} \quad & \frac{1}{m} \|(-Aw + eb + e)_+\|_1 + \\ & \frac{1}{k} \|(Bw - eb + e)_+\|_1 \end{aligned}$$

After Applying Lemma.

$$\begin{aligned} \min_{w, b, y, z} \quad & \frac{e^T y}{m} + \frac{e^T z}{k} \\ \text{Sub to:} \quad & Aw - eb + y \geq e \\ & -Bw + eb + z \geq e \\ & y, z \geq 0 \end{aligned} \quad \left. \vphantom{\begin{aligned} \min_{w, b, y, z} \quad & \frac{e^T y}{m} + \frac{e^T z}{k} \\ \text{Sub to:} \quad & Aw - eb + y \geq e \\ & -Bw + eb + z \geq e \\ & y, z \geq 0 \end{aligned}} \right\} \text{LPP.}$$

w and b are unrestricted.

LPP can be Efficiently Solved by using Simplex Algorithm.

Hence we can find whether A and B are Linearly separable or Not

Conclusion :

- A and B are linearly separable iff optimal value for the error minimizing LPP is ZERO . Hence there exists a Hyperplane $w.T * x = b$ which is Linear Separator
- A and B are NOT linearly separable then optimal value for the error minimizing LPP is NOT ZERO . Hence there exists NO Hyperplane $w.T * x = b$ which is Linear Separator

Q. Write error minimizing LPP for AND problem and check if the given Problem is Linearly Separable ?

AND:

x_1	x_2	0/1	Bipolar 0/1
0	0	0	-1
0	1	0	-1
1	0	0	-1
1	1	1	1

Solution:

An optimal solution for the LPP

$$w_1 = 2 \quad w_2 = 2 \quad b = 3$$

$$y_1 = 0 \quad z_1 = 0 \quad z_2 = 0 \\ z_3 = 0$$

Class +1

$$A = \begin{bmatrix} 1 & 1 \end{bmatrix}$$

$$m = 1$$

Class -1:

$$B = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \end{bmatrix}$$

$$K = 3$$

Optimization problem:

$$\text{Min } y_1 + \frac{1}{3} [z_1 + z_2 + z_3]$$

$$\text{Sub to: } 1 \cdot w_1 + 1 \cdot w_2 - b + y_1 \geq 1$$

$$-0 \cdot w_1 - 0 \cdot w_2 + b + z_1 \geq 1$$

$$-0 \cdot w_1 - 1 \cdot w_2 + b + z_2 \geq 1$$

$$-1 \cdot w_1 - 0 \cdot w_2 + b + z_3 \geq 1$$

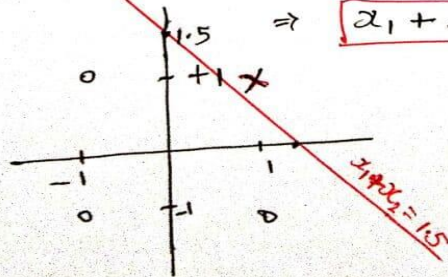
$$y_1, z_1, z_2, z_3 \geq 0$$

$$w_1, w_2, b \text{ unrestricted.}$$

→ The optimal value of the error minimizing LPP is zero, given AND problem is Linearly separable with Hyperplane.

$$w_1 x_1 + w_2 x_2 = b \Rightarrow 2x_1 + 2x_2 = 3$$

$$\Rightarrow x_1 + x_2 = 1.5$$



Q. Write error minimizing LPP for XOR problem and check if the given Problem is Linearly Separable ?

XOR problem:

x_1	x_2	o/p	Bipolar l/p
0	0	0	-1
0	1	1	1
1	0	1	1
1	1	0	-1

class +1

$$A = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

$$m = 2$$

class -1

$$B = \begin{bmatrix} 0 & 0 \\ 1 & 1 \end{bmatrix}$$

$$k = 2$$

optimization problem:

$$\text{Min } \frac{1}{2} [y_1 + y_2] + \frac{1}{2} [z_1 + z_2]$$

The solution for LPP.

$$w_1 = 2 \quad w_2 = 2 \quad b = 1$$

$$y_1 = 0 \quad y_2 = 0$$

$$z_1 = 0 \quad z_2 = 4$$

$$\text{sub to: } 0 \cdot w_1 + 1 \cdot w_2 - b + y_1 \geq 1$$

$$1 \cdot w_1 + 0 \cdot w_2 - b + y_2 \geq 1$$

$$-0 \cdot w_1 - 0 \cdot w_2 + b + z_1 \geq 1$$

$$-1 \cdot w_1 - 1 \cdot w_2 + b + z_2 \geq 1$$

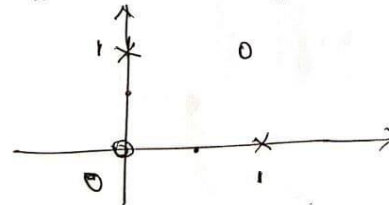
$$y_1, y_2, z_1, z_2 \geq 0$$

w_1, w_2, b are unrestricted.

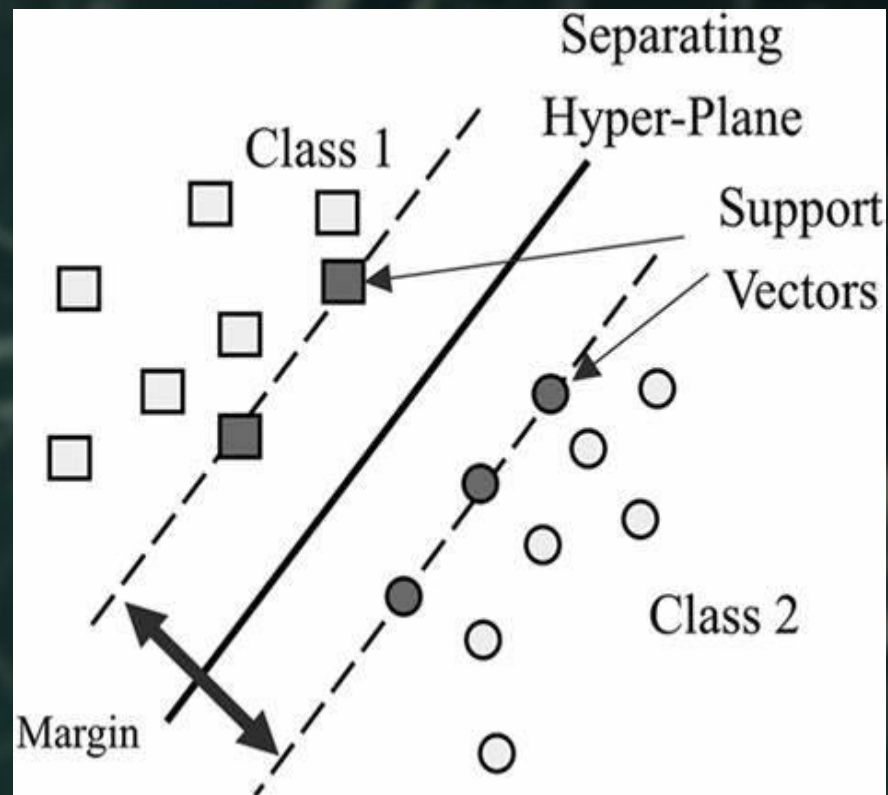
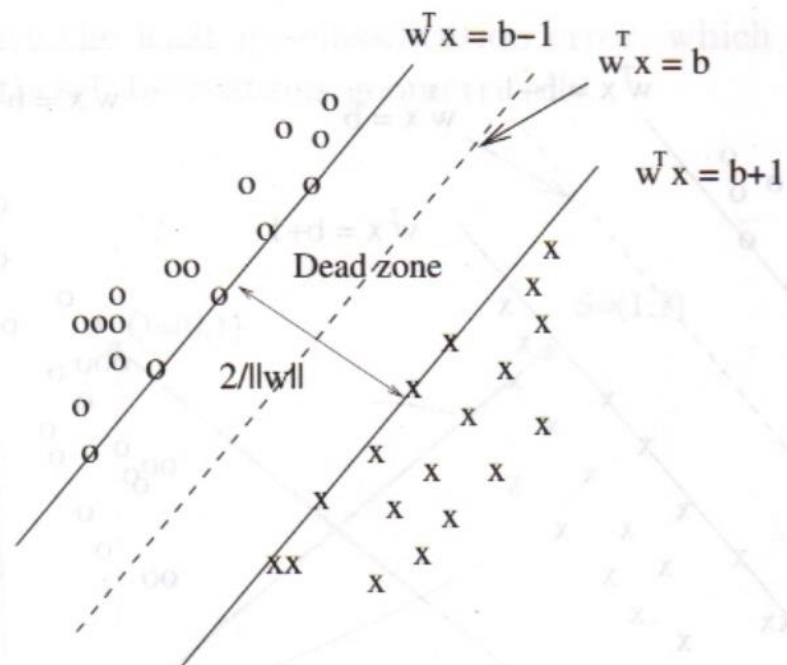
The optimal value for error minimizing LPP is 5.

optimal value $\neq 0$. [Hence There exists
No hyperplane]

$$2x_1 + 2x_2 = 1 \Rightarrow x_1 + x_2 = 0.5$$



- If A and B are linearly separable then there exists infinitely many hyperplanes
- So how should we choose a Hyperplane ?
- Is there any optimal Hyperplane ?



Canonical Hyperplane

Dead Zone

→ Canonical Hyperplane:

Let $A, B \subset \mathbb{R}^n$ be linearly separable.

Then the separating hyperplane $w^T x = b$ is called canonical hyperplane if it satisfies

$$Aw \geq eb + \epsilon$$

$$Bw \leq eb - \epsilon$$

→ Dead zone:

Let $A, B \subset \mathbb{R}^n$ be linearly separable and $w^T x = b$ be a canonical separating hyperplane.

Then there exists a ~~linear~~ Region

$\{x; (b-1) < w^T x < (b+1)\} \subset \mathbb{R}^n$ surrounding the

separating hyperplane $w^T x = b$, which is

void of points sets A and B .

This Region is called Dead zone for the separating hyperplane $w^T x = b$.

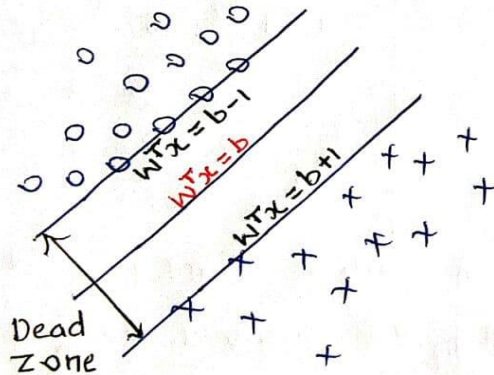
Margin

Margin:

Let $A, B \subset \mathbb{R}^n$ be linearly separable and $w^T x = b$ be canonical separating hyperplane.

Let $w^T x = b-1$ and $w^T x = b+1$ be the bounding hyperplanes which define the dead zone.

The distance between the bounds $w^T x = b-1$ and $w^T x = b+1$ is called the margin.



→ The perpendicular distance between $w^T x = b$ and $w^T x = b-1$ is $\frac{1}{\|w\|}$

→ so the perpendicular distance between $w^T x = b-1$ and $w^T x = b+1$ is

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

where $\|w\| = \sqrt{w_1^2 + \dots + w_n^2}$

[The distance of a point (x', y') to the plane $ax + by = 0$ is given by $\frac{|ax' + by'|}{\sqrt{a^2 + b^2}}$]

The distance of a point $x_i = (x_{i1}, x_{i2}, \dots, x_{in})$ to the hyperplane $w^T x + b = 0$ is given by

$$\begin{aligned} \text{dist}(x_i, \tilde{f}(x)) &= \frac{|w_1 x_{i1} + w_2 x_{i2} + \dots + w_n x_{in} + b|}{\sqrt{w_1^2 + w_2^2 + \dots + w_n^2}} \\ &= \frac{|w^T x_i + b|}{||w||} \end{aligned}$$

$w^T x + b \geq 1$ for all x in the positive class. Therefore

$\text{dist}(x, \tilde{f}(x)) \geq \frac{1}{\|w\|}$, for all x in the positive class. The points in the positive class, that lie closest to the separating hyperplane $w^T x + b = 0$ lies in the hyperplane

$$H_1 : w^T x + b = 1$$

Therefore $d^+ = \frac{1}{\|w\|}$.

In the negative class, for all points $w^T x + b \leq -1$. Therefore

$\text{dist}(x, \tilde{f}(x)) \geq \frac{1}{\|w\|}$ for all x in the negative class. The points in the negative class, that lie closest to the separating hyperplane $w^T x + b = 0$ lies in the hyperplane

$$H_2 : w^T x + b = -1$$

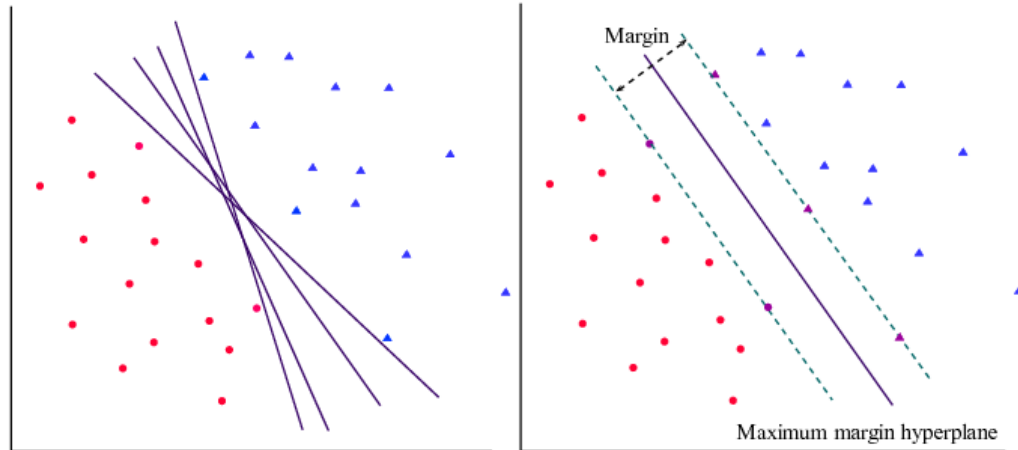
Therefore $d^- = \frac{1}{\|w\|}$.

Hence margin $\gamma = d^+ + d^- = \frac{2}{\|w\|}$

Optimal Separating Hyperplanes:

- Let sets A , B be linearly separable and $w.T * x = b$ be a canonical hyperplane .
- Then $w.T * x = b$ is called **Optimal Separating Hyperplane** if it's **Margin is Maximum** among all the Canonical Hyperplanes

eparating hyperplane, can achieve maximum margin.



- In Classification , our aim is not only to classify the available to classify the training data in accordance with the class labels
- But also to have good generalization on Unseen Data or Test Data
- The Larger the dead Zone , the smaller the misclassification Error will be.
- Therefore we would like to Maximize the dead Zone , which implies we would like to Maximize the Margin
- From the above Point of View , for the linearly Separable patterns , our aim is not only just finding the Hyperplane , But also for which the Margin is Maximum
- Ques : Will Error Minimizing LPP gives Optimal Separating Hyperplane ?
- So we should construct a Optimization problem in which Objective Function should contain Margin
- This Discussion Leads us to Hard Margin Classifier and Support Vector Machines

Hard Margin Classifier :

For the given patterns which are Linearly separable

→ Let $\{(x^{(i)}, y_i), i=1, 2, \dots, p\}$ be a set of finite Training samples (or) patterns where $x^{(i)} \in \mathbb{R}^n$ and $y_i \in \{-1, 1\}$.

→ Here $y_i \in \{-1, 1\}$ represents the Class Labels of the Patterns $x^{(i)}$, $(i=1, 2, \dots, p)$.

→ Let the given patterns be Linearly separable
Then There exists a $w \in \mathbb{R}^n$ and $b \in \mathbb{R}$, $\forall i=1, 2, \dots, p$

$$w^T x^{(i)} - b > 0 \quad \forall i \text{ having } y=1$$

$$w^T x^{(i)} - b < 0 \quad \forall i \text{ having } y=-1$$

By suitable scaling we can write above as

$$w^T x^{(i)} - b \geq 1 \quad \forall i \text{ having } y=1$$

$$w^T x^{(i)} - b \leq -1 \quad \forall i \text{ having } y=-1$$

writing in more compact form.

$$y_i (w^T x^{(i)} - b) \geq 1 \quad \forall i=1, 2, \dots, p.$$

→ Amongst all separating hyperplanes $w^T x = b$
we have to choose the one for which

margin $\frac{2}{\|w\|}$ is maximum, (or)

Equivalently $\frac{\|w\|}{2}$ is minimum.

→ How ever minimizing $\frac{\|w\|}{2}$ is equivalent to
minimizing $\frac{\|w\|^2}{2}$ because $\|w\|$ is monotonically
increasing for $w > 0$

∴ optimization problem:

$$\text{Min}_{(w,b)} \quad \frac{1}{2} w^T w$$

$$\text{sub to: } y_i (w^T x^{(i)} - b) \geq 1 \quad \forall i=1, 2, \dots, p.$$

→ The above problem is in the form of
Standard Quadratic Programming Problem [QPP],
which can be solved by standard QPP algorithm.

→ Once optimal solution of (w, b) is known,
The maximum margin classifier will be found

$$\boxed{w^T x = b} \rightarrow \text{Hard Margin classifier}$$

Difficulties in Solving Quadratic Programming Problem

- QPP has as many constraints as the number of patterns and hence extremely hard to solve for Large Data sets
- This QPP will make use of Kernel Methods, which again is difficult to Store kernel matrices.
- In Machine Learning there are special algorithms so that Not all Constraints are included at once
 - Chunking algorithm
 - Decomposition algorithm
 - SMO algorithm (Special Case of Decomposition)
- The above Techniques are also called as Active Set or Working Set approach

- In **Mathematical Programming** , Standard approach to handle Large number of constraints is to examine if **Dual problem** can be solved efficiently
- This analysis requires **KKT conditions**

- KKT conditions

min $1/2 \|w\|^2$

Sub : $y_i(\langle w, x_i \rangle + b) \geq 1$

$$L(w, b, \alpha) = \frac{1}{2} \langle w, w \rangle - \sum_{i=1}^N \alpha_i [y_i(\langle w, x_i \rangle + b) - 1]$$

where $\alpha_i \geq 0$ are the Lagrange multipliers.

The corresponding dual is found by differentiating with respect to w and b , imposing stationarity,

$$\frac{\partial L(w, b, \alpha)}{\partial w} = w - \sum_{i=1}^N y_i \alpha_i x_i = 0$$

$$\frac{\partial L(w, b, \alpha)}{\partial b} = \sum_{i=1}^N y_i \alpha_i = 0$$

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

For $\alpha_i \neq 0$
W exist

$$f(x) = \langle w, x \rangle = \sum_{i=1}^N y_i \alpha_i \langle \boxed{x_i}, x \rangle$$

Support Vectors:

- Only Those points for which x_i is influencing Weights 'W'
- KKT multipliers (Alpha) > 0
- So vector W is known once KKT multipliers are known.

KKT Complimentary Condition

Using Karush-Kuhn-Tucker complementarity conditions,

$$\alpha_i^* [y_i (\langle w, x_i \rangle + b) - 1] = 0, i = 1, 2, \dots, N$$

Therefore two cases

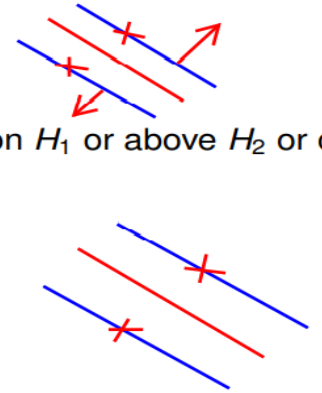
- $\alpha_i^* = 0$
 - $y_i (\langle w, x_i \rangle + b) - 1 \geq 0$
 - x_i lies either above H_1 or on H_1 or above H_2 or on H_2
- $\alpha_i^* \neq 0$
 - $y_i (\langle w, x_i \rangle + b) - 1 = 0$
 - $\boxed{x_i}$ lies either on H_1 or H_2

Support Vectors :

$\alpha \neq 0$,

α Contributes in weight W

Hence those x_i for $\alpha \neq 0$ are called Support Vectors



- The SVM algorithm clearly needs the knowledge of KKT multipliers
- Once KKT multipliers are known Support Vectors
 - Support Vectors:
 - Only Those points for which X_i is influencing Weights 'W'
 - KKT multipliers (Alpha) > 0
 - So vector W is known once KKT multipliers are known.
- Hence we get to know vector (w) and b
- From the concept of Non Linear Programming Duality , KKT multipliers are nothing but Dual Variables.
- Therefore we go for Wolfe Dual of the Problem to find KKT multipliers

$$\begin{aligned}
 \|w\|^2 &= \langle w, w \rangle \\
 &= \left\langle \sum_i \alpha_i y_i x_i, \sum_i \alpha_i y_i x_i \right\rangle
 \end{aligned}$$

$$\|w\|^2 = \sum_i \sum_j \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle$$

$$\begin{aligned}
 \sum_{i=1}^N \alpha_i [y_i (\langle w, x_i \rangle)] &= \sum_{i=1}^N \alpha_i [y_i (\langle \sum_{j=1}^N y_j \alpha_j x_j, x_i \rangle)] \\
 &= \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle
 \end{aligned}$$

$$L(w, b, \alpha) = \frac{1}{2} \langle w, w \rangle - \sum_{i=1}^N \alpha_i [y_i (\langle w, x_i \rangle + b) - 1]$$

$$\|w\|^{**2} = \sum_i \sum_j \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle$$

Dual Formulation

$$b \sum_i \alpha_i y_i = 0$$

$$\begin{aligned} \sum_{i=1}^N \alpha_i [y_i (\langle w, x_i \rangle)] &= \sum_{i=1}^N \alpha_i [y_i (\langle \sum_{j=1}^N y_j \alpha_j x_j, x_i \rangle)] \\ &= \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle \end{aligned}$$

Sub (8) and (9) into (5) to form the dual,

$$\begin{aligned} W(\alpha) &= \frac{1}{2} \sum_{i,j=1}^N y_i y_j \alpha_i \alpha_j \langle x_i, x_j \rangle - \sum_{i,j=1}^N y_i y_j \alpha_i \alpha_j \langle x_i, x_j \rangle + \sum_{i=1}^N \alpha_i \\ &= \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N y_i y_j \alpha_i \alpha_j \langle x_i, x_j \rangle \end{aligned}$$

Wolfe Dual :

- If we clearly observe dual problem , we can notice that there is only one constraint
- The objective Function of the Dual problem is Concave
- Therefore Dual problem is Easier to Solve Than Primal problem
- Once optimal alpha is known [w] and b can be computed.
- Hence separating hyperplane with Maximum Margin can be determined

Primal Problem : Number of Constraints = Number of Data Points

Dual Problem : Number of Constraints = 1 [Easy to Solve]

Hence the dual optimisation problem is

$$\text{maximise } W(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N y_i y_j \alpha_i \alpha_j \langle x_i, x_j \rangle$$


Only One Constraint Easy to Solve

$$\frac{\partial L(w, b, \alpha)}{\partial b} = \sum_{i=1}^N y_i \alpha_i = 0$$

$$\text{subject to } \sum_{i=1}^N y_i \alpha_i = 0$$

$$\alpha_i \geq 0, i = 1, 2, \dots, N$$

$L(w, \alpha, b)$ and $W(\alpha)$ arise from the same objective function but with different constraints and the solution is found by minimizing $L(w, \alpha, b)$ or by maximizing $W(\alpha)$.

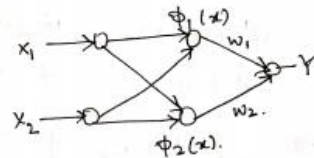
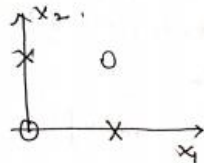
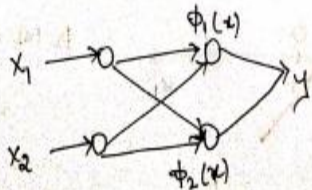
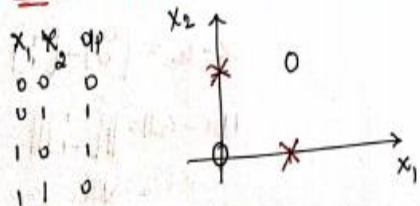
1. Gaussian function $\phi(x) = e^{-\frac{x^2}{2\sigma^2}}$ 

2. Multi quadratic functions $\phi(x) = \sqrt{x^2 + c^2}$ $c > 0$

3. Inverse multi quadratic fn $\phi(x) = \frac{1}{\sqrt{x^2 + c^2}}$ $c > 0$.

$$\phi_i(x) = e^{-\frac{1}{2\sigma^2} \|x - c_i\|^2} \quad i=1, 2, \dots, m.$$

Ex: XOR - problem. [linearly non separable]



$$\phi_1(x) = e^{-\|x - c_1\|^2} \quad \sigma^2 = 1$$

$$\phi_2(x) = e^{-\|x - c_2\|^2} \quad \sigma^2 = 1$$

where c_1 & c_2 are the center of the pattern clusters.

$$\text{let } c_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad c_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

x_1	x_2	$\phi_1(x)$	$\phi_2(x)$	Y
0	0	1	0.1	0
0	1	0.4	0.4	1
1	0	0.4	0.4	1
1	1	0.1	1	0

these centres
can be found out
in clustering
techniques.

$$x - c_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix} - \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\|x - c_1\| = \left\| \begin{bmatrix} 0 \\ 0 \end{bmatrix} \right\| = 0$$

$$\phi_1(x) = e^{-0} = 1$$

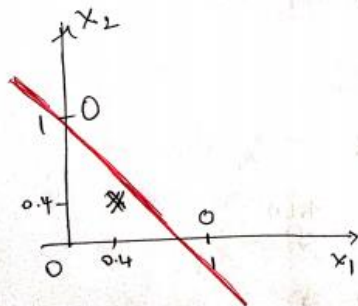
$$x - c_2 = \begin{bmatrix} 0 \\ 0 \end{bmatrix} - \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\|x - c_2\| = \left\| \begin{bmatrix} -1 \\ -1 \end{bmatrix} \right\|$$

$$= \sqrt{2}$$

$$\phi_2(x) = e^{-(\sqrt{2})^2}$$

$$= e^{-2}$$



Soft Margin Classifier :

For the given patterns which are **NOT** Linearly separable

→ Let $\{(x^{(i)}, y_i) \mid i = 1, 2, \dots, p\}$ be a finite training sample of patterns where $x^{(i)} \in \mathbb{R}^n$ and $y_i \in \{-1, 1\}$.

→ We consider the case where the patterns are NOT linearly separable

→ In this case Error minimizing problem will have a non-zero objective function value.

→ In the optimal solution not all the errors will be zero

→ Let the error variables will be denoted by ξ_i ($i = 1, 2, \dots, p$).

→ So our aim is to maximize margin $\frac{2}{\|w\|}$

and $\sum_{i=1}^p \xi_i$ is least.

→ \therefore minimize both $\frac{w^T w}{2}$ and $\sum_{i=1}^p \xi_i$

This will be multi objective optimization.

→ Since this is not possible, we consider

weighted sum

$$\min \frac{w^T w}{2} + c \cdot \sum_{i=1}^p \xi_i$$

Optimisation Problem

$$\min_{f \in \mathcal{F}, b \in \mathbb{R}} \frac{1}{2} \|f\|^2 + C \sum_{i=1}^N \xi_i$$

subject to

$$y_i(\langle f, k_{x_i} \rangle + b) - 1 + \xi_i \geq 0, \quad i = 1, \dots, N$$
$$\xi_i \geq 0.$$

- The objective function is quadratic and all the constraints are linear

The primal Lagrangian is,

$$L(f, b, \xi, \alpha) = \frac{\|f\|^2}{2} + C \sum_{i=1}^N \xi_i - \sum_i \alpha_i [y_i(\langle f, k_{x_i} \rangle + b) - 1 + \xi_i] - \sum_i \mu_i \xi_i$$

$$\frac{\partial L}{\partial f} = f - \sum_{i=1}^N \alpha_i y_i k_{x_i} = 0$$

$$f = \sum_{i=1}^N \alpha_i y_i k_{x_i}$$

$$\frac{\partial L}{\partial b} = \sum_{i=1}^N \alpha_i y_i = 0$$

For $i=1,2,\dots, N$,

$$\frac{\partial L(w, b, \alpha)}{\partial \xi_i} = C - \alpha_i - \mu_i = 0$$

$$\alpha_i \geq 0$$

$$\mu_i \geq 0$$

The KKT complimentary conditions are,

$$\alpha_i[y_i(\langle f, k_{x_i} \rangle + b) - 1 + \xi_i] = 0$$

$$\mu_i \xi_i = 0$$

$$C - \alpha_i = \mu_i \geq 0$$

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

Three cases

- $\alpha_i = 0$ $y_i(\langle f, k_{x_i} \rangle + b) \geq 1.$
- $0 < \alpha_i < C$ $y_i(\langle f, k_{x_i} \rangle + b) = 1$
- $\alpha_i = C$ $y_i(\langle f, k_{x_i} \rangle + b) \leq 1.$

Support vectors are those points for which $y_i(\langle f, k_{x_i} \rangle + b) \leq 1$.

- $y_i(\langle f, k_{x_i} \rangle + b) = 1$, k_{x_i} on H_1 or on H_2
- $0 \leq y_i(\langle f, k_{x_i} \rangle + b) < 1$ Between H_1 and H or H_2 and H or on H
- $y_i(\langle f, k_{x_i} \rangle + b) < 0$ Incorrect classification

Dual Optimisation

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

$$\text{subject to } \sum_{i=1}^N y_i \alpha_i = 0$$

$$C \geq \alpha_i \geq 0, i = 1, 2, \dots, N.$$

$$f = \sum_{i=1}^{sv} y_i \alpha_i^* k_{x_i}$$

where sv is the number of support vectors. Then

$$\begin{aligned} \tilde{f}(x) &= \langle f, k_x \rangle + b \\ &= \sum_{i=1}^{sv} y_i \alpha_i^* \langle k_{x_i}, k_x \rangle + b \\ &= \sum_{i=1}^{sv} y_i \alpha_i^* k(x_i, x) + b \end{aligned}$$

The background is a dark teal color with a complex pattern of faint, light-colored mathematical formulas and geometric diagrams. These include various equations, circles, lines, and points, creating a dense, intellectual atmosphere. The text "Thank You" is centered in a bright cyan color.

Thank You