# CS-671: Deep Learning and its Applications
## Lecture: 02 (a)
## Basics of Machine Learning: Linear Classification (I)

Aditya Nigam, Assistant Professor
School of Computing and Electrical Engineering (SCEE)
Indian Institute of Technology, Mandi
http://faculty.iitmandi.ac.in/ãditya/ aditya@iitmandi.ac.in

February - May, 2019

# Perceptron Algorithm :

It is a discriminative technique in which
* we try to learn the discriminant fn directly. We donot attempt
to learn the parameter distribution instead we will learn the
boundary explicitly. $\Longrightarrow$ Our Classifier will be using
discriminative features to learn discriminant
fn. itself.

* In other non-discriminant techniques (such as Bayes D. fn)
We learn the data distribution parameter and boundary.
was the result of that.

* Now we learn the boundary b/w 2 classes directly
and hence we need not to assume that data is following
certain distribution. Rather we assume data is linearly
Seperable or not.

Initially we assume linear discriminant fn. and perform LDA, as non-linear fn's are very difficult to learn and require more complex parameterization.
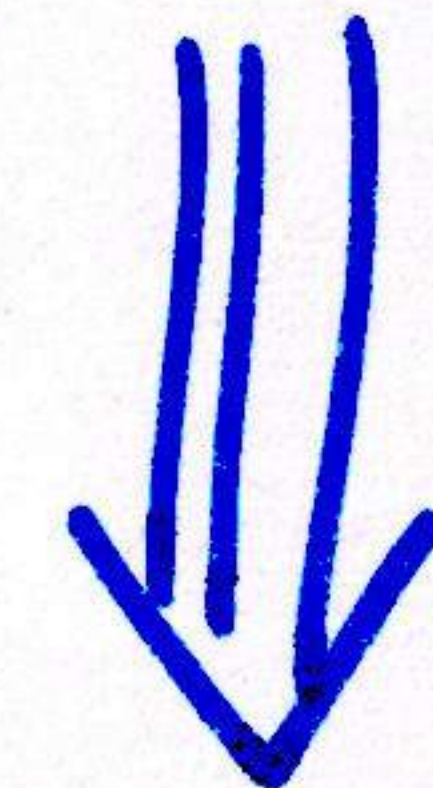
$$f(x) = \omega^T x + \omega_0 = 0$$

weight vector → ← data point → bias
[Column vector]

$$W = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ \vdots \\ w_d \end{bmatrix} \qquad x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_d \end{bmatrix}$$
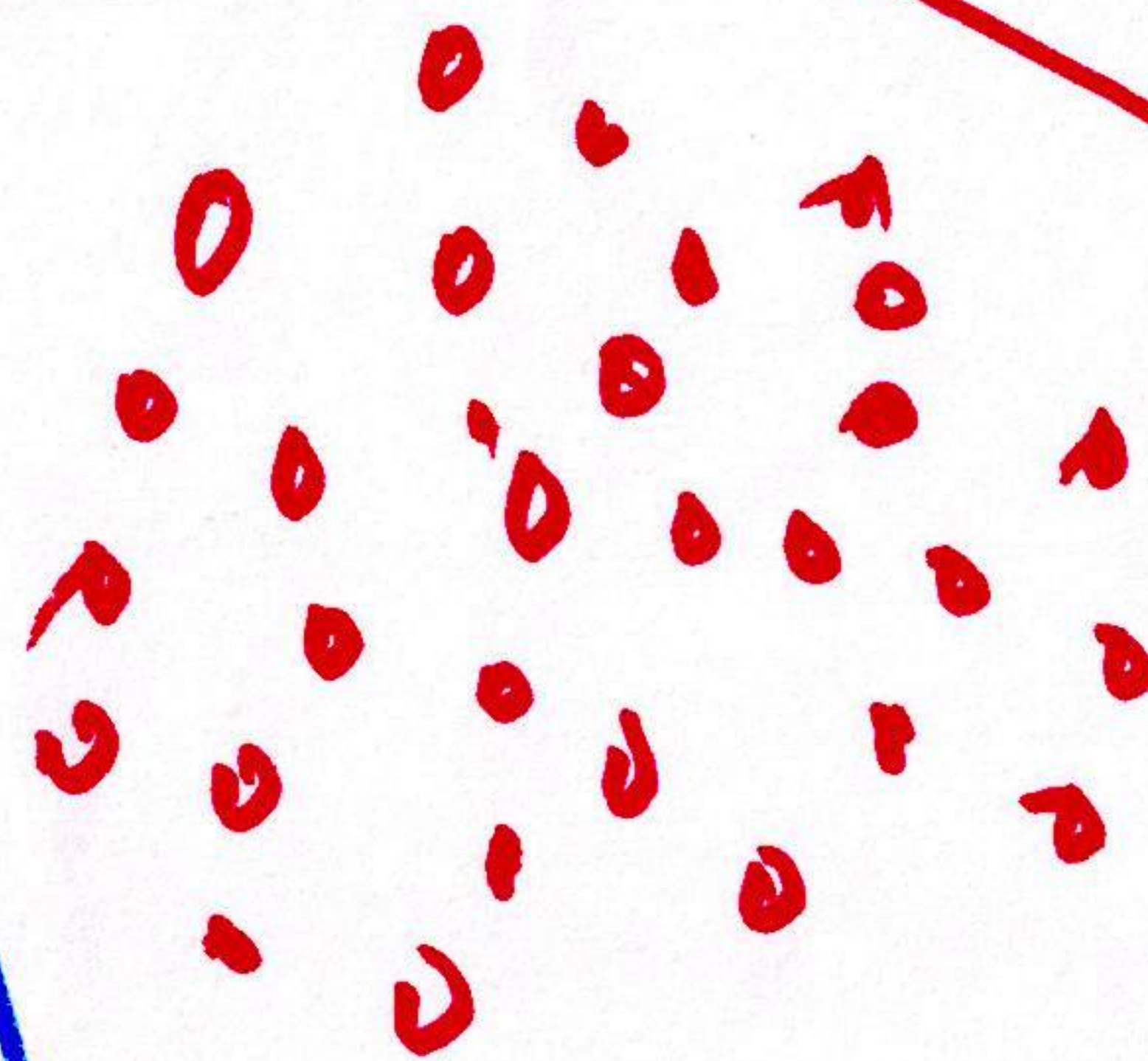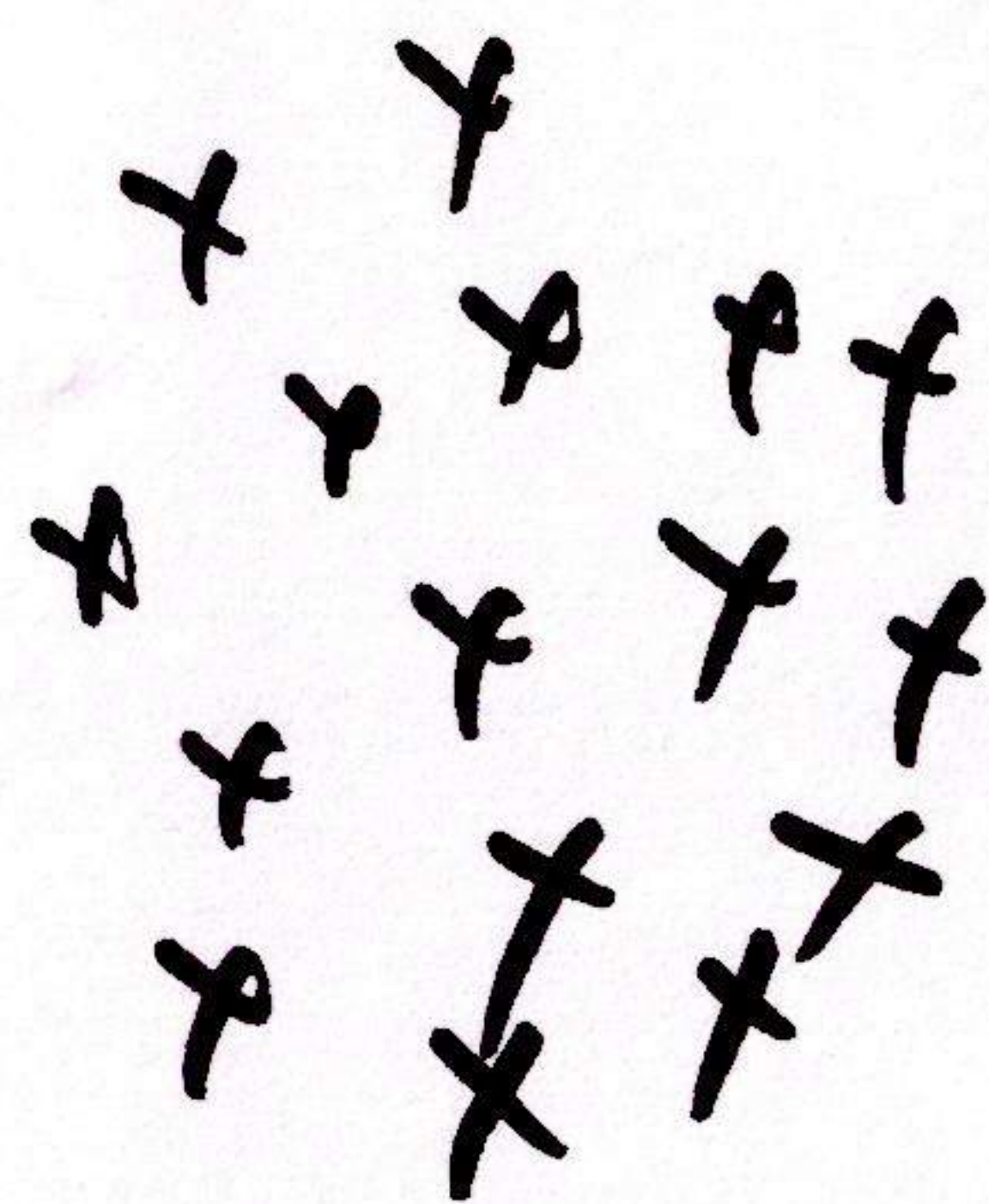
$$W_1 x_1 + W_2 x_2 + W_3 x_3 + \cdots W_0 = 0$$

$W^T x + W_0 < 0$

$W^T x + W_0 \geqslant 0$

$$a = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} \qquad z = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_d \end{bmatrix}$$

$$f(x) = \bar{a}^T z \qquad \begin{bmatrix} \text{Augumented} \\ \text{Representation} \end{bmatrix}$$

If data is linearly Seperable then in 2D space a line can be seen as a linear D.fn $f(x) = w_1 x_1 + w_2 x_2 + w_0$ that can seperate two classes, and we are looking to learn such fuction and its parameters. Hence basically we are looking to learn :

$$\begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \quad \text{and} \quad [w_0]$$

vector $\longrightarrow$     scalar $\longrightarrow$

weights/coefficient of the straight line.

$$w_1 x_1 + w_2 x_2 + w_0 = 0$$

$$x_2 = \boxed{-\frac{w_1}{w_2}} x_1 \boxed{-\frac{w_0}{w_2}}$$

$$\underline{y = mx + c}$$

$m$ [slope]          $c$ [intercept]

# TLN (Threshold logic neuron):

Linear neuron whose O/P is directed into a unit step or signum fn. It is adaptive as its weights are changing according to some fixed rule. __Widrow__ called it as adaptive linear element (Adaline)

Typically they use ← perceptron or LMS algorithm.



$$X_K = (x_0, x_1^K, x_2^K, \text{----} x_n^K)^T,$$

$$X_K \in \mathbb{R}^{n+1}$$

$$W_K = (w_0^K, w_1^K, w_2^K \text{----} w_n^K)^T$$

$$W_K \in \mathbb{R}^{n+1}$$

$K \equiv$ iteration number as these values are time dependent.

$X_K$: Pattern presented at $K^{th}$ iteration and $W_K$ is the neuronal weight vector.

The neuronal activation $\boxed{y_K = X_K^T W_K} \Longrightarrow$ Determins wheather Neuron $\boxed{\text{fires} - 1}$ or $\boxed{\text{not} - 0}$

Using this logic

Binary threshold neuronal Signal fn.

$$\delta(y_K) = \begin{cases} 1 & y_K > 0 \\ 0 & y_K < 0 \end{cases}$$

( $y_K = 0$ ; assumed as misclassification and weights are designed to avoid it.)

So, we have $C_0$ & $C_1$ (two classes) with $X_0$ $X_1$ = set of patterns of these classes. Hence Training set $\Longrightarrow X = X_0 \cup X_1$

We are looking to Compute weight vector $\boxed{W_S}$ s.t

$\forall X_K \in X_1$, $\delta(y_K) = 1$ and for all $X_K \in X_0$, $\delta(y_K) = 0$
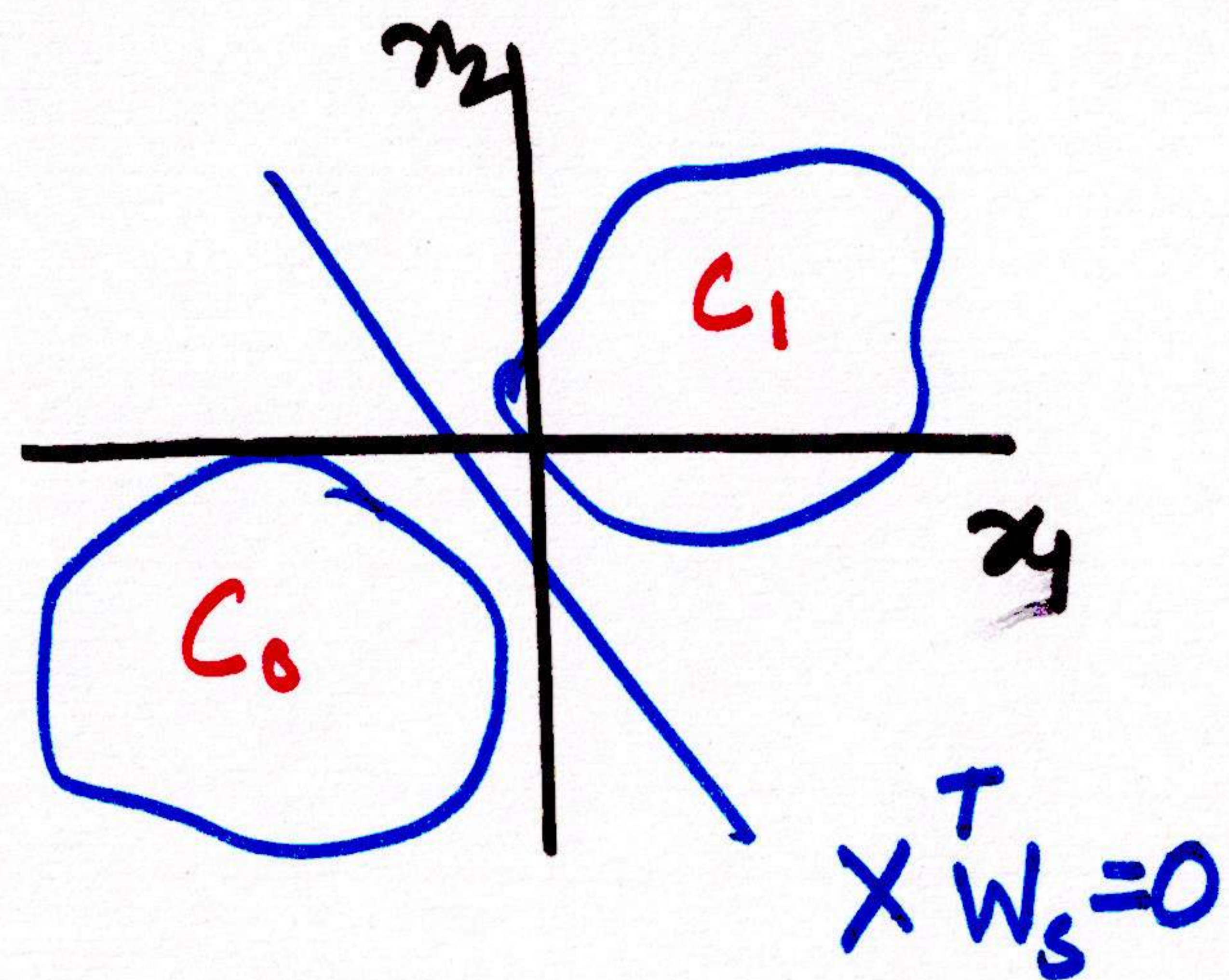
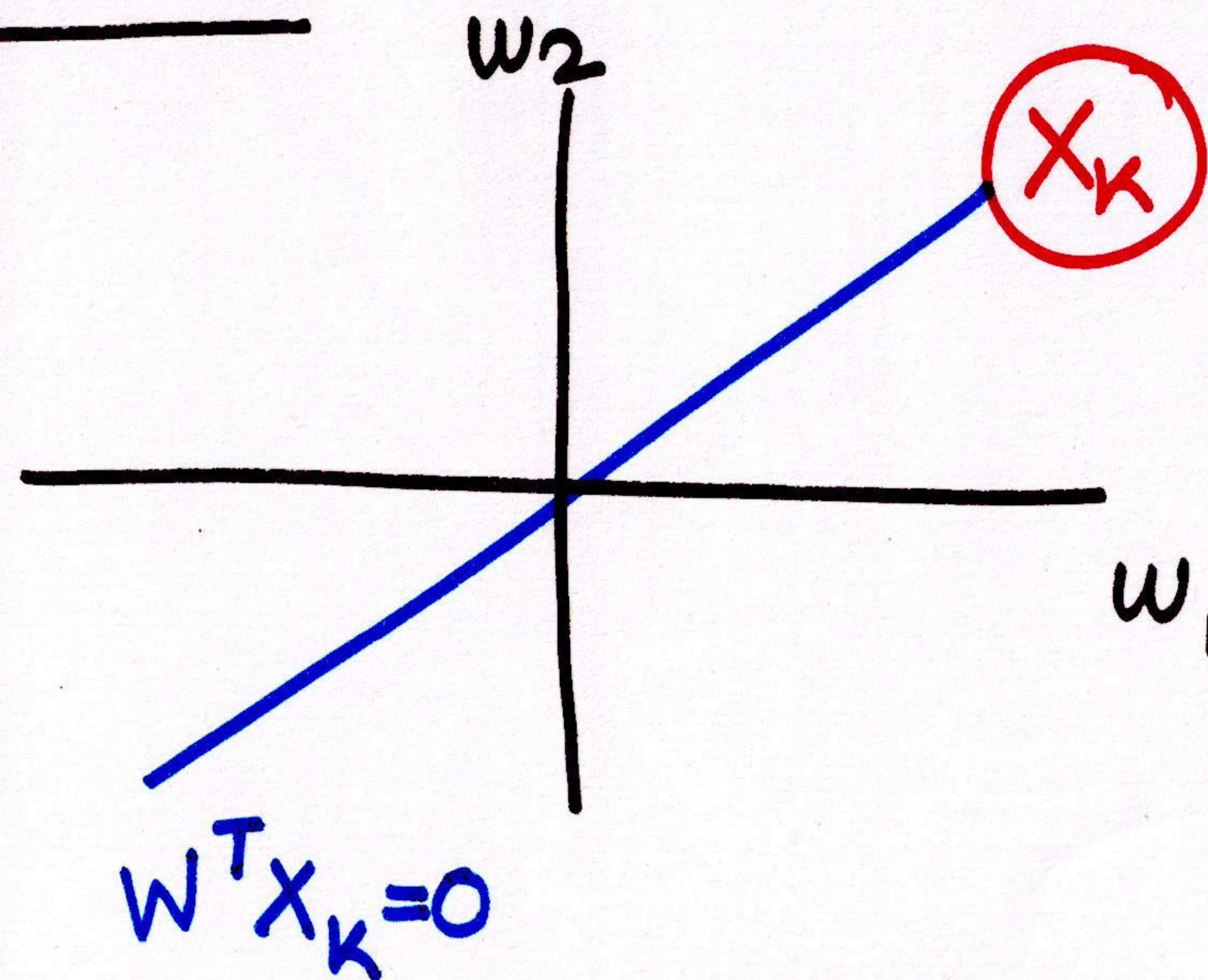$\Downarrow$ Innerproduct is +ve

$$\boxed{X_K^T W_S > 0}$$

$\Downarrow$ Innerproduct is -ve

$$\boxed{X_K^T W_S < 0}$$

# Pattern Space and Weight Space : '



$X^T W_s = 0$

## Pattern Space

$X^T W_s$ divides pattern space into two parts, one in which inner product is (+ve), other in which inner product is (-ve). In classes are linearly separable then $C_0$ & $C_1$ will be in 2 sides of this hyperplane.

$W^T X_k = 0$

## Weight space

$W^T X_k$ divides weight space into two parts. This hyperplane is the locus of all points (W) such that $\boxed{W^T X_k = 0}$.

for

For each (X_k) in pattern space there will be a corresponding h-plane in weight space, while for every point in weight space there is a corresponding h-plane in pattern space.

Consider 4 patterns divided into two pattern sets $X_1 = \{X_1, X_2\}$ ?
and $X_0 = \{X_3, X_4\}$. Using TLN classifier with neuronal signal ⓪
and neuronal signal ① for $X_1$. for $X_0$

(-ve inner product.)

Linear separability guarantees the existence of such a solution region.

(+ve inner product)



Solution region.

In order to design an automated weight update procedure that can search out a solution weight vector, starting with an arbitrary initial weight vector we have to:-

(*) Consider each pattern in turn to assess the correctness of the present classification.

(*) Subsequently adjust the weight vector to eliminate any classification error.

(*) As set of all solution vector forms a convex cone, the weight update procedure should terminate as soon as it penetrates the boundary of this cone.

# Perceptron Learning :

Our assumption is that data is of 2-classes and it is linearily seperable.

⊛ Starting with any initial (random) values of our weight vector $\bar{w}$ & $\bar{w}_0$.

⊛ Estimate the error wrt this $h$-plane. $(3)$.

⊛ Comeup with an update $\Delta w$, that can reduce this error.

This is an iterative procedure for error minimization.



- ve slope → weight has to be decreased

+ve slope → weight has to be increased.

Error Surface

Convex fn.

error

$w$

It looks for the direction that minimizes error.

$$\bar{w}(K+1) = \bar{w}(K) + \Delta\bar{w}(K)$$

Error → Mis classification

$$\left\{ \Delta\bar{w} \propto -\frac{\partial \varepsilon}{\partial w} \right\}$$

Classification is happening as

Classification $\begin{cases} \text{if } \bar{w}\bar{x} + w_0 \geqslant 0 \text{ then } \bar{x}_m \in C_1 \\ \\ \text{if } \bar{w}\bar{x} + w_0 < 0 \text{ then } \bar{x}_m \in C_2 \end{cases}$

$D_n$ : Have a set of misclassified Training examples.

$D_1$ and $D_2$ : Data of (+ve) and (-ve) class

$y_n$ : Class labels associated with $(x_n)$, $y_n \in \{+1, -1\}$

Augmented Vectors : $\bar{a} = \begin{bmatrix} w_0 \\ w_1 \\ | \\ w_d \end{bmatrix}$ and $z = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$

Class Assignment -:

$y_n = +1 \Rightarrow$ if $\quad w^T x_n + w_0 \geq 0 \Longrightarrow$ $\quad \overline{\bar{a}^T z_n \geq 0}$

$y_n = -1 \Rightarrow$ if $\quad w^T x_n + w_0 < 0 \Longrightarrow$ $\quad \overline{\bar{a}^T z_n < 0}$

$\Longrightarrow \quad y_n * (a^T z_n) > 0$ $\begin{cases} \to \text{TRUE for Correct classification} \\ \to \text{FALSE for Misclassification.} \end{cases}$

We have to modify our weights such that after each iteration $(D_m)$ reduces and we stops when it is empty.

$\Longrightarrow$ In order to show convergence one has to assume that linear sep. and weights are changed and error got reduced (atleast NOT increase) in any step.

# Perceptron Cost fn:

$$J_P = -\sum y_n (a^T z_n) \quad \text{(total error)}$$

$$z_m \in D_m \text{ (set of mis-classified samples)}$$

At any $k^{th}$ iteration:

$$\Delta \bar{a}(k) \propto -\frac{\partial J_P}{\partial w}$$

$\therefore$ $\Delta \bar{a}(k) = -\eta \frac{\partial J_P}{\partial a(k)} = -\eta \frac{\partial}{\partial a(k)} \left( -\sum_{x_m \in D_m(k)} y_n (\bar{a}(k))^T z_m \right)$

*n-plane parameter updation*

*Learning Rate*

$$= \eta \sum_{x_m \in D_m(k)} y_n \bar{z}_n$$

## Update Rule:

$$a(k+1) = a(k) + \Delta \bar{a}(k)$$

$\therefore$ $a(k+1) = a(k) + \eta \sum_{x_n \in D_m(k)} y_n z_n$ → Augmented misclassified example

*Learning rate*

*label*

Set of mis-classified examples at $k^{th}$ iteration.

**ALGO :** (1) Initilize augumented weight vector $a[0]$ for $0^{th}$ iteration. [12]

(2) At every (say $k^{th}$) iteration, Maintain a mis-classification set $D_m(k)$ for $\bar{a}(k)$ (n-plane).

(3) Update the h-plane parameter as:

$$a(k+1) = a(k) + \eta \sum_{z_n \in D_m(k)} y_n z_n$$

[ Batched Perceptron ]

h-plane updation only after full batch.

(4) Repeat Step 2 till $D_m(k)$ is empty.

# The perceptron learning requires only a single neuron and its Convergence proof states — that It will take only finite number of steps to Converge provided the data is linearly seperable. [PROOF in BOOK].

# Single sample perceptron :

It donot consider all misclassified examples, instead consider the current mis-classified example and update the h-plane.

## Instantaneous error →

$$J_P = -y_n (\bar{a}^T z_n)$$
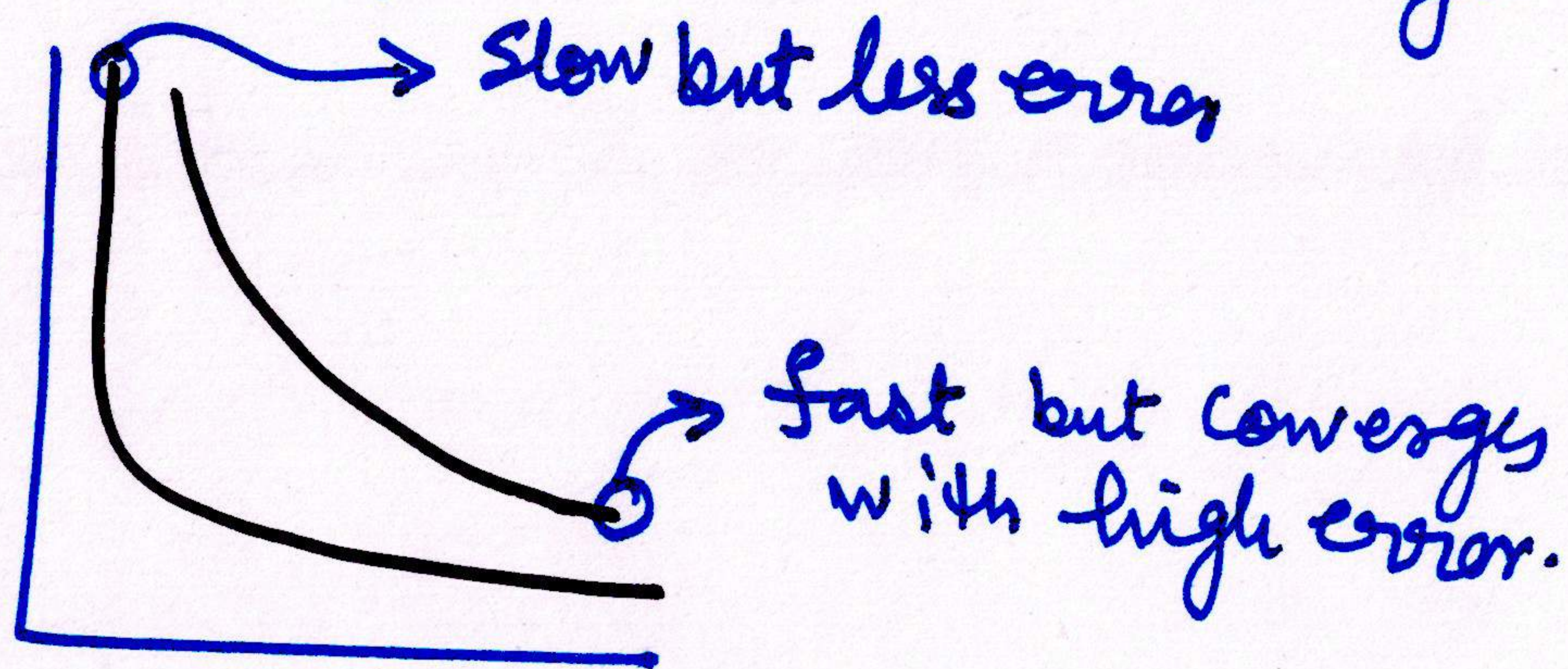
error wrt to a particular example.

1 - Initialize $a(0)$

2 - At any $K^{th}$ iteration : $a(K+1) = a(K) - \eta \dfrac{\partial J_P}{\partial a}$

$$a(K+1) = a(K) + \eta \, y_n z_n$$

3 - Repeat step 2 untill all $x_n$'s are correctly classified.

Batch update $\Rightarrow$ Slow Convergence | Single $\Rightarrow$ fast but may oscillate
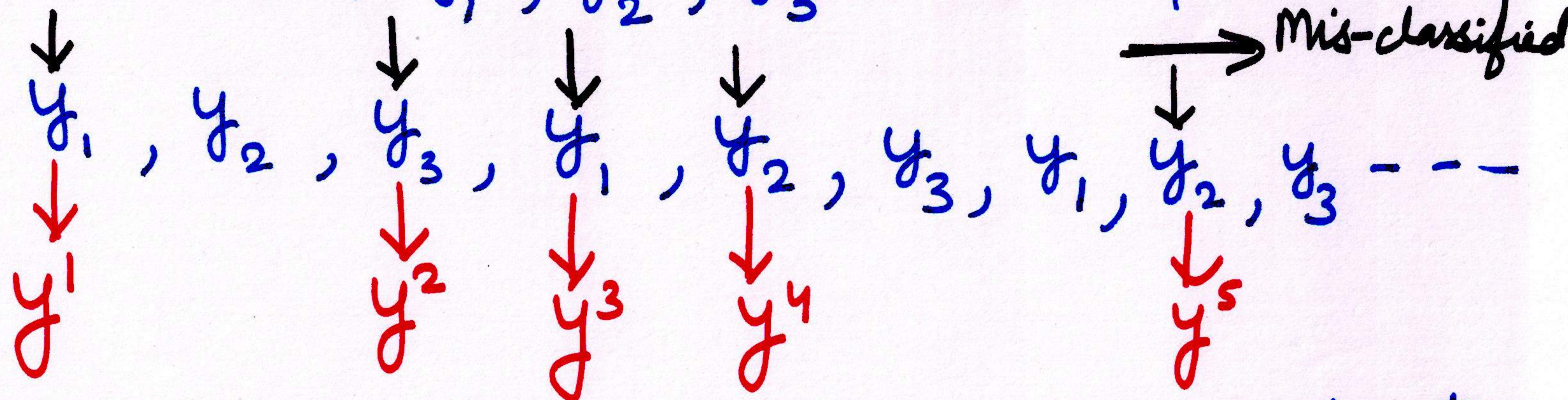
→ Slow but less error

→ fast but converges with high error.

We have seen fixed $\eta$.
One can start with high $\eta$ and then reduce $\eta$ as going on and fix it after few sets.

# Convergence Proof for Single-Sample Correction:

fixed increment ($\eta = 1$)

Assuming we have ③ samples, $y_1, y_2, y_3$ and are presented cyclically:

$$y_1, y_2, y_3, y_1, y_2, y_3, y_1, y_2, y_3 \,\text{---}$$

$$y^1 \qquad y^2 \quad y^3 \quad y^4 \qquad\qquad y^5$$

$\longrightarrow$ Mis-classified

Hence fixed-increment rule for generating a sequence of weight vectors

① $a(1)$ arbitrary ② $a(k+1) = a(k) + y^k$ $\quad k \geq 1$

Here $\boxed{a^t(k) y^k \leq 0}$ as mis-classified.

$\left( \begin{array}{c} \text{Update until} \\ \text{all samples are} \\ \text{correctly classified} \end{array} \right)$

Geometrically: ① $a(k)$ is mis-classified $y^k$ ; ② Hence $a(k)$ is not on the +ve side of $[a^t y \leq 0]$ h-plane.

③ Adding $\boxed{y^k}$ to $a(k)$ moves weight vector towards h-plane.

④ New inner product $\boxed{a^t(k+1) y^k}$ is larger than old inner product $a^t(k) y^k$ by amount $\| y^k \|^2$

Correction in good direction.

It terminates only when samples are linearly separable. [ we prove it indeed converge ]

If training samples are linearly separable, then sequence of weight vectors will terminate at a solution vector.

**Proof :** Let $\hat{a} \Rightarrow$ solution vector

(for sufficiently long) Solution vector

we try to show $\| a(K+1) - \hat{a} \| < \| a(K) - \hat{a} \|$ —①

$\circledast$ $\hat{a}^t y_i \geqslant 0, \forall i$ ⊖-2 $\circledast$ $\alpha \rightarrow$ +ve scale factor ③

$\therefore$ $a(K+1) - \alpha \hat{a} = (a(K) - \alpha \hat{a}) + y^K$ ④ ( from previous )

$\longrightarrow \| a(K+1) - \alpha \hat{a} \|^2 = \| a(K) - \alpha \hat{a} \|^2 + 2(a(K) - \alpha \hat{a})^t y^K + \| y^K \|^2$ ⑤

Using ②

$\| a(K+1) - \alpha \hat{a} \|^2 \leq \| a(K) - \alpha \hat{a} \|^2 - 2\alpha \hat{a}^t y^K + \| y^K \|^2$ ⑥

$\left. \begin{array}{l} \beta^2 = \max_i \| y_i \|^2 \\ \gamma = \min_i [\hat{a}^t y_i] > 0 \end{array} \right\}$ ⑦

if $\textcircled{\alpha}$ is very high it will dominate $\| y^K \|^2$. $\swarrow$ Strictly +ve

$$\therefore \ \|a(k+1) - \alpha\hat{a}\|^2 \leqslant \|a(k) - \alpha\hat{a}\|^2 - 2\alpha\gamma + \beta^2 \quad \text{⑧}$$

Choose: $\left[\alpha = \dfrac{\beta^2}{\gamma}\right]$

We get

$$\|a(k+1) - \alpha\hat{a}\|^2 \leqslant \|a(k) - \alpha\hat{a}\|^2 - \beta^2 \quad \text{⑨}$$

Hence the squared distance from $\boxed{a(k)}$ to $\boxed{\alpha\hat{a}}$ is reduced by at least $\boxed{\beta^2}$ at each correction $\text{⑩}$

After $\text{Ⓚ}$ corrections $\div$
Starting from $a(1)$

$$\|a(k+1) - \alpha\hat{a}\|^2 \leqslant \|a(1) - \alpha\hat{a}\|^2 - k\beta^2 \quad \text{⑪}$$

As above thing cannot be (-ve) $\rightarrow$ Sequence of correction must terminate $\text{⑫}$ after no more than $K_0$ corrections, where

$$\boxed{K_0 = \dfrac{\|a(1) - \alpha\hat{a}\|^2}{\beta^2}} \quad \text{⑬}$$