

## Lecture 2:

A linear classifier :

$$h(x; \theta, \theta_0) = \text{sign}(\theta^T x + \theta_0)$$

random

When we run our learning algorithm, we can see that while our training errors don't get more worse, the rate at which it improves is very slow and inefficient.  
 [Refer 8:13 to 18:08 Lecture - 2 YouTube]

## \* Perceptron Algorithm.

Perceptron ( $D_n; \tau$ )

Initialize  $\theta = [0 \ 0 \dots \ 0]^T$  [It has to have 'd' 0's]

Initialize  $\theta_0 = 0$

for  $t = 1$  to  $\tau$  // Represents an upper bound for time provided  
 changed = False.

for  $i = 1$  to  $n$  //  $n$  = size of data A. Pt. not on line and predict

(#i) if  $y^{(i)} (\theta^T x^{(i)} + \theta_0) \leq 0$  is wrong.

set  $\theta = \theta + y^{(i)} x^{(i)}$

B. Point is on line

Set  $\theta_0 = \theta_0 + y^{(i)}$ .

C. Initial Step.

changed = True.

if not changed:

break

Return  $\theta, \theta_0$

What does an update do?

$$y^{(i)} (\theta_{\text{updated}}^T x^{(i)} + \theta_0, \text{updated})$$

$$y^{(i)} ((\theta + (y^{(i)} x^{(i)})^T x^{(i)} + \theta_0 + y^{(i)}))$$

$$y^{(i)} ((\theta x^{(i)} + \theta_0)) + \underbrace{y^{(i)2} (x^{(i)T} x^{(i)} + 1)}_{\text{change}}$$

$$y^{(i)} (\theta x^{(i)} + \theta_0) + (||x^{(i)}||^2 + 1)$$

Might change the -ve value to +ve by adding

- The algorithm stops when it reaches 'D' - training error
- How does the classifier move when it has a misclassification?
- We're trying to get the angle right, so that we can correctly classify the particular data point

We hope:

$$|(\|x^{(i)}\|^2 + 1) > |y^{(i)}(\theta^T x^{(i)} + \theta_0)|$$

so that the entire value becomes more +ve and the if statement doesn't get executed.

It doesn't have to strictly decrease the error.

#### \* Classifier quality:

**Definition:** A training set  $D_n$  is linearly separable if there exist  $\theta, \theta_0$  such that, for every point index  $i \in \{1, \dots, n\}$ , we have

$$y^{(i)}(\theta^T x^{(i)} + \theta_0) > 0$$

The signed distance from a hyperplane defined by  $\theta, \theta_0$  to a point  $x^*$  is

$$\begin{aligned} & \text{projection of } x^* \text{ on } \theta - \text{signed distance of line to } x^* \\ &= \frac{\theta^T x^* - (-\theta_0)}{\|\theta\| \cdot \|\theta\|} = \frac{\theta^T x^* + \theta_0}{\|\theta\|} \end{aligned}$$

**Definition:** The margin of the labelled point  $(x^*, y^*)$  w.r.t. the hyperplane defined by  $\theta, \theta_0$  is:

$$y^* \left( \frac{\theta^T x^* + \theta_0}{\|\theta\|} \right)$$

**Definition:** The margin of the training set  $D_n$  w.r.t. the hyperplane defined by  $\theta, \theta_0$  is:

$$\min_{i \in \{1, \dots, n\}} y^{(i)} \left( \frac{\theta^T x^{(i)} + \theta_0}{\|\theta\|} \right)$$

\* If we get even one point wrong, the margin of  $D_n$  will be -ve, but if we get each one right, margin will be +ve.

#### \* Theorem: Perceptron Performance.

##### \* Assumptions:

A. Our hypothesis class = classifiers with separating hyperplanes that pass through the origin ( $\theta_0 = 0$ )

B. There exists  $\theta^*$  and  $\gamma$  such that  $\gamma > 0$  and, for every  $i \in \{1, \dots, n\}$ , we have

$$y^{(i)} \left( \frac{\theta^{*T} x^{(i)}}{\|\theta\|} \right) > \gamma$$

C. There exists  $R$  such that, for every  $i \in \{1, \dots, n\}$ , we have

$$\|x^{(i)}\| \leq R$$

**Conclusion:** Then the perceptron algo. will make at most  $(R/\gamma)^2$  updates to  $\theta$ . Once it goes through a pass of  $\theta$  without changes, the training error of its hypothesis will be zero.

\* why classifiers through the origin?

- if we're clever, we don't lose any flexibility

- classifier with offset.

$$x \in \mathbb{R}^d, \theta \in \mathbb{R}^d, \theta_0 \in \mathbb{R}$$

$$x: \theta^T x + \theta_0 = 0$$

- classifier w/out offset.

$$x_{\text{new}} \in \mathbb{R}^{d+1}, \theta_{\text{new}} \in \mathbb{R}^{d+1}$$

$$x_{\text{new}} = [x_1, x_2, \dots, x_d, 1], \theta_{\text{new}} = [\theta_1, \theta_2, \dots, \theta_d, \theta_0]$$

$$x_{\text{new}}, \theta: \theta^T x_{\text{new}} = 0$$

Problem: Typical real data sets aren't linearly separable

- Classification: Mapping to a discrete set.

- Regression: Mapping to continuous values

- Supervised Learning: Learn a mapping from features to labels

- Unsupervised Learning: No Labels; find patterns