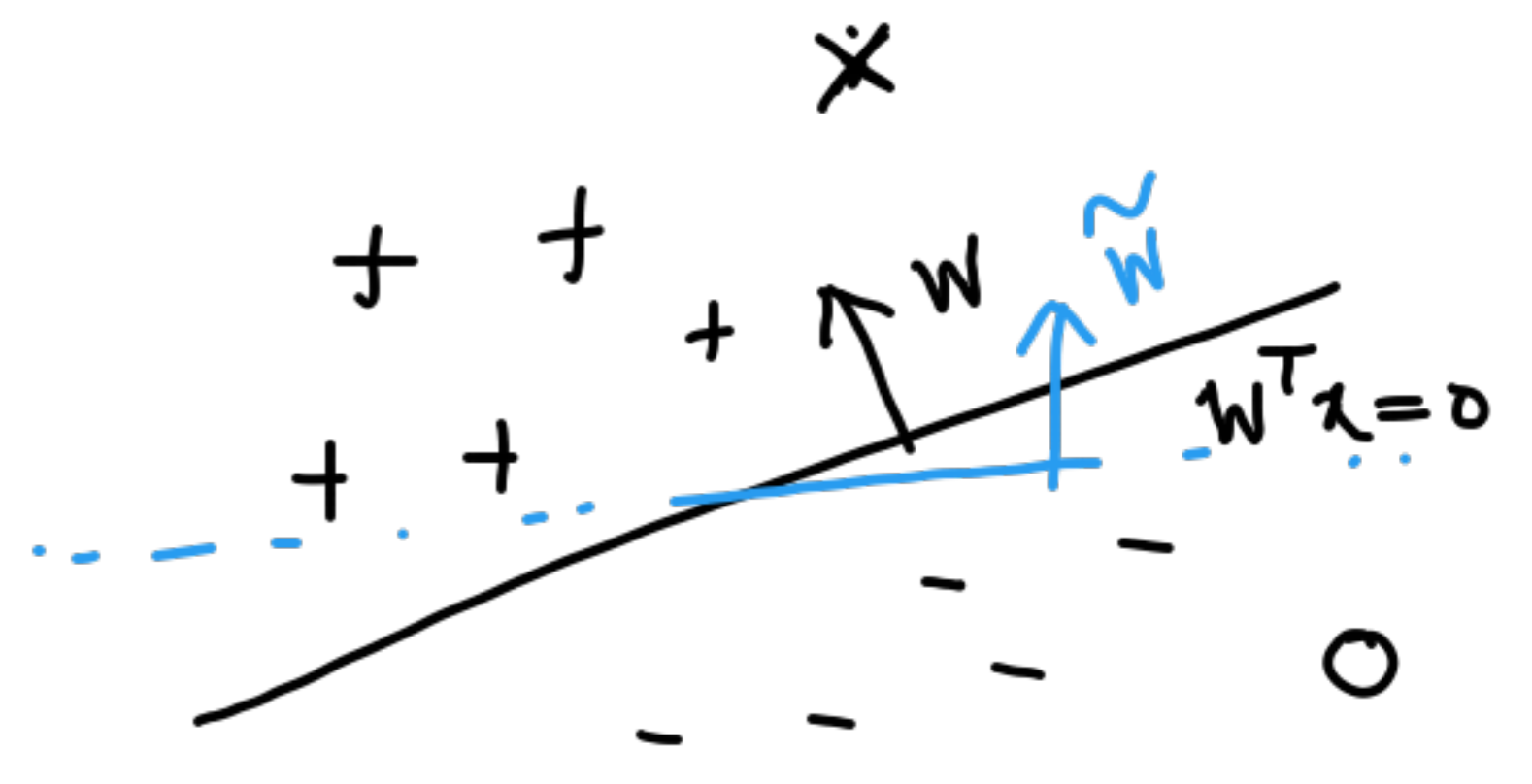# Lec 09: Perceptron : linear deterministic classification method

$$y_i = f(x_i) \in \{+1, -1\}$$



Goal : learn a weight vector w that <u>linearly</u> <u>separates</u> the training data points

$\begin{bmatrix} 1 \\ x_i \end{bmatrix} \; w^Tx$

$w^Tx_i + w_o \geqslant 0$ for $y_i = 1$
$\qquad\qquad < 0$ ow.

Remark: The decision boundary is not unique.

$$f_W(x) = sgn(W^T x)$$

Q: How does perceptron find W ?

1. Goes over the training examples $(x_i, y_i)$ one by one $\overset{W_0}{}$

2. Check if the current classifier $W_t$, i.e.
$$y_i = sgn(W_t^T x_i)$$

3. If correct — no update

4. If not — correct $W_{t+1}$
$$W_{t+1} \leftarrow W_t + y_i \, x_i$$
$$\uparrow \qquad = \qquad \uparrow_{\text{vector (d+1)}}$$

5. STOP if no update for a certain number of iteration.

# Algorithm: Perceptron

- Initialize $W_0$

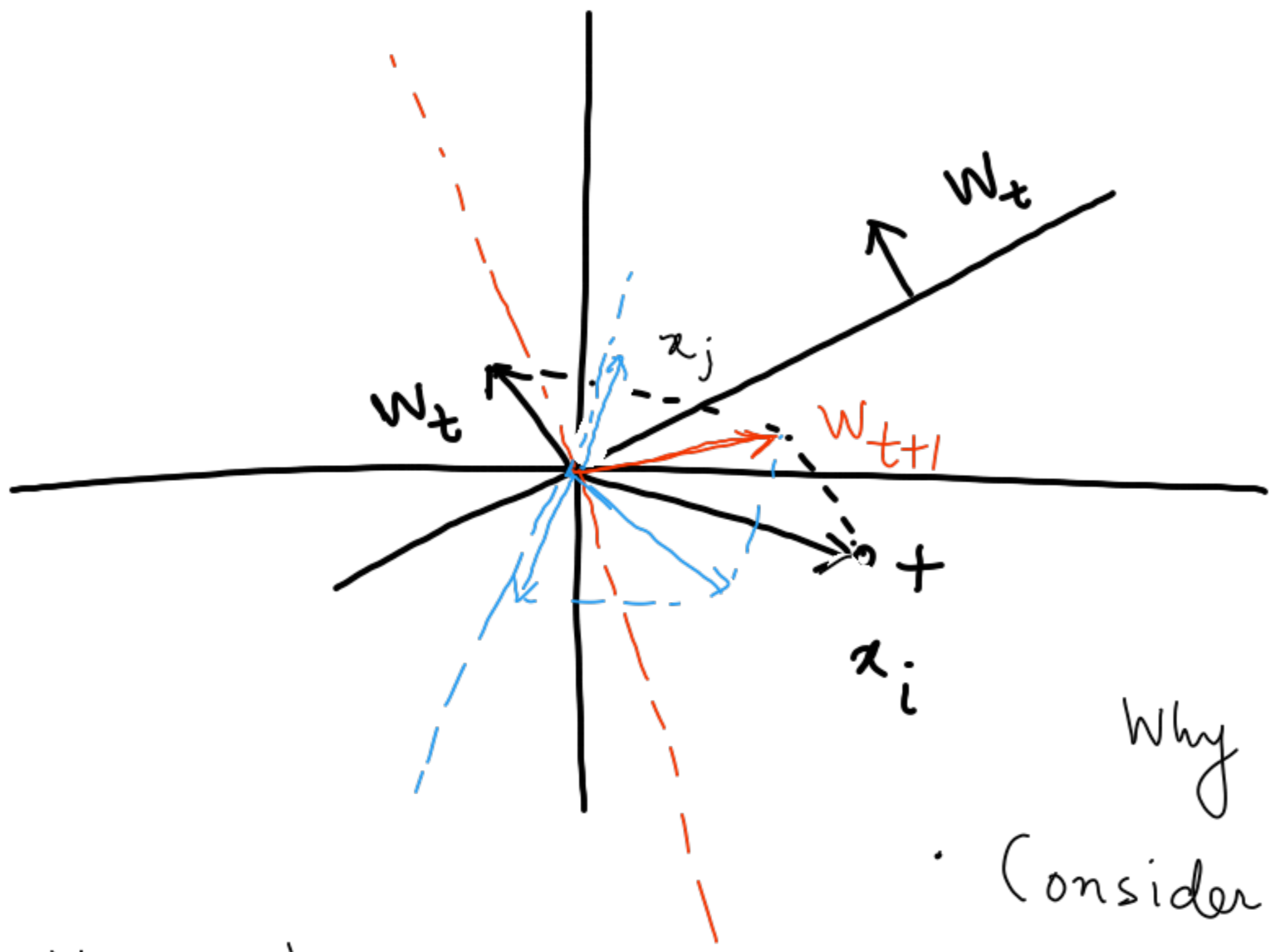- for $t = 0, 1, 2, \ldots, \text{max Rounds}:$         $y_i = \text{sgn}(W^T x_i)$

  Randomly choose a training example $(x_i, y_i)$

  if $y_i(W_t^T x_i) < 0$, then

  $$W_{t+1} \leftarrow W_t + y_i x_i$$

- STOP as before.

$$W_{t+1} \leftarrow W_t + \frac{1}{2}\left(y_i - \text{sgn}(W_t^T x_i)\right)x_i$$

$$W_{t+1} \leftarrow W_t + x_i$$

$$W_{t+2} \leftarrow W_{t+1} - x_j$$

Why is perceptron doing a meaningful update?

· Consider a misclassified example $(x_i, y_i)$

i.e. $\text{sgn}\left(W_{old}^T x_i\right) \neq y_i$

$$y_i\left(W_{new}^T x_i\right) = y_i\left(W_{old} + y_i x_i\right)^T x_i = y_i W_{old}^T x_i + \underline{\|x_i\|^2}$$

$$> y_i\left(W_{old}^T x_i\right)$$

$$W_{new} = W_{old} + y_i x_i$$

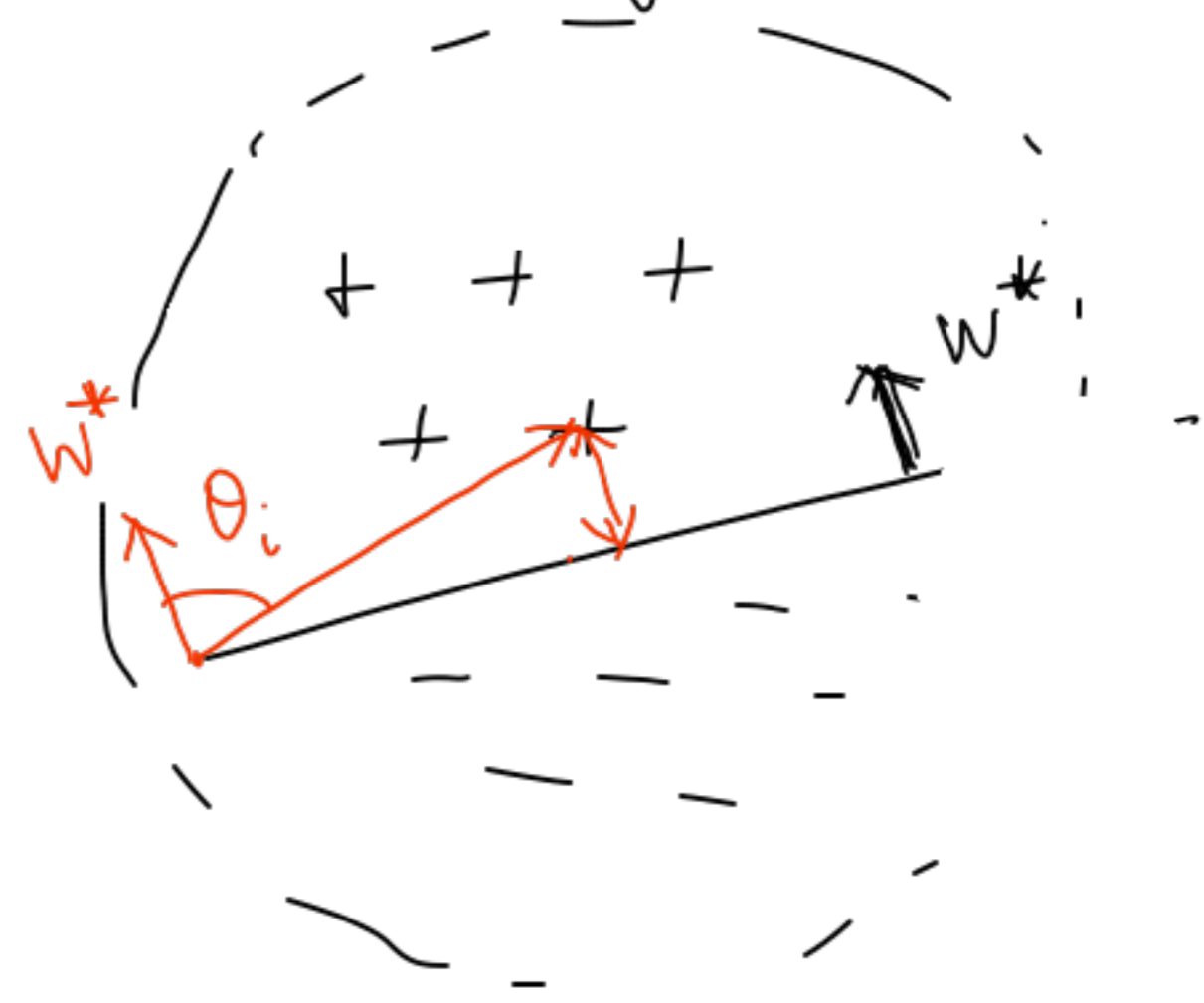Note: still not guaranteed to be correctly classified.

Summary : • Perceptron is a mistake-driven online learning algo.

• Guaranteed to converge for linearly separable training examples.

If the data is linearly separable,

$$\exists \, w^* \quad s.t. \quad y_i(w^{*T} x_i) \geq 0 \quad \forall \, i = 1, \ldots, n.$$

Assume $\|w^*\| = 1$, $\|x_i\| \leq 1$

Define, <span style="color:red">margin of separation</span> $\gamma = \min_i |w^{*T} x_i| = \|w^*\| \|x_i\| \cos \theta_i$

Theorem: If $\exists$ a unit vector $w^*$ s.t. $y_i \, w^{*T} x_i \geq \gamma \quad \forall (x_i, y_i) \in D$

Then the # of weight updates by perceptron is at most $\frac{1}{\gamma^2}$.

**Proof:** track two quantities ① $W_t^T w^*$ , ② $|W_t|_2^2$

① Claim: $W_t^T w^*$ on every update increases by at least $\gamma$

$$W_{t+1}^T w^* = (W_t + y_i x_i)^T w^*$$

$$= W_t^T w^* + \underbrace{y_i(w^{*T} x_i)}_{\geq \gamma} \geq W_t^T w^* + \gamma$$

② Claim: $\|W_t\|^2$ increases by at most 1.

$$\|w_{t+1}\|^2 = (W_t + y_i x_i)^T(W_t + y_i x_i) = \|w_t\|^2 + \underbrace{2 y_i w_t^T x_i}_{<0} + \underbrace{\|x_i\|^2}_{\leq 1} < \|w_t\|^2 + 1$$

Say $w_0 = 0$. After $k$ updates $\left.\begin{array}{c} W_{k+1}^T w^* \geq k\gamma \\ \|w_{k+1}\|^2 < k \end{array}\right\}$

$$\sqrt{k} > \|W_{k+1}\| \geq W_{k+1}^T w^* \geq k\gamma \implies k < \frac{1}{\gamma^2}$$

$$W_{k+1}^T w^* = \|W_{k+1}\| \underbrace{\|w^*\|}_{=1} \underbrace{\cos\theta}_{\leq 1}$$

finite number of mistakes if data is linearly separable.

Limitations:

① Not giving a rate of convergence

② The # of iterations can be large if $\gamma$ is small

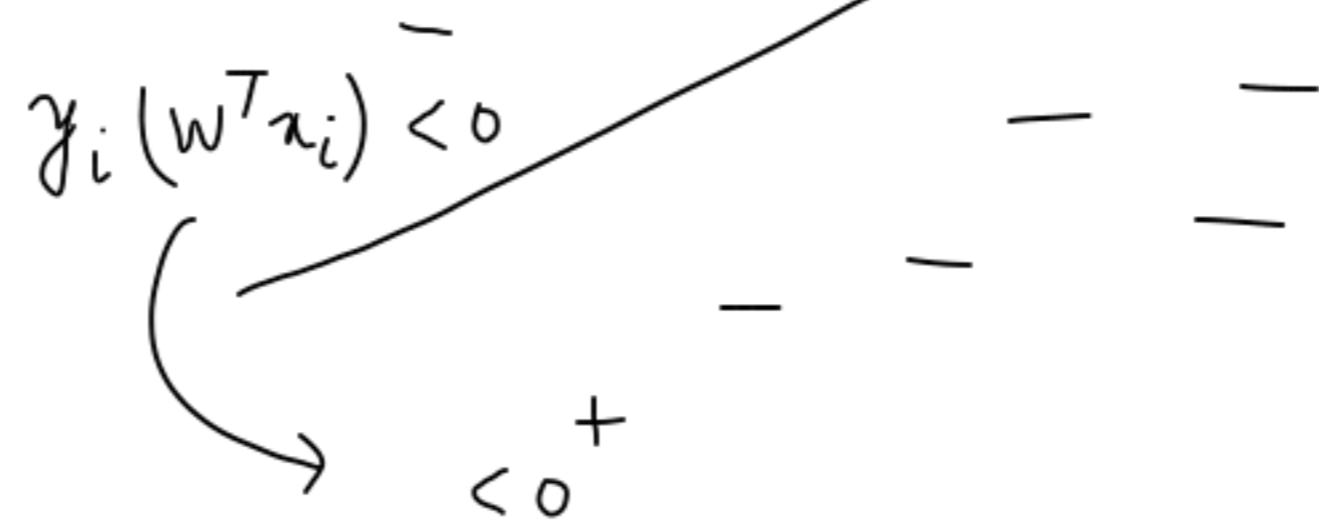③ May not converge if points are not linearly separable.

Find one such example of non convergence (HW)

$L + +$    $= -$

$= -$    $+ +$

The loss function view of perceptron

$y_i \cdot w^T x_i \ll 0$

maximize $y_i (w^T x_i)$

$y_i (w^T x_i) < 0$

$< 0$

$\Rightarrow \min_{w} \sum_{i \in \text{misclassified examples}} \left( - y_i (w^T x_i) \right)$

for any $i$

$y_i w^T x_i \geqslant 0 \rightarrow \text{loss} = 0$

$L_i (w, D) = \max \{ 0 , - y_i w^T x_i \}$

$< 0 \rightarrow \text{loss} = - y_i w^T x_i$

Hinge loss

$L(w, D) = \sum L_i (w, D)$

$y \, w^T x \longrightarrow$

Apply SGD :  $\longrightarrow$ randomly pick $i$ , compute $\nabla_w L_i (w, D)$

$$W_{t+1} \longleftarrow W_t - \nabla_w L_i (w, D)$$

$$-y_i w^T x_i$$

$$-y_i x_i$$

Hinge loss with SGD

is the perceptron algo. $\qquad W_t + y_i x_i$

# Decision Trees

| Fuel efficiency | cyl | disp | Origin | Year |
|---|---|---|---|---|
| good | 3 | low | | |
| | 4 | med | | |
| bad | 5 | | | |
| | 6 | high | | |

decision
tree

$x$