

LDA discussions:

① As a dimensionality reduction algorithm,

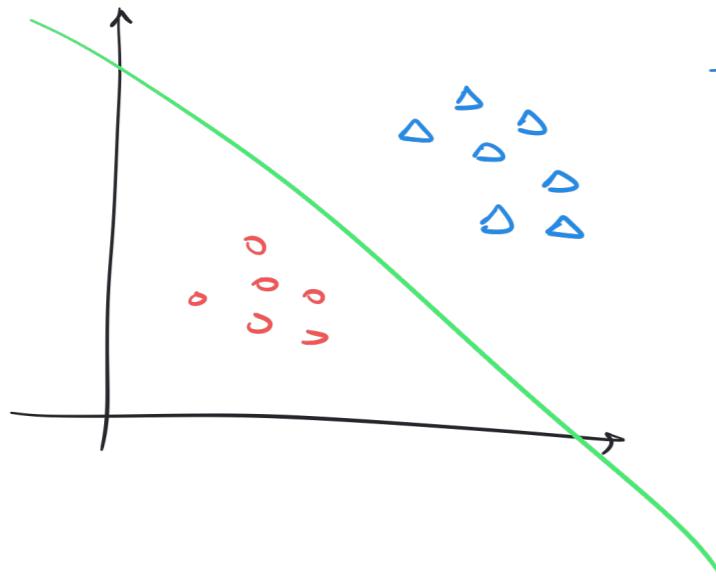
LDA projects up to

$$\min(C-1, D-1)$$

where $C \equiv \# \text{ classes}$

$D \equiv \text{dimensionality}$

② LDA assumes each class is Gaussian-distributed.



Least Squares
Regression

$$y = w^T x + w_0$$

Dataset: $\{(x_i, t_i)\}_{i=1}^N$, $x_i \in \mathbb{R}^d$, $t_i \in \mathbb{N}$

Mapper: $y = w^T x + w_0$

Objective: $J(w, w_0) = \sum_{i=1}^N (t_i - y_i)^2 = \sum_{i=1}^N \epsilon_i^2$

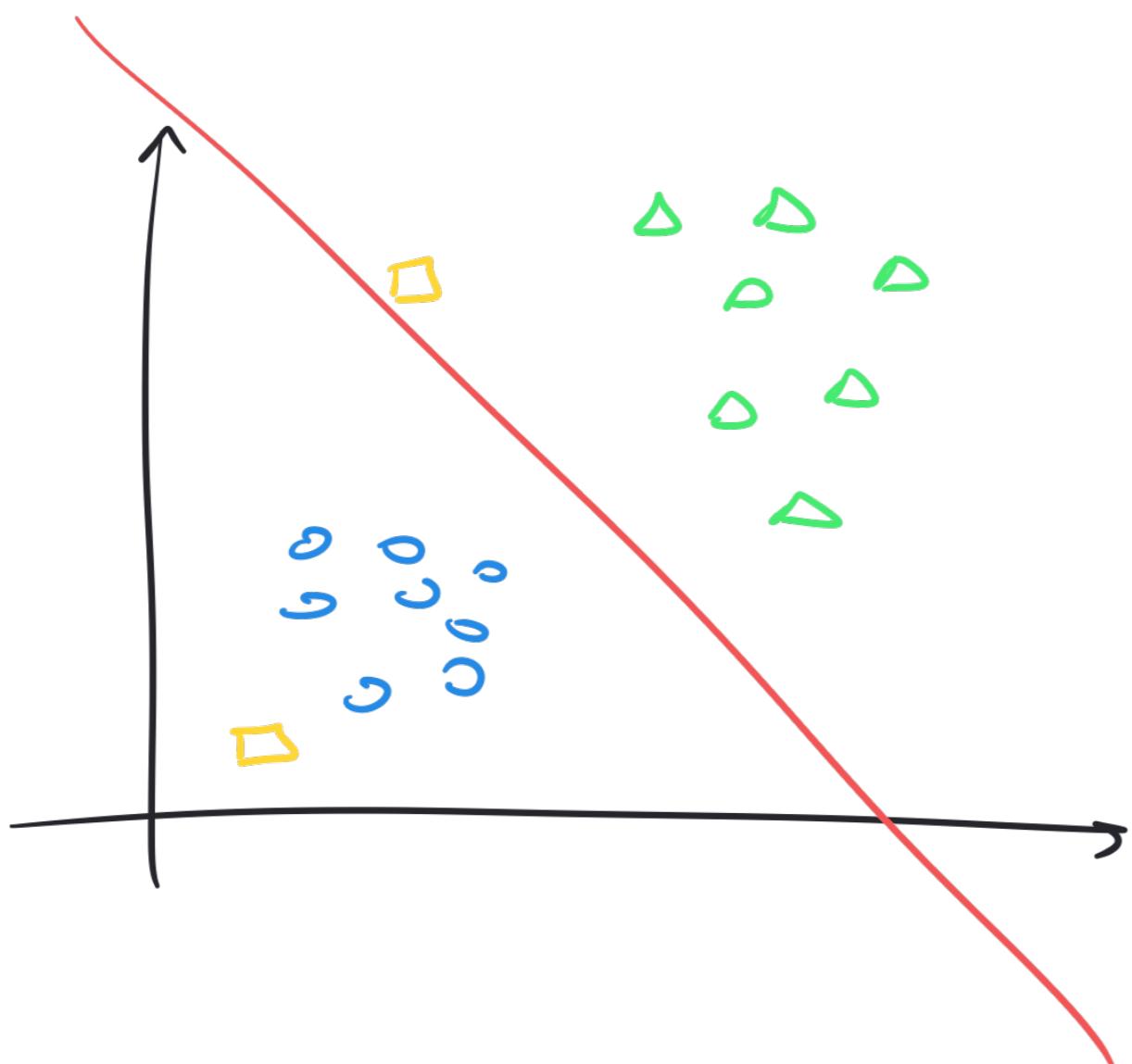
② the variance of estimator,

$E[(y - E[y])^2]$, is dependent on its mean, $E[y]$.

∴ Avoid using Least Squares as a classifier.

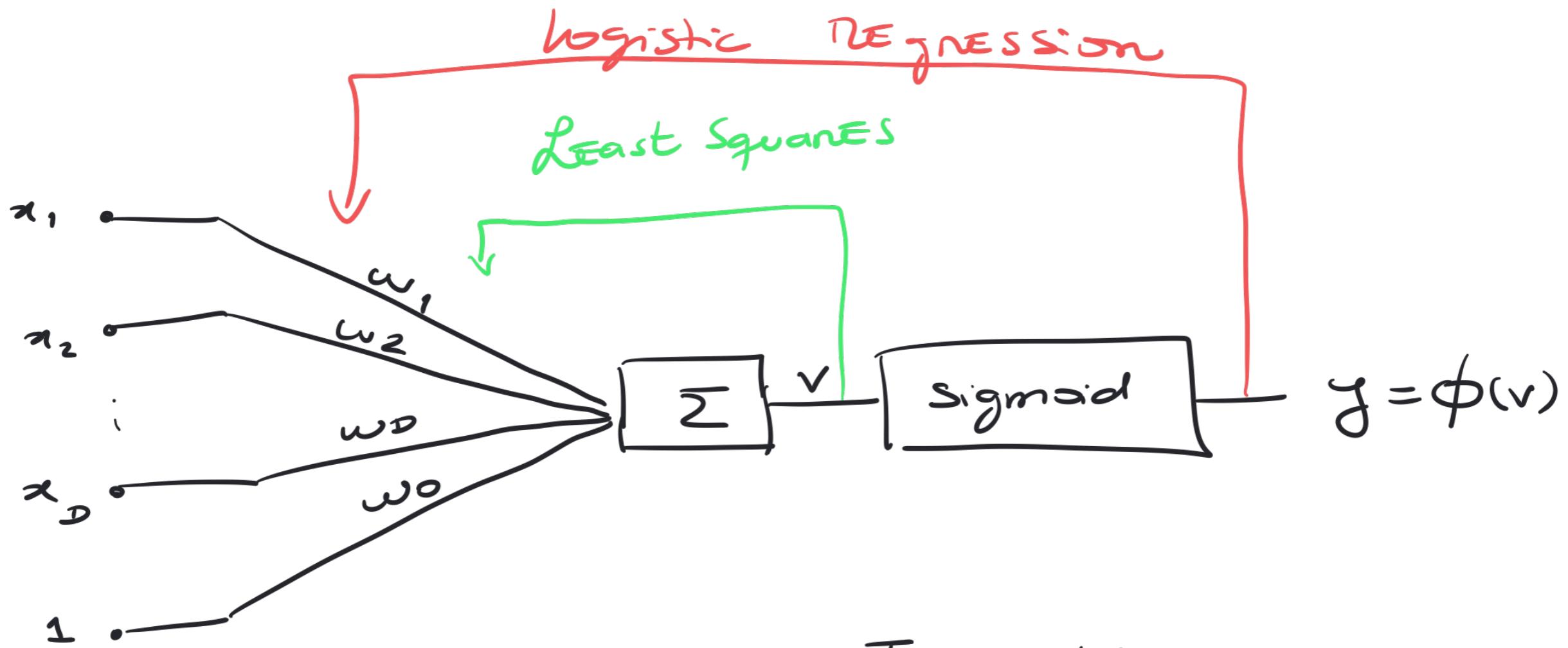
Logistic REGression

despite its name, it is a classifier.



The idea is to find the discriminant function such that we minimize uncertainty!

$$y = \phi(\omega^T x + w_0), \quad \phi(x) \equiv \text{logistic fct.}$$



where $v = \omega^T x + w_0$

and $\phi(x) = \frac{1}{1 + e^{-x}} \equiv \text{logistic fct}$

$y = \phi(v) \equiv \text{probabilities} \in [0, 1]$

Assume we only have 2 classes,

C_0 and $\overline{C_0}$:

$$\text{logit} \left(\underbrace{P(C_0|x)}_{\text{posterior prob.}} \right) = \ln \left(\frac{P(C_0|x)}{P(\overline{C_0}|x)} \right)$$

of C_0 for sample x

$$= \ln \left(\frac{P(C_0|x)}{1 - P(C_0|x)} \right) = \omega^T x + b$$

$$\Rightarrow \frac{P(C_0|x)}{1 - P(C_0|x)} = \exp(\omega^T x + b)$$

Applying exp.
fct

$$\Leftrightarrow P(C_0|x) = (1 - P(C_0|x)) \cdot \exp(\omega^T x + b)$$

$$P(C_0|x) = \frac{\exp(\underbrace{w^T x + b}_{=v})}{1 + \exp(\underbrace{w^T x + b}_{=v})}$$

$$= \frac{\exp(v)}{1 + \exp(v)}$$

$$= \frac{1}{1 + \exp(-v)} \equiv \text{logistic/ sigmoid function}$$

where $v = w^T x + b$

If $P(C_0|x) \underset{x \in C_0}{\underset{x \in \bar{C}_0}{\begin{matrix} \geq \\ < \end{matrix}}} P(\bar{C}_0|x)$

OBJECTIVE function for logistic Regression

TRAINING DATA : $\{(x_i, t_i)\}_{i=1}^N$, $t_i \in \{0, 1\}$

MAPPER : $y = \phi(w^T x + w_0)$ = probabilities

OBSERVED DATA LIKELIHOOD :

$$\begin{aligned}\hat{\mathcal{L}} &= P(y_1, y_2, \dots, y_N | x, w) \\ &= \prod_{i=1}^N P(y_i | x, w)\end{aligned}$$

$$\begin{array}{l} 1 - c_0 \\ 0 - \bar{c}_0 \end{array}$$

Example:

$$x = \{x_1, x_2, x_3\}$$

$$t = \{0, \underline{1}, \underline{1}\}$$

$$y_{w_1} = \{0.3, 0.8, 0.9\}$$

$$y_{w_2} = \{0.4, 0.6, 0.6\}$$

STEP 1

$$\mathcal{L}^o = \prod_{i=1}^N P(y_i | x_i, \omega)$$

$$= \prod_{i=1}^N \phi(v_i)^{t_i} \cdot (1 - \phi(v_i))^{1-t_i}$$

Example from previous

$$v_i = \omega^T \cdot x_i + \omega_0$$

$$\phi(v_i) = \frac{1}{1 + e^{-v_i}}$$

page:

$$\mathcal{L}_{w_1}^o = 0.8 \times 0.9 \times (1-0.3)$$

STEP 2

TAKE the OBSERVED log - likelihood:

$$\mathcal{L} = \ln(\mathcal{Z}^o)$$

$$= \ln \left(\prod_{i=1}^n \phi(v_i)^{t_i} \cdot (1-\phi(v_i))^{1-t_i} \right)$$

$$= \sum_{i=1}^n [t_i \cdot \ln(\phi(v_i)) + (1-t_i) \cdot \ln(1-\phi(v_i))]$$

Goal: $\arg \max_{\{\omega, \omega_0\}} \mathcal{L} = \arg \min_{\{\omega, \omega_0\}} -\mathcal{L}$

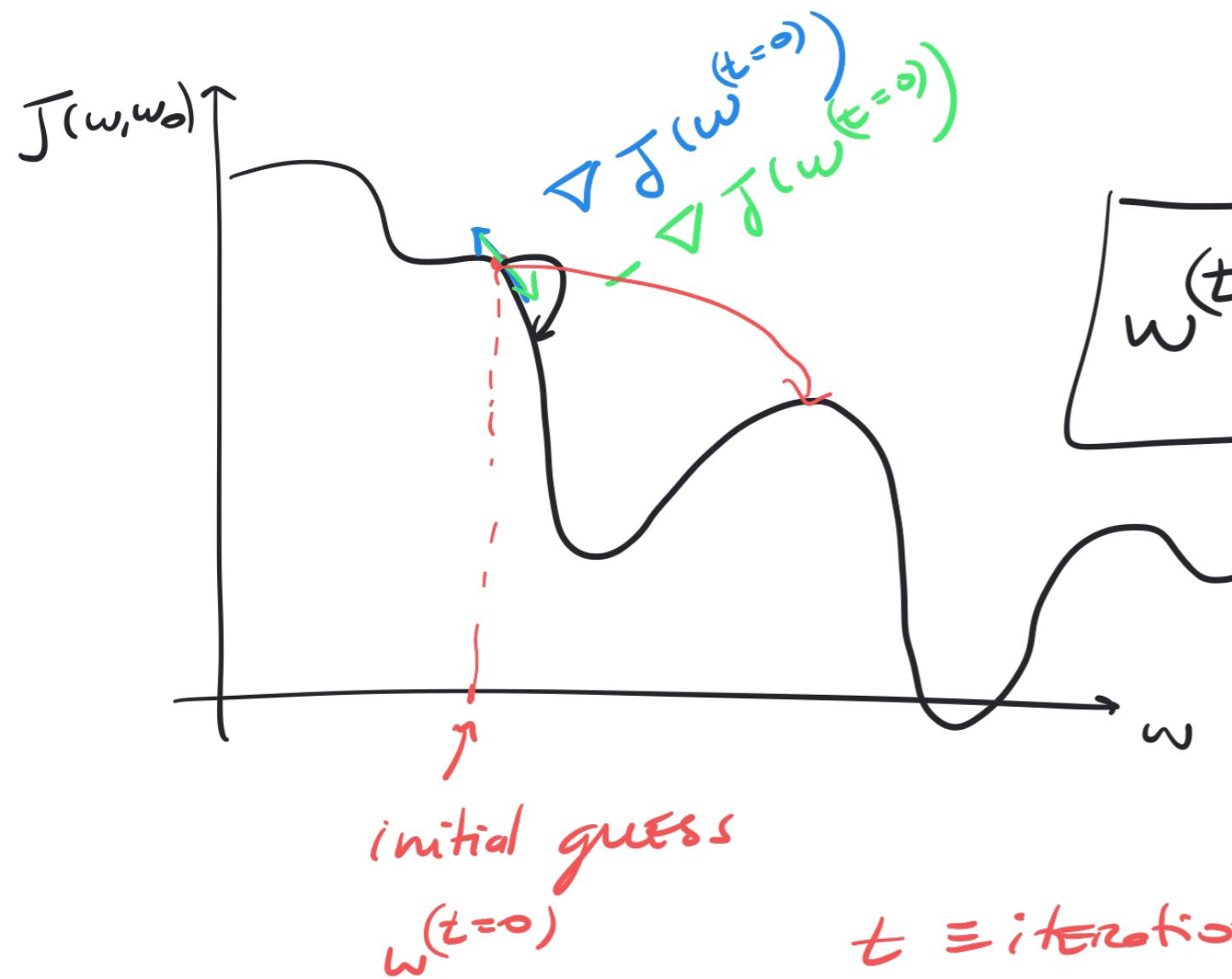
As a minimization problem ,

we optimize the

Cross - Entropy:

$$J(w, w_0) = \sum_{i=1}^n [-t_i \cdot \ln(\phi(v_i)) - (1-t_i) \cdot \ln(1-\phi(v_i))]$$

this is a non-convex function , so we
use search algorithm to find
a solution \leftarrow Gradient descent



$\gamma \equiv \text{LEARNING RATE ("ETA")}$

$$w^{(t+1)} \leftarrow w^{(t)} - \gamma \nabla J(w^{(t)})$$

γ is a hyperparameter

STABILITY RANGE

$$0 < \gamma < 1$$

Solution converges to a local minima
not necessarily global Δc of initialization.

$$J(\omega, \omega_0) = \sum_{i=1}^n [-t_i \cdot \ln(\phi(v_i)) - (1-t_i) \cdot \ln(1-\phi(v_i))]$$

where $v_i = \omega^T x_i + \omega_0$

$$\textcircled{1} \quad \frac{\partial J(\omega, \omega_0)}{\partial \omega} = 0$$

$$\Leftrightarrow \sum_{i=1}^n \left[-t_i \cdot \underbrace{\frac{\partial \ln(\phi(v_i))}{\partial \phi(v_i)}}_{= \frac{1}{\phi(v_i)}} \cdot \underbrace{\frac{\partial \phi(v_i)}{\partial v_i}}_{= \phi'(v_i)} \cdot \underbrace{\frac{\partial v_i}{\partial \omega}}_{= x_i} \right.$$

$$\left. - (1-t_i) \cdot \underbrace{\frac{\partial \ln(1-\phi(v_i))}{\partial \phi(v_i)}}_{= \frac{-1}{1-\phi(v_i)}} \cdot \underbrace{\frac{\partial \phi(v_i)}{\partial v_i}}_{= \phi'(v_i)} \cdot \underbrace{\frac{\partial v_i}{\partial \omega}}_{= x_i} \right]$$

$$\frac{\partial \bar{J}}{\partial w} = \sum_{i=1}^n \left[\frac{t_i}{\phi(v_i)} - \frac{1-t_i}{1-\phi(v_i)} \right] \cdot \phi'(v_i) \cdot x_i$$

where $\phi'(x) = \phi(x) \cdot (1-\phi(x))$

reducing to same denominator:

$$\frac{\partial \bar{J}}{\partial w} = \sum_{i=1}^n (t_i - \phi(v_i)) \cdot x_i$$

Similarly,

$$\frac{\partial \bar{J}}{\partial w_0} = \sum_{i=1}^n (t_i - \phi(v_i))$$

Pseudo - code:

Input: DATA \mathcal{D} , target t , learning rate γ .

$t = 0$ (iteration)

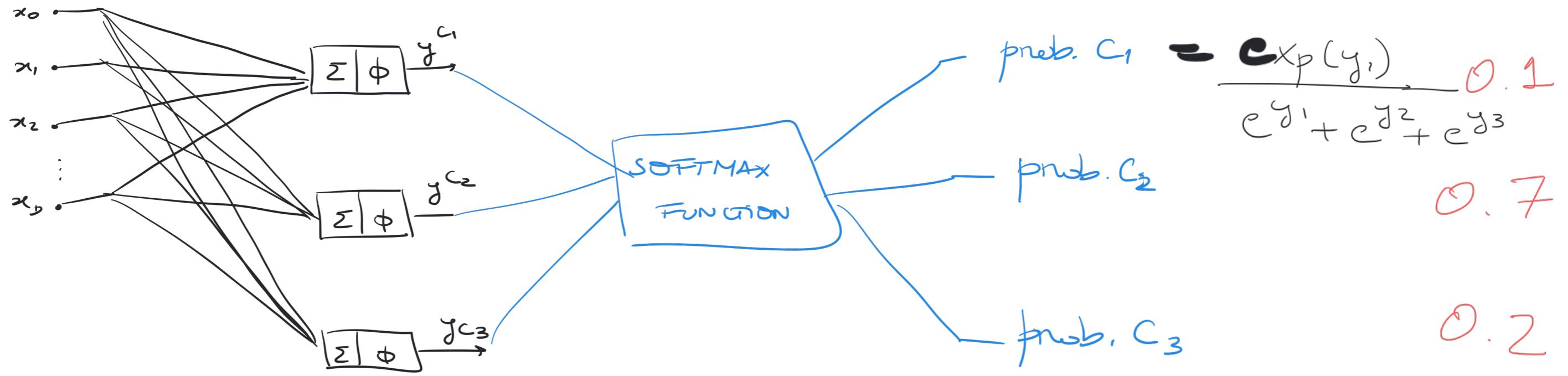
① Initialize parameters $w^{(t)}, w_0^{(t)}$

② Update the parameters

$$w^{(t+1)} \leftarrow w^{(t)} - \gamma \cdot \frac{\partial J(w, w_0)}{\partial w^{(t)}}$$

$$w_0^{(t+1)} \leftarrow w_0^{(t)} - \gamma \cdot \frac{\partial J(w, w_0)}{\partial w_0^{(t)}}$$

③ $t < t + 1$. Iterate until convergence criterion is met.



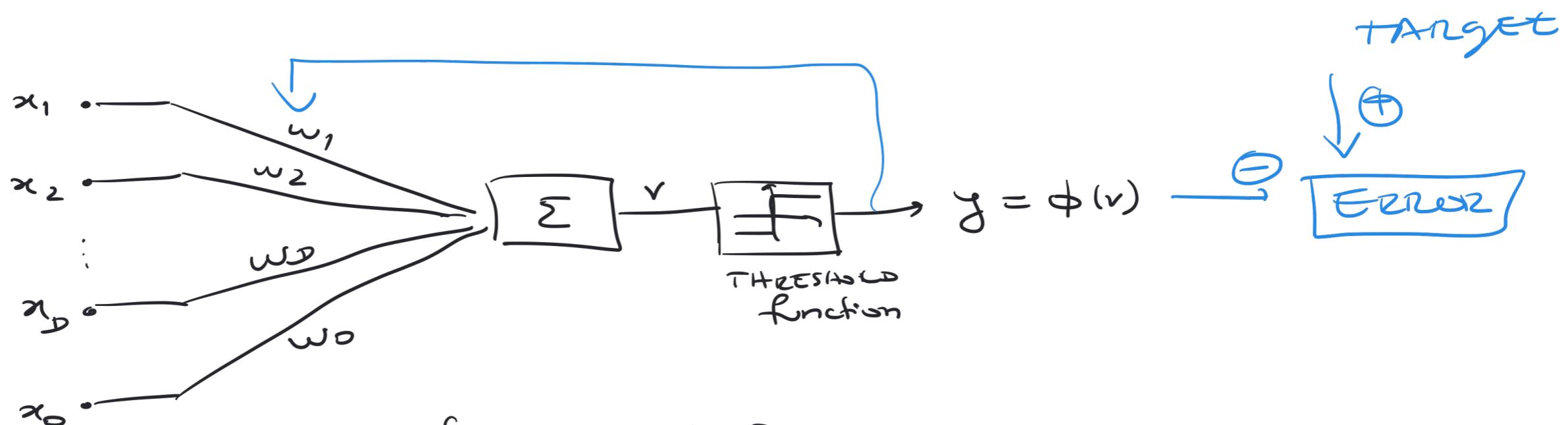
SOFTMAX function :

$$\frac{\exp(y_j(\alpha))}{\sum_j \exp(y_j(\alpha))} = p_j \quad \sum = 1$$

The Perceptron, 1957, Rosenblatt

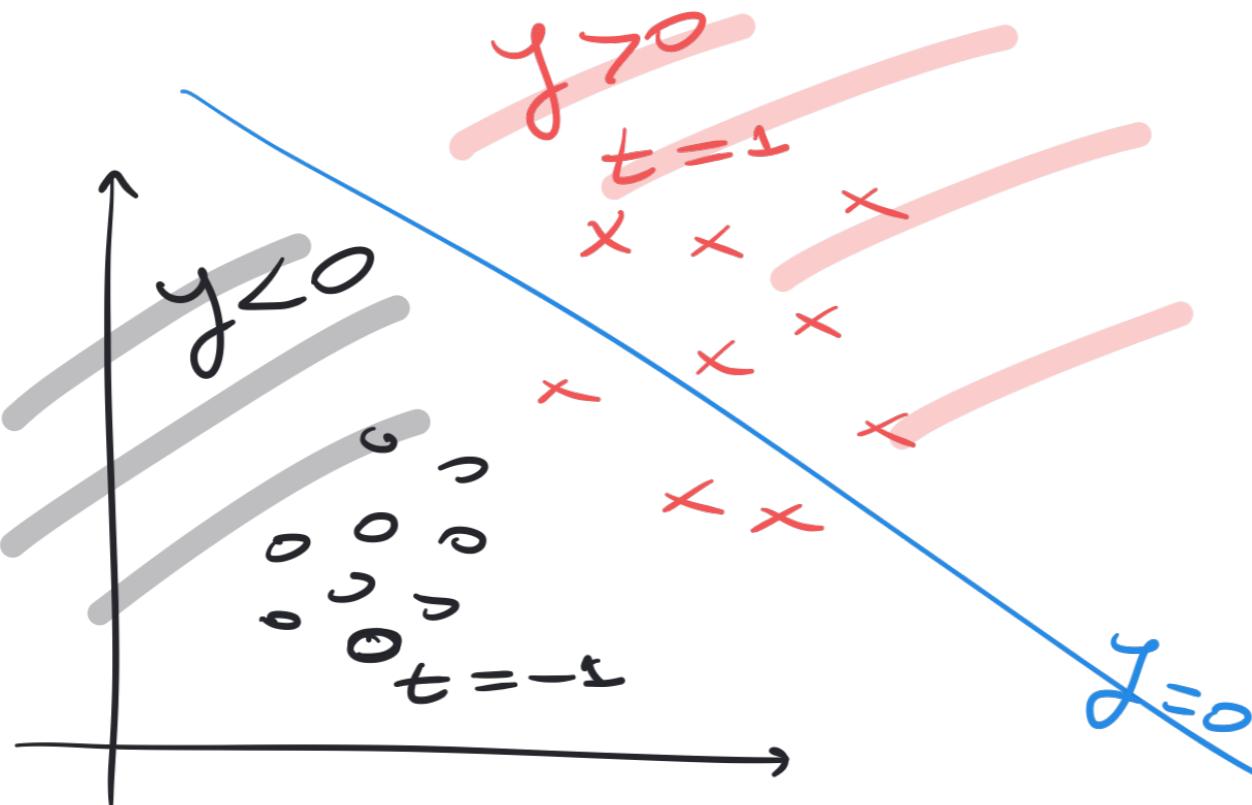
Linear binary classifier

↳ discriminative



$$\phi(x) = \begin{cases} 1, & x \geq 0 \\ -1, & x < 0 \end{cases}$$

$$y = \phi(\omega^T x + w_0)$$



\circ - class -1

\times - class 1

initial discriminant function

NOTE that ALL
MISCLASSIFIED points satisfy:

$$t_n \cdot y(x_n) < 0$$

and all correctly classified points

$$t_n \cdot y(x_n) > 0$$

Objective function

uses only misclassified points.

and it iterates through new values
for w and w_0 until no misclassified
points exist.

$$E_p(w, w_0) = - \sum_{n \in M} t_n \cdot (w^T \cdot x_n + w_0)$$

$M \equiv$ set of misclassified points

$$\arg \min_{\{w, w_0\}} E_p(w, w_0)$$