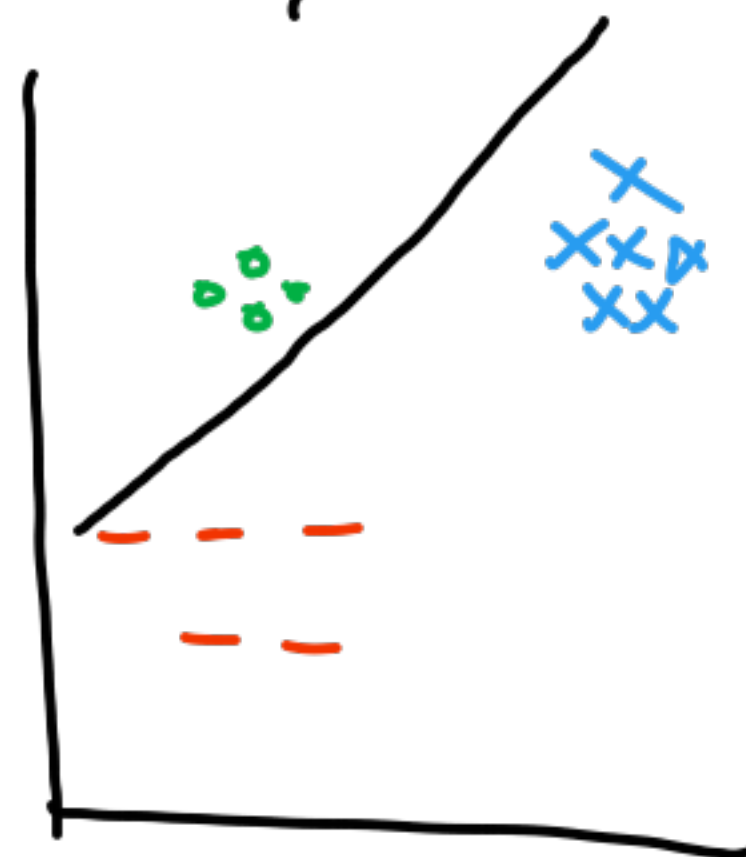
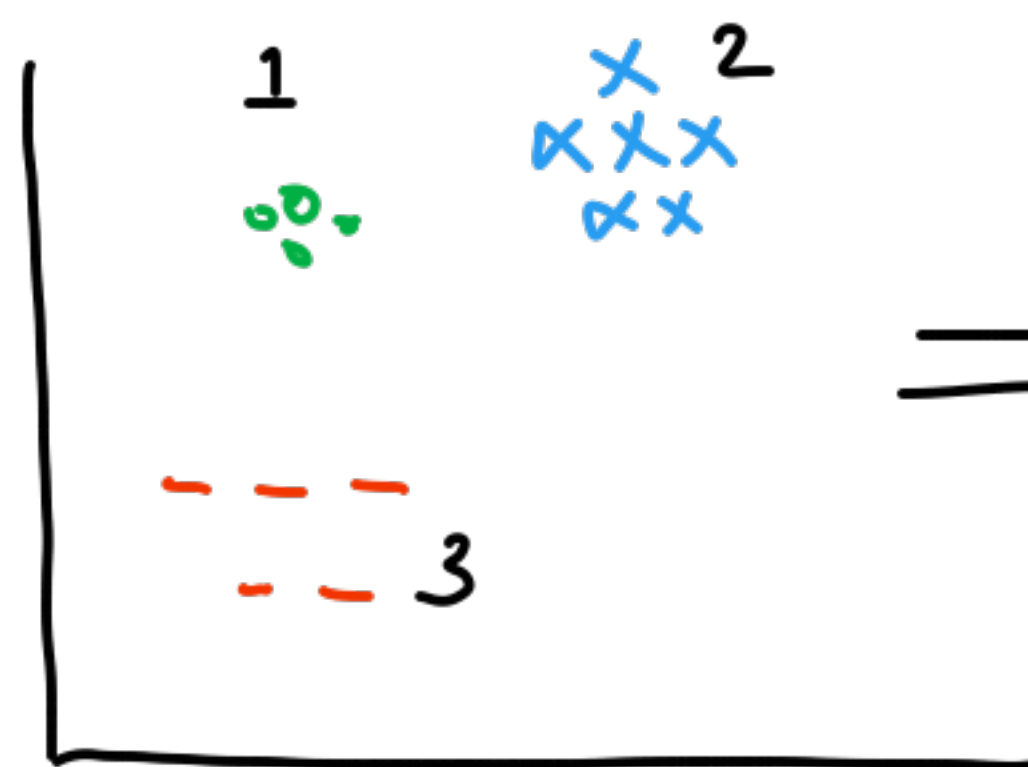


CS725

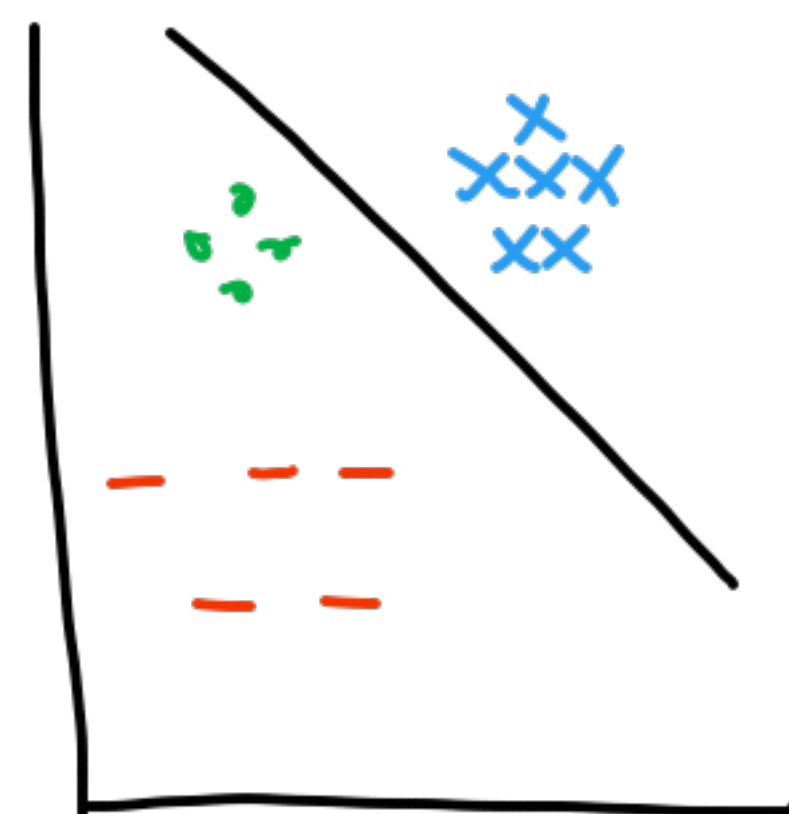
# Logistic Regression with multiple classes

(A) One vs. rest

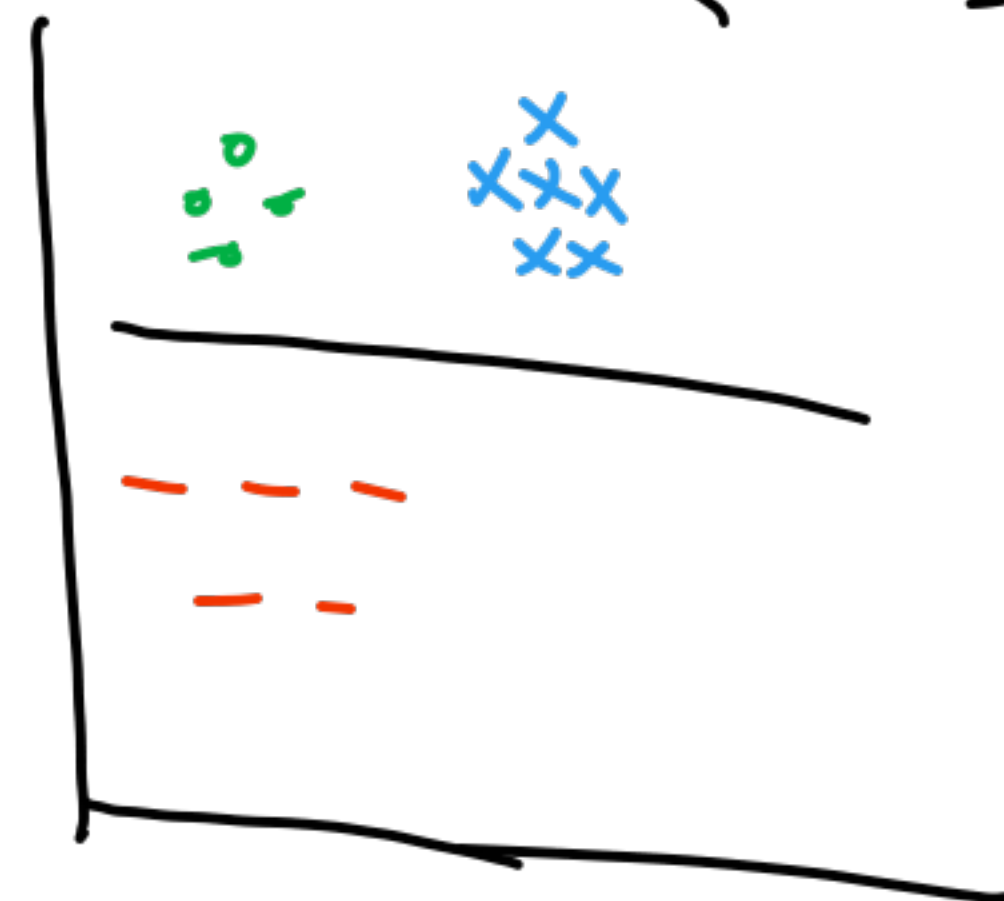
Pick the label with the largest probability



$P(y=1|x)$



$P(y=2|x)$



$P(y=3|x)$

## ② SOFTMAX REGRESSION (MULTINOMIAL LOGISTIC REGRESSION)

Given  $y_i \in \{1, 2, \dots, K\}$ , for a point  $x$ , we want a predictor function

$$h_w(x) = \begin{bmatrix} P(y=1|x;w) \\ P(y=2|x;w) \\ \vdots \\ P(y=K|x;w) \end{bmatrix} \longrightarrow \textcircled{1}$$

Softmax regression assumes weight vectors, one per class, to yield scores for  $K$  classes;

$$\begin{bmatrix} w_1^T x \\ w_2^T x \\ \vdots \\ w_K^T x \end{bmatrix} \longrightarrow \textcircled{2}$$

$w_1, w_2, \dots, w_K$ , where  $w_j \in \mathbb{R}^d$   
if  $j = 1, \dots, K$

How do we go from ② to ①? SOFTMAX TRANSFORMATION

What is softmax? If  $V = [v_1, v_2, \dots, v_k]$ ,  $\text{softmax}(V) = \frac{1}{\sum_k e^{v_k}} \begin{bmatrix} e^{v_1} \\ e^{v_2} \\ \vdots \\ e^{v_k} \end{bmatrix}$   
 $v_j \in \mathbb{R}$

$$h_w(x) = \begin{bmatrix} P(y=1|x;w) \\ P(y=2|x;w) \\ \vdots \\ P(y=k|x;w) \end{bmatrix} = \frac{1}{\sum_j e^{w_j^T x}} \begin{bmatrix} e^{w_1^T x} \\ e^{w_2^T x} \\ \vdots \\ e^{w_k^T x} \end{bmatrix} \xrightarrow{\text{NORMALIZATION FACTOR}} \textcircled{3}$$

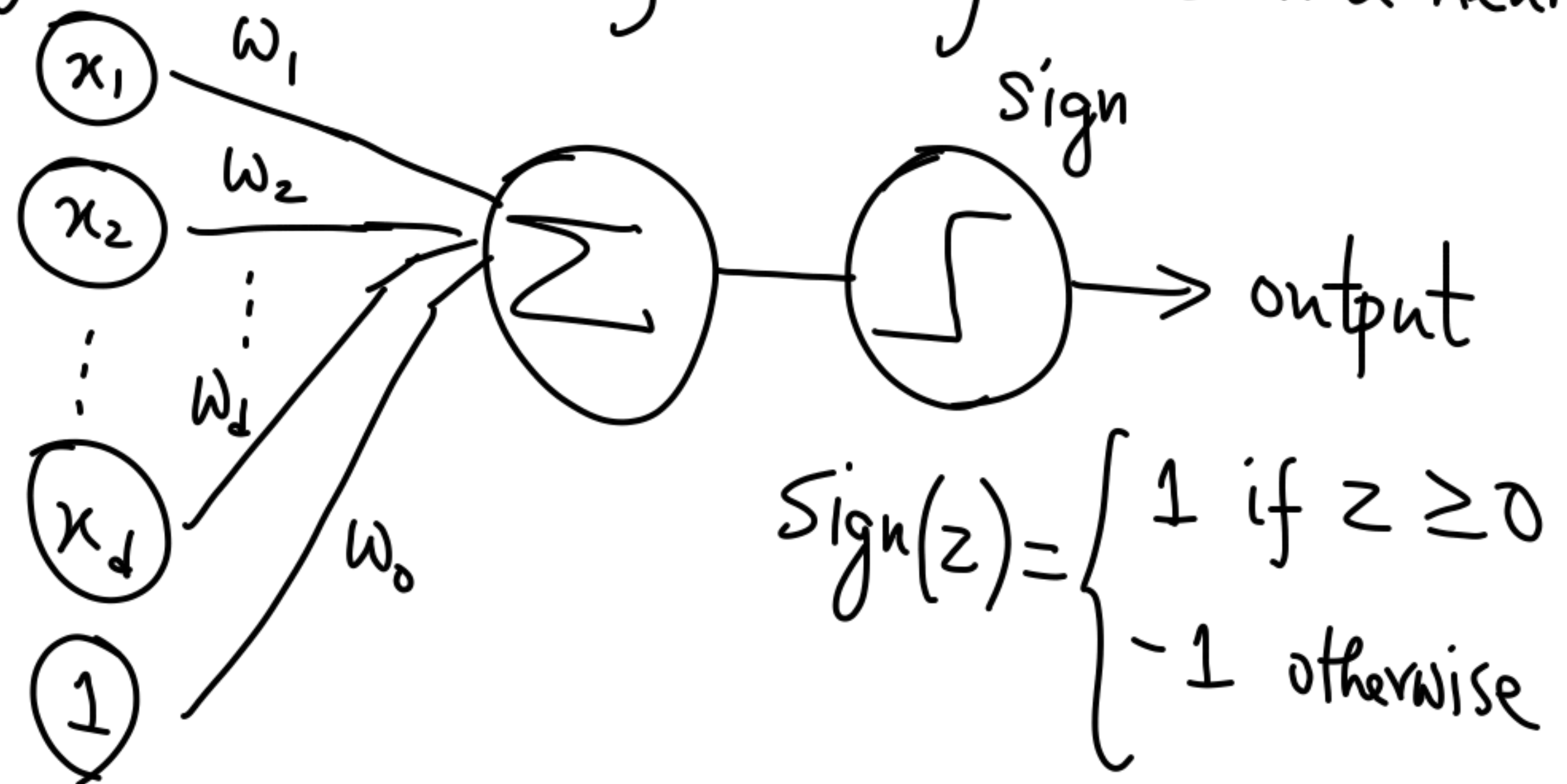
Loss function remains cross-entropy loss with the predicted distribution derived as in (3) and the reference distribution being a one-hot vector over  $K$  classes.

# PERCEPTRON MODEL OF CLASSIFICATION

Linear model of classification. The perceptron model is

(A) the basic unit of neural learning. Loosely based on a neuron.

ROSENBLATT'S  
PERCEPTRON  
MODEL



(B) Online Learner  $\Rightarrow$  that is, it acts on examples one after the other

(C) Error-driven or mistake-driven algorithm  $\Rightarrow$  make a weight update only when the learner makes a mistake



# PERCEPTRON ALGORITHM

Inputs: Set of training points  $\mathcal{D}$ , max # of iterations  $T$ , initialize the wt vector  $w$  [ Let  $w \leftarrow \vec{0}$  (all zero's vector) ]

for  $t = 1, \dots, T$

for all  $(x, y) \in \mathcal{D}, x \in \mathbb{R}^d, y \in \{-1, 1\}$

Can also write  
if  $(y w^T x \leq 0)$

Predict  $\hat{y} = \text{sign}(w^T x)$   
if  $(\hat{y} \neq y)$

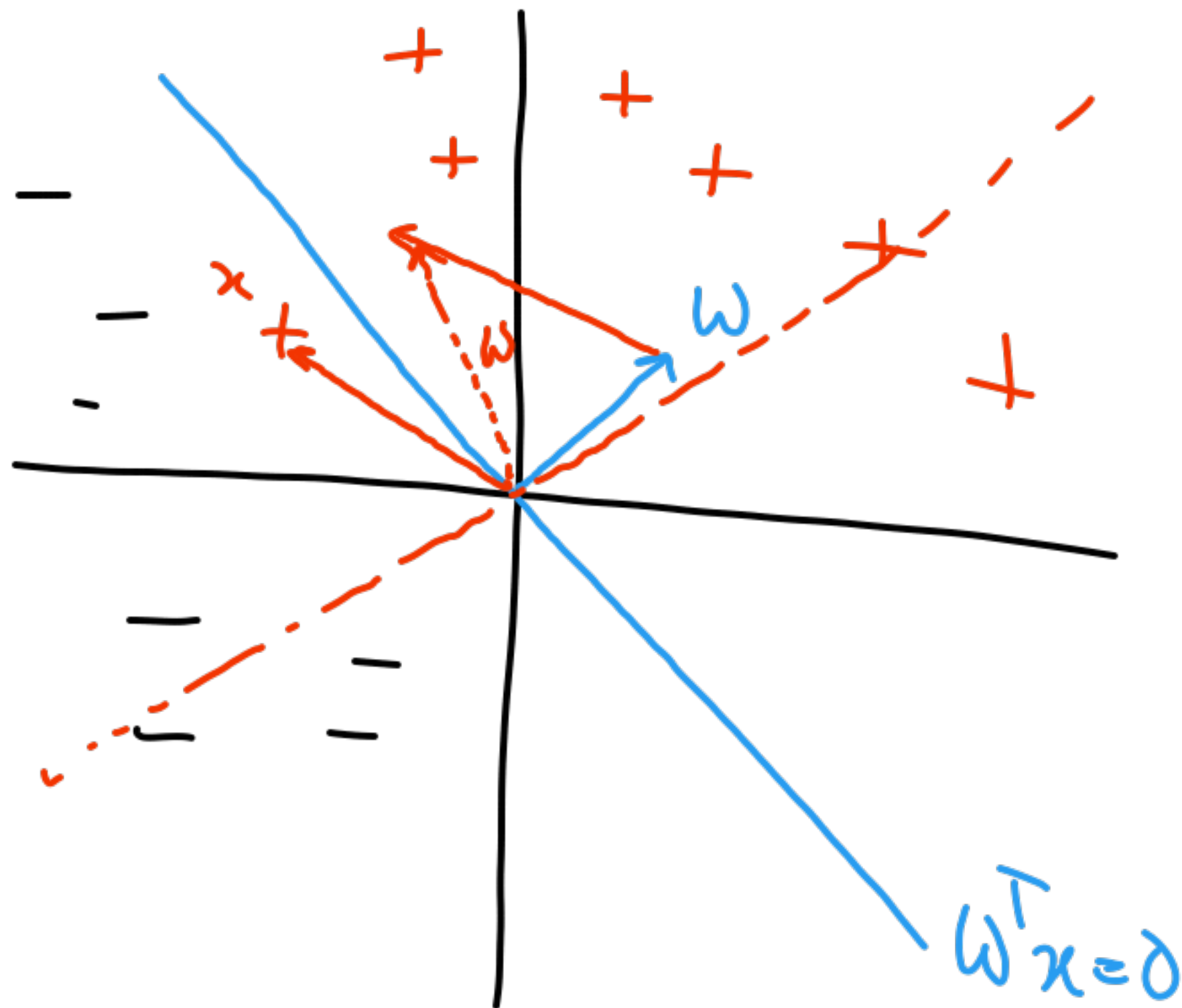
$w \leftarrow w + \eta y x$

$$w = \begin{cases} w + x & \text{if } y = 1 \\ w - x & \text{if } y = -1 \end{cases}$$

return  $w$

# Geometric Intuition for the Perceptron

Moving  $w$  towards or away from a misclassified example depending on whether its true label is 1 or -1, respectively





Does the perceptron guarantee to be more correct on a current misclassified example?

Ⓐ Yes, more correct but not guaranteed to give the correct label with a single update

Consider a misclassified example  $(x, y) \in \mathcal{D}$

$$w_{\text{new}} \leftarrow w_{\text{old}} + yx$$

$$y w_{\text{new}}^T x = y (w_{\text{old}} + yx)^T x = y w_{\text{old}}^T x + y^2 \|x\|_2^2 > y w_{\text{old}}^T x$$

We know that  $yw_{old}^T x \leq 0$ , and  $yw_{new}^T x > yw_{old}^T x$

$\Rightarrow$  thus guaranteeing that we are more correct on  $(x, y)$

What is the perceptron optimizing?

Minimize the misclassified examples

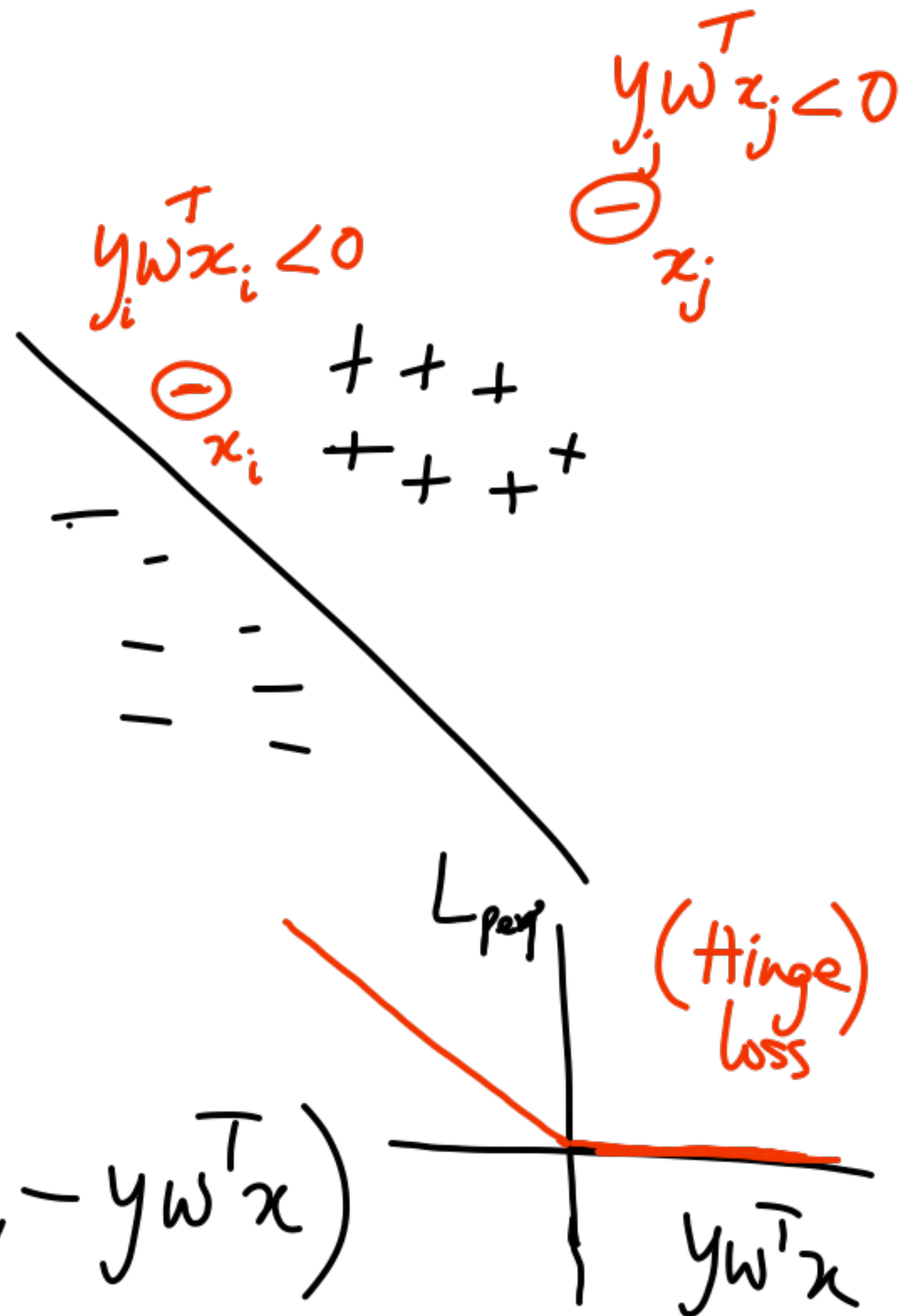
Let  $\mathcal{D}'$  be the set of misclassified examples, then

$$\text{we want } \min \sum_{(x,y) \in \mathcal{D}'} -y w^T x$$

Perceptron loss for one example  $(x,y) \Rightarrow \max(0, -y w^T x)$

$$\text{Perceptron loss for } \mathcal{D} \Rightarrow L_{\text{per}}(w; \mathcal{D}) = \sum_i \max(0, -y_i w^T x_i)$$

(Gradient descent on  $L_{\text{per}}$  yields  $w \leftarrow w + yx$  if  $y w^T x \leq 0$ )



The perceptron is guaranteed to converge (i.e., make no more mistakes) on linearly separable data