

Lecture 2

Linear Classifiers

Professor Slobodan Ribarić
University of Zagreb, Faculty of Electrical Engineering and
Computing (FER)
Croatia

GIAN course, March 2018

Linear decision functions

A general linear decision function:

$d(\mathbf{x}) = w_1x_1 + w_2x_2 + \dots + w_nx_n + w_{n+1}$ – linear combinations of feature (pattern) vector's components

$w_i ; i = 1, 2, \dots, n$ – weights

w_{n+1} – bias, threshold weight

$\mathbf{w}_0 = [w_1, w_2, \dots, w_n]^T$ – weight vector

$\mathbf{x} = [x_1, x_2, \dots, x_n]^T$

$$d(\mathbf{x}) = \mathbf{w}_0^T \mathbf{x} + w_{n+1}$$

$\mathbf{x} = [x_1, x_2, \dots, x_n, 1]^T$ and $\mathbf{w} = [w_1, w_2, \dots, w_n, w_{n+1}]^T$ – augmented vectors

$$d(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$$

Linear decision function for feature vector $\mathbf{x} = [x_1, x_2, \dots, x_n]^T$ is:

- for $n = 2$ line
- for $n = 3$ plane
- for $n > 3$ hyperplane

Linear decision functions for $M = 2$ classes (ω_1 and ω_2)

Boundary between two subspaces which correspond to classes ω_1 and ω_2 :

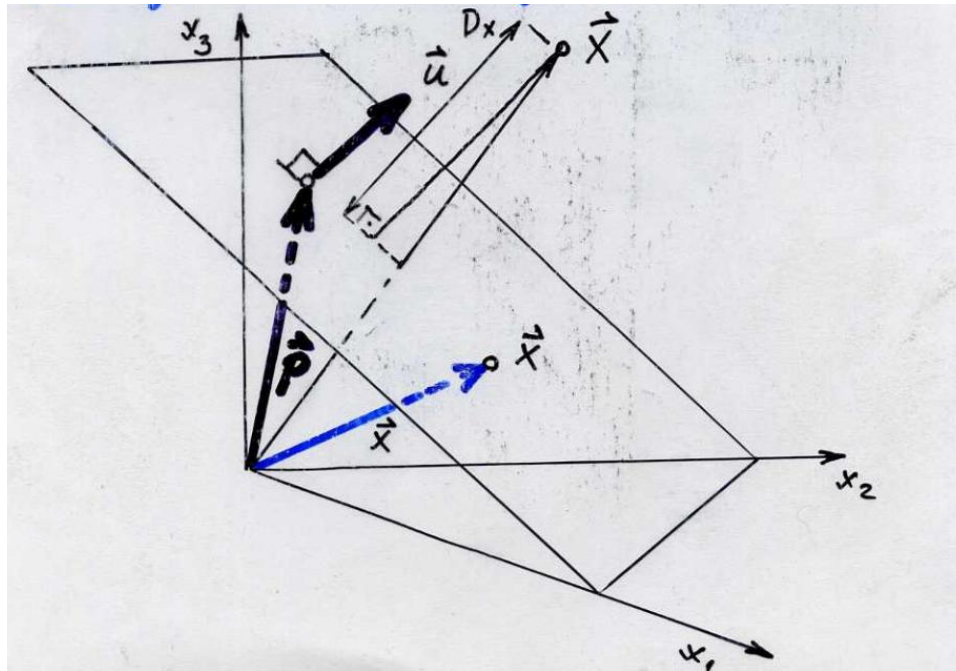
$$d(\mathbf{x}) = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + w_{n+1} = 0$$

A decision function $d(\mathbf{x})$ is assumed to have the property:

$$\begin{aligned} d(\mathbf{x}) = \mathbf{w}^T \mathbf{x} &> 0 \text{ if } \mathbf{x} \in \omega_1 \text{ and} \\ d(\mathbf{x}) = \mathbf{w}^T \mathbf{x} &< 0 \text{ if } \mathbf{x} \in \omega_2 \end{aligned}$$

If $d(\mathbf{x}) = 0$, \mathbf{x} lies on hyperplane and its class membership is not defined

Geometrical interpretation of linear decision function



\mathbf{u} – a unit vector normal to the hyperplane at some point \mathbf{p} and oriented to the positive side of the hyperplane

Equation of the hyperplane:

$$\mathbf{u}^T (\mathbf{x} - \mathbf{p}) = 0$$

$$\mathbf{u}^T \mathbf{x} = \mathbf{u}^T \mathbf{p}$$

$$\mathbf{u}^T \mathbf{x} = \mathbf{u}^T \mathbf{p} \quad (1)$$

$$d(\mathbf{x}) = \mathbf{w}_0^T \mathbf{x} + w_{n+1} = 0$$

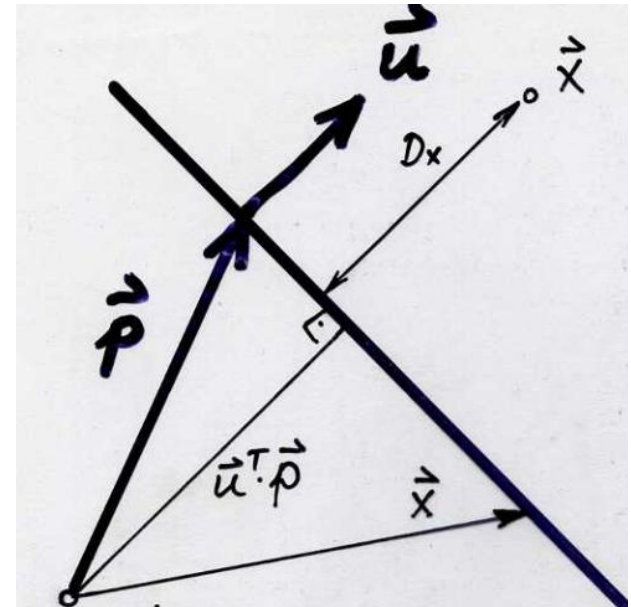
$$\mathbf{w}_0^T \mathbf{x} = -w_{n+1} / \|\mathbf{w}_0\|$$

$$\mathbf{w}_0^T \mathbf{x} / \|\mathbf{w}_0\| = -w_{n+1} / \|\mathbf{w}_0\| \quad (2)$$

By comparing Eq. (1) and (2) follows:

$$\mathbf{u} = \mathbf{w}_0 / \|\mathbf{w}_0\|$$

$$\mathbf{u}^T \mathbf{p} = -w_{n+1} / \|\mathbf{w}_0\|$$

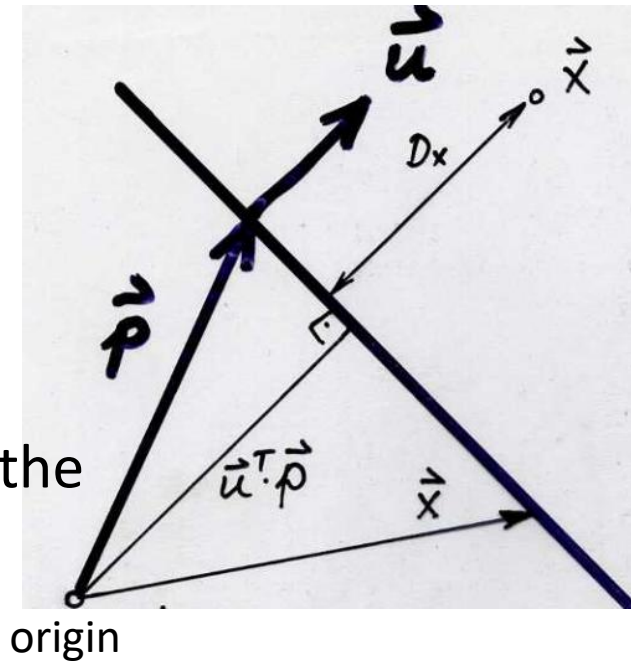


$$\|\mathbf{w}_0\| = \sqrt{w_1^2 + w_2^2 + \dots + w_n^2}$$

$$\mathbf{u} = \mathbf{w}_0 / \|\mathbf{w}_0\|$$

$$\mathbf{u}^T \mathbf{p} = -w_{n+1} / \|\mathbf{w}_0\|$$

Absolute value of $\mathbf{u}^T \mathbf{p}$ represents the normal distance D_u from the origin to the hyperplane



$\mathbf{u} = \mathbf{w}_0 / \|\mathbf{w}_0\|$ defines the orientation of the hyperplane

If any component of \mathbf{u} is zero, the hyperplane is parallel to the coordinate axis which corresponds to that component

- The normal distance D_x from hyperplane and an arbitrary pattern vector \mathbf{x} :

$$D_x = \text{abs}(\mathbf{u}^T \mathbf{x} - \mathbf{u}^T \mathbf{p})$$

$$D_x = \text{abs}(\mathbf{w}_0^T \mathbf{x} / \|\mathbf{w}_0\| + w_{n+1} / \|\mathbf{w}_0\|)$$

Multicategory /multiclass/ case: $M > 2$

- there is more than one way to devise multicategory classifiers employing linear decision functions

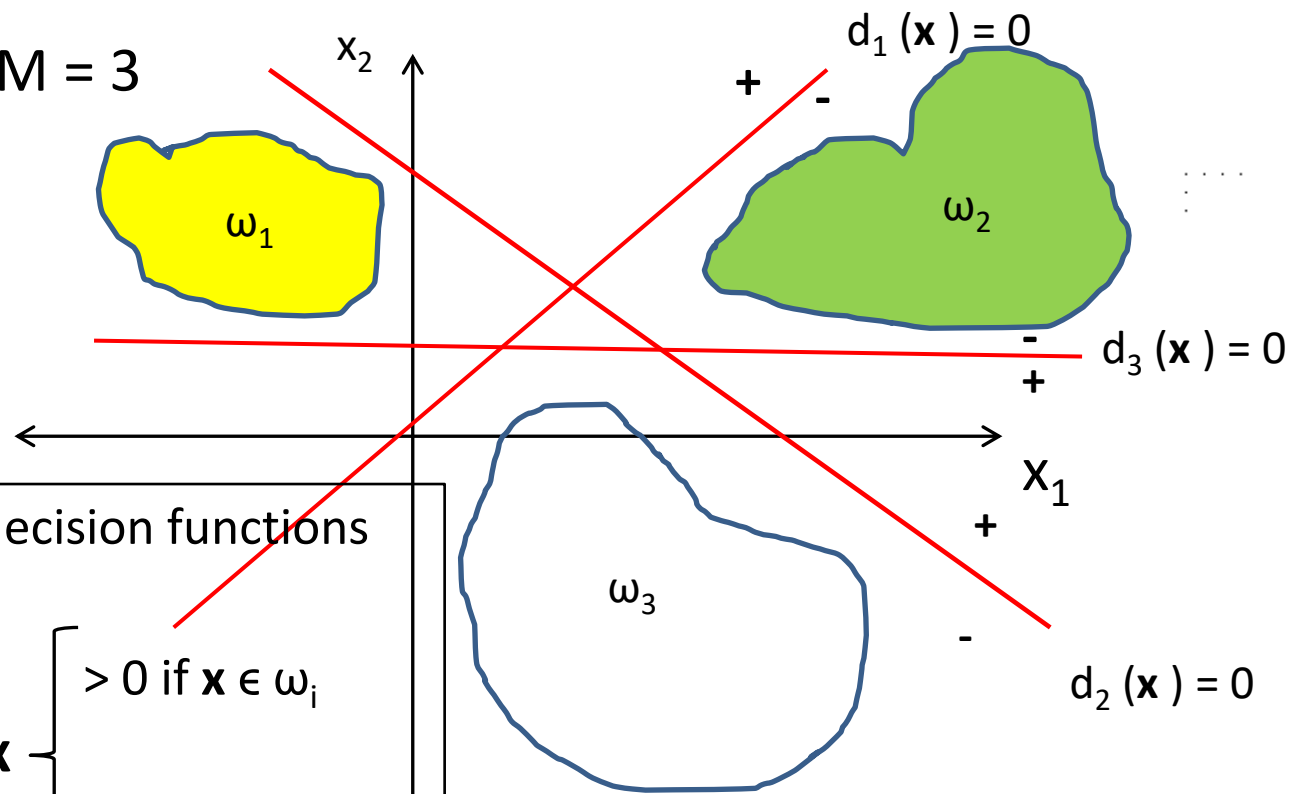
For example – reduce the problem to M two-classes problems where the i -th problem is solved by linear decision function that separates patterns assigned to ω_i from those not assigned to ω_i .

Three cases!

Case 1:

Each class is separable from other classes by a single decision surface.

Example: $M = 3$



There are M decision functions

$$d_i(\mathbf{x}) = \mathbf{w}^T \mathbf{x} \begin{cases} > 0 \text{ if } \mathbf{x} \in \omega_i \\ < 0 \text{ otherwise} \end{cases}$$

Case 2:

Each class is separable from every other individual class by a distinct decision surface – the classes are pairwise separable.

There are $M(M-1)/2$ decision surfaces;

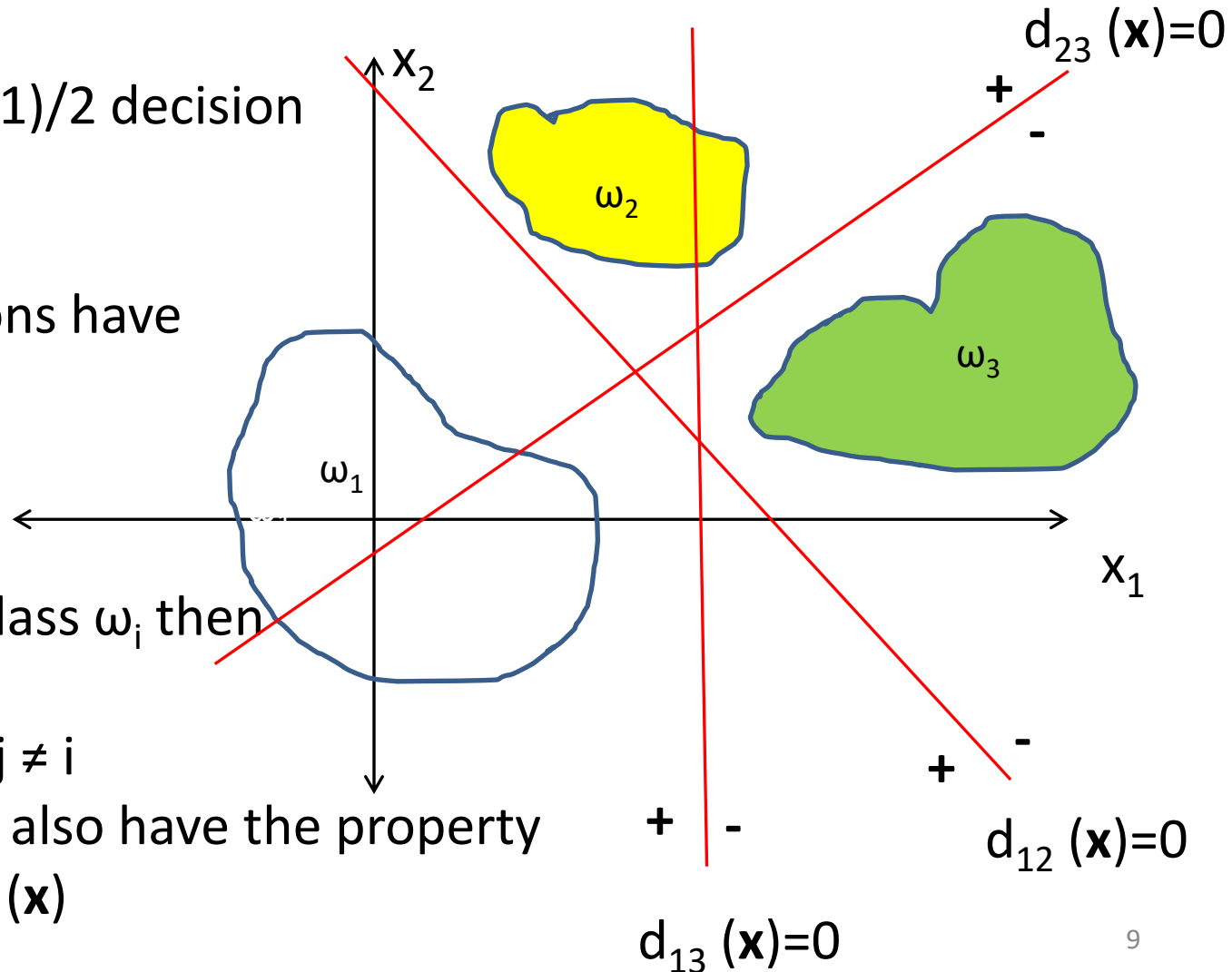
Decision functions have the form:

$$d_{ij}(\mathbf{x}) = \mathbf{w}_{ij}^T \mathbf{x}$$

If \mathbf{x} belongs to class ω_i then

$$d_{ij}(\mathbf{x}) > 0 \text{ for all } j \neq i$$

These functions also have the property that $d_{ij}(\mathbf{x}) = -d_{ji}(\mathbf{x})$



Case 3:

There exist M decision functions $d_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x}$, $k = 1, 2, \dots, M$

If \mathbf{x} belongs to class ω_i then

$d_i(\mathbf{x}) > d_j(\mathbf{x})$ for all $j \neq i$

Example:

$$d_1(\mathbf{x}) = -x_1 + x_2$$

$$d_2(\mathbf{x}) = x_1 + x_2 - 1$$

$$d_3(\mathbf{x}) = -x_2$$

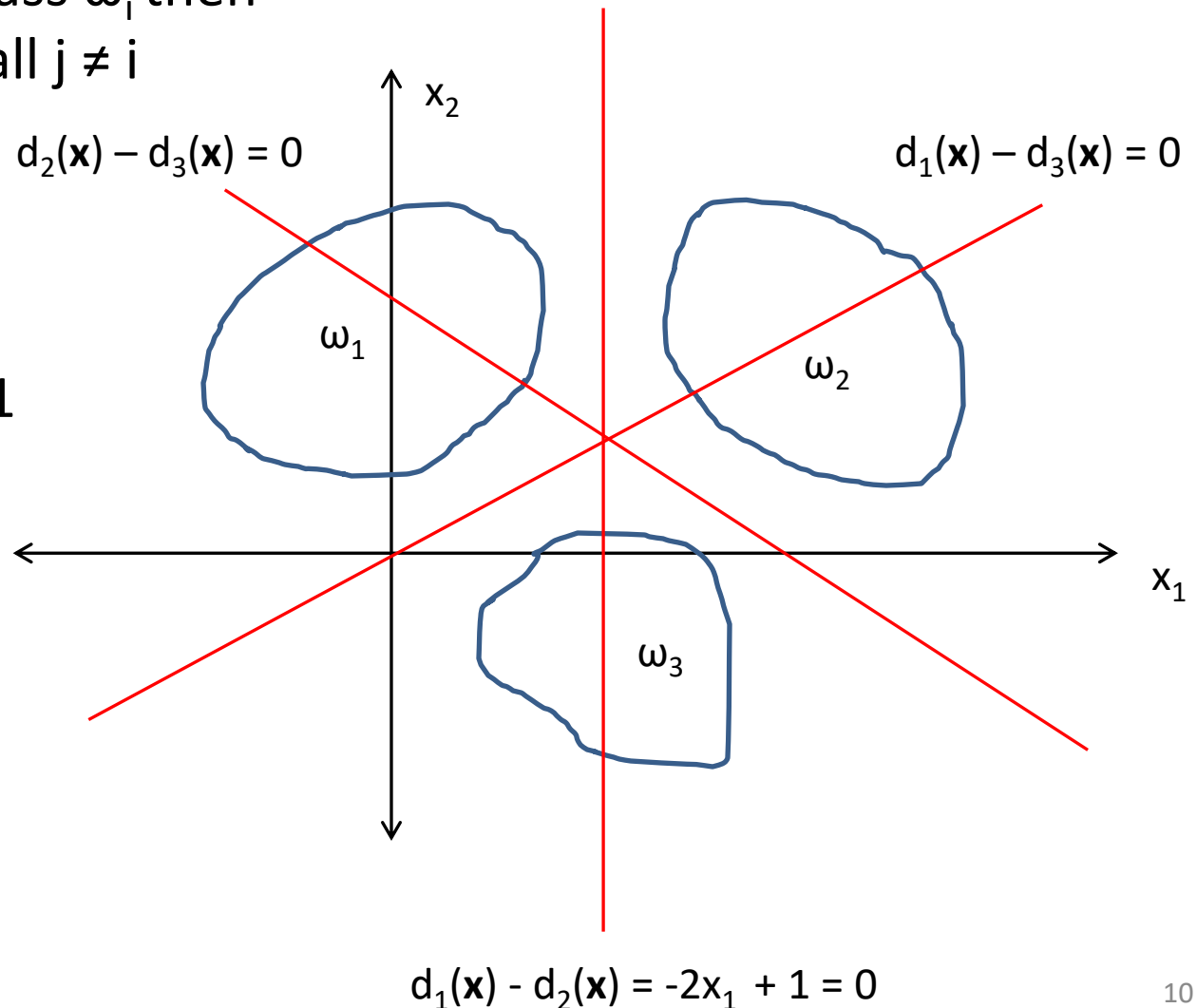
$$\mathbf{x} = [2, 4]^T$$

$$d_1(\mathbf{x}) = 2$$

$$d_2(\mathbf{x}) = 5$$

$$d_3(\mathbf{x}) = -4$$

$$\mathbf{x} \in \omega_2$$



Case 3:

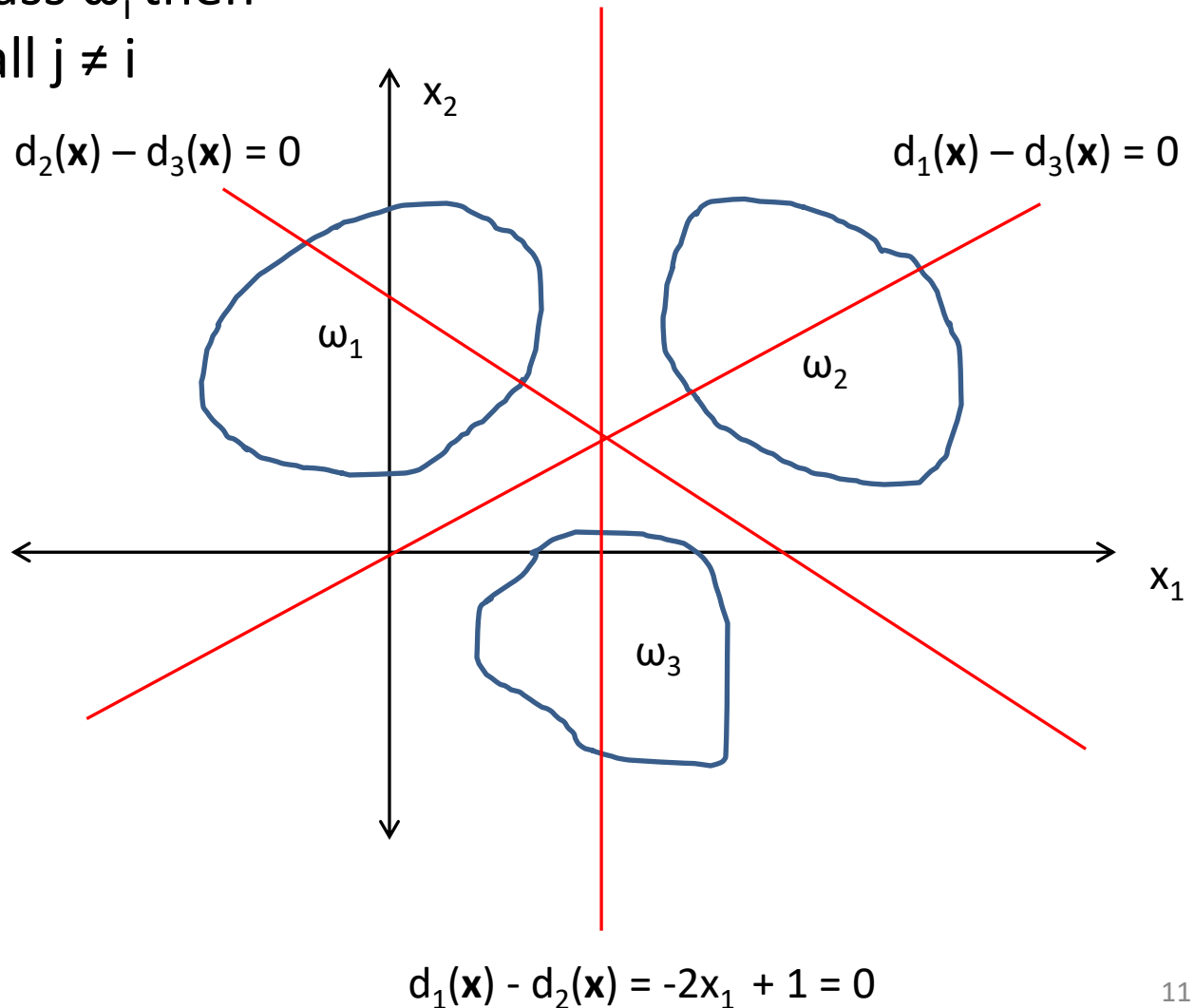
There exist M decision functions $d_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x}$, $k = 1, 2, \dots, M$

If \mathbf{x} belongs to class ω_i then

$d_i(\mathbf{x}) > d_j(\mathbf{x})$ for all $j \neq i$

Case 3 is a special instance of Case 2 since we may define:

$$\begin{aligned} d_{ij}(\mathbf{x}) &= d_i(\mathbf{x}) - d_j(\mathbf{x}) \\ &= (\mathbf{w}_i - \mathbf{w}_j)^T \mathbf{x} \\ &= \mathbf{w}_{ij}^T \mathbf{x} \end{aligned}$$



Learning algorithms for linear classifiers (deterministic approach)

Assumption: All feature vectors from available classes can be classified correctly using **linear classifier**

Advantages of linear classifiers: simplicity and computational attractiveness

$$d(\mathbf{x}) = w_1x_1 + w_2x_2 + \dots + w_nx_n + w_{n+1}$$

$$d(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$$

Problem: determination of the components of the weight vector

$$\mathbf{w} = [w_1, w_2, \dots, w_n, w_{n+1}]^T$$

Learning of decision function for $M = 2$ classes

$$d(\mathbf{x}) = \mathbf{w}^T \mathbf{x} > 0 \text{ if } \mathbf{x} \in \omega_1 \text{ and}$$

$$d(\mathbf{x}) = \mathbf{w}^T \mathbf{x} < 0 \text{ if } \mathbf{x} \in \omega_2$$

Learning of decision function for $M = 2$ classes (cont.)

A set of N training pattern vectors:

$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$ - labeled augmented pattern vectors

$$\mathbf{x}_i = [x_1, x_2, \dots, x_n, 1]^T; i = 1, 2, \dots, N$$

$$d(\mathbf{x}) = \mathbf{w}^T \mathbf{x} > 0 \text{ if } \mathbf{x} \in \omega_1 \text{ and}$$

$$d(\mathbf{x}) = \mathbf{w}^T \mathbf{x} < 0 \text{ if } \mathbf{x} \in \omega_2$$

If all pattern vectors of ω_2 are multiplied by -1 we obtain:

$$\mathbf{w}^T \mathbf{x} < 0 / (-1)$$

$$\mathbf{w}^T \mathbf{x} > 0$$

Now we have equivalent condition $\mathbf{w}^T \mathbf{x} > 0$ for all pattern vectors

Redefinition of the problem:

We are looking for the components of the weight vector $\mathbf{w} = [w_1, w_2, \dots, w_n, w_{n+1}]^T$ which have to satisfy unique condition:

$$\mathbf{w}^T \mathbf{x} > 0$$

for all pattern vectors from training set.

We are looking for the components of the weight vector $\mathbf{w} = [w_1, w_2, \dots, w_n, w_{n+1}]^T$ which have to satisfy unique condition:

$$\mathbf{w}^T \mathbf{x} > 0$$

for all pattern vectors from training set.

The above condition we can write in a form:

$$\mathbf{X} \mathbf{w} > \mathbf{0}, \text{ where}$$

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \mathbf{x}_3^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix} \quad \text{and } \mathbf{0} \text{ is the zero vector.}$$

Example:

$$\mathbf{x}_1 = [0, 0]^T, \mathbf{x}_2 = [0, 1]^T, \mathbf{x}_3 = [1, 0]^T, \mathbf{x}_4 = [1, 1]^T; \mathbf{x}_1, \mathbf{x}_2 \in \omega_1$$
$$\mathbf{x}_3, \mathbf{x}_4 \in \omega_2$$

Augmented pattern vectors are:

$$\mathbf{x}_1 = [0, 0, 1]^T, \mathbf{x}_2 = [0, 1, 1]^T, \mathbf{x}_3 = [1, 0, 1]^T, \mathbf{x}_4 = [1, 1, 1]^T$$

Multiply \mathbf{x}_3 and \mathbf{x}_4 by (-1) : $\mathbf{x}_3 = [-1, 0, -1]^T$ and $\mathbf{x}_4 = [-1, -1, -1]^T$

$$\mathbf{X} = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 1 \\ -1 & 0 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$
$$\begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 1 \\ -1 & 0 & -1 \\ -1 & -1 & -1 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} > \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

The weight vector **w** which satisfies the system of inequalities

$$\mathbf{X} \mathbf{w} > \mathbf{0}$$

is separating vector.

If **w** exists – classes ω_1 and ω_2 are (linear) separable.

Note: All pattern vectors from ω_2 are multiplied by -1.

The gradient technique

In general, the gradient of function $f(\mathbf{y})$ with respect to the vector

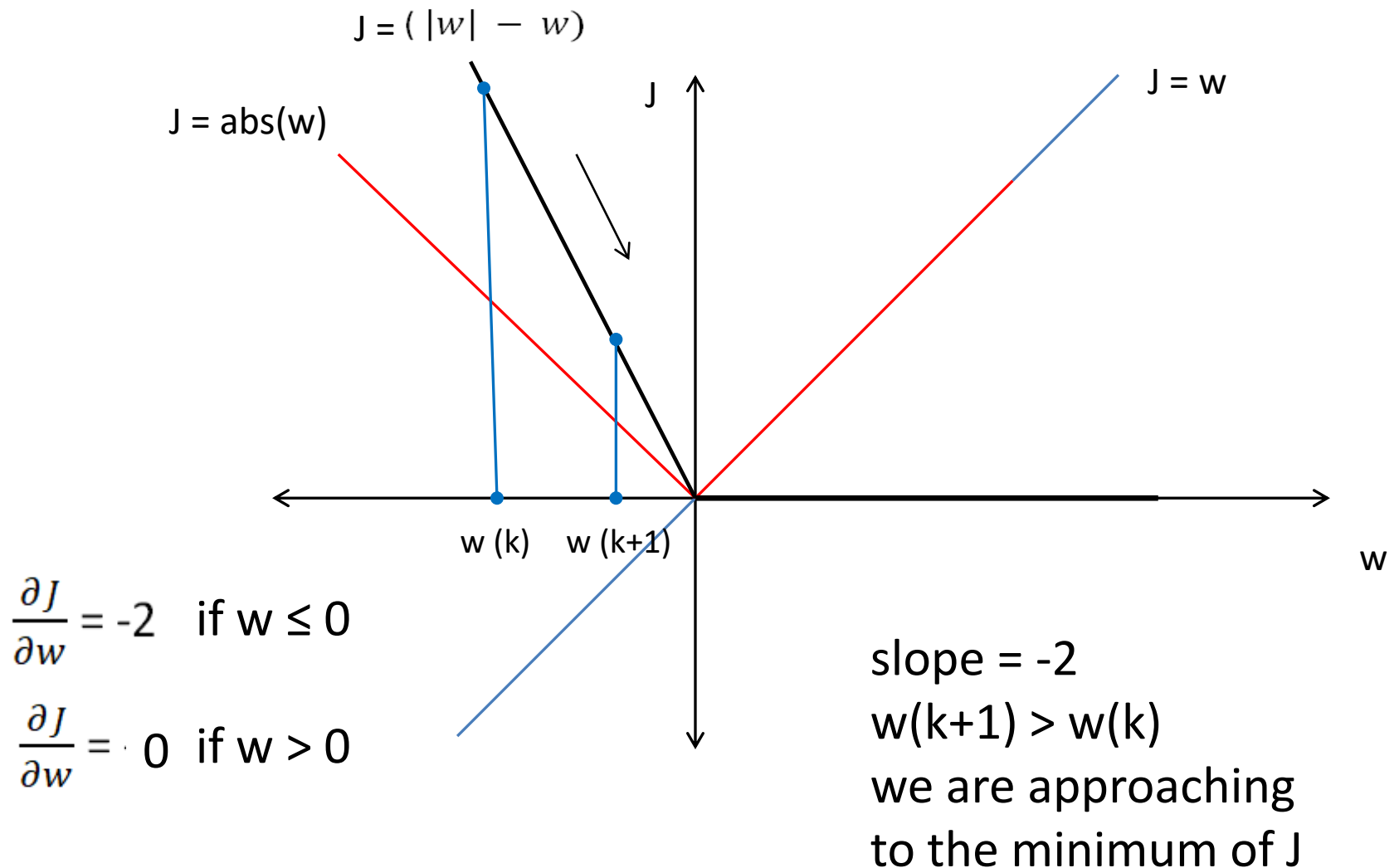
$\mathbf{y} = [y_1, y_2, \dots, y_n]^T$ is defined as:

$$\text{grad } f(\mathbf{y}) = \frac{df(\mathbf{y})}{d\mathbf{y}} = \begin{pmatrix} \frac{df}{dy_1} \\ \frac{df}{dy_2} \\ \vdots \\ \frac{df}{dy_n} \end{pmatrix}$$

- $\text{grad } f(\mathbf{y})$: Gradient of a scalar function of a vector argument is a vector
- each component of the gradient gives the rate of change of the function in the direction of that component
- the positive of gradient points in direction of the maximum rate of increase of the function f when the argument increases
- the negative of the gradient points in the direction of the maximum rate of decrease of the function f
- the above properties can be used for finding the minimum (or maximum) of a function f

Example:

Consider the function $J(w, 1) = (|w| - w)$



Basic idea - select the proper function $J(\mathbf{w}, \mathbf{x})$ /criterion function/ which achieves the minimum when the condition $\mathbf{w}^T \mathbf{x}_i > 0$ for all \mathbf{x}_i , $i = 1, 2, \dots, N$, is satisfied!

- Incrementing \mathbf{w} in the direction of negative gradient of $J(\mathbf{w}, \mathbf{x})$, in order to seek the minimum of the function.

Gradient decent algorithm

$\mathbf{w}(k)$ – value of \mathbf{w} at the k th step

$$\mathbf{w}(k + 1) = \mathbf{w}(k) - c \left(\frac{\partial J(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w}} \right)_{\mathbf{w} = \mathbf{w}(k)}$$

where $\mathbf{w}(k + 1)$ represents the new value of \mathbf{w} , and $c > 0$ dictates the magnitude of the correction.

Perceptron algorithm

Let us define the criterion function:

$$J(\mathbf{w}, \mathbf{x}) = \frac{1}{2} (|\mathbf{w}^T \mathbf{x}| - \mathbf{w}^T \mathbf{x})$$

The partial derivation of J with respect to \mathbf{w} :

$$\frac{\partial}{\partial \mathbf{w}} J(\mathbf{w}, \mathbf{x}) = \frac{1}{2} (\mathbf{x} \operatorname{sgn}(\mathbf{w}^T \mathbf{x}) - \mathbf{x})$$

where:

$$\operatorname{sgn}(x) \begin{cases} 1 & \text{if } x > 0 \\ -1 & \text{if } x \leq 0 \end{cases}$$

Some rules for partial derivations:

$$\frac{d (\mathbf{x}^T \mathbf{A})}{d\mathbf{x}} = \mathbf{A}$$

$$\frac{d}{d\mathbf{x}} (\mathbf{x}^T) = \mathbf{I}$$

$$\frac{d}{d\mathbf{x}} (\mathbf{x}^T \mathbf{a}) = \frac{d}{d\mathbf{x}} (\mathbf{a}^T \mathbf{x}) = \mathbf{a}$$

$$\frac{d}{d\mathbf{X}} (\mathbf{a}^T \mathbf{X} \mathbf{b}) = \mathbf{a} \mathbf{b}^T$$

$$\frac{d}{d\mathbf{X}} (\mathbf{a}^T \mathbf{X} \mathbf{a}) = \frac{d}{d\mathbf{X}} (\mathbf{a}^T \mathbf{X}^T \mathbf{a}) = \mathbf{a} \mathbf{a}^T$$

$$\frac{d}{d\mathbf{x}} (\mathbf{x}^T \mathbf{C} \mathbf{x}) = (\mathbf{C} + \mathbf{C}^T) \mathbf{x}^T$$

$$\text{if } \mathbf{C} = \mathbf{C}^T \text{ then } \frac{d}{d\mathbf{x}} (\mathbf{x}^T \mathbf{C} \mathbf{x}) = 2 \mathbf{C} \mathbf{x}$$

$$\frac{d}{d\mathbf{x}} (\mathbf{x}^T \mathbf{x}) = 2 \mathbf{x}$$

$$\frac{d}{d\mathbf{x}} (\mathbf{A} \mathbf{x} + \mathbf{b})^T (\mathbf{D} \mathbf{x} + \mathbf{e}) = \mathbf{A}^T (\mathbf{D} \mathbf{x} + \mathbf{e}) + \mathbf{D}^T (\mathbf{A} \mathbf{x} + \mathbf{b})$$

$$\text{sgn}(\mathbf{w}^T \mathbf{x}) \begin{cases} 1 & \text{if } \mathbf{w}^T \mathbf{x} > 0 \\ -1 & \text{if } \mathbf{w}^T \mathbf{x} \leq 0 \end{cases} \quad \frac{d}{d\mathbf{x}} (\mathbf{x}^T \mathbf{a}) = \frac{d}{d\mathbf{x}} (\mathbf{a}^T \mathbf{x}) = \mathbf{a}$$

$$\frac{\partial}{\partial \mathbf{w}} J(\mathbf{w}, \mathbf{x}) = \frac{1}{2} (\mathbf{x} \text{sgn}(\mathbf{w}^T \mathbf{x}) - \mathbf{x}) \quad (1)$$

Substituting Eq. (1) into $\mathbf{w}(k+1) = \mathbf{w}(k) - c \left(\frac{\partial J(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w}} \right)_{\mathbf{w} = \mathbf{w}(k)}$

yields:

$$\mathbf{w}(k+1) = \mathbf{w}(k) - c \left(\frac{1}{2} (\mathbf{x}(k) \text{sgn}(\mathbf{w}^T \mathbf{x}(k)) - \mathbf{x}(k)) \right)$$

$$\mathbf{w}(k+1) = \mathbf{w}(k) + \frac{c}{2} (\mathbf{x}(k) - \mathbf{x}(k) \text{sgn}(\mathbf{w}^T \mathbf{x}(k)))$$


$$\mathbf{w}(k + 1) = \mathbf{w}(k) - c \left(\frac{1}{2} (\mathbf{x}(k) \operatorname{sgn}(\mathbf{w}^T \mathbf{x}(k)) - \mathbf{x}(k)) \right)$$

$$\mathbf{w}(k + 1) = \mathbf{w}(k) + \frac{c}{2} [\mathbf{x}(k) - \mathbf{x}(k) \operatorname{sgn}(\mathbf{w}^T \mathbf{x}(k))]$$

$$\mathbf{w}(k + 1) = \mathbf{w}(k) + c \begin{cases} \mathbf{0} & \text{if } \mathbf{w}^T(k) \mathbf{x}(k) > 0 \\ \mathbf{x}(k) & \text{if } \mathbf{w}^T(k) \mathbf{x}(k) \leq 0 \end{cases}$$

where $c > 0$ and $\mathbf{w}(1)$ is arbitrary.


The criterion function:

$$\frac{d}{dx}(uv) = u \frac{dv}{dx} + v \frac{du}{dx}$$


$$J(\mathbf{w}, \mathbf{x}) = \frac{1}{4\mathbf{x}^T \mathbf{x}} (|\mathbf{w}^T \mathbf{x}|^2 - |\mathbf{w}^T \mathbf{x}| \mathbf{w}^T \mathbf{x})$$

The partial derivative of above function:

$$\frac{\partial J(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w}} = \frac{1}{4\mathbf{x}^T \mathbf{x}} [2|\mathbf{w}^T \mathbf{x}| \mathbf{x} \operatorname{sgn}(\mathbf{w}^T \mathbf{x}) - (|\mathbf{w}^T \mathbf{x}| \mathbf{x} + (\mathbf{w}^T \mathbf{x}) \mathbf{x} \operatorname{sgn}(\mathbf{w}^T \mathbf{x}))]$$

$$(\mathbf{w}^T \mathbf{x}) \mathbf{x} \operatorname{sgn}(\mathbf{w}^T \mathbf{x}) = |\mathbf{w}^T \mathbf{x}| \mathbf{x}$$


$$\frac{\partial J(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w}} = \frac{1}{4\mathbf{x}^T \mathbf{x}} [2|\mathbf{w}^T \mathbf{x}| \mathbf{x} \operatorname{sgn}(\mathbf{w}^T \mathbf{x}) - 2|\mathbf{w}^T \mathbf{x}| \mathbf{x}]$$

$$\frac{\partial J(\mathbf{w}, \mathbf{x})}{\partial \mathbf{w}} = \frac{1}{2\mathbf{x}^T \mathbf{x}} [|\mathbf{w}^T \mathbf{x}| \mathbf{x} \operatorname{sgn}(\mathbf{w}^T \mathbf{x}) - |\mathbf{w}^T \mathbf{x}| \mathbf{x}]$$

$$\mathbf{w}(k+1) = \mathbf{w}(k) + \lambda \frac{|\mathbf{w}^T(k) \mathbf{x}(k)|}{2\mathbf{x}(k)^T \mathbf{x}(k)} [\mathbf{x}(k) - \mathbf{x}(k) \operatorname{sgn}(\mathbf{w}(k)^T \mathbf{x}(k))]$$

The fraction-correction algorithm:

$$\mathbf{w}(k+1) = \mathbf{w}(k) + \lambda \frac{|\mathbf{w}^T(k) \mathbf{x}(k)|}{\mathbf{x}(k)^T \mathbf{x}(k)} \begin{cases} \mathbf{0} & \text{if } \mathbf{w}^T(k) \mathbf{x}(k) > 0 \\ \mathbf{x}(k) & \text{if } \mathbf{w}^T(k) \mathbf{x}(k) \leq 0 \end{cases}$$

λ – correction factor; $0 < \lambda < 2$

$\mathbf{w}(1) \neq \mathbf{0}$

Variations of perceptron algorithm

$$\mathbf{w}(k+1) = \mathbf{w}(k) + c \begin{cases} \mathbf{0} & \text{if } \mathbf{w}^T(k)\mathbf{x}(k) > 0 \\ \mathbf{x}(k) & \text{if } \mathbf{w}^T(k)\mathbf{x}(k) \leq 0 \end{cases}$$

- i) Fixed-increment algorithm
- ii) Algorithm with absolute correction
- iii) Fractional-correction algorithm

- i) $c > 0$
- ii) If $\mathbf{w}(k)^T \mathbf{x}(k) \leq 0$ then select c that:

$$\mathbf{w}^T(k+1) \mathbf{x}(k) = [\mathbf{w}(k) + c\mathbf{x}(k)]^T \mathbf{x}(k) > 0$$

$$\text{integer value } c > \frac{|\mathbf{w}^T(k) \mathbf{x}(k)|}{\mathbf{x}^T(k) \mathbf{x}(k)}$$

iii)

$$|\mathbf{w}^T(k) \mathbf{x}(k) - \mathbf{w}^T(k+1) \mathbf{x}(k)| = \lambda |\mathbf{w}^T(k) \mathbf{x}(k)|$$

$$\mathbf{w}(k+1) = \mathbf{w}(k) + c \mathbf{x}(k)$$

$$|\mathbf{w}^T(k) \mathbf{x}(k) - (\mathbf{w}^T(k) + c \mathbf{x}^T(k)) \mathbf{x}(k)| = \lambda |\mathbf{w}^T(k) \mathbf{x}(k)|$$

$$c = \lambda \frac{|\mathbf{w}^T(k) \mathbf{x}(k)|}{\mathbf{x}^T(k) \mathbf{x}(k)}$$

Example:

-Training set: $\omega_1 = \{(0, 0)^T, (0, 1)^T\}$
 $\omega_2 = \{(1, 0)^T, (1, 1)^T\}$

$c = 1$

Augmented pattern vectors:

$$\omega_1 = \{(0, 0, 1)^T, (0, 1, 1)^T\}$$

$$\omega_2 = \{(1, 0, 1)^T, (1, 1, 1)^T\}$$

Pattern vectors from ω_2 have to be multiplied by (-1)

$$\omega_2 = \{(-1, 0, -1)^T, (-1, -1, -1)^T\}$$

Let us suppose that vectors are linear separable by $d(\mathbf{x})$:

$$d(\mathbf{x}) = \mathbf{w}^T \mathbf{x} = w_1 x_1 + w_2 x_2 + w_3$$

Let us select $\mathbf{w}(1) = (-1, 0, 0)^T$

1.Step

$$\mathbf{w}^T(1) \mathbf{x}(1) = (-1, 0, 0) \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} = 0,$$

$$\mathbf{w}(2) = \mathbf{w}(1) + \mathbf{x}(1) = \begin{pmatrix} -1 \\ 0 \\ 1 \end{pmatrix}$$

2. Step

$$\mathbf{w}^T(2) \mathbf{x}(2) = (-1, 0, 1) \begin{pmatrix} 0 \\ 1 \\ 1 \end{pmatrix} = 1$$

$$\mathbf{w}^T(2) \mathbf{x}(2) > 0 \quad \dots \quad \mathbf{w}(3) = \mathbf{w}(2)$$

3. Step

$$\mathbf{w}^T(3) \mathbf{x}(3) = (-1, 0, 1) \begin{pmatrix} -1 \\ 0 \\ -1 \end{pmatrix} = 0$$

3. Step (cont.)

$$\mathbf{w}(4) = \mathbf{w}(3) + \mathbf{x}(3)$$

$$\mathbf{w}(4) = \begin{pmatrix} -1 \\ 0 \\ 1 \end{pmatrix} + \begin{pmatrix} -1 \\ 0 \\ -1 \end{pmatrix} = \begin{pmatrix} -2 \\ 0 \\ 0 \end{pmatrix}$$

4. Step

$$\mathbf{w}^T(4) \mathbf{x}(4) = (-2, 0, 0) \begin{pmatrix} -1 \\ -1 \\ -1 \end{pmatrix} = 2$$

$$\mathbf{w}(5) = \mathbf{w}(4)$$

There were corrections on the weight vector in the 1. and 3. steps!

Solution has been obtained ONLY when algorithm yields a complete **ERROR-FREE** iteration through all patterns

The second iteration through all pattern vectors:

$$x(5) = x(1); x(6) = x(2); x(7) = x(3); x(8) = x(4)$$

5. Step

$$\mathbf{w}^T(5) \mathbf{x}(5) = (-2, 0, 0) \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} = 0$$

$$\mathbf{w}(6) = \mathbf{w}(5) + \mathbf{x}(5) = \begin{pmatrix} -2 \\ 0 \\ 0 \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} = \begin{pmatrix} -2 \\ 0 \\ 1 \end{pmatrix}$$

6. Step

$$\mathbf{w}^T(6) \mathbf{x}(6) = (-2, 0, 1) \begin{pmatrix} 0 \\ 1 \\ 1 \end{pmatrix} = 1$$

$$\mathbf{w}(7) = \mathbf{w}(6)$$

7. Step

$$\mathbf{w}^T(7) \mathbf{x}(7) = (-2, 0, 1) \begin{pmatrix} -1 \\ 0 \\ -1 \end{pmatrix} = 1$$

$$\mathbf{w}(8) = \mathbf{w}(7)$$

8. Step

$$\mathbf{w}^T(8) \mathbf{x}(8) = (-2, 0, 1) \begin{pmatrix} -1 \\ -1 \\ -1 \end{pmatrix} = 1$$

Two error occurred in this iteration!

$$x(9) = x(1); x(10) = x(2); x(11) = x(3); x(12) = x(4)$$

Third iteration through all vectors

$$\mathbf{x}(9) = \mathbf{x}(1); \mathbf{x}(10) = \mathbf{x}(2); \mathbf{x}(11) = \mathbf{x}(3); \mathbf{x}(12) = \mathbf{x}(4)$$

9. Step

$$\mathbf{w}^T(9) \mathbf{x}(9) = (-2, 0, 1) \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} = 1$$

$$\mathbf{w}(10) = \mathbf{w}(9)$$

10. Step

$$\mathbf{w}^T(10) \mathbf{x}(10) = (-2, 0, 1) \begin{pmatrix} 0 \\ 1 \\ 1 \end{pmatrix} = 1$$

$$\mathbf{w}(11) = \mathbf{w}(10)$$

11. Step

$$\mathbf{w}^T(11) \mathbf{x}(11) = (-2, 0, 1) \begin{pmatrix} -1 \\ 0 \\ -1 \end{pmatrix} = 1$$

$$\mathbf{w}(12) = \mathbf{w}(11)$$

12. Step

$$\mathbf{w}^T(12) \mathbf{x}(12) = (-2, 0, 1) \begin{pmatrix} -1 \\ -1 \\ -1 \end{pmatrix} = 1$$

All pattern vectors are classified correctly - Error-free

Solution:

$$\mathbf{w} = (-2, 0, 1)^T$$

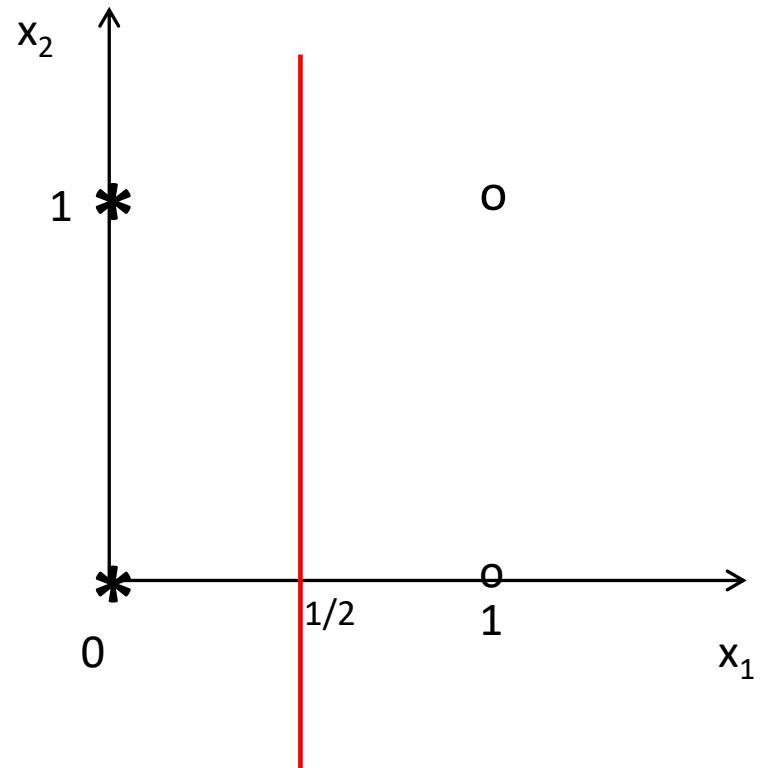
Decision function:

$$d(\mathbf{x}) = \mathbf{w}^T \mathbf{x} = -2x_1 + 1 = 0$$

$$x_1 = 1/2$$

$$* \in \omega_1$$

$$o \in \omega_2$$



Generalized perceptron algorithm (Multicategory classification)

Training method for $M > 2$

Case 3:

$\omega_1, \omega_2, \dots, \omega_M$

If $\mathbf{x} \in \omega_i$ then $d_i(\mathbf{x}) > d_j(\mathbf{x})$ for all $j \neq i$;

At k -th iterative step during training, a pattern vector $\mathbf{x}(k)$ belonging to **class ω_i** is presented to the machine;

Evaluate M decision functions

$$d_j(\mathbf{x}(k)) = \mathbf{w}_j^T \mathbf{x}(k), j = 1, 2, \dots, M$$

Then, if

$$d_i(\mathbf{x}(k)) > d_j(\mathbf{x}(k)) \quad j = 1, 2, \dots, M; j \neq i;$$

the weight vectors are not adjusted:

$$\mathbf{w}_j(k+1) = \mathbf{w}_j(k) \quad j = 1, 2, \dots, M$$

If $d_i(\mathbf{x}(k)) \leq d_L(\mathbf{x}(k))$

The following weight adjustments are made:

$$\mathbf{w}_i(k+1) = \mathbf{w}_i(k) + c\mathbf{x}(k)$$

$$\mathbf{w}_L(k+1) = \mathbf{w}_L(k) - c\mathbf{x}(k)$$

$$\mathbf{w}_j(k+1) = \mathbf{w}_j(k) \text{ for } j=1, 2, \dots, M; j \neq i;$$
$$j \neq L$$

c is positive constant.

If the classes are separable under *Case 3*, the algorithm converges in the finite number of iterations.

Example:

$$M = 3$$

Each class contains a single pattern:

$$\omega_1 = \{(0, 0)^T\}$$

$$\omega_2 = \{(1, 1)^T\}$$

$$\omega_3 = \{(-1, 1)^T\}$$

0. Augment the pattern vectors:

$$\omega_1 = \{(0, 0, 1)^T\}$$

$$\omega_2 = \{(1, 1, 1)^T\}$$

$$\omega_3 = \{(-1, 1, 1)^T\}$$

Note that none of pattern vectors is multiplied by -1.

$$c = 1; \mathbf{w}_1(1) = \mathbf{w}_2(1) = \mathbf{w}_3(1) = (0, 0, 0)^T$$

Apply the **Generalized perceptron algorithm!**

- Find subspaces of feature space which correspond to each class!

Solution:

$$d_1(\mathbf{x}) = -2x_2$$

$$d_2(\mathbf{x}) = 2x_1 - 2$$

$$d_3(\mathbf{x}) = -2x_1 - 2$$

Let us check the solution:

$$\mathbf{x} = (0, 0)^T \in \omega_1$$

$$d_1(\mathbf{x}) = 0$$

$$d_2(\mathbf{x}) = -2$$

$$d_3(\mathbf{x}) = -2$$

$$\mathbf{x} = (1, 1)^T \in \omega_2$$

$$d_1(\mathbf{x}) = -2$$

$$d_2(\mathbf{x}) = 0$$

$$d_3(\mathbf{x}) = -4$$

$$\mathbf{x} = (-1, 1)^T \in \omega_3$$

$$d_1(\mathbf{x}) = -2$$

$$d_2(\mathbf{x}) = -4$$

$$d_3(\mathbf{x}) = 0$$

LMSE – Least-Mean-Square-Error Algorithm

Ho-Kashyap algorithm

- The perceptron algorithm and its variations converge when the classes are separable by linear decision functions
- In non-separable situations these algorithms oscillate
- It is not possible to precompute the number of steps required for convergence in linear separable situation
- **Ho-Kashyap** algorithm indicates that the classes are not separable!

Instead of finding weight vector \mathbf{w} such that $\mathbf{X}\mathbf{w} > \mathbf{0}$ is satisfied, Ho-Kashyap algorithm searches for vectors \mathbf{w} and \mathbf{b} such that:

$$\mathbf{X}\mathbf{w} = \mathbf{b}$$

where the components of $\mathbf{b} = (b_1, b_2, \dots, b_N)^T$ are **all positive**.

N – number of pattern vectors in training set

The criterion function:

$$J(\mathbf{w}, \mathbf{x}, \mathbf{b}) = \frac{1}{2} \sum_{j=1}^N (\mathbf{w}^T \mathbf{x}_j - b_j)^2 = \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{b}\|^2 = \frac{1}{2} (\mathbf{X}\mathbf{w} - \mathbf{b})^T (\mathbf{X}\mathbf{w} - \mathbf{b})$$

where $\|\mathbf{X}\mathbf{w} - \mathbf{b}\|$ is the magnitude of the vector $(\mathbf{X}\mathbf{w} - \mathbf{b})$.

The function $J(\mathbf{w}, \mathbf{x}, \mathbf{b})$ achieves its minimum whenever

$$\mathbf{X}\mathbf{w} = \mathbf{b}$$

- Both variables **w** and **b** can be used in the minimization procedure;
We expect that above can improve the convergence rate of the algorithm;

- Function $J(\mathbf{w}, \mathbf{x}, \mathbf{b})$ will be minimized with respect to **w** and **b**;

- Gradients:

$$\frac{\partial J}{\partial \mathbf{w}} \text{ and } \frac{\partial J}{\partial \mathbf{b}}$$

$$\frac{\partial J}{\partial \mathbf{w}}$$

$$J(\mathbf{w}, \mathbf{x}, \mathbf{b}) = \frac{1}{2} (\mathbf{X}\mathbf{w} - \mathbf{b})^T (\mathbf{X}\mathbf{w} - \mathbf{b})$$

$$\frac{d}{d\mathbf{x}} (\mathbf{x}^T \mathbf{x}) = 2 \mathbf{x}$$

$$\frac{d}{d\mathbf{x}} (\mathbf{A}\mathbf{x} + \mathbf{b})^T (\mathbf{D}\mathbf{x} + \mathbf{e}) = \mathbf{A}^T (\mathbf{D}\mathbf{x} + \mathbf{e}) + \mathbf{D}^T (\mathbf{A}\mathbf{x} + \mathbf{b})$$

$$J(\mathbf{w}, \mathbf{x}, \mathbf{b}) = \frac{1}{2} (\mathbf{X}\mathbf{w} - \mathbf{b})^T (\mathbf{X}\mathbf{w} - \mathbf{b})$$

$$\cdot \frac{d}{d\mathbf{x}} (\mathbf{A}\mathbf{x} + \mathbf{b})^T (\mathbf{D}\mathbf{x} + \mathbf{e}) = \mathbf{A}^T (\mathbf{D}\mathbf{x} + \mathbf{e}) + \mathbf{D}^T (\mathbf{A}\mathbf{x} + \mathbf{b})$$

$$\mathbf{A} = \mathbf{X}, \mathbf{e} = -\mathbf{b}, \mathbf{D} = \mathbf{X}, \mathbf{b} = -\mathbf{b}$$

$$\frac{\partial J}{\partial \mathbf{w}} = \frac{1}{2} (\mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{b}) + \mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{b}))$$

$$\frac{\partial J}{\partial \mathbf{w}} = \mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{b})$$

$$\frac{\partial J}{\partial \mathbf{b}}$$

$$J(\mathbf{w}, \mathbf{x}, \mathbf{b}) = \frac{1}{2} (\mathbf{X}\mathbf{w} - \mathbf{b})^T (\mathbf{X}\mathbf{w} - \mathbf{b})$$

$$\frac{\partial J(\mathbf{w}, \mathbf{x}, \mathbf{b})}{\partial \mathbf{b}} = \frac{\partial J(\mathbf{w}, \mathbf{x}, \mathbf{b})}{\partial \mathbf{b}} \left\{ \frac{1}{2} [(\mathbf{X}\mathbf{w})^T \mathbf{X}\mathbf{w} - \mathbf{b}^T \mathbf{X}\mathbf{w} - (\mathbf{X}\mathbf{w})^T \mathbf{b} + \mathbf{b}^T \mathbf{b}] \right\}$$

$$\frac{d}{d\mathbf{x}} (\mathbf{x}^T \mathbf{a}) = \frac{d}{d\mathbf{x}} (\mathbf{a}^T \mathbf{x}) = \mathbf{a}$$

$$\frac{\partial J(\mathbf{w}, \mathbf{x}, \mathbf{b})}{\partial \mathbf{b}} = \frac{1}{2} [-\mathbf{X}\mathbf{w} - \mathbf{X}\mathbf{w} + 2\mathbf{b}]$$

$$\frac{\partial}{\partial \mathbf{b}} J(\mathbf{w}, \mathbf{x}, \mathbf{b}) = -(\mathbf{X}\mathbf{w} - \mathbf{b})$$

Since \mathbf{w} is not constrained in any way, and we can set

$$\frac{\partial J}{\partial \mathbf{w}} = 0$$

$$\mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{b}) = 0$$

$$\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{b}$$

$$\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{b} \ / \ (\mathbf{X}^T \mathbf{X})^{-1}$$

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{b}$$

$(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = \mathbf{X}^\#$ - generalized inverse matrix

$$\boxed{\mathbf{w} = \mathbf{X}^\# \mathbf{b}} \longrightarrow \mathbf{w}(k+1) = \mathbf{X}^\# \mathbf{b} (k+1)$$

- A vector $\mathbf{b} = (b_1, b_2, \dots, b_n)^T$ is a positive vector – all components of \mathbf{b} are constrained to be **positive**
- This vector must be varied in a such a manner as never to violate this constraint:

$$\mathbf{b}(k+1) = \mathbf{b}(k) + \delta \mathbf{b}(k)$$

$$(*) \quad \delta b_i(k) = \begin{cases} 2c[\mathbf{X}\mathbf{w}(k) - \mathbf{b}(k)]_i & \text{if } [\mathbf{X}\mathbf{w}(k) - \mathbf{b}(k)]_i > 0 \\ 0 & \text{if } [\mathbf{X}\mathbf{w}(k) - \mathbf{b}(k)]_i \leq 0 \end{cases}$$

k denotes the iteration index, i denotes the index of the vector components, c is positive correction increment

- Equation (*) may be written in vector form:

$$\delta \mathbf{b}(k) = c[\mathbf{X}\mathbf{w}(k) - \mathbf{b}(k) + |\mathbf{X}\mathbf{w}(k) - \mathbf{b}(k)|]$$

- where $|\mathbf{X}\mathbf{w}(k) - \mathbf{b}(k)|$ indicates the absolute value of each component of the vector $[\mathbf{X}\mathbf{w}(k) - \mathbf{b}(k)]$

$$\mathbf{w} = \mathbf{X}^\# \mathbf{b}$$

$$\delta \mathbf{b}(k) = c[\mathbf{X}\mathbf{w}(k) - \mathbf{b}(k) + |\mathbf{X}\mathbf{w}(k) - \mathbf{b}(k)|]$$

$$\mathbf{w}(k+1) = \mathbf{X}^\# \mathbf{b}(k+1)$$

$$\mathbf{w}(k+1) = \mathbf{X}^\# [\mathbf{b}(k) + \delta \mathbf{b}(k)]$$

$$\mathbf{w} = \mathbf{X}^\# \mathbf{b}$$

$$\mathbf{w}(k+1) = \mathbf{X}^\# \mathbf{b}(k) + \mathbf{X}^\# \delta \mathbf{b}(k)$$

$$\mathbf{w}(k+1) = \mathbf{w}(k) + \mathbf{X}^\# \delta \mathbf{b}(k)$$

- let us denote

$$\mathbf{e}(k) = \mathbf{X}\mathbf{w}(k) - \mathbf{b}(k)$$

- we have **the following algorithm**:

$$\mathbf{w}(1) = \mathbf{X}^\# \mathbf{b}(1) , \quad b_i > 0$$

$$\mathbf{e}(k) = \mathbf{X}\mathbf{w}(k) - \mathbf{b}(k)$$

$$\delta \mathbf{b}(k) = c[\mathbf{X}\mathbf{w}(k) - \mathbf{b}(k) + |\mathbf{X}\mathbf{w}(k) - \mathbf{b}(k)|]$$

$$\delta \mathbf{b}(k) = c[\mathbf{e}(k) + |\mathbf{e}(k)|]$$

Where $|\mathbf{e}(k)|$ denotes the vector whose components are the absolute value of the component $\mathbf{e}(k)$

$$\mathbf{w}(k+1) = \mathbf{w}(k) + c\mathbf{X}^\# [\mathbf{e}(k) + |\mathbf{e}(k)|]$$

$$\mathbf{b}(k+1) = \mathbf{b}(k) + c[\mathbf{e}(k) + |\mathbf{e}(k)|]$$

- $\mathbf{w}(k+1)$ can be also calculated:

$$\mathbf{w}(k+1) = \mathbf{X}^\# \mathbf{b}(k+1)$$

Algorithm:

$$\mathbf{w}(1) = \mathbf{X}^\# \mathbf{b}(1) \quad \text{vector } \mathbf{b}(1) - \text{arbitrary but such that } b_i > 0 \\ i = 1, 2, \dots, n$$

$$\mathbf{w}(k+1) = \mathbf{w}(k) + c\mathbf{X}^\# [\mathbf{e}(k) + |\mathbf{e}(k)|]$$

$$\mathbf{e}(k) = \mathbf{X}\mathbf{w}(k) - \mathbf{b}(k)$$

- when the inequalities $\mathbf{X}\mathbf{w} > \mathbf{0}$ have solution the algorithm converges for $0 < c \leq 1$

- if *all* the components of $\mathbf{e}(k)$ cease to be positive (*but not all zero*) at any iteration step, this indicates that the classes are not linear separable

- When $\mathbf{e}(k) = \mathbf{0} \rightarrow$ the $\mathbf{w}(k)$ is a solution

$$\mathbf{e}(k) = \mathbf{X}\mathbf{w}(k) - \mathbf{b}(k)$$

$$\mathbf{X} \mathbf{w} = \mathbf{b}$$

Example:

-Training set: $\omega_1 = \{(0, 0)^T, (0, 1)^T\}$

$$\omega_2 = \{(1, 0)^T, (1, 1)^T\}$$

Augmented pattern vectors:

$$\omega_1 = \{(0, 0, 1)^T, (0, 1, 1)^T\}$$

$$\omega_2 = \{(1, 0, 1)^T, (1, 1, 1)^T\}$$

Pattern vectors from ω_2 have to be **multiplied** by (-1)

$$\omega_2 = \{(-1, 0, -1)^T, (-1, -1, -1)^T\}$$

$$\mathbf{b}(1) = (1, 1, 1, 1)^T \text{ and } c = 1$$

Form the matrix \mathbf{X} :

$$\begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 1 \\ -1 & 0 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$

- Generalized inverse $\mathbf{X}^\# = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$

$$\mathbf{X}^\# = \frac{1}{2} \begin{bmatrix} -1 & -1 & -1 & -1 \\ -1 & 1 & 1 & -1 \\ 3/2 & 1/2 & -1/2 & 1/2 \end{bmatrix}$$

1. Step

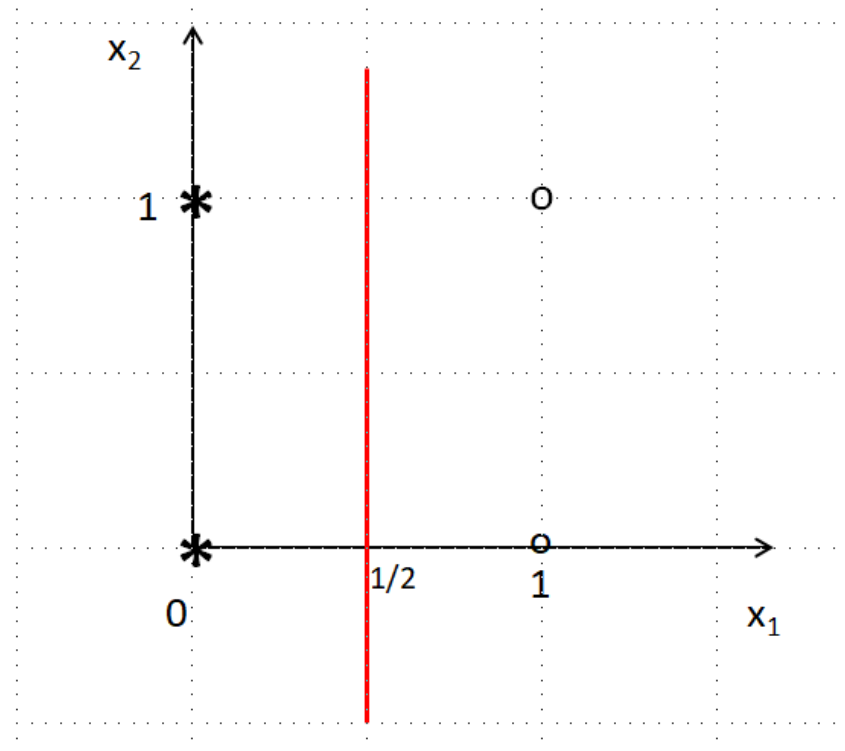
$$\mathbf{w}(1) = \mathbf{X}^\# \mathbf{b}(1) = \frac{1}{2} \begin{bmatrix} -1 & -1 & -1 & -1 \\ -1 & 1 & 1 & -1 \\ 3/2 & 1/2 & -1/2 & 1/2 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} -2 \\ 0 \\ 1 \end{bmatrix}$$

$$\mathbf{e}(1) = \mathbf{X} \mathbf{w}(1) - \mathbf{b}(1) = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 1 \\ -1 & 0 & -1 \\ -1 & -1 & -1 \end{bmatrix} \begin{bmatrix} -2 \\ 0 \\ 1 \end{bmatrix} - \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$\mathbf{e}(1) = \mathbf{0}$ $\mathbf{w}(1)$ is solution!

$$\mathbf{w}(1) = (-2, 0, 1)^T$$

$$d(\mathbf{x}) = -2x_1 + 1$$



Example:

-Training set: $\omega_1 = \{(0, 0)^\top, (1, 1)^\top\}$

$$\omega_2 = \{(0, 1)^\top, (1, 0)^\top\}$$

Augmented pattern vectors:

$$\omega_1 = \{(0, 0, 1)^\top, (1, 1, 1)^\top\}$$

$$\omega_2 = \{(0, 1, 1)^\top, (1, 0, 1)^\top\}$$

Pattern vectors from ω_2 have to be **multiplied** by (-1)

$$\omega_2 = \{(0, -1, -1)^\top, (-1, 0, -1)^\top\}$$

$$\mathbf{b}(1) = (1, 1, 1, 1)^\top \text{ and } c = 1$$

$$\mathbf{X} = \begin{bmatrix} 0 & 0 & 1 \\ 1 & 1 & 1 \\ 0 & -1 & -1 \\ -1 & 0 & -1 \end{bmatrix}$$

Generalized inverse $\mathbf{X}^\# = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$

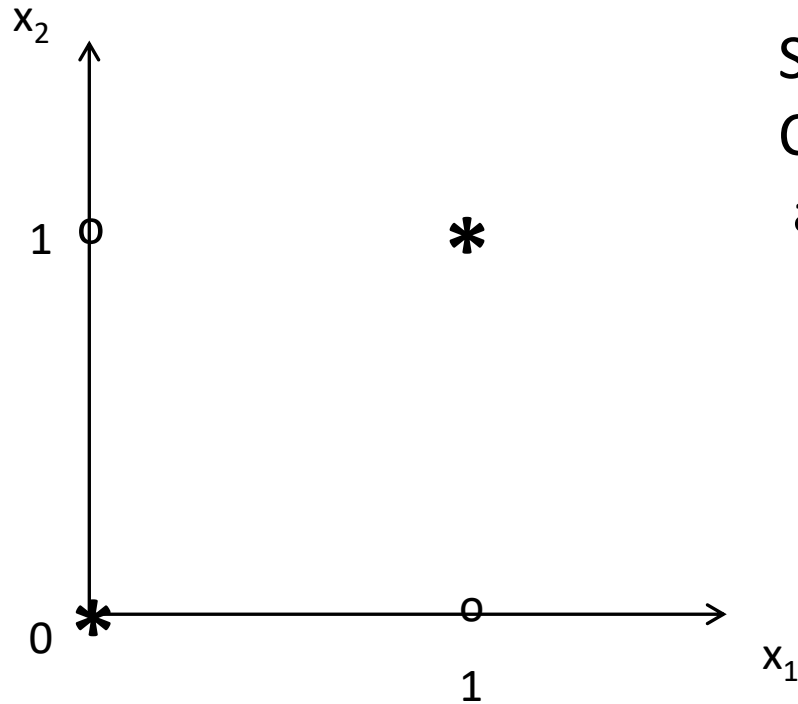
$$\mathbf{X}^\# = \frac{1}{2} \begin{bmatrix} -1 & 1 & 1 & -1 \\ -1 & 1 & -1 & 1 \\ 3/2 & -1/2 & -1/2 & -1/2 \end{bmatrix}$$

1. Step

$$\mathbf{w}(1) = \mathbf{X}^\# \mathbf{b}(1) = \frac{1}{2} \begin{bmatrix} -1 & 1 & 1 & -1 \\ -1 & 1 & -1 & 1 \\ 3/2 & -1/2 & -1/2 & -1/2 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\mathbf{e}(1) = \mathbf{X} \mathbf{w}(1) - \mathbf{b}(1) = \begin{bmatrix} 0 & 0 & 1 \\ 1 & 1 & 1 \\ 0 & -1 & -1 \\ -1 & 0 & -1 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} - \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

$\mathbf{e}(1) = (-1, -1, -1, -1)^T$ – negative vector indicates that $\mathbf{X} \mathbf{w} > \mathbf{0}$
has no solution



So-called XOR problem
Classes ω_1 and ω_2
are not linearly separable

* $\in \omega_1$

o $\in \omega_2$