

Machine Learning for Intelligent Systems

Lecture 6: Linear Classifiers and Perceptron

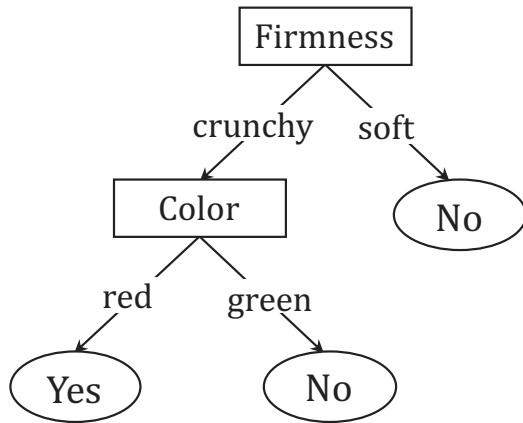
Reading: UML 9.1

Instructors: Nika Haghtalab (this time) and Thorsten Joachims

Hypothesis Spaces

$(\text{color} = \text{red}) \wedge (\text{size} = \text{small})$

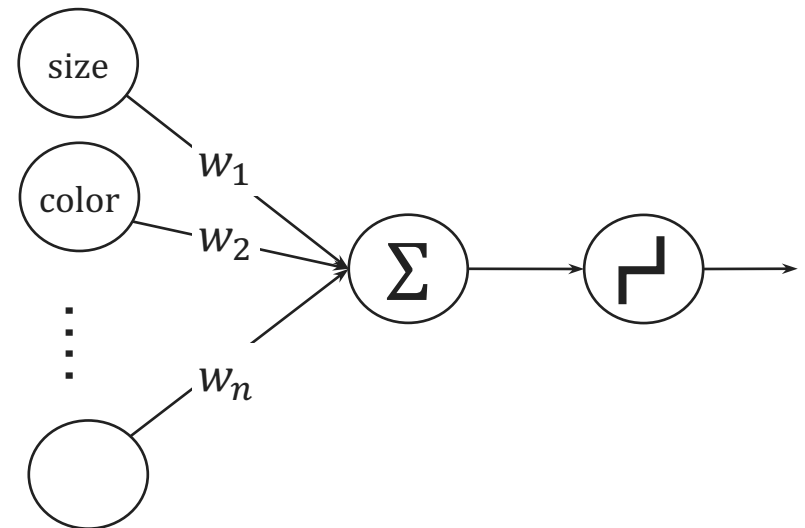
AND of feature-values



Decision Trees

$(\text{color} = \text{red}) \vee (\text{size} = \text{small})$

OR of feature-values



Linear Classifiers

Encoding in Euclidean Space

Represent apples as vectors in \mathbb{R}^d .

- Old:

$$X = \{A, B\} \times \{\text{red, green}\} \times \{\text{large, medium, small}\} \times \{\text{crunchy, soft}\}.$$

- New:

$$\rightarrow X \subseteq \mathbb{R}^4.$$

$$x_1 \text{ farm: } A \rightarrow 1, B \rightarrow -1$$

$$x_2 \text{ color: } \text{red} \rightarrow 1, \text{green} \rightarrow -1$$

$$x_3 \text{ size: } \text{large} \rightarrow 1, \text{medium} \rightarrow 0, \text{small} \rightarrow -1$$

$$x_4 \text{ firmness: } \text{crunchy} \rightarrow 1, \text{soft} \rightarrow -1.$$

$$\rightarrow Y = \{-1, +1\}: \text{Tasty} \rightarrow +1, \text{Not Tasty} \rightarrow -1$$

Reuters Business News text classification:

- 9947 keywords (more accurately, word “stems”)
- $X = \{0,1\}^{9947}$, where $x_i = 1$ if the keyword i appears in document.
- $Y = \{-1, +1\}$.

Linear Classifiers

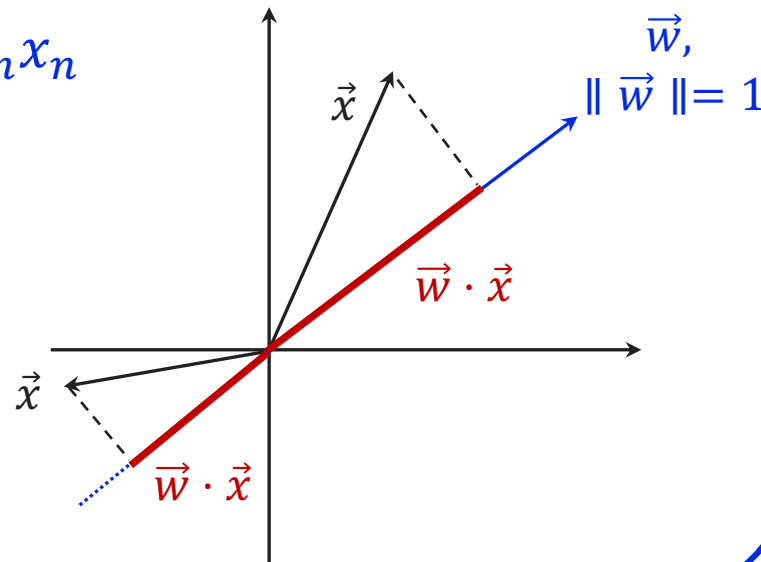
For a vector $\vec{w} \in \mathbb{R}^d$ and $b \in \mathbb{R}$, the hypothesis $h_{\vec{w},b}: \mathbb{R}^d \rightarrow \mathbb{R}$ defined below is called a **linear classifier/linear predictor/halfspace**,

$$h_{\vec{w},b}(\vec{x}) = \text{sign}(\vec{w} \cdot \vec{x} + b) = \begin{cases} +1 & \vec{w} \cdot \vec{x} + b > 0 \\ -1 & \vec{w} \cdot \vec{x} + b \leq 0 \end{cases}$$

Recall: Dot products

For two vectors: $\vec{w} = (w_1, w_2, w_3, \dots, w_n)$ and $\vec{x} = (x_1, x_2, x_3, \dots, x_n)$.

- $\vec{w} \cdot \vec{x} = w_1x_1 + w_2x_2 + w_3x_3 + \dots + w_nx_n$
- $\vec{w} \cdot \vec{x}$ is the (signed) length of the projection of \vec{x} on unit vector \vec{w} .

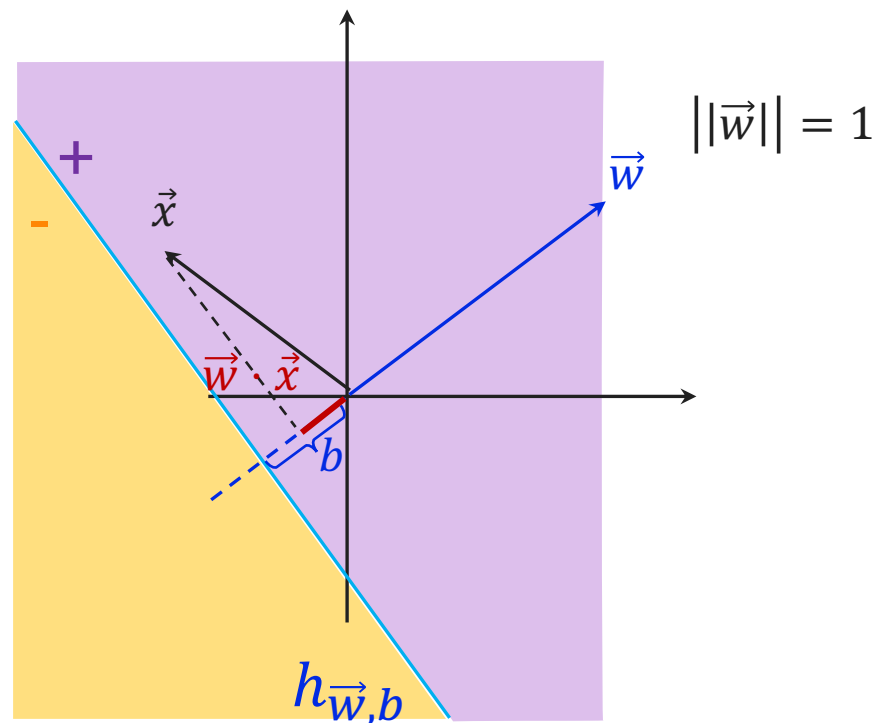
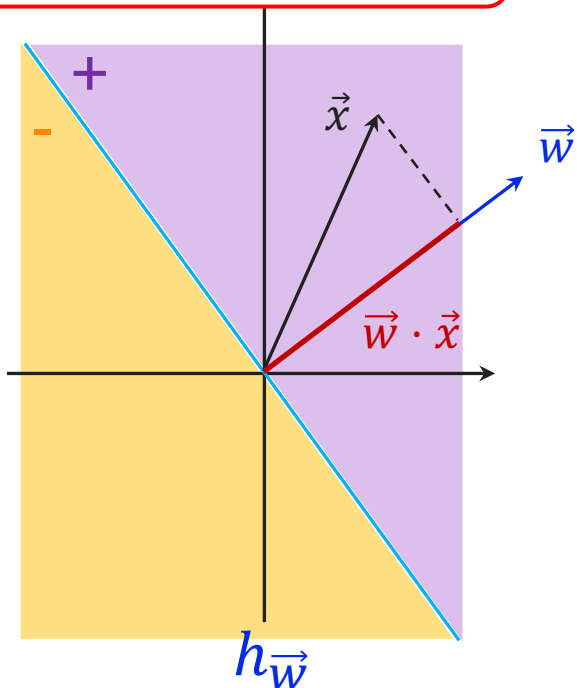


Linear Classifiers

For a vector $\vec{w} \in \mathbb{R}^d$ and $b \in \mathbb{R}$, the hypothesis $h_{\vec{w},b}: \mathbb{R}^d \rightarrow \mathbb{R}$ defined below is called a **linear classifier/linear predictor/halfspace**

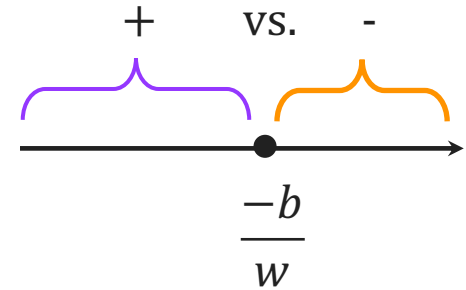
$$h_{\vec{w},b}(\vec{x}) = \text{sign}(\vec{w} \cdot \vec{x} + b) = \begin{cases} +1 & \vec{w} \cdot \vec{x} + b > 0 \\ -1 & \vec{w} \cdot \vec{x} + b \leq 0 \end{cases}$$

Homogenous linear classifier: $b = 0$

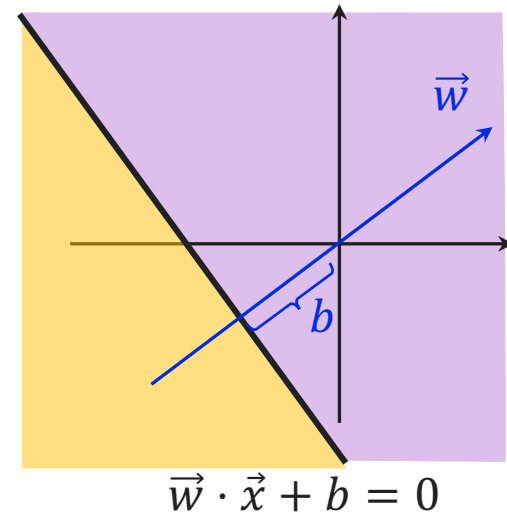


Linear Classifiers in all dimensions

- One-dimension: $h_{w,b}(x) = \text{sign}(wx + b)$
→ Decision boundary: point



- Two-dimension: $h_{\vec{w},b}(x) = \text{sign}(\vec{w} \cdot \vec{x} + b)$
→ Decision boundary: line



- d-dimension: $h_{\vec{w},b}(x) = \text{sign}(\vec{w} \cdot \vec{x} + b)$
→ Decision boundary: hyperplane $\vec{w} \cdot \vec{x} + b = 0$

Representational Power

Assume that x_1, x_2 take are binary values 0, 1. Represent the following using linear thresholds.

- $x_1 \wedge x_2$

- $x_1 \vee x_2$

- $x_1 \oplus x_2$

→ \oplus represent XOR, where $x_1 \oplus x_2 = 1$ when exactly one of x_1 and x_2 is set to 1.

Homogenous vs. Non-homogenous

Any d -dimensional learning problem for **non-homogenous linear classifiers** has a **homogenous** form in $(d+1)$ dimension.

Non-Homogenous $HS^d = \{h_{\vec{w},b} \mid \vec{w} \in \mathbb{R}^d, b \in \mathbb{R}\}$	Homogenous $HS_{homogenous}^{d+1} = \{h_{\vec{w}'} \mid \vec{w}' \in \mathbb{R}^{d+1}\}$
\vec{x}	$\vec{x}' = (\vec{x}, +1)$
\vec{w}, b	$\vec{w}' = (\vec{w}, b)$
$\vec{w} \cdot \vec{x} + b$	$\vec{w}' \cdot \vec{x}' = \vec{w} \cdot \vec{x} + b$

Without loss of generality, focus on **homogenous linear classifiers**.

Find a consistent classifier

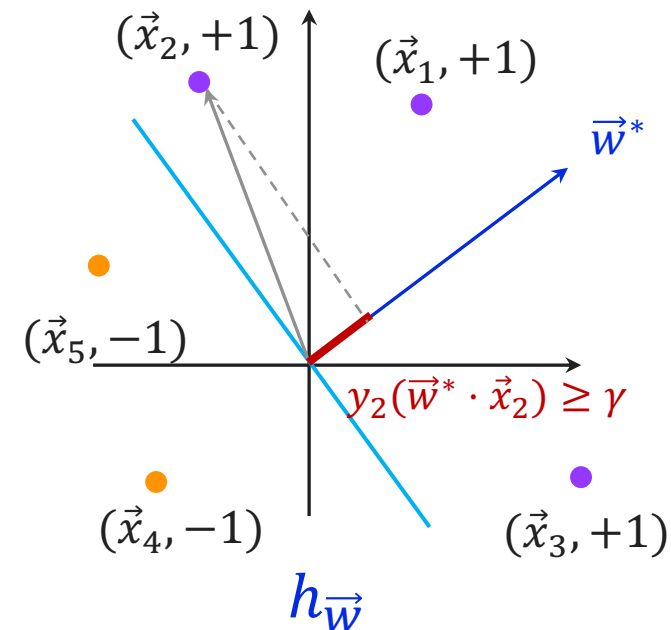
If there is a homogeneous linear classifier that is consistent with $\{(\vec{x}_1, y_1), (\vec{x}_2, y_2), \dots, (\vec{x}_m, y_m)\}$, how can we find it?

Unit vector \vec{w}^* is such that for all (\vec{x}_i, y_i) , $y_i(\vec{w}^* \cdot \vec{x}_i) \geq \gamma > 0$.

On the correct side,
with wiggle room

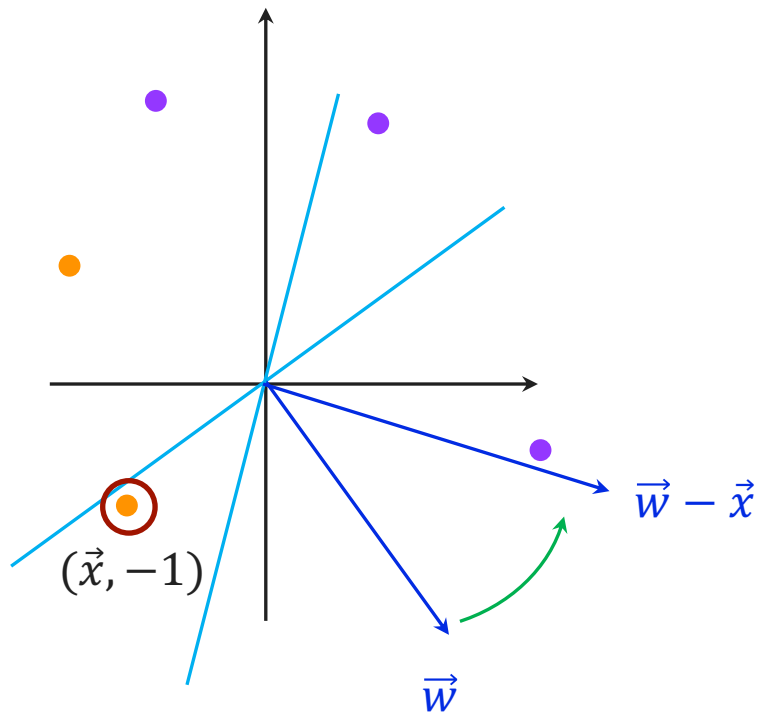
We want to find a \vec{w} such that $y_i(\vec{w} \cdot \vec{x}_i) > 0$ for all (\vec{x}_i, y_i) .

Can be done with a linear program

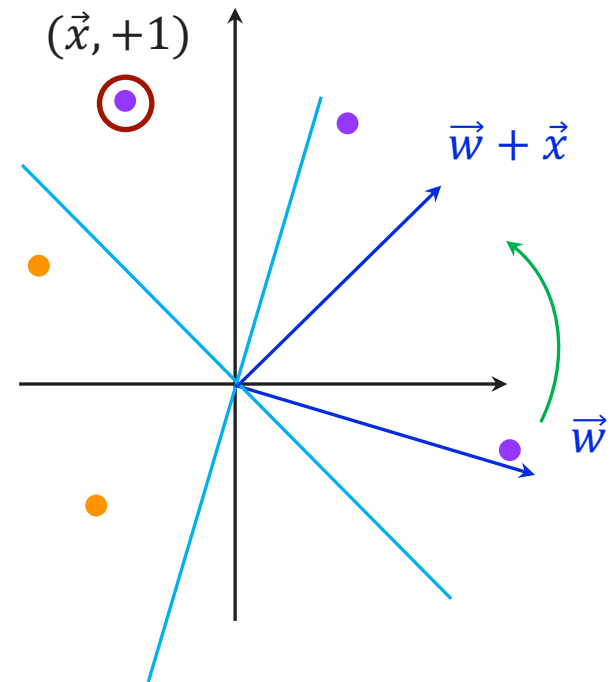


Improving a linear classifier

Start with a guess and improve it.



Move away from negative misclassified points



Move towards positive misclassified points

Perceptron (homogeneous & batch)

Input: Training data set $\{(\vec{x}_1, y_1), (\vec{x}_2, y_2), \dots, (\vec{x}_m, y_m)\}$

Initialize $\vec{w}^{(0)} = (0, \dots, 0)$, $t = 0$

While there is $i \in [m]$, such that $y_i (\underbrace{\vec{w}^{(t)} \cdot \vec{x}_i}_{\text{misclassified}}) \leq 0$ then,

- $\vec{w}^{(t+1)} = \vec{w}^{(t)} + y_i \vec{x}_i$ $\begin{cases} \vec{w}^{(t)} + \vec{x}_i & \text{for positive instances} \\ \vec{w}^{(t)} - \vec{x}_i & \text{for negative instances} \end{cases}$
- $t \leftarrow t + 1$

End While

Output $\vec{w}^{(t)}$



Frank Rosenblatt
@ Cornell!

Example: Reuters Text Classification

