

# Perceptron Algo Continued

Recall:

- linear decision function  $f(x) = w \cdot x$  (no  $b$  for simplicity)
- decision boundary:  $\{x: f(x) = 0\}$  (hyperplane through origin)
- sample points  $x_1, \dots, x_n \in \mathbb{R}^d$ , class labels  $y_1, \dots, y_n = \pm 1$
- goal: find weights such that  $y_i x_i \cdot w \geq 0$  (constraints are satisfied, all points classified properly)
- goal, revised: find  $w$  that minimizes  $R(w) = \sum_{i \in V} -y_i x_i \cdot w$ ,  $V = \text{set of } i \text{ where } y_i x_i \cdot w < 0$

$w$  is in a different space than our hyperplane  
"w space"

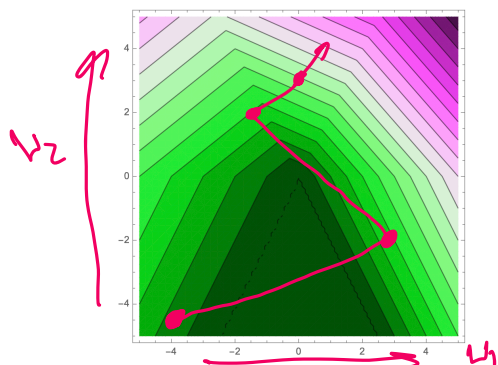
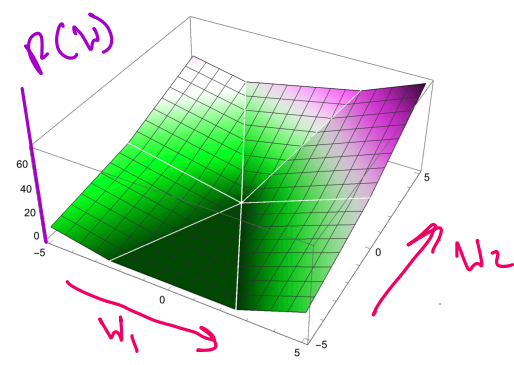
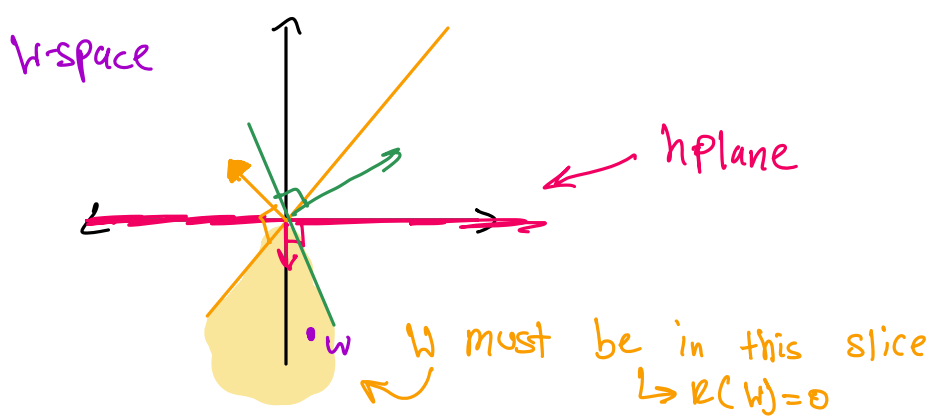
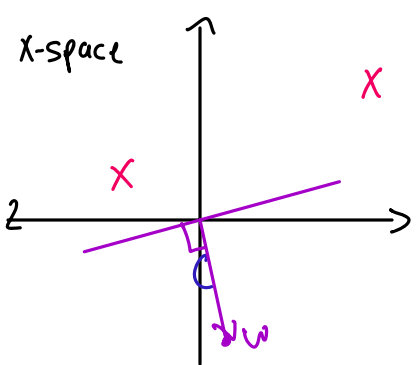
Objects in  $x$ -space transform to objects in  $w$ -space

$x$ -space	$w$ -space
hyperplane: $\{z: w \cdot z = 0\}$	point: $w$ <span style="color: green;">normal vector</span>
Point: $x$	hyperplane: $\{z: x \cdot z = 0\}$

Consider point  $x$  on  $h$ -plane  $\{z: w \cdot z = 0\} \iff w \cdot x = 0$

$\iff$  point  $w$  on  $h$ -plane  $\{z: x \cdot z = 0\}$

If we want  $x_i \cdot w \geq 0 \rightarrow$  in  $x$ -space,  $x$  should be on same side as  $w$  if  $x$  is a positive sample  
in  $w$ -space,  $w$  should be on same side of  $h$ -plane as  $x$ ,  $x$  is the normal vector



Optimization Algorithm: **gradient descent** on  $R$

no gradient @  
 $w = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

- Given a starting point  $w$ , find gradient of  $R$  w.r.t.  $w$  (direction of steepest ascent)
- Take step in opposite direction

Recall:  $\nabla R(w) = \begin{bmatrix} \frac{\partial R}{\partial w_1} \\ \vdots \\ \frac{\partial R}{\partial w_d} \end{bmatrix}$  and  $\nabla_w (z \cdot w) = \begin{bmatrix} z_1 \\ \vdots \\ z_d \end{bmatrix} = \vec{z}$

$$\rightarrow \nabla R(w) = \sum_{i \in V} \nabla -y_i x_i \cdot w = \sum_{i \in V} -y_i x_i$$

At any point  $w$ , walk downhill in direction of steepest descent ( $-\nabla R(w)$ )

- doesn't work if not linearly separable
- 1) arbitrary  $w \neq \vec{0}$  (good choice is any  $y_i x_i$ )  
 $\rightarrow$  will classify at least 1 point correctly
  - 2) while  $R(w) > 0$   
 $V = \text{Set of indices } i \text{ s.t. } y_i x_i < 0$   
 $w = w + \xi \sum_{i \in V} y_i x_i$  (where  $\sum_{i \in V} y_i x_i = -\nabla R(w)$ )
  - 3) return  $w$

Newton's method?

$\rightarrow$  not on piecewise linear function

$\xi > 0$  is the **step size / learning rate**

PROBLEM: slow algo. each step takes  $O(nd)$  time

2nd optimization algo: **stochastic gradient descent**

Idea: in each step, pick **one** misclassified point  $x_i$ ;

do gradient descent on  $L(x_i \cdot w, y_i)$

$\swarrow$  dim of sample points  
= # features

Called the **perceptron algorithm**. each step takes  $O(d)$  time

While some  $y_i x_i \cdot w < 0$ :

$$w = w + \xi y_i x_i$$

return  $w$

\* SGD can't always replace vanilla gradient descent

What if separating h-plane doesn't pass through origin?

↳ add a fictitious dimension

$$\rightarrow f(x) = w \cdot x + d = \begin{bmatrix} w_1 & w_2 & d \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ 1 \end{bmatrix}$$

now: Sample points are in  $\mathbb{R}^{d+1}$ ,  $x_{d+1} = 1$

→ run perceptron algo in  $(d+1)$ -dimensional space

online algo: If new data comes along while training, you can incorporate them

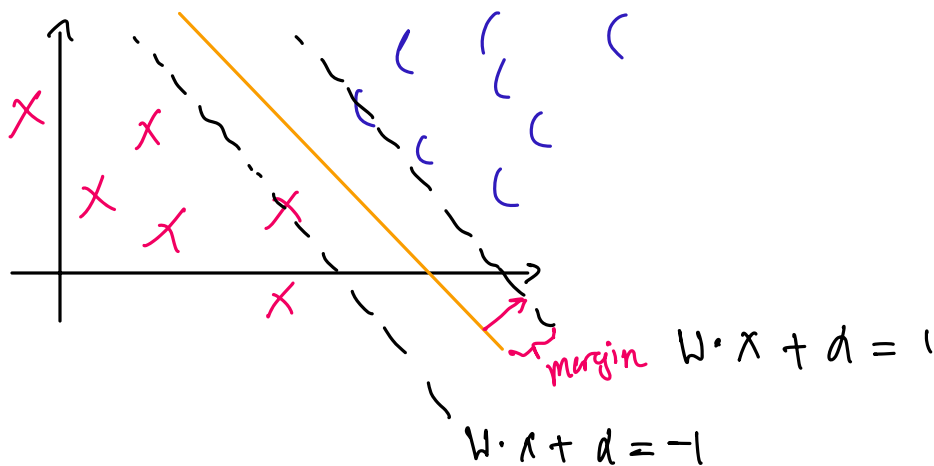
Perceptron Convergence Thm.

$$O\left(\frac{r^2}{\gamma^2}\right) \quad r = \max \|x_i\| \quad \text{"radius of data"}$$

$$\gamma = \text{max margin}$$

Maximum Margin Classifiers

The margin of a linear classifier is the distance from the decision boundary to the nearest sample point



enforce some constraints:

✓ 0 in perceptron algo.

$$y_i [W \cdot x_i + d] \geq 1 \quad \text{for } i \in [1 \dots n]$$

↖ makes it impossible for  $W = \vec{0}$

Recall: If  $\|W\| = 1$ , then signed distance from h-plane to  $x_i$  is  $W \cdot x_i + d$

for  $\|\vec{w}\| \neq 1$ : Signed distance is  $\frac{w}{\|\vec{w}\|} \cdot x_i + \frac{d}{\|\vec{w}\|}$

→ margin is  $\min_i \frac{1}{\|\vec{w}\|} \cdot \underbrace{|w \cdot x_i + d|}_{\geq 1} \geq \frac{1}{\|\vec{w}\|}$

→ maximize margin  $\frac{1}{\|\vec{w}\|} = \text{minimizing } \|\vec{w}\|$  Subject to  $y_i [w \cdot x_i + d] \geq 1$

optimization

problem: find  $w$  &  $d$  that minimize  $\|\vec{w}\|^2$  s.t.  $y_i [x_i \cdot w + d] \geq 1 \forall i \in [1..n]$  ✓  $\|\vec{w}\|^2$  is smooth, even at origin

called a **quadratic program** in  $\mathbb{R}^{d+1}$  - (since  $d$  brings 1 dimension)  
-  $n$  constraints

One unique solution  $^*$  if data is linearly separable

Solution: **maximum margin classifier / hard-margin SVM**

- At the optimal solution, the margin is exactly  $\frac{1}{\|\vec{w}\|}$  when  $\|\vec{w}\|$  is maximized
- there is a slab w/ width  $\frac{2}{\|\vec{w}\|}$  w/ no data points

