<u>Perceptron</u> : for binary classification.
  - Algorithm for learning <mark>half spaces.</mark>

  • Script $C$ = {half spaces}

definition of half space : $f(x) = \text{sign}\left(\sum_{i=1}^{n} w_i \cdot x_i - \theta\right)$ ⟵ threshold.

dot product

vector $w$ = unknown.
$f(x) \in \{-1, +1\}$
$x_i \in \mathbb{R}^n$.

A "complicated" algorithm for learning halfspace :

$(x^{(1)}, +1) \longrightarrow w_1 x_1^{(1)} + \ldots + w_n x_n^{(1)} \geq \theta$
        ↑
      label

$(x^{(2)}, -1) \longrightarrow w_1 x_1^{(2)} + \ldots + w_n x_n^{(2)} < \theta$.

       ⋮                        ⋮

  $m$ training examples $\longrightarrow$ $m$ inequalities.

we can use linear programming to find
weight vector $w$ consistent w/ training set.
  • LP is very expensive, use perceptron.

Perceptron :

1) Initially, $w^0 = (0, \cdots, 0)$ or unit $w^0 = \left(\frac{1}{\sqrt{n}}, \cdots, \frac{1}{\sqrt{n}}\right)$

2) Learner has $w$ as its state.

3) Teacher presents challenge :

$$\overset{\longleftarrow}{X \in \mathbb{R}^n} \quad \text{Teacher}$$

4) Learner responds with $\text{sign}(w \cdot x)$, which is current state.

$$\text{Learner} \xrightarrow{\quad \text{sign}(w \cdot x) \quad}$$

5) If mistake is made :

    Case 1 : $X$ was truly a negative example :
$$W_{new} = W_{old} - X .$$

    Case 2 : $X$ was truly a positive example :
$$W_{new} = W_{old} + X .$$

    Equivalent way to update rule ( applies only when mistake).
$$W_{new} = W_{old} + y \cdot \vec{x}$$
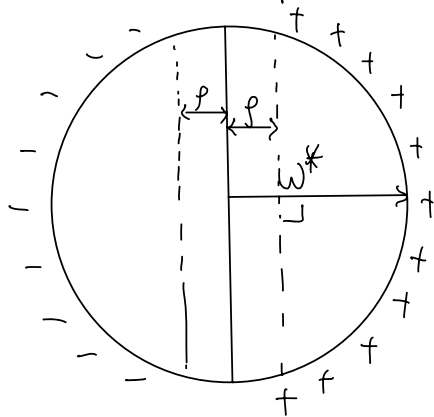$$\underset{\nwarrow}{} \text{label} \in \{-1, +1\}$$

Assumptions :

· Assume $\exists \, w^*$, which is true unknown weight vector, which has norm 1, $\|w^*\| = 1$.

· Assume $X$ has norm 1, $\|X\| = 1$.

· $\theta = 0$.

<u>Main Assumption</u> : There exists a margin $\rho$ where all points are at least distance $\rho$ from $W^*$.



All +/- points have distance $\geq \rho$ from half space.

$$\|x\| = \|W^*\| = 1 \quad (\text{assumption})$$

Equivalently, $|\langle X, W^* \rangle| \geq \rho$

$\underbrace{\qquad\qquad\qquad\qquad}$

"margin Assumption."

$\rho$ = "cushion" of the classifier.

<u>Main Theorem</u> : "Perceptron convergence theorem" ;

The mistake bound of the perceptron algorithm is $O\left(\frac{1}{\rho^2}\right)$ mistakes.

proof :

Recall update step : $W_{new} = W_{old} + y \cdot X$

Let's say: $w$ is current state of learner,
$w^*$ is true normal to half space.

Claim 1 : on every mistake, $w \cdot w^*$ increases by at least $\rho$.

claim 2 : $\|w\|^2$ increases after every mistake by at most 1.

Question : How to obtain $O\left(\frac{1}{\rho^2}\right)$ mistake bound given claims 1 and 2?

Let $t = $ # mistakes we've made at some point during execution.

$$\underbrace{t \cdot \rho \leq w \cdot w^*}_{\text{Claim 1}} \leq \|w\| \|w^*\| \underset{\text{Claim 2}}{\overset{\leq \sqrt{t} \text{ since } \underline{\|w\|^2 \leq t}}{\longrightarrow}} 1$$

↑ Cauchy-Schwartz inequality

$$\Rightarrow t \rho \leq \sqrt{t}$$

$$\Rightarrow \boxed{t \leq \frac{1}{\rho^2}}$$

Proving the claims:

**Claim 1:** $w \cdot w^*$ increases on every mistake by $\geq \rho$.

$w_{new} = w_{old} + y\vec{x}$ (only when there is a mistake).

$$\Rightarrow w_{new} \cdot w^* = (w_{old} + y \cdot \vec{x}) \cdot \vec{w}^*$$

$$= (w_{old} \cdot w^* + \underbrace{y \cdot \vec{x} \cdot \vec{w}^*}_{})$$

−1 since updates only when mistake.

⊖ since mistake is made, and we know the absolute value $\geq \rho$ $(\leq -\rho)$ per the margin assumption

$$\underbrace{\phantom{xxxxxxxxxxxxxxxxxxxxx}}_{\geq \rho}$$

proving claim 2:

$\|w\|^2$ increases by ==at most 1== on every mistake.

$$\|W_{new}\|^2 = \|W_{old} + y \cdot X\|^2 = \|W_{old}\|^2 + 2 \cdot y \underbrace{\langle \vec{X}, W_{old} \rangle}_{} + \underbrace{\|X\|^2}_{1}$$

negative, bc

$\langle \vec{X} \cdot W_{old} \rangle$ has different sign from $y$ since mistake is made, therefore their product = negative.

$$\underbrace{\qquad\qquad}_{\leq 1} \qquad \blacksquare$$

Look back at some assumptions:

• $\theta = 0$ : Add a new feature, call it $X_{n+1}$

$(X_1, \cdots X_n) \longrightarrow (X_1, \cdots, X_{n+1})$
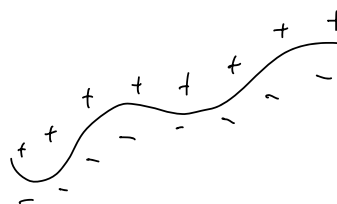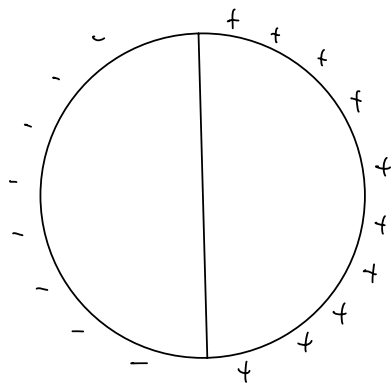
$\hookrightarrow$ always set to 1.

• $\|x\| = 1$. If $\|x\| = R$. $\Rightarrow$ M.B. $O\left(\frac{R^2}{\rho^2}\right)$

# Perceptron Learning: Polynomial Threshold Functions (PTFs)

Definition: $f = \text{sign}\left(p(x)\right)$

$p$ is a multivariate polynomial of degree $d$.

How can we use perceptron to learn this function class?



Higher degree polynomial

feature map:

$$(x_1, \cdots, x_n) \rightarrow X_1^2, X_1 X_2, \cdots, X_n^2 = n^2 \text{ new variables.}$$

$\underbrace{\phantom{(x_1, \cdots, x_n)}}_{\text{degree 2}}$

Call them $Y_1, \cdots, Y_{N=n^2}$, $f = \text{sign}\left(\sum_{i=1}^{N} w_i Y_i\right)$

---

Learning PTFs of degree $d$ is equivalent to learning halfspaces in $n^d$ dimensions.

---

↳ can run perceptron to learn this half-space in higher dimension.

• Runtime: just computing the feature map takes time $n^d$

• What is the margin in this $n^d$ dimension space?
   May be costly.

we can save on the running time, using something called the kernel trick!

# Perceptron Learning: Kernel Functions:

$X \longleftarrow$ input

☆ Denote $\varphi(x)$ to be the image of $x$ in the feature space.

$$X \in R^n \xrightarrow[\text{function}]{\text{kernel}} \varphi(x) \in R^{n^d}$$

$K(X^1, X^2)$ outputs $\langle \varphi(X^1), \varphi(X^2) \rangle$

$\Bigg($ and let's assume $K(X^1, X^2)$ is easy to compute.
kernel function $\qquad \underset{\text{kernel function.}}{\uparrow}$

Kernel Perception (want to work in $R^{n^d}$ )

$W = 0^{n^d} = \langle 0$ vector of length $n^d \rangle$

Let's Assume we make a mistake on first try (point $X^1$)

$$W_{new} = W_{old} + y\,\varphi(x)$$

$\qquad\qquad\quad \underset{\substack{\text{0 vector bc} \\ \text{first mistake.}}}{\downarrow} \qquad \text{in } n^d \text{ space}$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad X^{[2]}$ is the new point.

we need to evaluate $W_{new} \odot \varphi(X^{(2)})$

$$= \langle y \cdot \varphi(x^1), \varphi(x^2) \rangle \qquad (1)$$

$$= y \cdot \underline{K(x^{(1)}, x^{(2)})}.$$

$\quad \underset{\text{Scalar}}{\uparrow} \quad$ easy to compute.

Note: $K(X^1, X^2) = \langle \varphi(X^1), \varphi(X^2) \rangle$

therefore (1) is $y \cdot K(X^1, X^2)$

$$W_{t+1} = \sum_{i=1}^{t} y^i \cdot \varphi(x^i) \quad \in \mathbb{R}^{n^d}$$

If need to compute $\langle W_{t+1} , \varphi(x^{t+1}) \rangle$:

$$\sum_{i=1}^{t} y^i \cdot \underbrace{\langle \varphi(x^i) , \varphi(x^{t+1}) \rangle}_{K(x^i, x^{t+1})}$$

efficiency computable.

So, we're able to simulate the perceptron algorithm in a much higher dimensional space with just low dimensional vectors and this kernel function!

## Example of Kernel Function:

Consider degree-2 polynomial threshold function:

$$\varphi(x_1, \ldots, x_n) = (x_1 x_1 , x_1 x_2 , \ldots \ x_{n-1} x_n , x_n^2)$$

$$K(x, z) = \underbrace{\langle \varphi(x) , \varphi(z) \rangle}_{\text{inner product}} =$$
$(x, z \in \mathbb{R}^n)$

$$= (x_1^2 z_1^2 + x_1 x_2 z_1 z_2 + \cdots + x_n^2 z_n^2)$$

$$= \sum_{i, j} x_i x_j z_i z_j = \left( \sum_{i=1}^{n} x_i z_i \right) \cdot \left( \sum_{j=1}^{n} x_j z_j \right)$$

$$= \boxed{\underbrace{(x \cdot z)^2}_{\text{inner product}} = K(x, z)}$$

Example of Kernel methods.

## Other Kernels

$$K(x, z) = (x \cdot z + C)^2$$

$$\varphi(x) = (x_1^2, \cdots, x_n^2, \sqrt{2C} \, x_1, \cdots, \sqrt{2C} \, x_n, C)$$

Gaussian Kernels = $K(x, z) = \underbrace{e^{-\|x-z\|_2^2}}_{\text{radial basis kernel.}}$