

# Lecture 17: Convergence Proof of Perceptron Algo, Kernel perceptron, Logistic Regression (Begin)

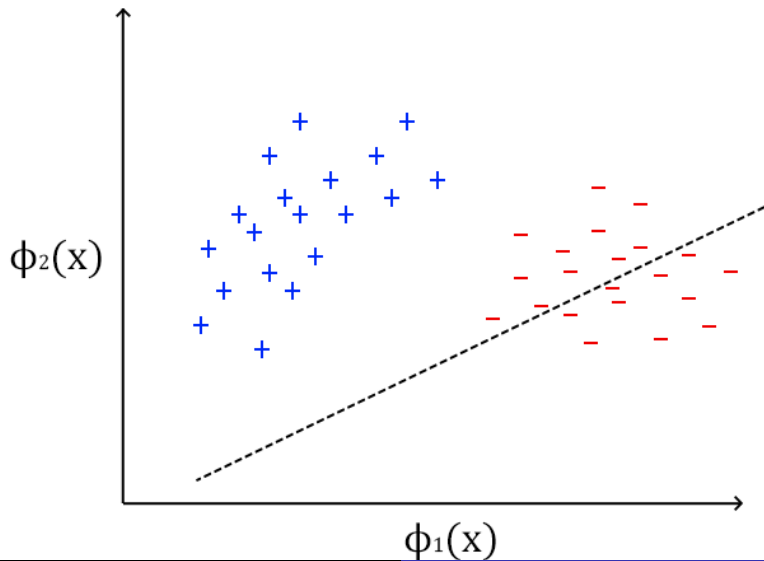
Instructor: Prof. Ganesh Ramakrishnan

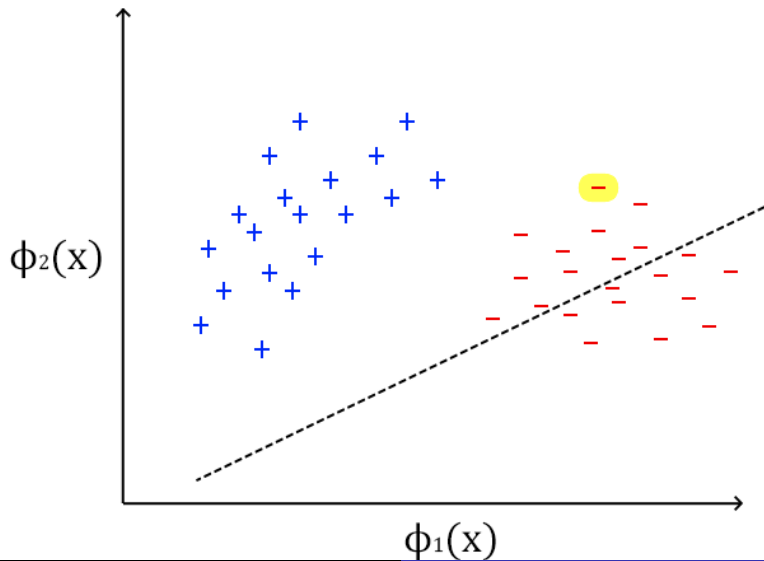
# Recap: Perceptron Update Rule

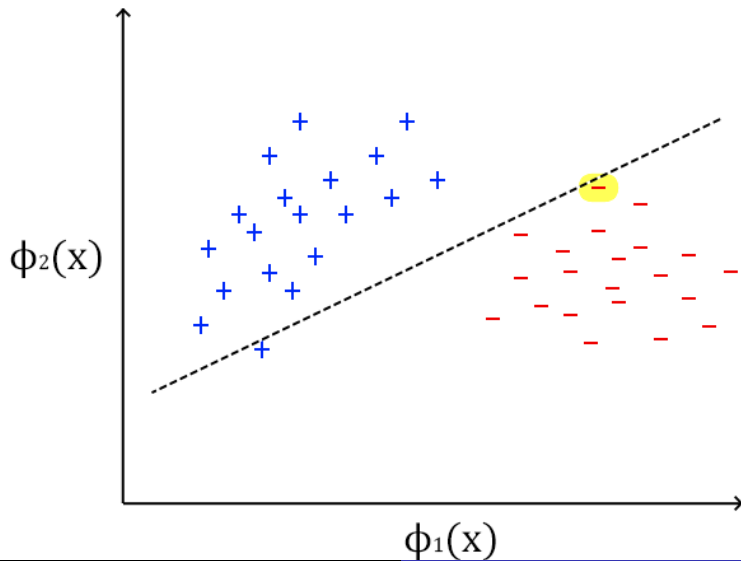
- Perceptron works for two classes ( $y = \pm 1$ ). A point is misclassified if  $y\mathbf{w}^T(\phi(\mathbf{x})) < 0$  (error is on those  $\mathbf{x}$  for which the "unsigned distance" turns out be "signed" !!)
- Perceptron Algorithm:
  - INITIALIZE:  $\mathbf{w} = \text{ones}()$
  - REPEAT: for each  $\langle \mathbf{x}, y \rangle$ 
    - If  $y\mathbf{w}^T\phi(\mathbf{x}) < 0$
    - then,  $\mathbf{w} = \mathbf{w} + \eta\phi(\mathbf{x}) \cdot y$
    - endif
- Intuition: } → By design, won't converge if dataset is NOT linearly separable! In practice you run this for some max # iterations

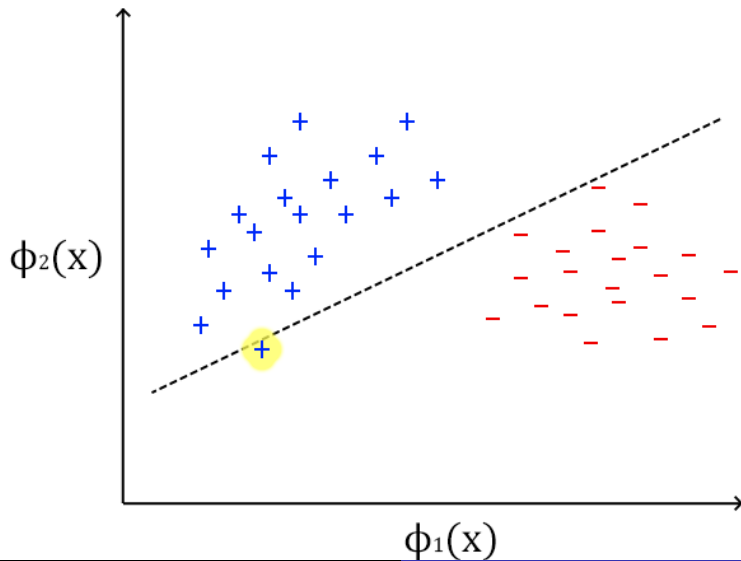
$$\begin{aligned} \underline{y(\mathbf{w}^{(k+1)})^T\phi(\mathbf{x})} &= y(\mathbf{w}^k + \eta y\phi^T(\mathbf{x}))\phi(\mathbf{x}) \\ &= y(\mathbf{w}^k)^T\phi(\mathbf{x}) + \eta y^2\|\phi(\mathbf{w})\|^2 \\ &\quad \text{unsigned distance becomes increasingly less neg.} \leftarrow \textcircled{>} \underline{y(\mathbf{w}^k)^T\phi(\mathbf{x})} \end{aligned}$$

Since  $y(\mathbf{w}^k)^T\phi(\mathbf{x}) \leq 0$ , we have  $y(\mathbf{w}^{(k+1)})^T\phi(\mathbf{x}) > y(\mathbf{w}^k)^T\phi(\mathbf{x}) \Rightarrow$  more hope that this point is classified correctly now.

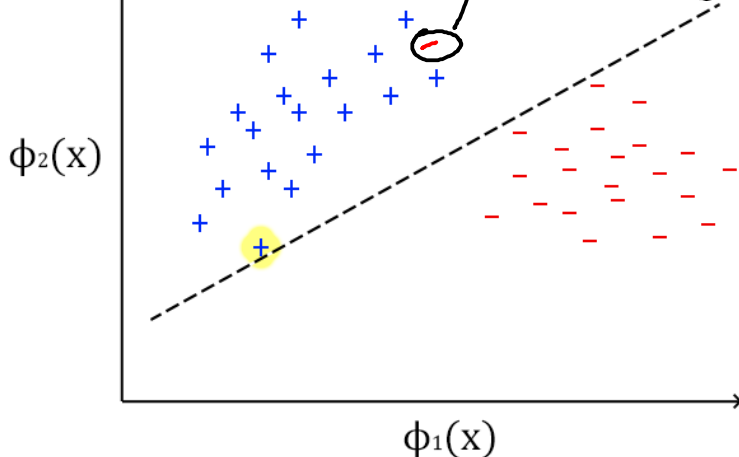








Especially in the presence of such noise,  
it is a good idea to alternate between  $+$  and  $-$   
points to reasonably ensure that "noise" pts are  
(possibly) left alone.



Efficiently computing # of misclassified pts

requires  $\mathbb{R}^m$   $\left\{ \begin{bmatrix} \phi(x^{(1)}) y^{(1)} \\ \phi(x^{(2)}) y^{(2)} \\ \phi(x^{(m)}) y^{(m)} \end{bmatrix} \begin{bmatrix} \omega^k \end{bmatrix} \right\} \mathbb{R}^n$

$\mathbb{R}^n$

#multiplications =  $mn$  → Can be empirically reduced using decompositions on  $\Phi y$  matrix.

#additions =  $m(n-1)$

Please use inbuilt libraries for matrix vector mult in scipy & numpy & tensorflow!



# Convergence of Perceptron Algorithm

# Perceptron Update Rule: Further analysis

(An intuitive sufficient condition)

- Formally,:- If  $\exists$  an optimal separating hyperplane with parameters  $\mathbf{w}^*$  such that,

$$\forall (\mathbf{x}, y), y\phi^T(\mathbf{x})\mathbf{w}^* \geq 0$$

then the perceptron algorithm converges.

**Proof:-** We want to show that

gap(k) =  $\|\mathbf{w}^{(k)} - \beta \mathbf{w}^*\|_2^2$  *At least  $\Theta^2$  (independent of k) reduction in each step!*

$$\lim_{k \rightarrow \infty} \|\mathbf{w}^{(k+1)} - \rho \mathbf{w}^*\|^2 = 0 \quad (1)$$

$\text{gap}(k+1) < \text{gap}(k) - \Theta^2 \rightarrow (3)$

(If this happens for some constant  $\rho$ , we are fine.)

$$\|\mathbf{w}^{(k+1)} - \beta \mathbf{w}^*\|_2^2 = \|\mathbf{w}^{(k)} + \underbrace{y\phi(\mathbf{x})}_{\text{independent of } k} - \beta \mathbf{w}^*\|_2^2 \rightarrow (1)$$

In our discussion, let  $\eta = 1$

Note:  $(\mathbf{x}, y)$  is such that  $y(\mathbf{w}^{(k)})^T \phi(\mathbf{x}) < 0 \rightarrow (2)$

# Perceptron Update Rule: Further analysis

- **Formally**,:- If  $\exists$  an optimal separating hyperplane with parameters  $\mathbf{w}^*$  such that,

$$\forall (\mathbf{x}, y), y\phi^T(\mathbf{x})\mathbf{w}^* \geq 0$$

then the perceptron algorithm converges.

**Proof**:- We want to show that

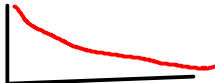
$$\lim_{k \rightarrow \infty} \|\mathbf{w}^{(k+1)} - \rho \mathbf{w}^*\|^2 = 0 \quad (1)$$

(If this happens for some constant  $\rho$ , we are fine.)

$$\|\mathbf{w}^{(k+1)} - \rho \mathbf{w}^*\|^2 = \|\mathbf{w}^k - \rho \mathbf{w}^*\|^2 + \underbrace{\|y\phi(\mathbf{x})\|^2 + 2y(\mathbf{w}^k - \rho \mathbf{w}^*)^T \phi(\mathbf{x})}_{\text{Desire to be } \leq -\Theta^2} \quad (2)$$

$\underbrace{\|\mathbf{w}^{(k+1)} - \rho \mathbf{w}^*\|^2}_{\text{gap}(k+1)} = \underbrace{\|\mathbf{w}^k - \rho \mathbf{w}^*\|^2}_{\text{gap}(k)} + \dots$   
Substitute  $\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + y\phi(\mathbf{x})$

Desire to be  $\leq -\Theta^2$   
 $1/k > 0$  but no constant bounds it



# Perceptron Update Rule: Further analysis

- **Formally**,:- If  $\exists$  an optimal separating hyperplane with parameters  $\mathbf{w}^*$  such that,

$$\forall (\mathbf{x}, y), y\phi^T(\mathbf{x})\mathbf{w}^* \geq 0$$

then the perceptron algorithm converges.

**Proof**:- We want to show that

$$\lim_{k \rightarrow \infty} \|\mathbf{w}^{(k+1)} - \rho\mathbf{w}^*\|^2 = 0 \quad (1)$$

(If this happens for some constant  $\rho$ , we are fine.)

$$\|\mathbf{w}^{(k+1)} - \rho\mathbf{w}^*\|^2 = \|\mathbf{w}^k - \rho\mathbf{w}^*\|^2 + \|y\phi(\mathbf{x})\|^2 + \underbrace{2y(\mathbf{w}^k - \rho\mathbf{w}^*)^T \phi(\mathbf{x})}_{\text{blue wavy line}} \quad (2)$$

- For convergence of perceptron, we need L.H.S. to be less than R.H.S. at every step, although by some small but non-zero value (with  $\theta \neq 0$ )

$$\|\mathbf{w}^{(k+1)} - \rho\mathbf{w}^*\|^2 \leq \|\mathbf{w}^k - \rho\mathbf{w}^*\|^2 - \underbrace{\theta^2}_{\text{blue wavy line}} \quad (3)$$

# Perceptron Update Rule: Further analysis

- **Need** that  $\|w^{(k+1)} - \rho w^*\|^2$  reduces by at least  $\theta^2$  at every iteration.

we know  $\rightarrow \|w^{(k+1)} - \rho w^*\|^2 = \|w^k - \rho w^*\|^2 + 2y(w^k - \rho w^*)^T \phi(x) + \|y\phi(x)\|^2$   
we need  $\rightarrow \|w^{(k+1)} - \rho w^*\|^2 \leq \|w^k - \rho w^*\|^2 - \theta^2$  (4)

- Based on (2) and (4), we **need** to find  $\theta$  such that,

$$-\theta^2 \geq 2y(w^k - \rho w^*)^T \phi(x) + \|y\phi(x)\|^2$$

If I show such a  $\theta^2$  exists, I know (4) will hold!

# Perceptron Update Rule: Further analysis

- **Need** that  $\|\mathbf{w}^{(k+1)} - \rho\mathbf{w}^*\|^2$  reduces by at least  $\theta^2$  at every iteration.

$$\|\mathbf{w}^{(k+1)} - \rho\mathbf{w}^*\|^2 \leq \|\mathbf{w}^k - \rho\mathbf{w}^*\|^2 - \theta^2 \quad (4)$$

- Based on (2) and (4), we **need** to find  $\theta$  such that,

$$\|\phi(\mathbf{x})\|^2 + 2y(\mathbf{w}^k - \rho\mathbf{w}^*)^T \phi(\mathbf{x}) \leq -\theta^2$$

$$(\|y\phi(\mathbf{x})\|^2 = \|\phi(\mathbf{x})\|^2 \text{ since } y = \pm 1)$$

- The number of iterations would be:  $O\left(\frac{\|\mathbf{w}^{(0)} - \rho\mathbf{w}^*\|^2}{\theta^2}\right)$
- Tutorial 6, Problem 4 has more concerning the number of iterations. But first we will discuss how convergence holds in the first place!

$O\left(\frac{r^2}{\theta^2}\right)$  independent of arbit  $\rho$  &  $\mathbf{w}^*$

# Perceptron Update Rule: Further analysis

$$\delta = -\min_{\mathbf{x} \in \mathcal{D}} 2y\mathbf{w}^{*T}\phi(\mathbf{x}) = \max_{\mathbf{x} \in \mathcal{D}} -2y\mathbf{w}^{*T}\phi(\mathbf{x}) \geq y\hat{\mathbf{w}}^T\phi(\mathbf{x}') \quad \forall \mathbf{x}', y'$$

- Observations:-

- ①  $y(\mathbf{w}^k)^T\phi(\mathbf{x}) < 0$  ( $\because \mathbf{x}$  was misclassified)

- ②  $\Gamma^2 = \max_{\mathbf{x} \in \mathcal{D}} \|\phi(\mathbf{x})\|^2$

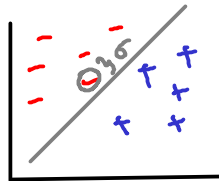
- ③  $\delta = -\min_{\mathbf{x} \in \mathcal{D}} 2y\mathbf{w}^{*T}\phi(\mathbf{x})$

(if  $\sigma$  = unsigned dist of closest pt in  $\mathcal{D}$  to  $\hat{\mathbf{w}}$ , then  $\delta = -2\sigma$ )

- Here, negative margin  $\delta = -2y\mathbf{w}^{*T}\phi(\hat{\mathbf{x}})$  is the negative of unsigned distance of closest point  $\hat{\mathbf{x}}$  from separating hyperplane :

$$\hat{\mathbf{x}} = \operatorname{argmax}_{\mathbf{x} \in \mathcal{D}} -2y\mathbf{w}^{*T}\phi(\mathbf{x}) = \operatorname{argmin}_{\mathbf{x} \in \mathcal{D}} y\mathbf{w}^{*T}\phi(\mathbf{x})$$

- Since the data is linearly separable,



Margin = 2 \* distance of closest point  
(will revisit for support vector classification)

# Perceptron Update Rule: Further analysis

- **Observations:-**

- 1  $y(\mathbf{w}^k)^T \phi(\mathbf{x}) < 0$  ( $\because \mathbf{x}$  was misclassified)

- 2  $\Gamma^2 = \max_{\mathbf{x} \in \mathcal{D}} \|\phi(\mathbf{x})\|^2$

- 3  $\delta = -\min_{\mathbf{x} \in \mathcal{D}} 2y\mathbf{w}^{*T} \phi(\mathbf{x})$

$$\geq -2y\mathbf{w}^{*T} \phi(\mathbf{x}) \quad \forall \mathbf{x}, y$$

$$-\theta^2 \geq 2y(\mathbf{w}^k - \rho\mathbf{w}^*)^T \phi(\mathbf{x}) + \|\phi(\mathbf{x})\|^2$$

need  $\rho$  s.t. this is negative

- Here, negative margin  $\delta = -2y\mathbf{w}^{*T} \phi(\hat{\mathbf{x}})$  is the negative of unsigned distance of closest point  $\hat{\mathbf{x}}$  from separating hyperplane :

$$\hat{\mathbf{x}} = \operatorname{argmax}_{\mathbf{x} \in \mathcal{D}} -2y\mathbf{w}^{*T} \phi(\mathbf{x}) = \operatorname{argmin}_{\mathbf{x} \in \mathcal{D}} y\mathbf{w}^{*T} \phi(\mathbf{x})$$

- Since the data is linearly separable,  $\hat{y}\mathbf{w}^{*T} \phi(\hat{\mathbf{x}}) \geq 0$ , so,  $\delta \leq 0$ . Consequently, we **need** from  $\rho$  that:

$$2y(\mathbf{w}^k - \rho\mathbf{w}^*)^T \phi(\mathbf{x}) + \|\phi(\mathbf{x})\|^2 < \underline{-2\rho y(\mathbf{w}^*)^T \phi(\mathbf{x})} + \|\phi(\mathbf{x})\|^2$$

$$\underline{y(\mathbf{w}^k)^T \phi(\mathbf{x})} < 0 \quad < \rho\delta + \Gamma^2 \leq 0$$



# Perceptron Update Rule: Further analysis

- **Observations:-**

- ①  $y(\mathbf{w}^k)^T \phi(\mathbf{x}) < 0$  ( $\because \mathbf{x}$  was misclassified)

- ②  $\Gamma^2 = \max_{\mathbf{x} \in \mathcal{D}} \|\phi(\mathbf{x})\|^2$

- ③  $\delta = -\min_{\mathbf{x} \in \mathcal{D}} 2y\mathbf{w}^*{}^T \phi(\mathbf{x})$

- Here, negative margin  $\delta = -2y\mathbf{w}^*{}^T \phi(\hat{\mathbf{x}})$  is the negative of unsigned distance of closest point  $\hat{\mathbf{x}}$  from separating hyperplane :

$$\hat{\mathbf{x}} = \operatorname{argmax}_{\mathbf{x} \in \mathcal{D}} -2y\mathbf{w}^*{}^T \phi(\mathbf{x}) = \operatorname{argmin}_{\mathbf{x} \in \mathcal{D}} y\mathbf{w}^*{}^T \phi(\mathbf{x})$$

- Since the data is linearly separable,  $\hat{y}\mathbf{w}^*{}^T \phi(\hat{\mathbf{x}}) \geq 0$ , so,  $\delta \leq 0$ . Consequently, we **need** from  $\rho$  that:

$$0 \leq \|\mathbf{w}^{(k+1)} - \rho\mathbf{w}^*\|^2 < \underbrace{\|\mathbf{w}^k - \rho\mathbf{w}^*\|^2 + \Gamma^2 + \rho\delta}$$

## Perceptron Update Rule: Further analysis

- Since,  $\mathbf{w}^{*T} \phi(\hat{\mathbf{x}}) \geq 0$ , so,  $\delta \leq 0$ . Consequently, restating what we **need** from  $\rho$  is:

$$0 \leq \|\mathbf{w}^{(k+1)} - \rho \mathbf{w}^*\|^2 < \|\mathbf{w}^k - \rho \mathbf{w}^*\|^2 + \underline{\Gamma^2 + \rho \delta}$$

to look like

$$0 \leq \|\mathbf{w}^{(k+1)} - \rho \mathbf{w}^*\|^2 < \|\mathbf{w}^k - \rho \mathbf{w}^*\|^2 - \theta^2$$

- Taking,  $\rho =$

$$\rho \delta + \Gamma^2 \leq 0 \Rightarrow \rho \delta \leq -\Gamma^2 \Rightarrow \rho \geq \frac{-\Gamma^2}{\delta}$$

$$(\because \delta < 0)$$

$$\rho = \frac{-2\Gamma^2}{\delta} \text{ is fine!}$$

# Perceptron Update Rule: Further analysis

- Since,  $\mathbf{w}^{*T} \phi(\hat{\mathbf{x}}) \geq 0$ , so,  $\delta \leq 0$ . Consequently, restating what we **need** from  $\rho$  is:

$$0 \leq \|\mathbf{w}^{(k+1)} - \rho \mathbf{w}^*\|^2 < \|\mathbf{w}^k - \rho \mathbf{w}^*\|^2 + \Gamma^2 + \rho \delta$$

to look like

$$0 \leq \|\mathbf{w}^{(k+1)} - \rho \mathbf{w}^*\|^2 < \|\mathbf{w}^k - \rho \mathbf{w}^*\|^2 + \theta^2$$

- Taking,  $\rho = \frac{2\Gamma^2}{-\delta}$ ,  $0 \leq \|\mathbf{w}^{(k+1)} - \rho \mathbf{w}^*\|^2 \leq \|\mathbf{w}^k - \rho \mathbf{w}^*\|^2 - \Gamma^2$

- Thus, we get,  $\Gamma^2 = \theta^2$ , that we were looking for in eq.(3).

$\therefore \|\mathbf{w}^{(k+1)} - \rho \mathbf{w}^*\|^2$  decreases by at least  $\Gamma^2$  at every iteration.

- Summarily:  $\mathbf{w}^k$  converges to  $\rho \mathbf{w}^*$  by making a minimum  $\theta^2$  decrement at each step.
- Thus, for  $k \rightarrow \infty$ ,  $\|\mathbf{w}^k - \rho \mathbf{w}^*\| \rightarrow 0$ . This proves convergence.

## Number of iterations: (Tutorial 6, Problem 4)

- A statement on number of iterations for convergence:  
If  $\|\mathbf{w}^*\| = 1$  and if there exists  $\delta > 0$  such that for all  $i = 1, \dots, n$ ,  
 $y_i(\mathbf{w}^*)^T \phi(\mathbf{x}_i) \geq \delta$  and  $\|\phi(\mathbf{x}_i)\|^2 \leq \Gamma^2$  then the perceptron algorithm will make  
atmost  $\frac{\Gamma^2}{\delta^2}$  errors (that is take atmost  $\frac{\Gamma^2}{\delta^2}$  iterations to converge)