

# Lecture 15: Kernel perceptron, Neural Networks, SVMs etc

Instructor: Prof. Ganesh Ramakrishnan

# Binary Classification using Perceptron

# Perceptron Classifier

- Consider a binary classification problem:  $f(x) \in \{-1, +1\}$
  - Assuming linearly separability, is there a learning rule that converges in finite time?
- In general for multiclass*  
 $f(x) \in \{c_1, c_2 \dots c_k\}$

2 approaches to multiclass

Wrapper around binary.

Principally extend the binary classifier

1-vs-rest      1-vs-1

for each  $C_i$       for each pair  $(C_i, C_j)$

Wrapping based on single vs multilabeled

Effectiveness

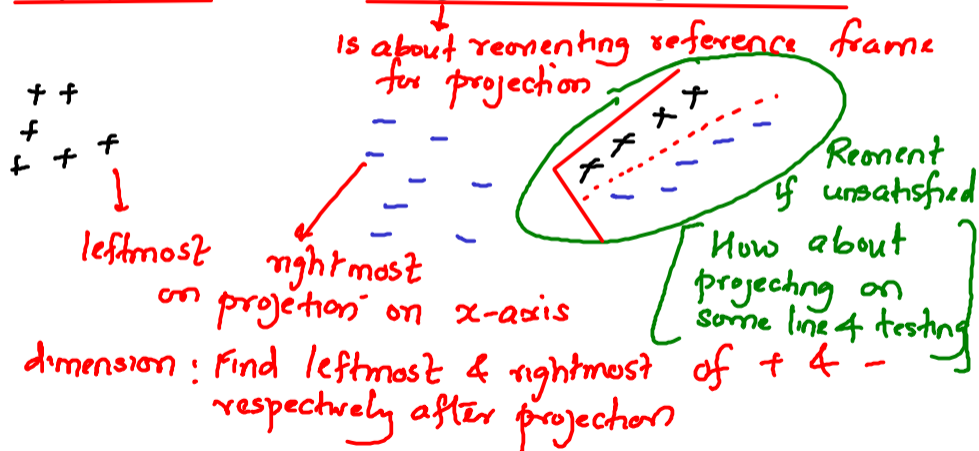
Efficiency

→ measured using "real world" performance measures

→ measured in terms of how easy it is to train models

# Perceptron Classifier

- Consider a binary classification problem:  $f(\mathbf{x}) \in \{-1, +1\}$
- Assuming linearly separability, is there a learning rule that converges in finite time?



## Perceptron Classifier

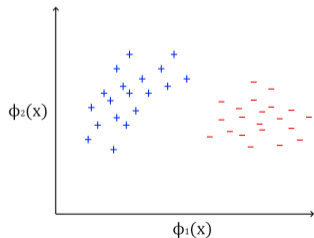
- Consider a binary classification problem:  $f(\mathbf{x}) \in \{-1, +1\}$
- Assuming linear separability, is there a learning rule that converges in finite time?
- Naive Idea: Perform linear regression by constraining  $y \in \{+1, -1\}$ .

$$\omega_{\text{ridge}} = (\Phi^T \Phi + \lambda \mathbf{I})^{-1} \Phi^T \mathbf{y}$$

Does not respect that  $y = -500$  &  $y = -1$  should  
be treated similarly &  
so should  $y = 500$  &  $y = +1$

# Perceptron Classifier

- Consider a binary classification problem:  $f(\mathbf{x}) \in \{-1, +1\}$
- Assuming linear separability, is there a learning rule that converges in finite time?
- Naive Idea: Perform linear regression by constraining  $y \in \{+1, -1\}$ .
- Can we do better? What is ideal?
- Desirable: Any new (unseen) input pattern similar to a seen pattern is **classified** correctly



$\phi(x_u)$

$\phi(x_s)$

$\text{sgn}(\phi^T(x)\omega + b)$

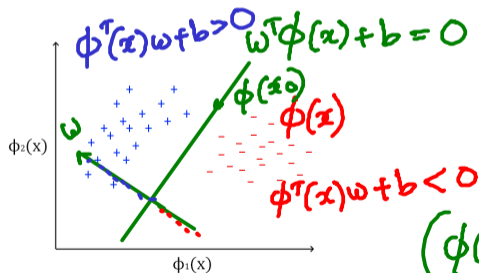
Want happen

if  $y = +1$  &  $y = -1$   
are treated to  
be more similar  
than  $y = +1$  &  $y = +500!$

# Perceptron Classifier

Note:  $\phi^T(x_0)w + b = 0 \Rightarrow \phi^T(x_0)w = -b$   
 for any  $\phi(x_0)$  on hyperplane

- Consider a binary classification problem:  $f(x) \in \{-1, +1\}$
- Assuming linear separability, is there a learning rule that converges in finite time?
- Naive Idea: Perform linear regression by constraining  $y \in \{+1, -1\}$ .
- Can we do better? What is ideal?
- Desirable: Any new (unseen) input pattern similar to a seen pattern is **classified** correctly



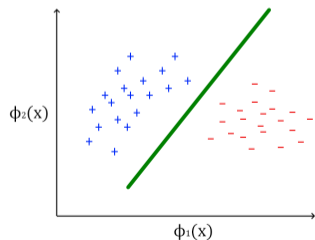
sgn(distance of point from hyperplane computed along  $w$ ) =  $y$

Linear Classification?

$$(\phi(x) - \phi(x_0))^T w = \phi^T(x)w - \phi^T(x_0)w = \phi^T(x)w + b$$

# Perceptron Classifier

- Consider a binary classification problem:  $f(\mathbf{x}) \in \{-1, +1\}$
- Assuming linear separability, is there a learning rule that converges in finite time?
- Naive Idea: Perform linear regression by constraining  $y \in \{+1, -1\}$ .
- Can we do better? What is ideal?
- Desirable: Any new (unseen) input pattern similar to a seen pattern is **classified** correctly



Signed distance of pt  
from separating hyperplane

Linear Classification?

$$\left. \begin{aligned} w^\top \phi(x) + b &\geq 0 \text{ for +ve points } (y = +1) \\ w^\top \phi(x) + b &< 0 \text{ for -ve points } (y = -1) \end{aligned} \right\}$$

$$w, \phi \in \mathbb{R}^m$$

$$y(w^\top \phi(x) + b) \geq 0$$



## Perceptron Classifier: Setting up Notation

- Often,  $b$  is indirectly captured by including it in  $\mathbf{w}$ , and using a  $\phi$  as:  $\phi_{aug} = [\phi, 1]$
- Thus,  $\mathbf{w}^T \phi(\mathbf{x})$

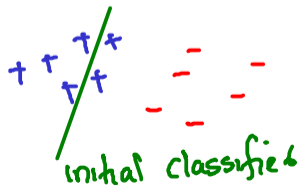
$$= \begin{bmatrix} w_1 & w_2 & w_3 & \dots & w_m & \underline{b} \end{bmatrix} \begin{bmatrix} \phi_1 \\ \phi_2 \\ \phi_3 \\ \vdots \\ \phi_m \\ \underline{1} \end{bmatrix}$$

*hide b within  $\mathbf{w}_{aug}$*

- $\mathbf{w}^T \phi(\mathbf{x}) = 0$  is the separating hyperplane.

# Perceptron Intuition

- 1 Go over all the existing examples, whose class is known, and check their classification with the current weight vector
- 2 If correct, continue
- 3 If not, marginally correct the weights
  - ▶ By adding to the weights a quantity that is proportional to the product of the input pattern with the desired output  $y = \pm 1$



such that unsigned distance of that misclassified point increases

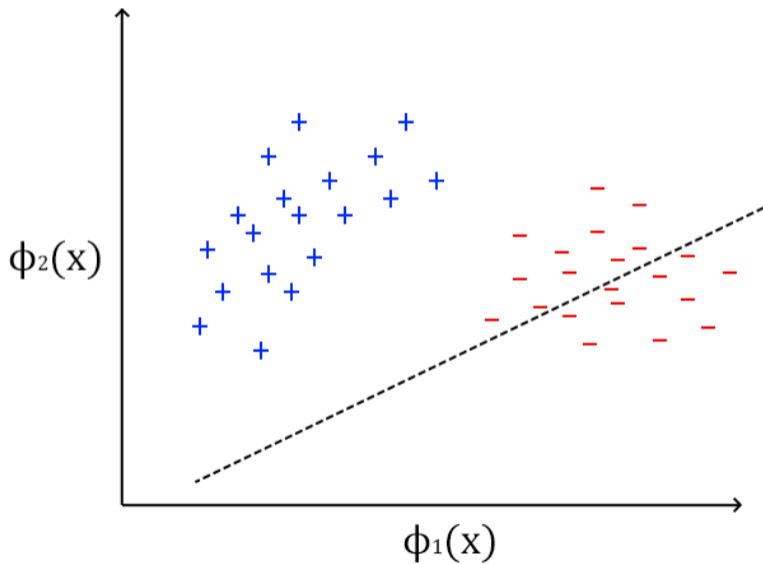
$$y_i (\mathbf{w}^T \phi(\mathbf{x}_i) + b) = \text{unsigned distance}$$

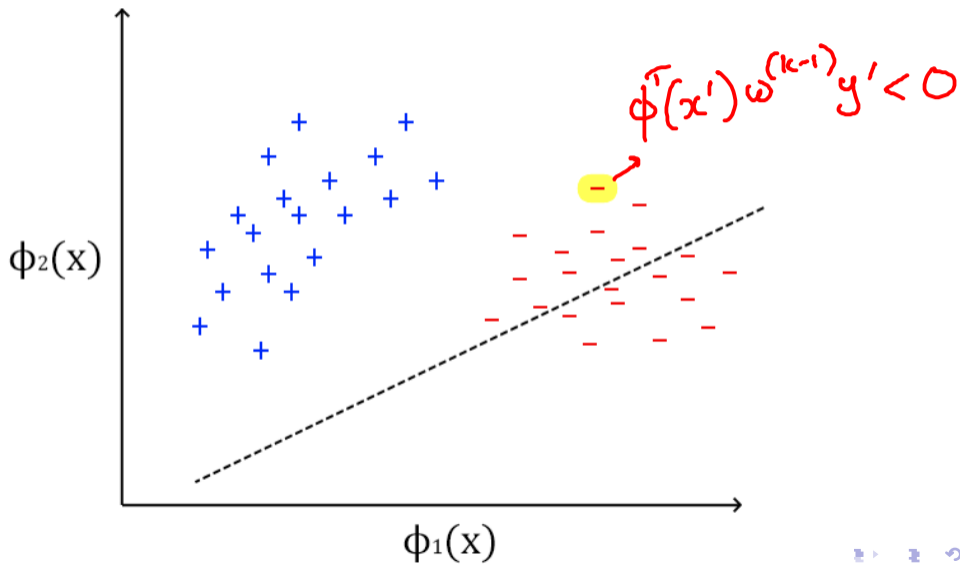
# Perceptron Update Rule

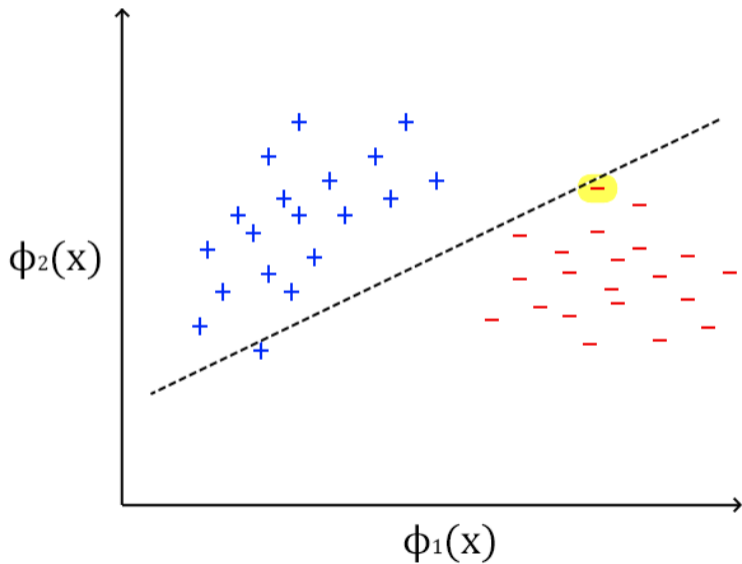
- Start with some weight vector  $\mathbf{w}^{(0)}$ , and for  $k = 0, 1, 2, 3, \dots, n$  (for every example), do:  
 $\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \Gamma \phi(\mathbf{x}')$
- where  $\mathbf{x}'$  s.t.  $\mathbf{x}'$  is misclassified by  $(\mathbf{w}^{(k)})^\top \phi(\mathbf{x})$   
i.e.  $y'(\mathbf{w}^{(k)})^\top \phi(\mathbf{x}') < 0$

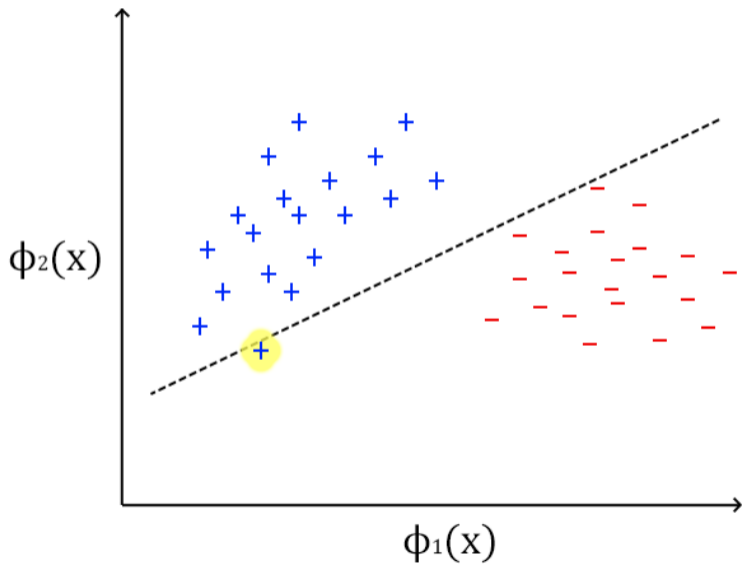
$$\text{Find: } \mathbf{x}' \text{ s.t. } (\mathbf{w}^{(k-1)})^\top \phi(\mathbf{x}') y < 0$$
$$\mathbf{w}^{(k)} = \mathbf{w}^{(k-1)} + \Gamma \phi(\mathbf{x}')$$

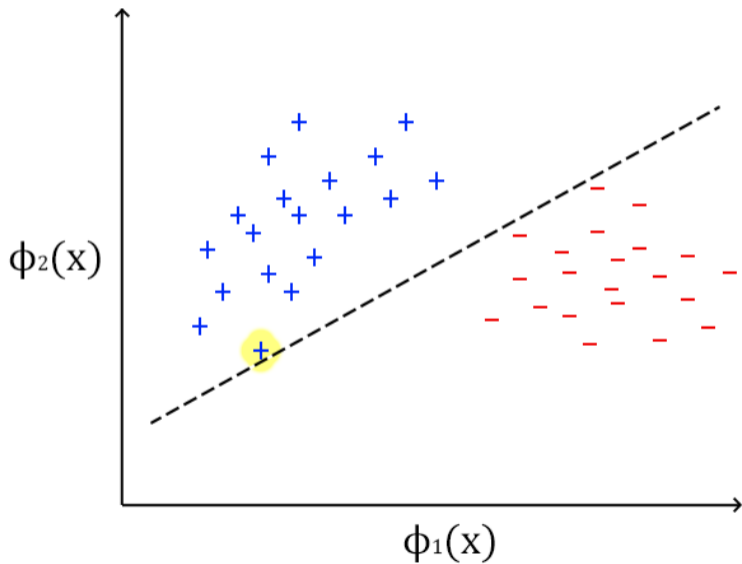
What happens  
to the  
unsigned  
distance of  $\mathbf{x}'$   
after this  
update?











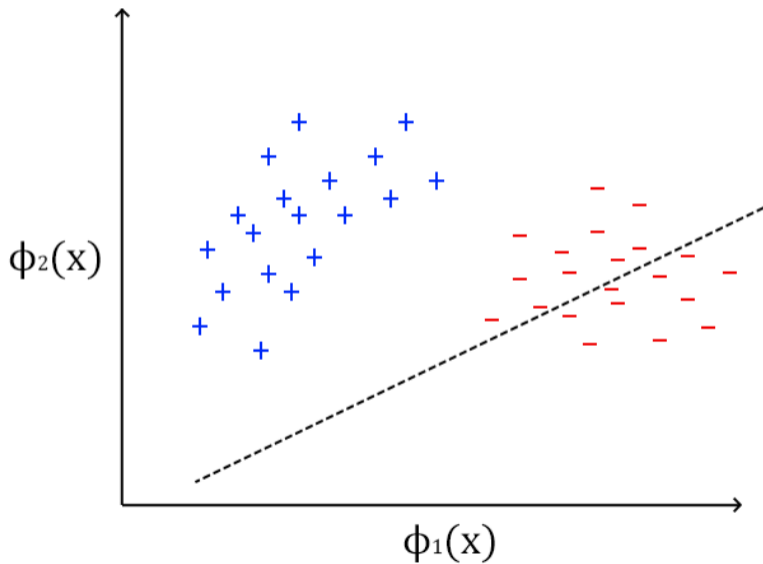


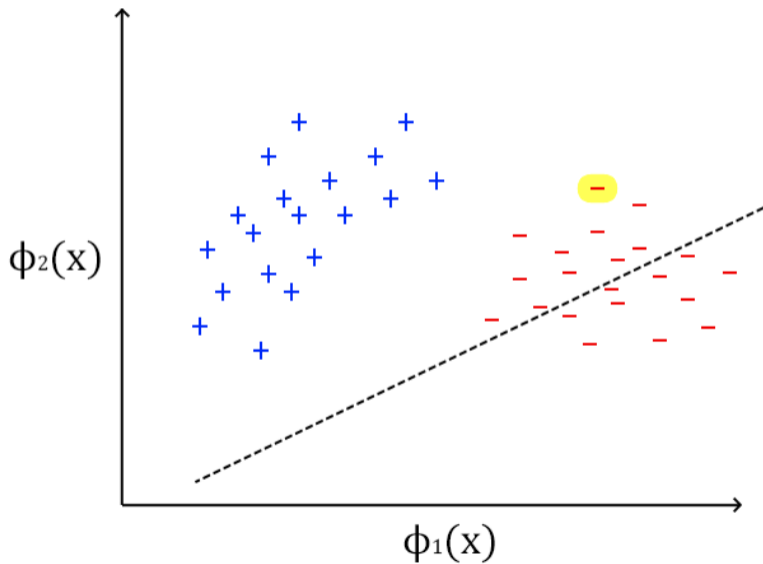
## Some Notes about the Perceptron Algorithm

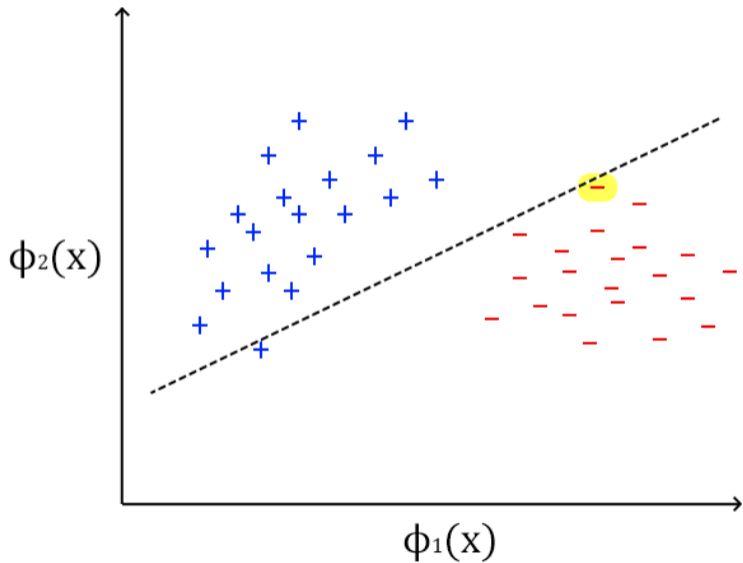
- Perceptron does not find the *best* separating hyperplane, it finds *any* separating hyperplane.
- In case the initial  $\mathbf{w}$  does not classify all the examples, the separating hyperplane corresponding to the final  $\mathbf{w}^*$  will often pass through an example.
- The separating hyperplane does not provide enough breathing space – this is what SVMs address!

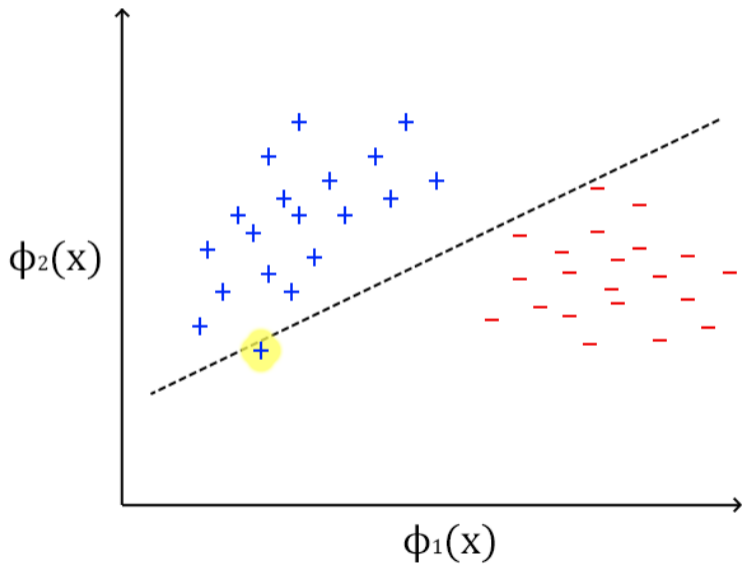
## Perceptron Update Rule

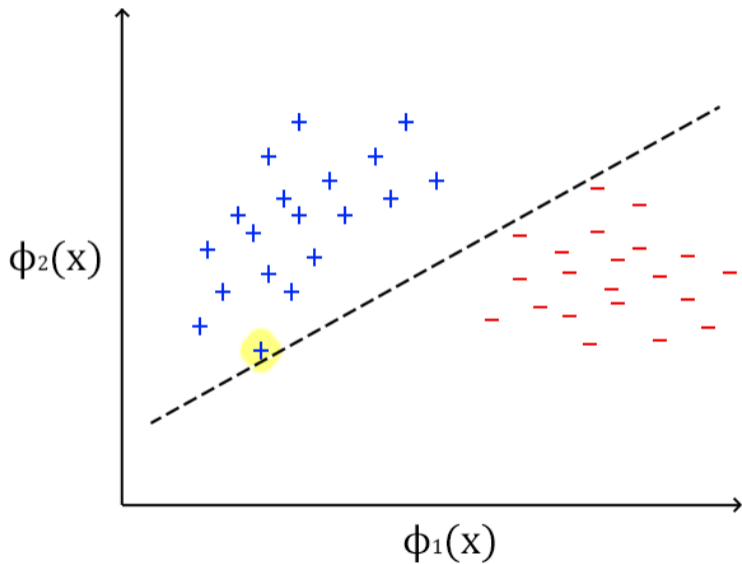
- Start with some weight vector  $w^{(0)}$ , and for  $k = 1, 2, 3, \dots, n$  (for every example), do:  
 $w^{(k)} = w^{(k-1)} + \Gamma \phi(x')$
- where  $x'$  s.t.  $x'$  is misclassified by  $(w^{(k)})^\top \phi(x)$   
i.e.  $y'(w^{(k)})^\top \phi(x') < 0$









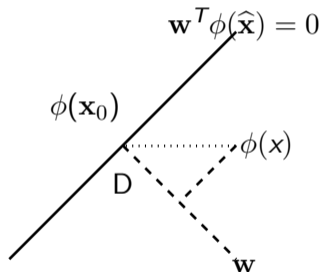


## Perceptron Update Rule: Further analysis

- Explicitly account for signed distance of (misclassified) points from the hyperplane

$$\mathbf{w}^T \phi(\hat{\mathbf{x}}) = 0$$

Distance from hyperplane can be calculated as follows:



$$\text{perpendicular distance} = y\mathbf{w}^T(\phi(\mathbf{x}_0) - \phi(\mathbf{x}))$$

$$\text{Since } \mathbf{w}^T(\phi(\mathbf{x}_0)) = 0 \text{ we get distance} = -y\mathbf{w}^T(\phi(\mathbf{x}))$$



## Perceptron Update Rule: Further analysis

- Perceptron works for two classes ( $y = \pm 1$ ). A point is misclassified if  $y\mathbf{w}^T(\phi(\mathbf{x})) < 0$

- 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,  $\mathbf{w} = \mathbf{w} + \Phi(\mathbf{x}) \cdot y \lambda$
  - endif

if  $y' = -1$  &  $\phi^T(x^i) \omega^{(k-1)} > 0$   
 then  $\phi^T(x^i) \omega^k = \phi^T(x^i) \omega^{(k-1)}$   
 + negative qty

& vice versa

- Intuition:

$$\phi^T(x^i) \omega^{(k)} = \phi^T(x^i) (\omega^{(k-1)} + \lambda \phi(x^i) y^i)$$

$$= \phi^T(x^i) \omega^{(k-1)} + \lambda \phi^T(x^i) \phi(x^i) y^i$$

$$= \phi^T(x^i) \omega^{(k-1)} + \lambda \|\phi(x^i)\|^2 y^i$$

Unsigned distance  $\equiv y^i \phi^T(x^i) \omega^{(k)} = y^i \phi^T(x^i) \omega^{(k-1)} + \lambda \|\phi(x^i)\|^2 (y^i)^2 \geq 0$

## Perceptron Update Rule: Further analysis

- Perceptron works for two classes ( $y = \pm 1$ ). A point is misclassified if  $y\mathbf{w}^T(\phi(\mathbf{x})) < 0$
- 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,  $\mathbf{w} = \mathbf{w} + \Phi(\mathbf{x}) \cdot y$
    - ★ endif

- **Intuition:**

$$\begin{aligned}y(\mathbf{w}^{(k+1)})^T\phi(\mathbf{x}) &= y(\mathbf{w}^k + y\phi^T(\mathbf{x})\phi(\mathbf{x})) \\ &= y(\mathbf{w}^k)^T\phi(\mathbf{x}) + y^2\|\phi(\mathbf{w})\|^2 \\ &> y(\mathbf{w}^k)^T\phi(\mathbf{x})\end{aligned}$$

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.

## Perceptron Update Rule: Further analysis

- Tries to minimize the error function  $E$  (sum of unsigned distances from hyperplane to misclassified points) over misclassified examples

$$E = - \sum_{(\mathbf{x}, y) \in \mathcal{M}} y \mathbf{w}^T \phi(\mathbf{x})$$

where  $\mathcal{M} \subseteq \mathcal{D}$  is the set of misclassified examples.

- **Gradient Descent (Batch Perceptron) Algorithm**  $\nabla_{\mathbf{w}} E = - \sum_{(\mathbf{x}, y) \in \mathcal{M}} y \phi(\mathbf{x})$

$$\begin{aligned} \mathbf{w}^{(k+1)} &= \mathbf{w}^k - \eta \nabla_{\mathbf{w}} E \\ &= \mathbf{w}^k + \eta \sum_{(\mathbf{x}, y) \in \mathcal{M}} y \phi(\mathbf{x}) \end{aligned}$$

Less meaningful to take entire  $\mathcal{M}$  since the fate of misclassified points is interdependent

## Perceptron Update Rule: Further analysis

- Batch update considers all misclassified points simultaneously

$$\begin{aligned}\mathbf{w}^{(k+1)} &= \mathbf{w}^k - \eta \nabla_{\mathbf{w}} E \\ &= \mathbf{w}^k + \eta \sum_{(\mathbf{x}, y) \in \mathcal{M}} y \phi(\mathbf{x})\end{aligned}$$

- Perceptron update  $\Rightarrow$  Stochastic Gradient Descent:

$$\nabla_{\mathbf{w}} E = - \sum_{(\mathbf{x}, y) \in \mathcal{M}} y \phi(\mathbf{x}) = - \sum_{(\mathbf{x}, y) \in \mathcal{M}} \nabla_{\mathbf{w}} E(\mathbf{x}) \text{ s.t. } E(\mathbf{x}) = -y\mathbf{w}^T \phi(\mathbf{x})$$

$$\begin{aligned}\mathbf{w}^{(k+1)} &= \mathbf{w}^k - \eta \nabla_{\mathbf{w}} E(\mathbf{x}) \\ &= \mathbf{w}^k + \eta y \phi(\mathbf{x})\end{aligned} \quad (\text{for any } (\mathbf{x}, y) \in \mathcal{M})$$

## 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.)

