



Machine Learning 1

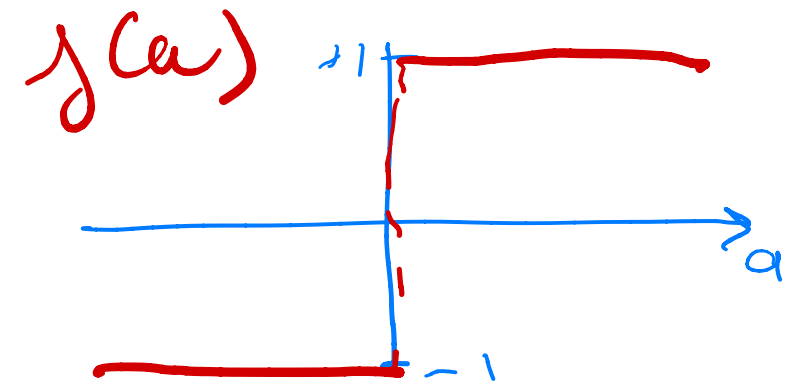
Lecture 6.5 - Supervised Learning
Classification - Discriminative Models - The
Perceptron

Erik Bekkers

(Bishop 4.1.7)



The Perceptron Algorithm



Input: $\underline{x} \in \mathbb{R}^D$

targets: $t \in \{C_1, C_2\}, \rightarrow t \in \{-1, 1\}$

2 classes

Prediction: $y(\mathbf{x}) = f(\mathbf{w}^T \phi(\mathbf{x}))$ $f(a) = \begin{cases} 1 & , a \geq 0 \\ -1 & , a < 0 \end{cases}$

Class decisions: assign \mathbf{x} to class C_1 if $\dots \dots \dots \underline{w}^T \underline{\phi}_n > 0$ (and C_{-1} if $\underline{w}^T \underline{\phi}_n < 0$)

For correct classification: find \mathbf{w} such that for all (\mathbf{x}_n, t_n) :

$$\underline{w}^T \underline{\phi}_n t_n \geq 0$$

Perceptron criterion: $E_P(\mathbf{w}) = - \sum_{n \in \mathcal{M}} w^T \phi(\mathbf{x}_n) t_n$
 \leftarrow all misclassified

$$\mathcal{M} : \{ n : \underline{w}^T \underline{\phi}_n t_n < 0 \}$$

Perceptron: Stochastic Gradient Descent

$$\diamond E_P(\mathbf{w}) = - \sum_{n \in \mathcal{M}} w^T \phi(\mathbf{x}_n) t_n$$

$$= \sum_{n \in \mathcal{M}} E_n(\mathbf{w})$$

$E_n(\mathbf{w}) = -\mathbf{w}^T \phi_n t_n$

◆ Stochastic Gradient Descent (SGD).
For each misclassified \mathbf{x}_n :

$$\begin{aligned} \mathbf{w}^{(\tau+1)} &= \mathbf{w}^{(\tau)} - \eta \nabla E_n(\mathbf{w}) \\ &= \mathbf{w}^{(\tau)} + \gamma (\phi_n t_n) \end{aligned}$$

◆ If \mathbf{X} is linearly separable, then
perceptron SGD will converge

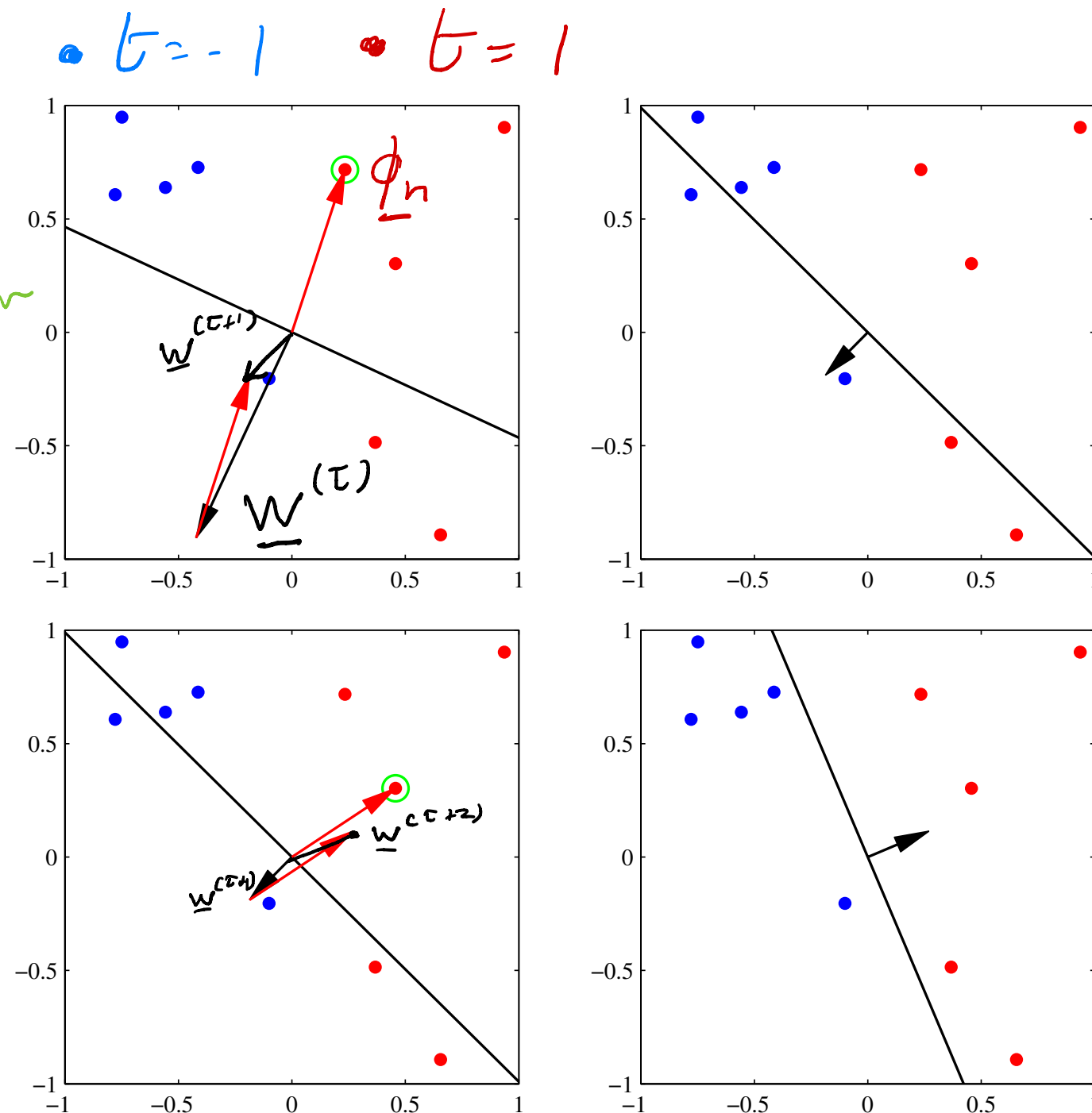


Figure: for \mathbf{x}_n in C_1 : add $\phi(\mathbf{x}_n)$ to \mathbf{w} , for \mathbf{x}_n in C_2 : subtract $\phi(\mathbf{x}_n)$ from \mathbf{w} . SGD for perceptron criterion (Bishop 4.7)

Problems: Perceptron

- ▶ Perceptron only works for 2 classes
- ▶ There might be many solutions depending on the initialization of \mathbf{w} and on the order in which data is presented in SGD
- ▶ If dataset is not linearly separable, the perceptron algorithm will not converge.
- ▶ Based on linear combination of fixed basis functions.