



Class: Machine Learning

Neural Networks: Perceptron

Instructor: Matteo Leonetti

Learning outcomes



UNIVERSITY OF LEEDS

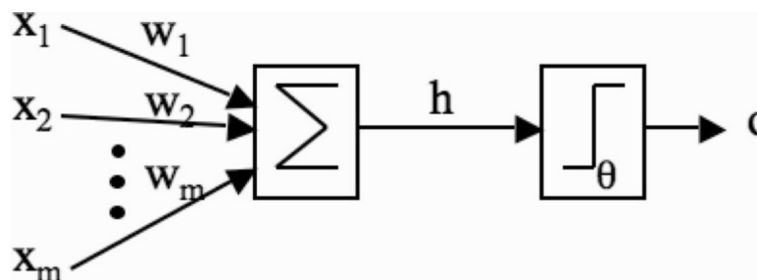
- Define an appropriate error function for the perceptron.
- Derive the corresponding update algorithm.
- Describe the difference between gradient descent and stochastic gradient descent.

Recap

We want to apply gradient descent:

$$x_{t+1} = x_t - \eta \nabla f(x_t)$$

To the parameters of a perceptron:



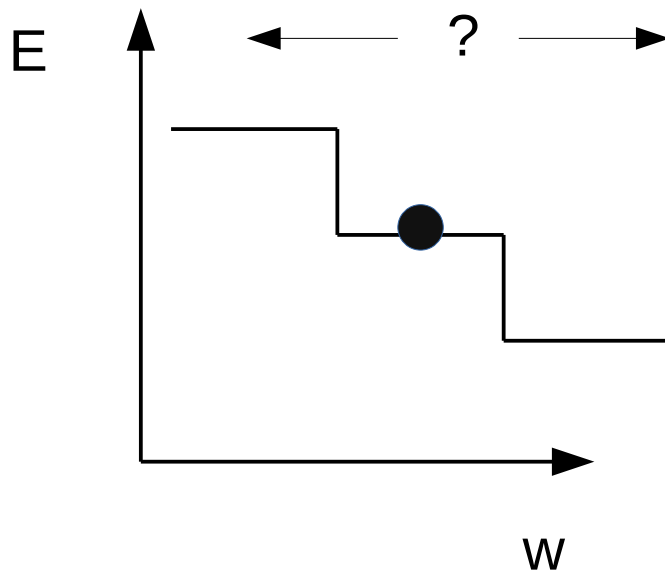
So as to minimise an error (or loss) function, such as:

$$E(\mathbf{X}) = \sum_{\vec{x}_n \in \mathbf{X}} |y_n - t_n|$$

Number of mistakes

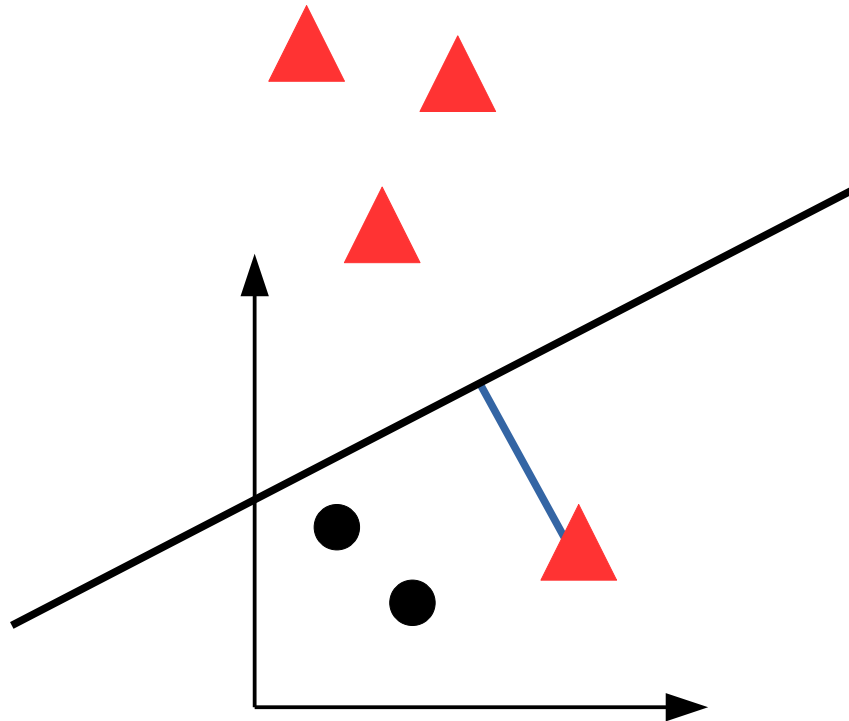
$$E(\mathbf{X}) = \sum_{\vec{x}_n \in \mathbf{X}} |y_n - t_n|$$

Number of mistakes on the dataset. Piecewise constant \rightarrow no gradient.



There is no local information
on the direction of
improvement

Towards a better error function



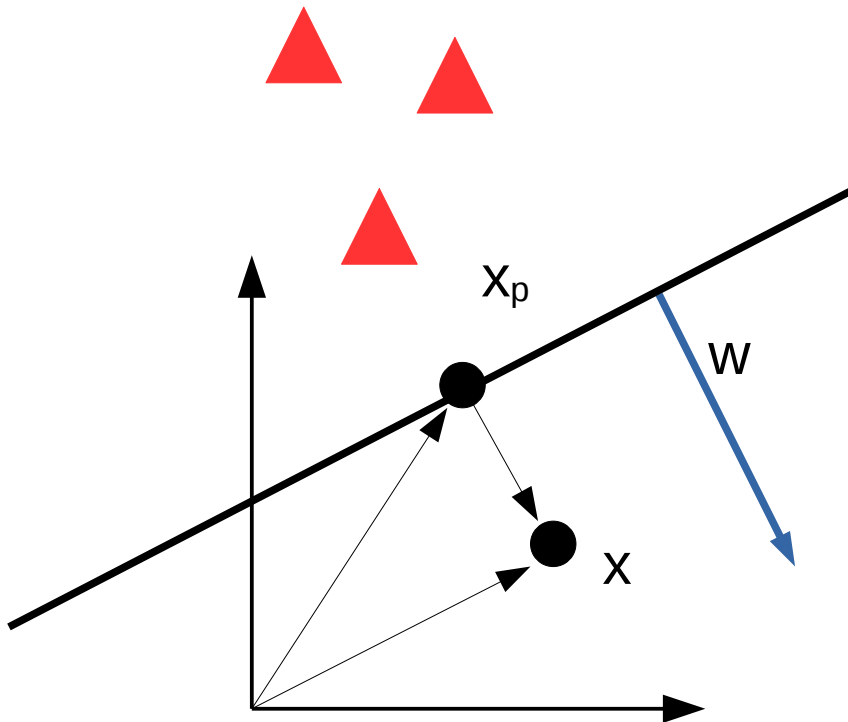
For each misclassified point, we would like to know not only that they are on the wrong side, but also **by how much**.

Towards a better error function



UNIVERSITY OF LEEDS

$$h_w(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = 0$$



Distance to the hyperplane

$$\mathbf{x} = \mathbf{x}_p + d \frac{\mathbf{w}}{\|\mathbf{w}\|}$$

$$\begin{aligned} h_w(\mathbf{x}) &= \mathbf{w}^T \left(\mathbf{x}_p + d \frac{\mathbf{w}}{\|\mathbf{w}\|} \right) + w_0 \\ &= \cancel{\mathbf{w}^T \mathbf{x}_p} + w_0 + d \frac{\mathbf{w}^T \mathbf{w}}{\|\mathbf{w}\|} = d \|\mathbf{w}\| \end{aligned}$$

Recall that:

$$\mathbf{w}^T \mathbf{w} = w_1^2 + w_2^2 + \dots + w_n^2 = \|\mathbf{w}\|^2$$

Towards a better error function

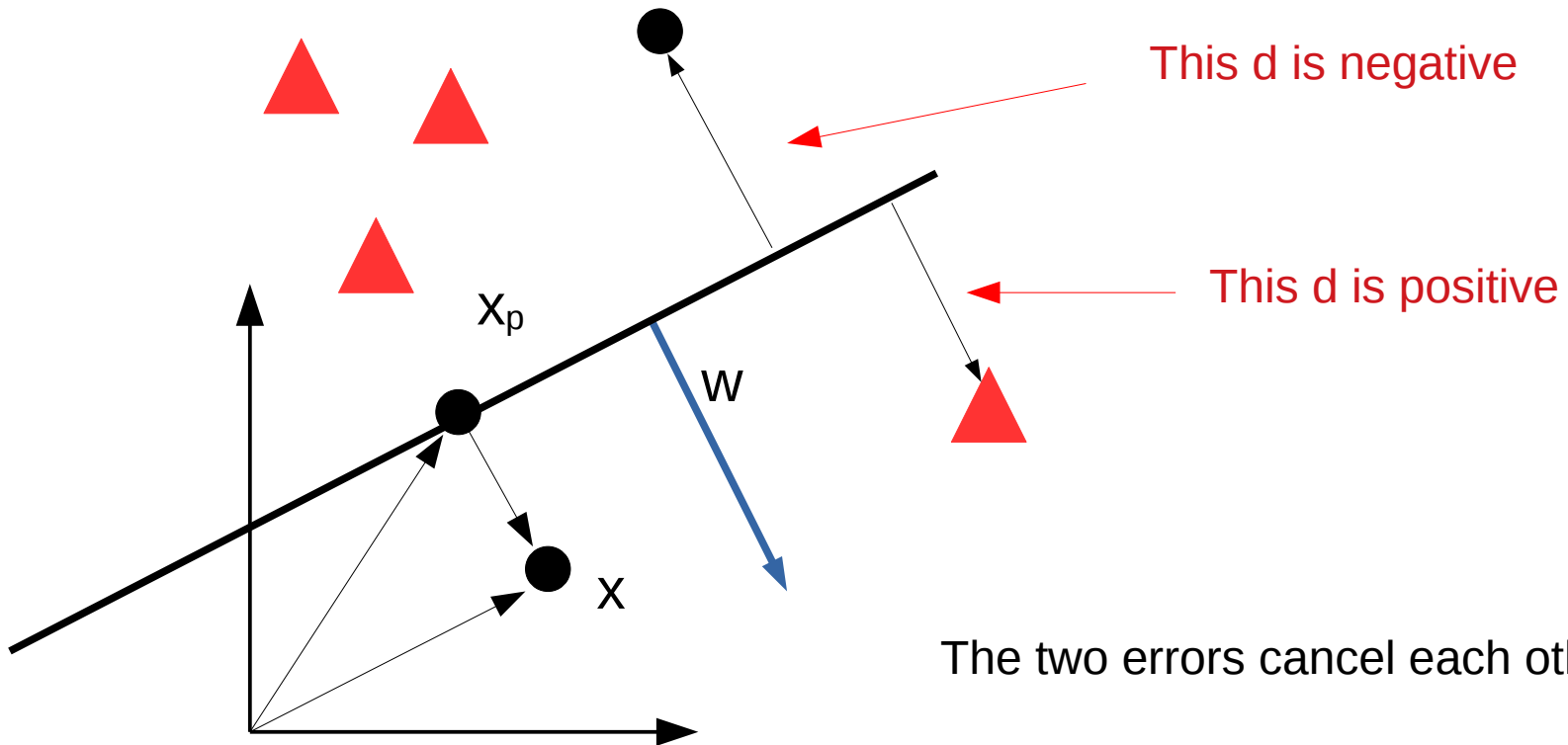
$$h_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = d \|\mathbf{w}\|$$

$$E(\mathbf{X}) = \sum_{\mathbf{x}_n \in \mathbf{X}} (\mathbf{w}^T \mathbf{x}_n + w_0)$$

Is this a good error?

Towards a better error function

$$\mathbf{x} = \mathbf{x}_p + d \frac{\mathbf{w}}{\|\mathbf{w}\|}$$



The perceptron criterion

$$h_w(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = 0 \quad \text{apply the bias input}$$

if $\mathbf{w}^T \mathbf{x} > 0$ then $y = 1$ In case of mistake: $t = 0$ $(y - t) = 1$

if $\mathbf{w}^T \mathbf{x} \leq 0$ then $y = 0$ In case of mistake: $t = 1$ $(y - t) = -1$

Therefore, if mistake: $\mathbf{w}^T \mathbf{x} (y - t) > 0$

$$E(\mathbf{X}) = \sum_{\mathbf{x}_n \in X} |y_n - t_n|$$

Number of mistakes on the dataset.
Piecewise constant \rightarrow gradient
useless.

$$E_p(\mathbf{X}) = \sum_{\mathbf{x}_n \in X} \mathbf{w}^T \mathbf{x}_n (y_n - t_n)$$

Proportional to distance of
misclassified points from surface.
 \rightarrow gradient ok.

What is the derivative of

$$y = 2x \quad ?$$

Given the perceptron error (below), what is the gradient with respect to \mathbf{w} ?

$$E_p(\mathbf{X}) = \mathbf{w}^T \mathbf{x} (y - t) = (w_0 x_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n) (y - t)$$

Solution



UNIVERSITY OF LEEDS

$$E_p(\mathbf{x}) = \mathbf{w}^T \mathbf{x}(y-t) = w_0 x_0(y-t) + w_1 x_1(y-t) + \dots + w_m x_m(y-t)$$

$$\nabla E_p(\mathbf{x}) = \begin{bmatrix} \frac{\partial}{\partial w_0} E_p(\mathbf{x}) \\ \frac{\partial}{\partial w_1} E_p(\mathbf{x}) \\ \frac{\partial}{\partial w_2} E_p(\mathbf{x}) \\ \dots \\ \frac{\partial}{\partial w_n} E_p(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} x_0(y-t) \\ x_1(y-t) \\ x_2(y-t) \\ \dots \\ x_n(y-t) \end{bmatrix}$$

Gradient descent

$$\nabla E_p(\mathbf{X}) = \sum_{\mathbf{x}_n \in \mathbf{X}} \mathbf{x}_n (y_n - t_n)$$

Recall that gradient descent does the following update:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \eta \nabla f(\mathbf{w}_k)$$

Which leads us to the update rule for the perceptron:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \eta \sum_{\mathbf{x}_n \in \mathbf{X}} \mathbf{x}_n (y_n - t_n)$$

$$E_p(\mathbf{X}) = \frac{1}{N} \sum_{\mathbf{x}_n \in X} \mathbf{w}^T \mathbf{x}_n (y_n - t_n) = \mathbf{E}[\mathbf{w}^T \mathbf{x}_n (y_n - t_n)]$$

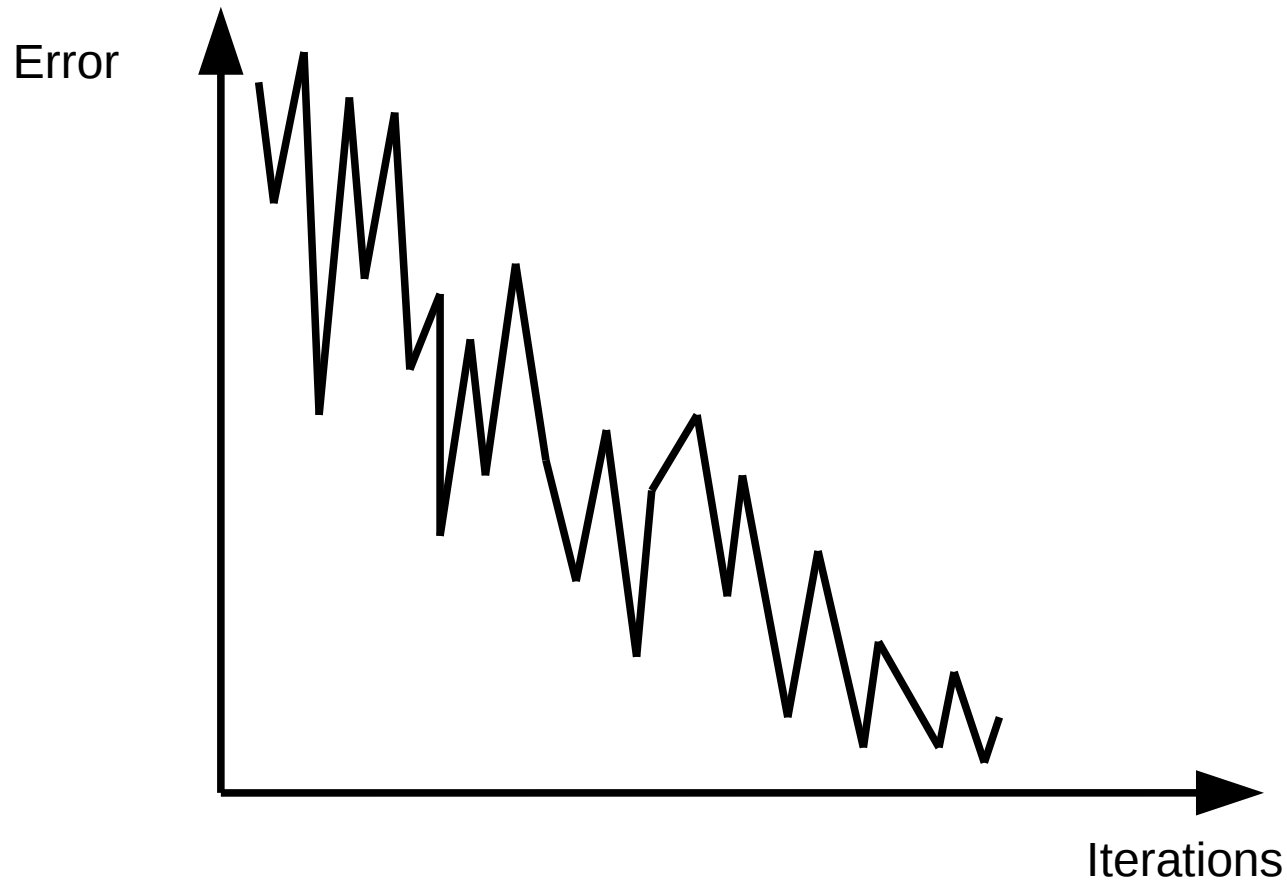
Gradient:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \eta \frac{1}{N} \sum_{\mathbf{x}_n \in X} \mathbf{x}_n (y_n - t_n)$$

Stochastic gradient descent:

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \eta \mathbf{x} (y - t)$$

Stochastic gradient descent





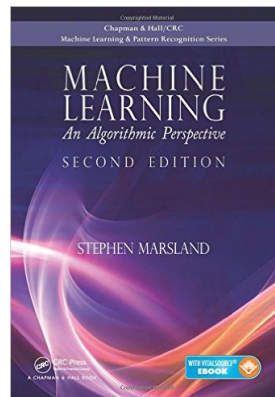
Conclusion

Learning outcomes

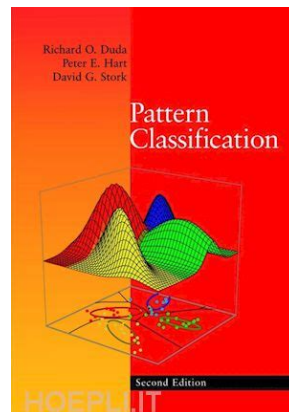


UNIVERSITY OF LEEDS

- Define an appropriate error function for the perceptron.
- Derive the corresponding update algorithm.
- Describe the difference between gradient descent and stochastic gradient descent.



Section 3.4



Book in Minerva
in “ Online Course Readings Folder”

Section 5.2.1, 5.4. and 5.5
(without convergence proof)