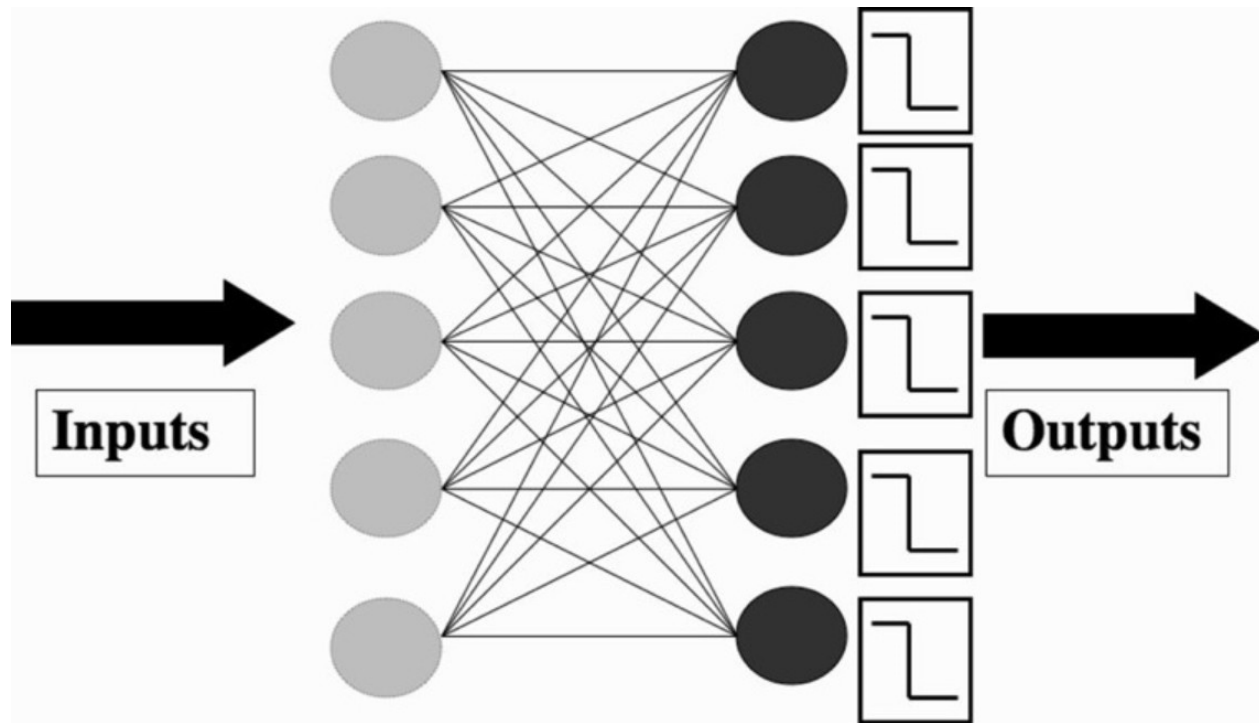


Machine Learning

Perceptron

Jian Liu

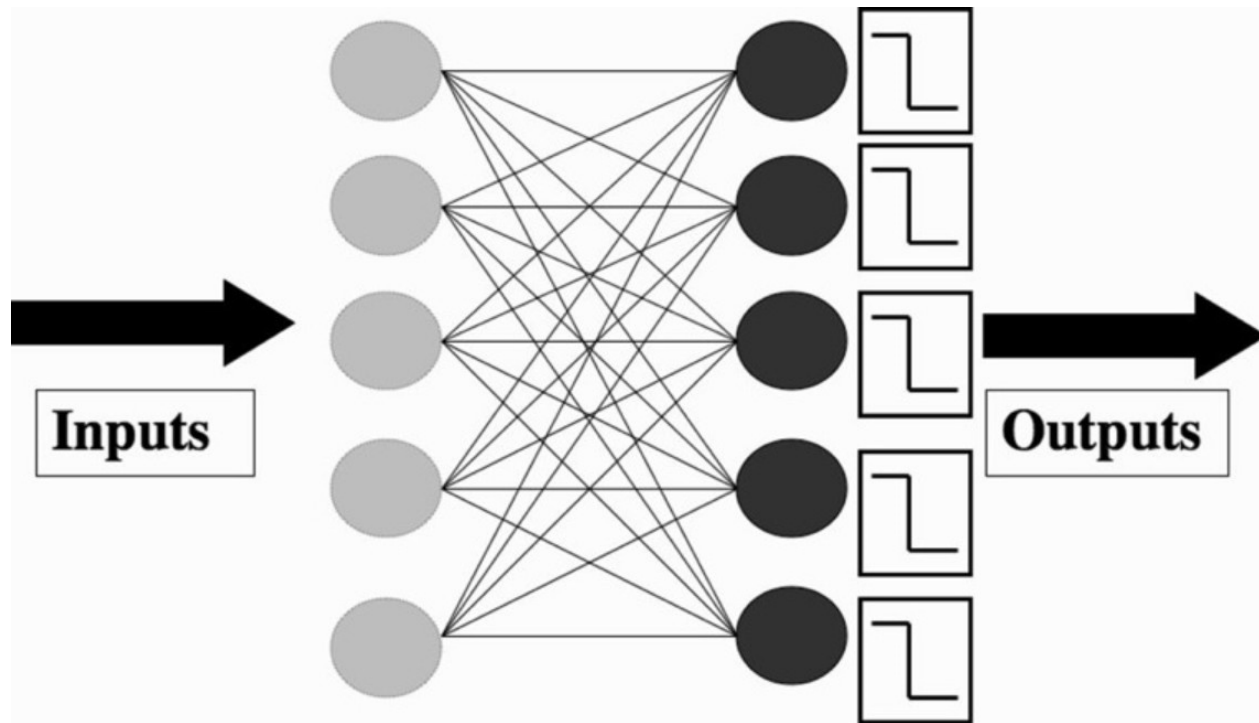


Machine Learning

Perceptron

Jian Liu

Part 2: Error Functions



Training the perceptron or training neural network models

Learning happens through **optimisation**.

We define an error function, and then an optimisation algorithm finds the parameters that obtain the minimum error.

For example, the error function is the total number of mistakes:

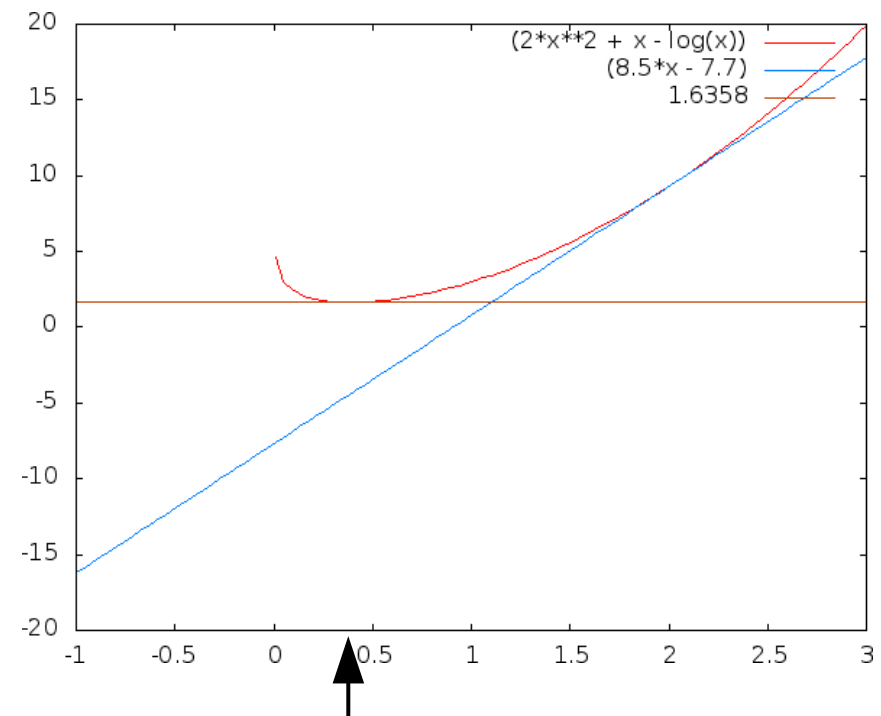
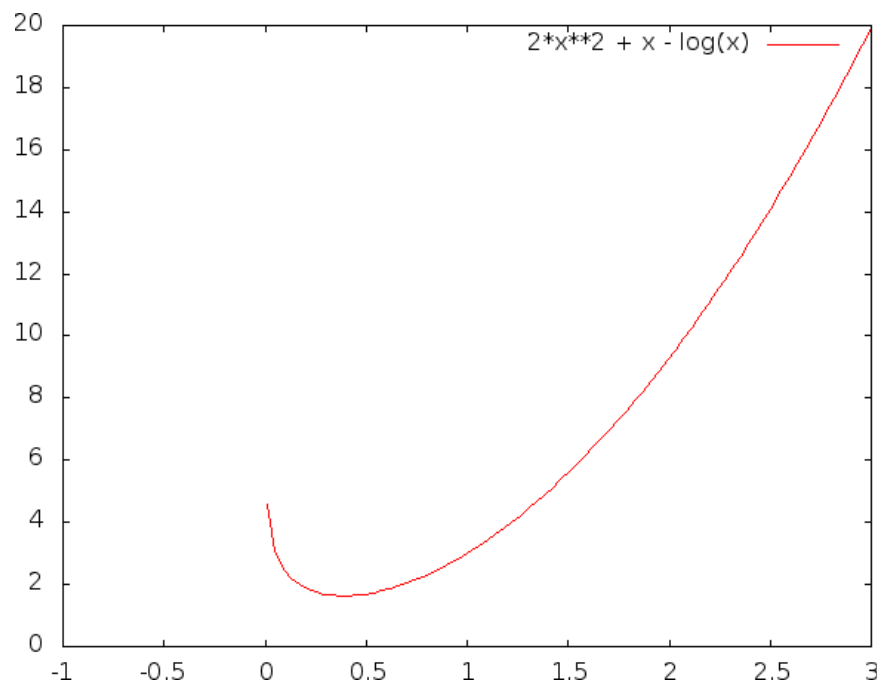
$$E(\mathbf{X}) = \sum_{\vec{x}_n \in \mathbf{X}} |y_n - t_n|$$

Where y_n is the output of the perceptron on point n , and $t_n \in \{0, 1\}$ is the desired class, and \mathbf{X} is the dataset.

Elements of Local Optimisation

Goal

Find the minimum point of a given function:



The minimum is at 0.39

Local Methods



Gradient descent

First order: gradient descent

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

step parameter

Second order: Newton's method

$$f(x_n + \Delta x) \approx f(x_n) + f'(x_n)\Delta x + \frac{1}{2}f''(x_n)\Delta x^2$$

Taylor's
expansion

$$\frac{\partial}{\partial \Delta x} f(x_n + \Delta x) = f'(x_n) + f''(x_n)\Delta x = 0$$

Optimal step

$$\Delta x = \frac{-f'(x_n)}{f''(x_n)}$$

Many dimensions: $x_{t+1} = x_t - H^{-1}|_{x_n} \nabla f$

Question

The current point is $\langle 1, 0 \rangle$
compute the next point following gradient
descent on the function $f(x, y) = x^3 + 2y^2 - y$ with
step size 0.1.

Question

We want to compute: $x_{t+1} = \langle 1, 0 \rangle - 0.1 \nabla f(x_t)$

$\nabla f = \langle 3x^2, 4y-1 \rangle$ Evaluated in $\langle 1, 0 \rangle$ is $\langle 3, -1 \rangle$

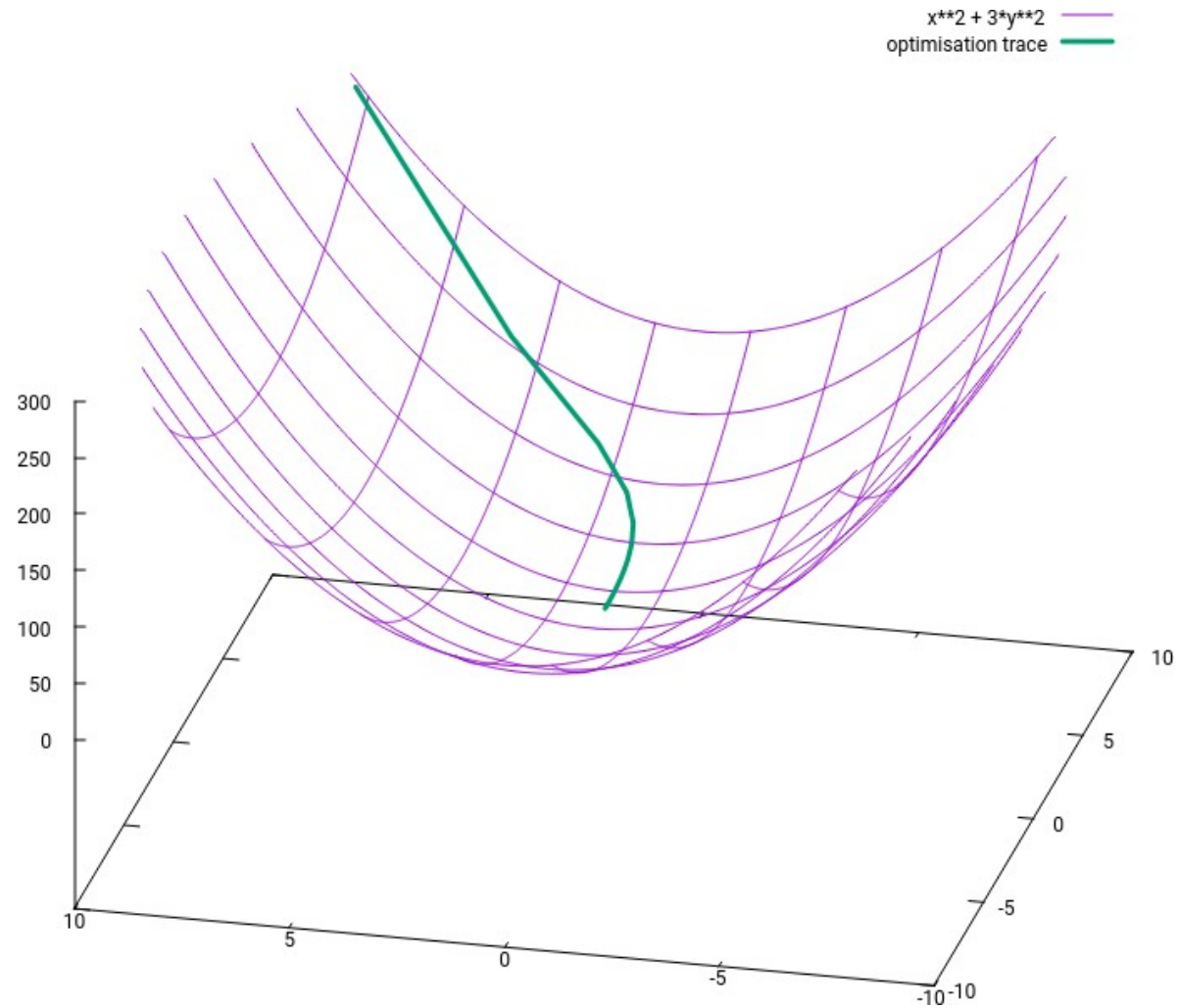
$$x_{t+1} = \langle 1, 0 \rangle - 0.1 \cdot \langle 3, -1 \rangle = \langle 0.7, 0.1 \rangle$$

$$f(1, 0) = 1$$

Our solution has improved!

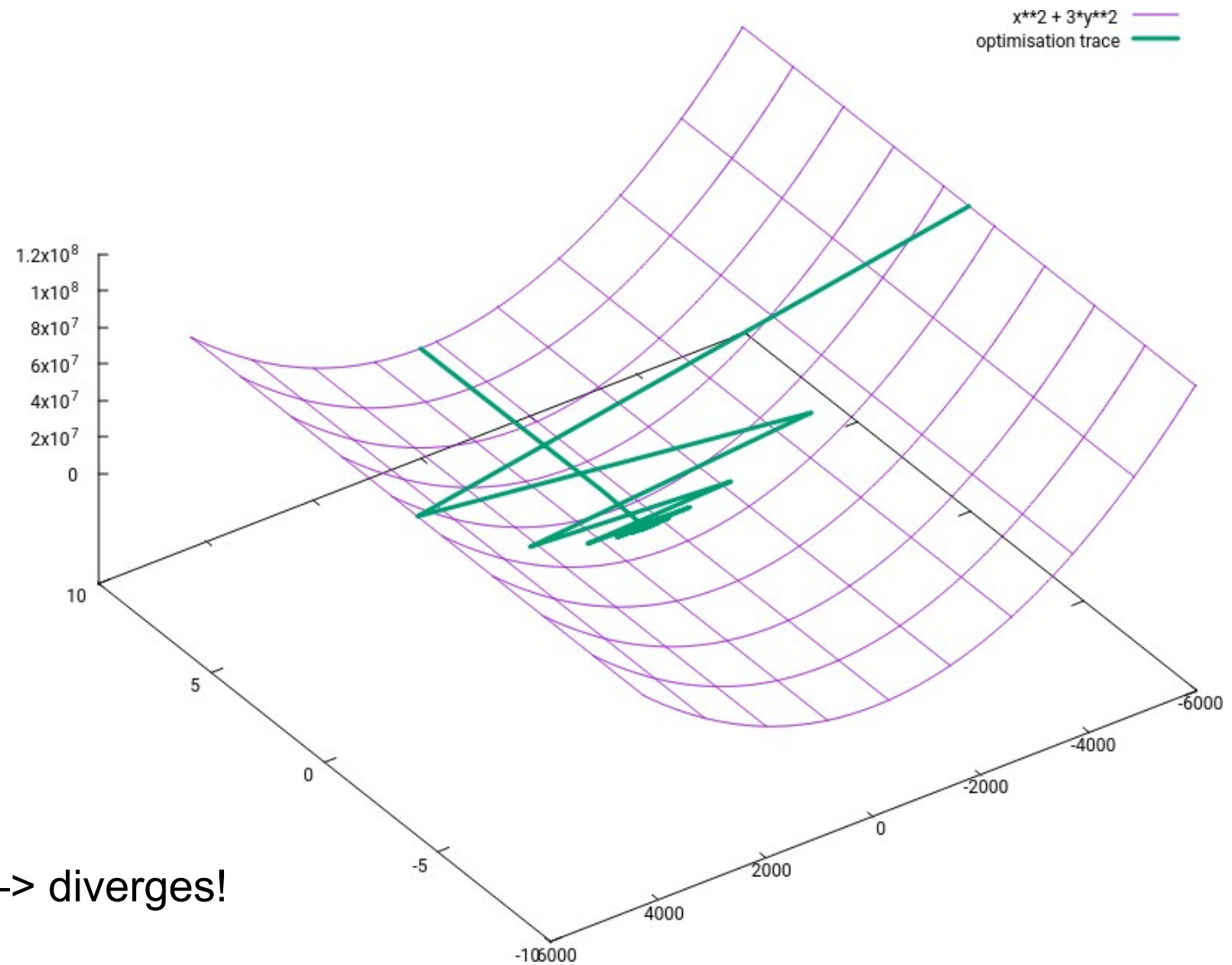
$$f(0.7, 0.1) = 0.263$$

In 3D



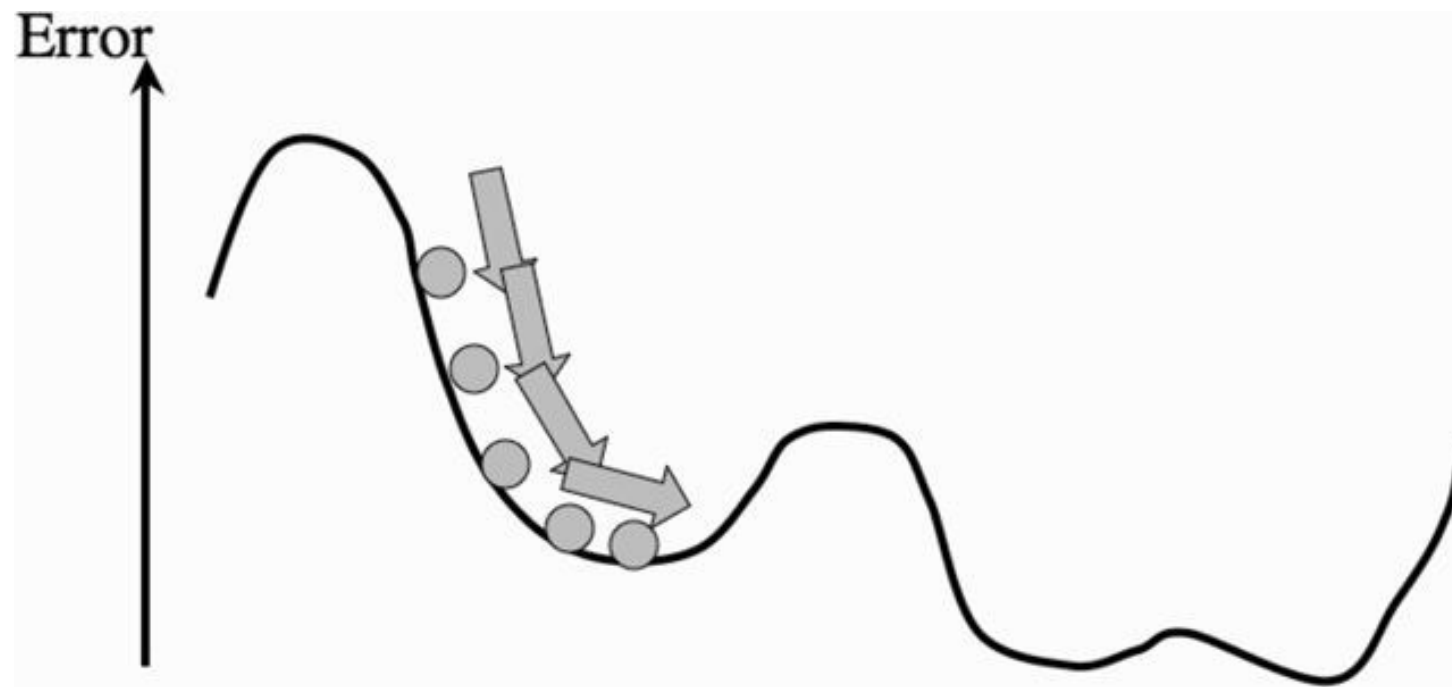
Step size: 0.1

In 3D



Step size: 0.14 → diverges!

Local Minima

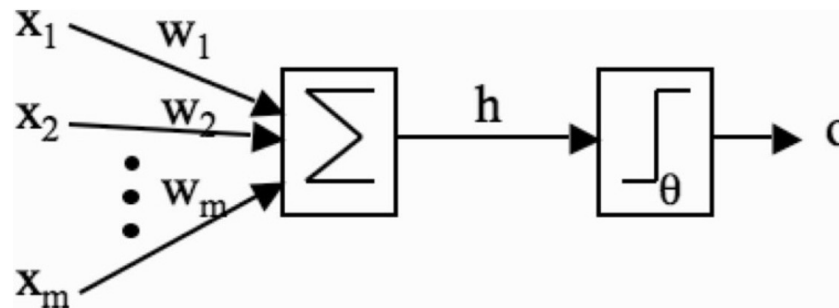


Training the perceptron

We want to apply gradient descent:

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

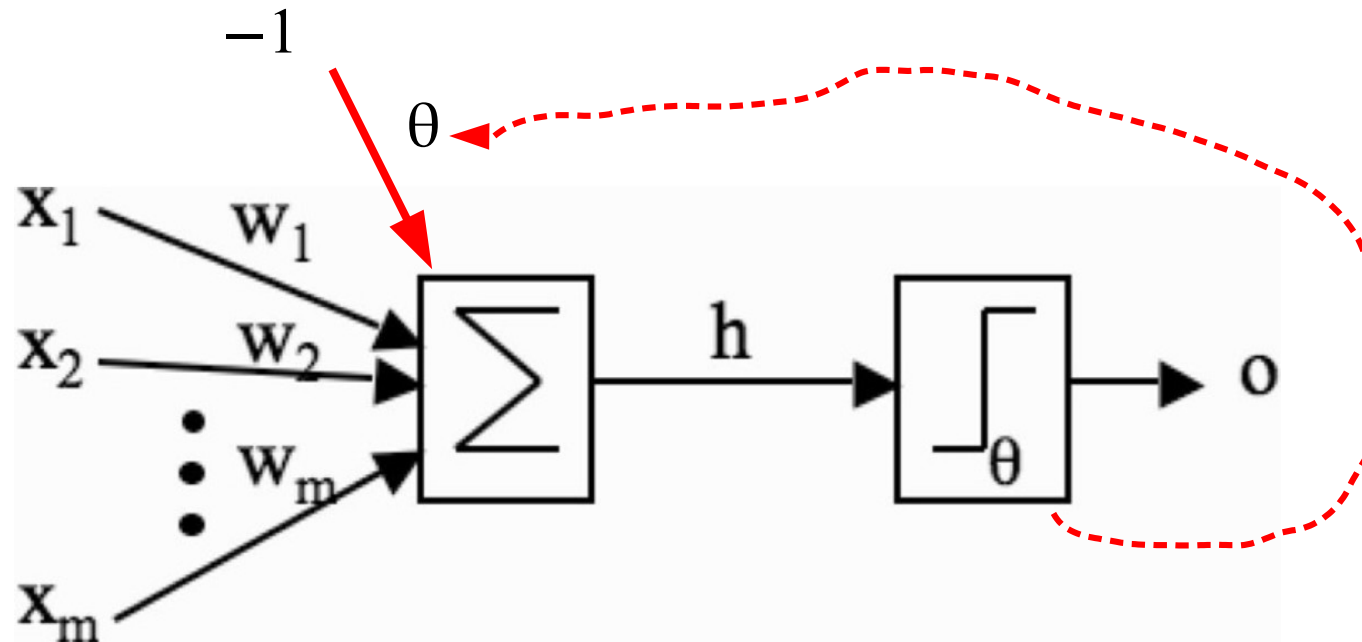
To the parameters of a perceptron:



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

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

Bias input



$$h_w(\mathbf{x}) = \sum_i w_i x_i = \mathbf{w} \cdot \mathbf{x}$$

$$o(h_w) = \begin{cases} 1 & \text{if } h_w > \theta \\ 0 & \text{if } h_w \leq \theta \end{cases}$$

$$w_1 x_1 + w_2 x_2 + \dots + w_n x_n - 1 \cdot \theta > 0$$



$$h_w(\mathbf{x}) = \sum w_i x_i - \theta_i$$

$$\mathbf{x}_{new} = \langle \mathbf{x}, -1 \rangle \quad \mathbf{w}_{new} = \langle \mathbf{w}, \theta \rangle$$

$$h_w(\mathbf{x}_{new}) = \mathbf{w}_{new} \cdot \mathbf{x}_{new}$$

$$o(h) = \begin{cases} 1 & \text{if } h_w > 0 \\ 0 & \text{if } h_w \leq 0 \end{cases}$$

Question

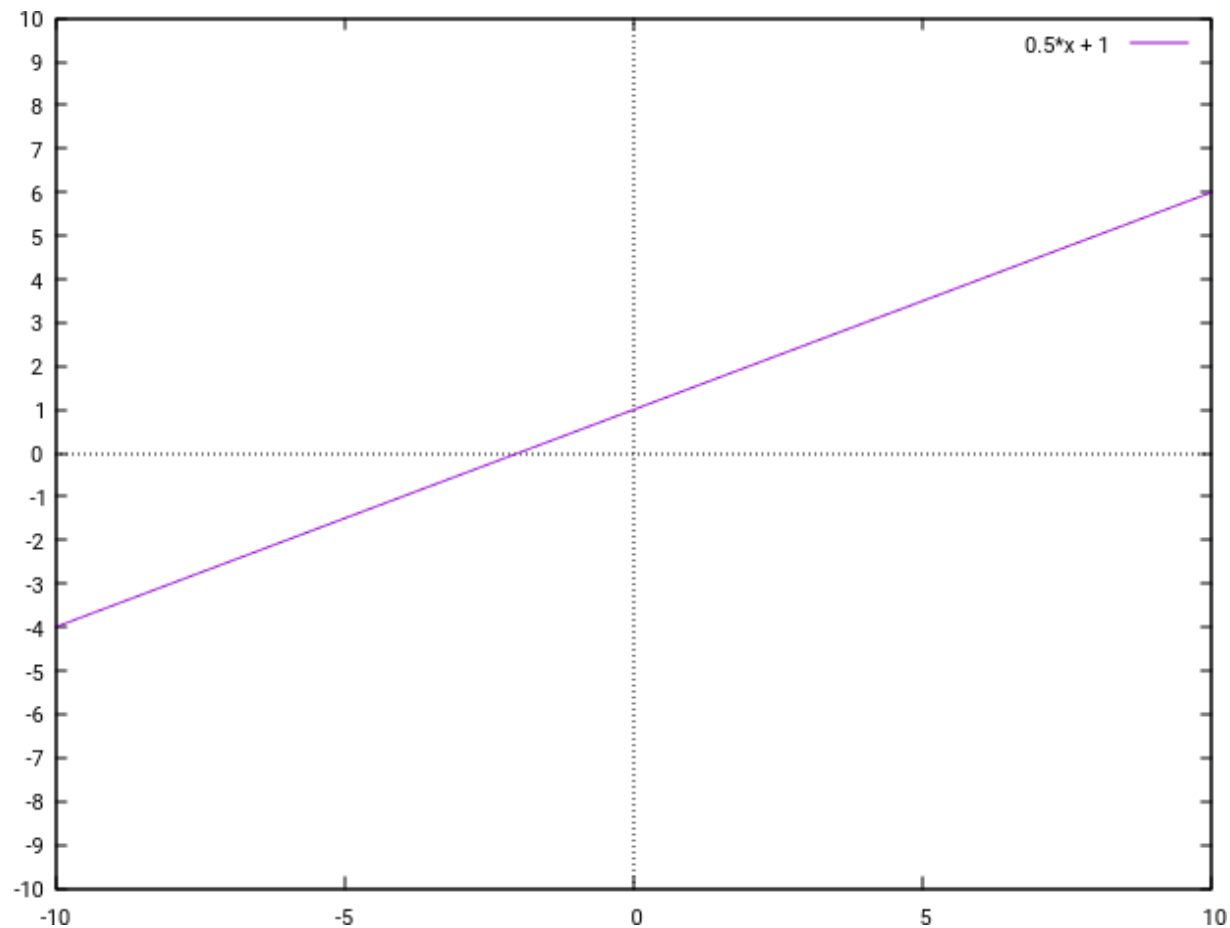
The decision boundary of the perceptron is the function below:

$$f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} = 0$$

Plot the following function: $\frac{1}{2}x - y + 1 = 0$

Solution

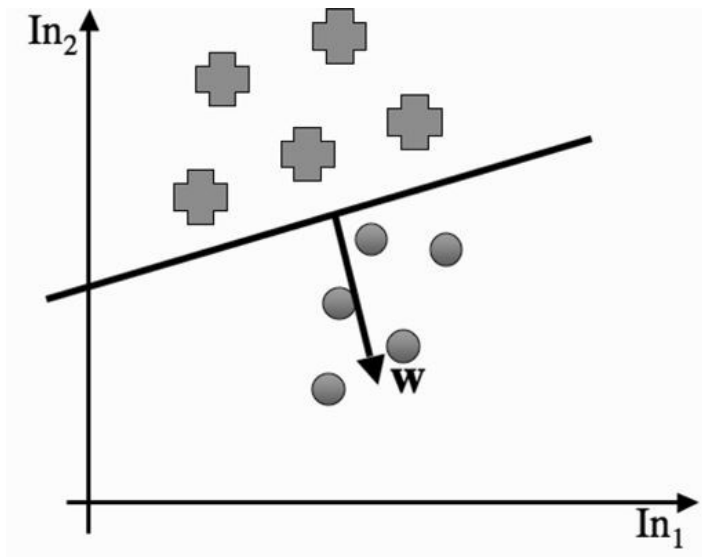
Plot the following function: $\frac{1}{2}x - y + 1 = 0$



Linear separability

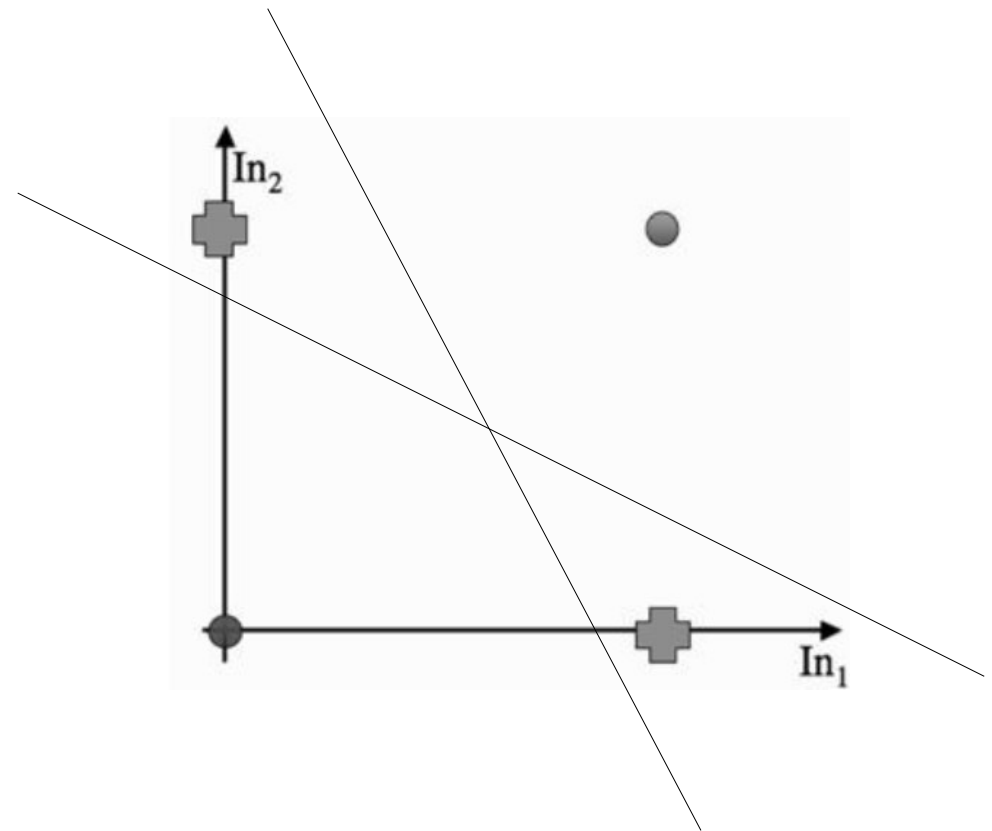
We have established that the decision boundary is a hyperplane.

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



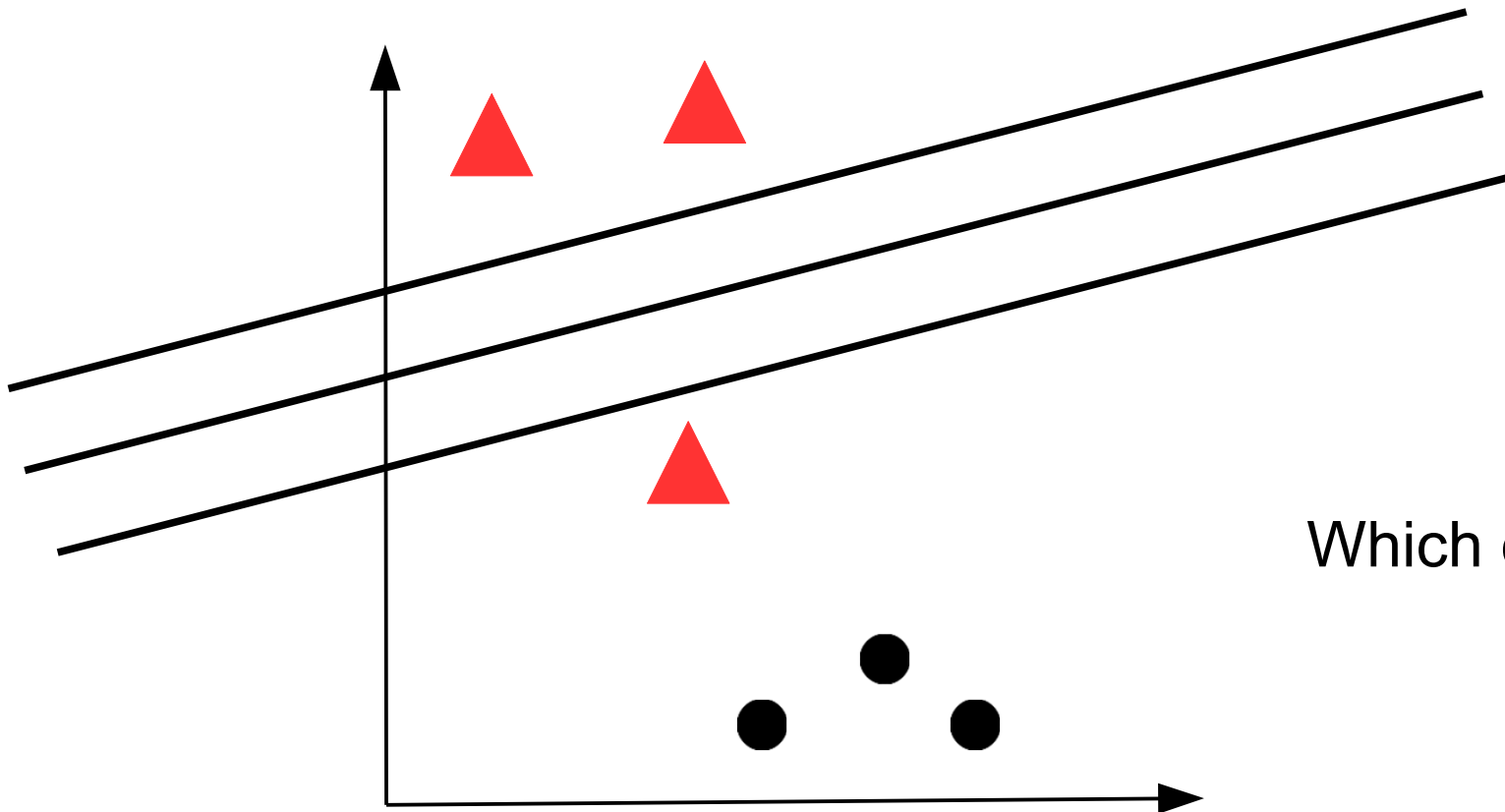
XOR

Not linearly separable!



Number of mistakes as Error

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

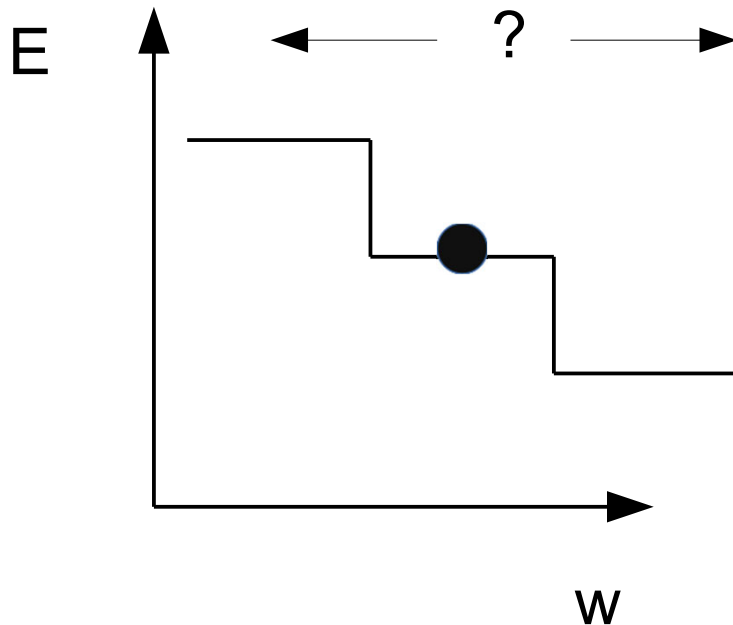


Which one is *better*?

Number of mistakes

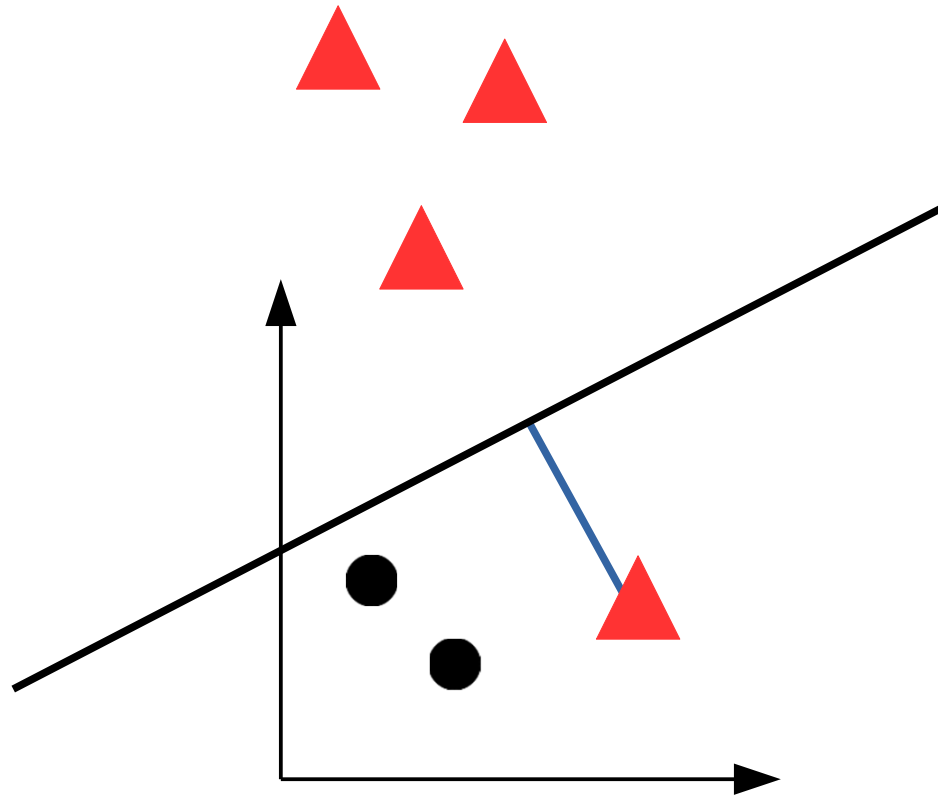
$$E(\mathbf{X}) = \sum_{\vec{x}_n \in \mathbf{X}} |y_n - t_n|$$

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



There is no local information
on the direction of
improvement

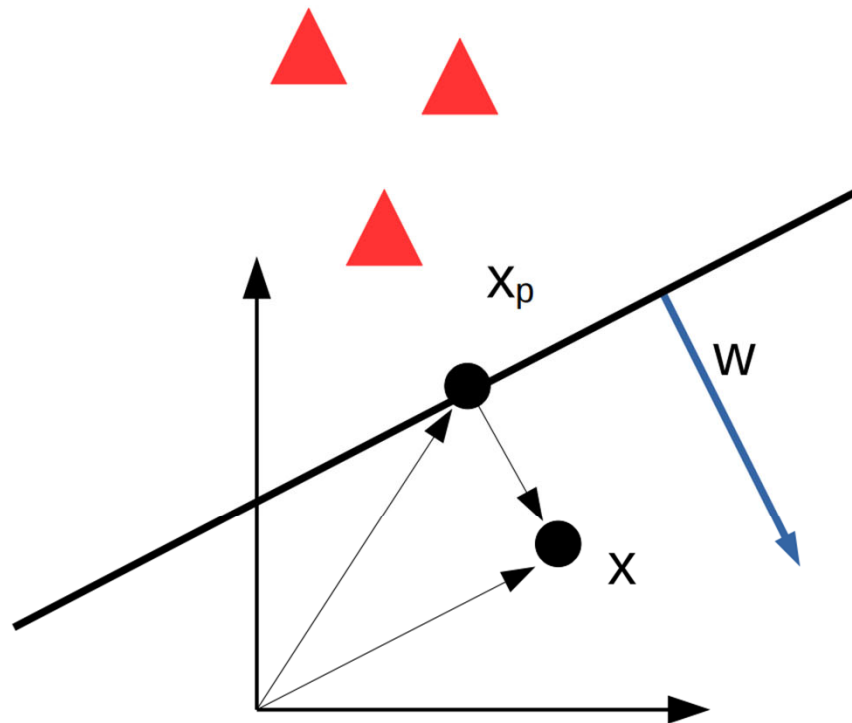
Towards a better error function



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

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



Distance to the hyperplane

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

$$h_w(\mathbf{x}) = \mathbf{w} \left(\mathbf{x}_p + d \frac{\mathbf{w}}{\|\mathbf{w}\|} \right) + w_0$$

$$= \cancel{\mathbf{w} \mathbf{x}_p} + w_0 + d \frac{\mathbf{w}^T \mathbf{w}}{\|\mathbf{w}\|} = d \|\mathbf{w}\|$$

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

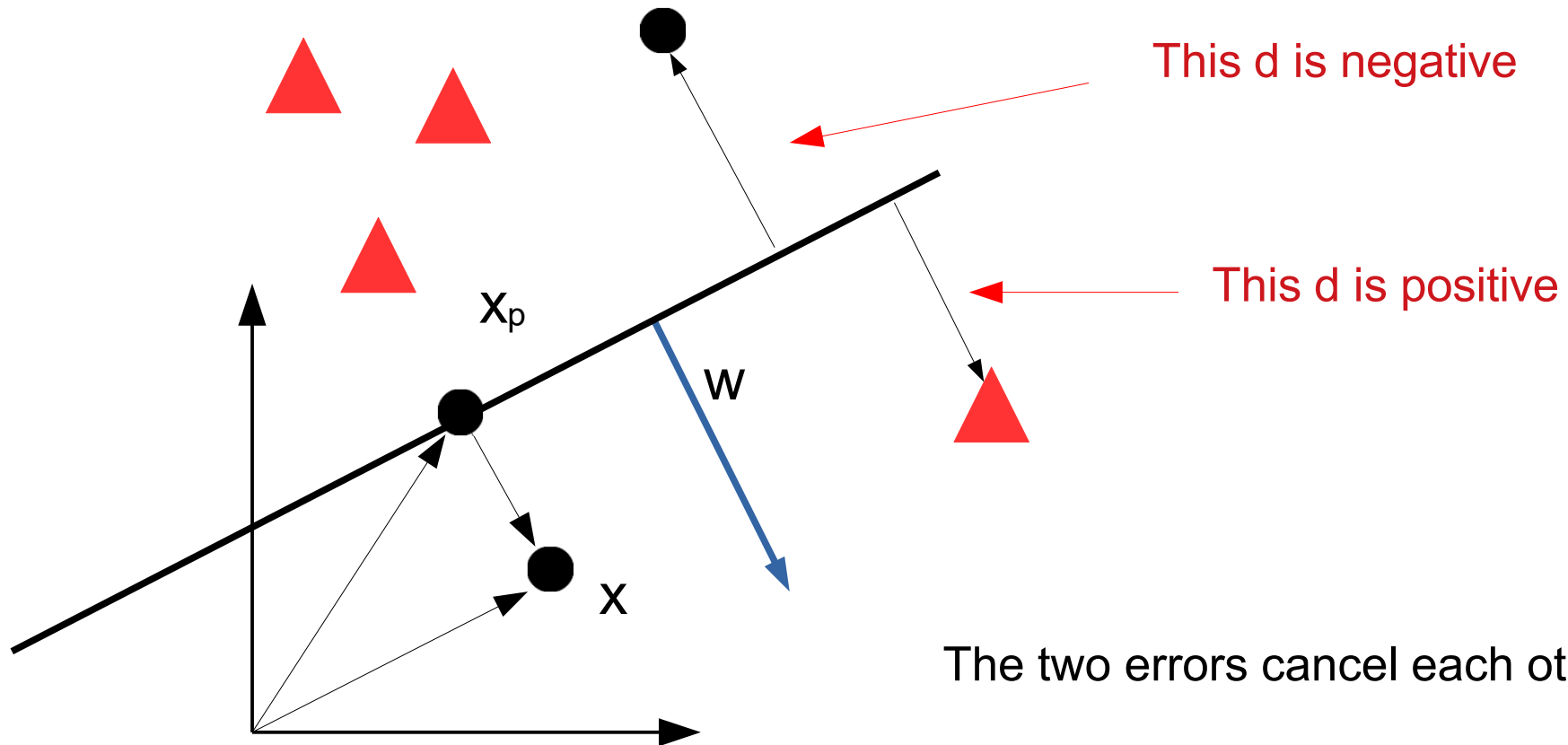
$$h_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = d \|\mathbf{w}\|$$

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

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

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

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

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

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

Question

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

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

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

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:

$$w_{k+1} = w_k - \eta \nabla f(w_k)$$

Which leads us to the update rule for the perceptron:

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

Stochastic gradient descent

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

Gradient:

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

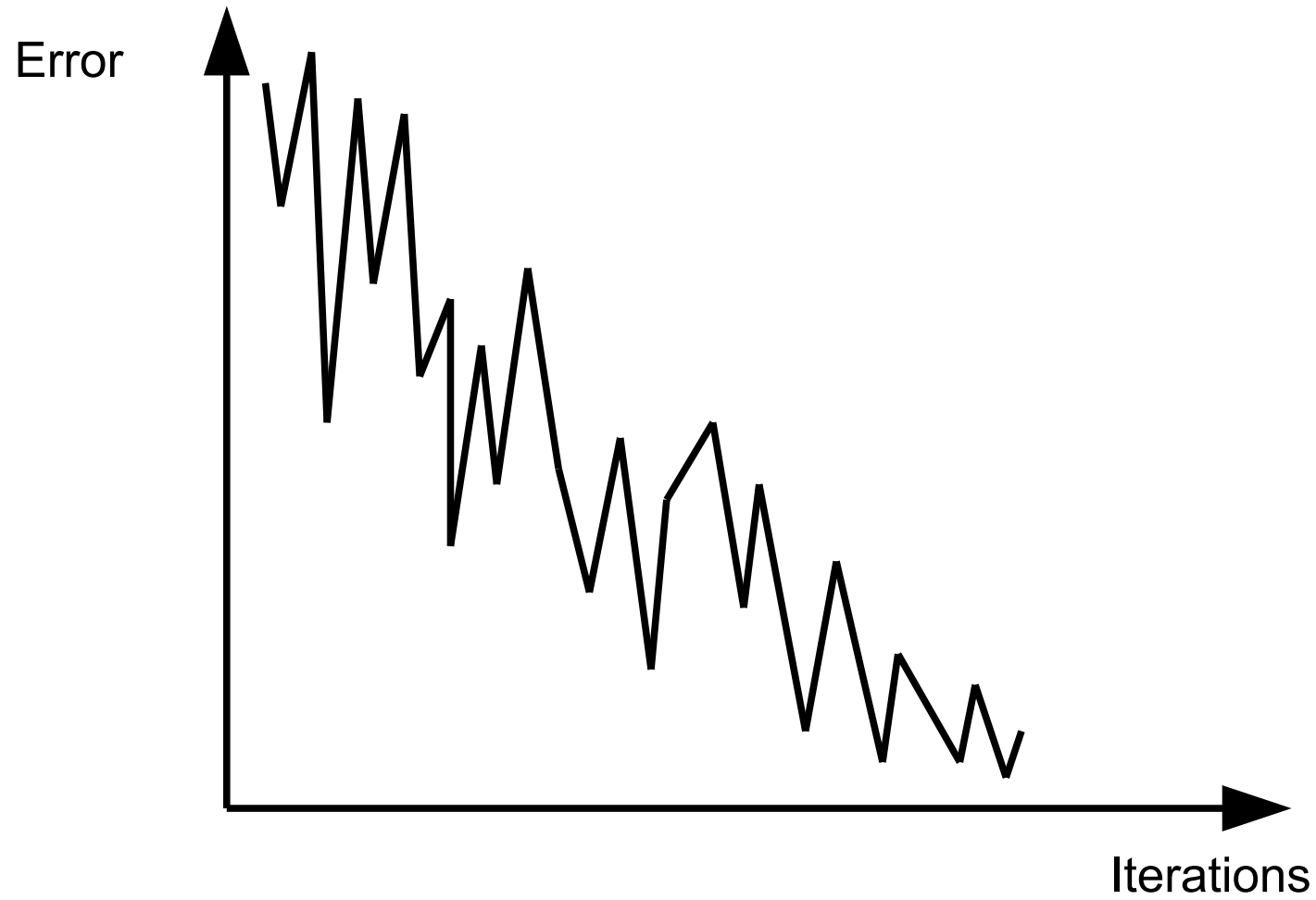
Stochastic gradient descent:

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

GD: you have to run through ALL the samples in your training set to do a single update for a parameter in a particular iteration

SGD: you use ONLY ONE or SUBSET of training sample from your training set to do the update for a parameter in a particular iteration. If you use SUBSET, it is called Minibatch Stochastic gradient Descent.

Stochastic gradient descent



Summary

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

