

Artificial neural networks

The Delta rule

- We wish to fit $y(\mathbf{x}) = g(\mathbf{w}^\top \mathbf{x})$ to a set of learning examples $\{(\mathbf{x}_1, t_1), \dots, (\mathbf{x}_N, t_N)\}$, where $\mathbf{x}_n \in \mathbb{R}^d, t_n \in \mathbb{R}$
- Define the (empirical) **mean-square error** of the network as:

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (t_n - y(\mathbf{x}_n))^2 = \frac{1}{2} \sum_{n=1}^N \left[t_n - g \left(\sum_{i=1}^d w_i x_{n,i} + w_0 \right) \right]^2$$

Artificial neural networks

The Delta rule

Let $f : \mathbb{R}^r \rightarrow \mathbb{R}$ differentiable; we wish to minimize it by making changes in its variables. Then the increment in each variable is proportional to the corresponding derivative: $x_i(t+1) := x_i(t) + \Delta x_i(t)$, with

$$\Delta x_i(t) = -\alpha \left. \frac{\partial f}{\partial x_i} \right|_{\mathbf{x}=\mathbf{x}(t)}, \quad \alpha > 0, i = 1, \dots, r$$

Illustration ($r = 1$). Let $f(x) = 3x^2 + x - 4$ and take $\alpha = 0.05$. We have $f'(x) = 6x + 1$. Then $x(0) = 1, x(1) = x(0) - \alpha f'(1) = 1 - 0.05 \cdot 7 = 0.65, x(2) = 0.65 - 0.05 \cdot 4.9 = 0.405, \dots$. We find $\lim_{i \rightarrow \infty} x(i) = -\frac{1}{6}$.

Artificial neural networks

The Delta rule

In our case, the function to be minimized is the empirical error and the variables are the weights w of the network:

$$\Delta w_j(t) = -\alpha \left. \frac{\partial E(w)}{\partial w_j} \right|_{w=w(t)}, \quad \alpha > 0, \quad j = 0, \dots, d$$

$$\frac{\partial E(w)}{\partial w_j} = - \sum_{n=1}^N (t_n - y(x_n)) g'(w^\top x_n) x_{n,j}$$

- $t_n - y(x_n)$ is called the **delta**
- $w^\top x_n$ is called the **net input**

Artificial neural networks

The Delta rule

$$\longrightarrow \Delta w_j(t) = \alpha \sum_{n=1}^N (t_n - y(x_n)) g'(w(t)^T x_n) x_{n,j}$$

When g is the identity, we get the α -LMS Learning Rule:

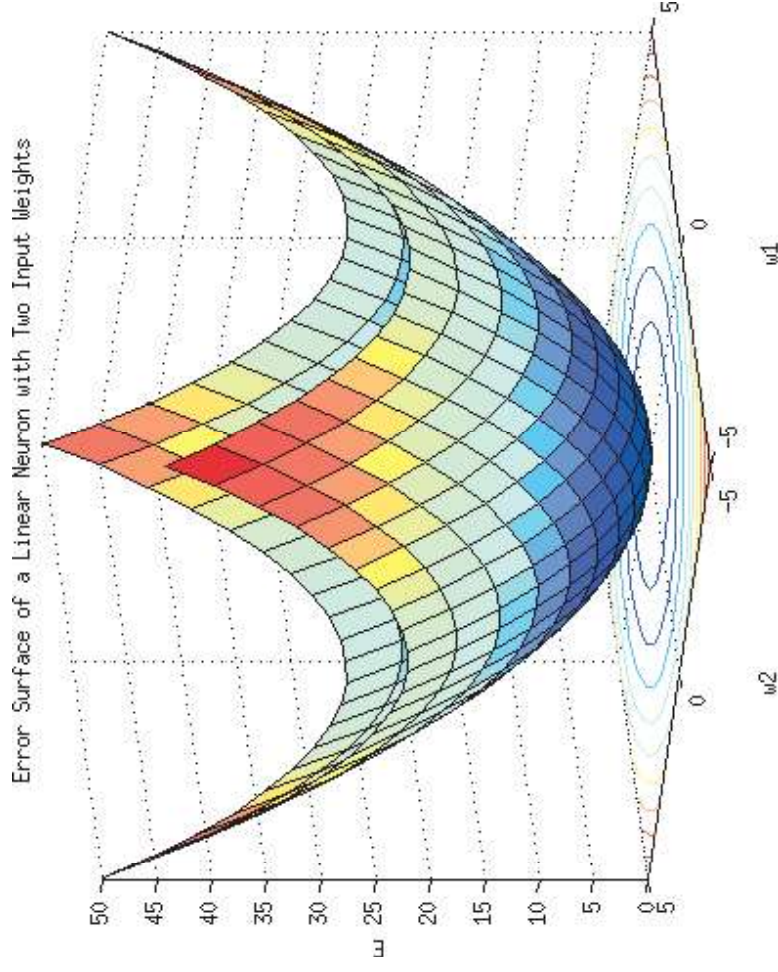
$$\Delta w_j(t) = \alpha \sum_{n=1}^N (t_n - y(x_n)) x_{n,j} = \alpha \sum_{n=1}^N (t_n - w^T x_n) x_{n,j}$$

-
- This technique represents a **linear regressor** where the regression coefficients are estimated iteratively (probably the most analyzed and applied **learning rule**)
 - This is a form of learning (because of the adaptation to the data) but it is **not incremental**: we need all the examples from the beginning (“**batch**” learning)

Artificial neural networks

The Delta rule

The function $E(w)$ is convex in w : it defines a convex hyper-paraboloidal surface with a single **global minimum** w^*



Artificial neural networks

The Delta rule

1. The constant α controls the stability and speed of convergence. If chosen sufficiently small, the gradient descent procedure asymptotically converges towards \mathbf{w}^* , regardless of the initial vector $\mathbf{w}(0)$
2. A sufficient condition is $0 < \alpha < \frac{2}{\lambda_{max}}$, where λ_{max} is the largest eigenvalue of the input auto-correlation matrix $\mathbb{E}[\mathbf{x}\mathbf{x}^\top] \approx \mathbf{X}\mathbf{X}^\top$

In practice, one may use $\alpha < \frac{2}{\sum_{n=1}^N \|\mathbf{x}_n\|^2}$, since $\lambda_{max} < \text{Tr}(\mathbf{X}\mathbf{X}^\top) \approx \sum_{n=1}^N \|\mathbf{x}_n\|^2$

Artificial neural networks

The Delta rule

In the **on-line** version, we also begin with $w(0)$ arbitrary and apply:

$$\Delta w_j(t) = \alpha_t (t_{n(t)} - y(x_{n(t)})) x_{n(t),j}$$

- At each step t , the example $n(t)$ is drawn at random from $\{1, \dots, N\}$
- It can be shown that if $\sum_{t \geq 0} \alpha_t = \infty$ and $\sum_{t \geq 0} \alpha_t^2 < \infty$, then $w(t)$ converges to the **global minimum** w^* in the mean square sense:

$$\lim_{t \rightarrow \infty} \|w(t) - w^*\|^2 = 0$$

One such procedure is $\alpha_t = \frac{\alpha}{t+1}$, with $\alpha > 0$