

Learning Processes :

Learning is a process by which the free parameters of a NN are adapted through a process of stimulation by the environment in which the network is embedded. The type of learning is determined by the manner in which the parameter changes take place.

Five basic learning rules :

- ① Error-Correction learning
- ② Memory-based learning
- ③ Hebbian Learning.
- ④ Competitive learning &
- ⑤ Boltzmann Learning.

Error-Correction learning is rooted in optimum Filtering.

* Memory-based learning operates by memorizing the training data explicitly.

* Hebbian & Competitive learning are both inspired by neurobiological considerations.

* Boltzmann is based on statistical mechanics.

Error-Correction Learning:

Consider the simple case of a neuron k constituting the only computational node in the output layer of a feedforward neural network as shown in the fig.

Neuron k is driven by a signal vector $x(n)$ produced by one or more layers of hidden neurons, which are themselves driven by an input vector (stimulus) applied to the source nodes (ie input layer) of the neural network.

n denotes discrete time

The output signal of neuron k is denoted by $y_k(n)$.

The output signal is compared to a desired response or target output, denoted by $d_k(n)$.

An error signal denoted by $e_k(n)$ is produced & given as

$$e_k(n) = d_k(n) - y_k(n) \quad \text{--- (1)}$$

The corrective adjustments are designed to make the output signal $y_k(n)$ come closer to the desired response $d_k(n)$ in a step-by-step manner. This is done by minimizing the cost function or index of performance, $E(n)$

$$E(n) = \frac{1}{2} e_k^2(n) \quad \text{--- (2)}$$

$E(n)$ is the instantaneous value of the error energy. The step-by-step adjustments to the synaptic weights of neuron k are continued until the system reaches a steady state. At that point the learning process is terminated.

The learning process is referred to as error-correction learning.

$E(n)$ the cost function leads to a learning rule commonly referred to as the delta rule or Widrow-Hoff rule.

Let $w_{kj}(n)$ denote the value of synaptic weight w_{kj} of neuron k excited by element $x_j(n)$ of the signal vector $x(n)$ at time step n .

According to the delta rule, the adjustment $\Delta w_{kj}(n)$ applied to the synaptic weight w_{kj} at time step n is defined by

$$\Delta w_{kj}(n) = \eta e_k(n) x_j(n) \quad \text{--- (3)}$$

η is a positive constant that determines the rate of learning

$\eta \rightarrow$ learning-rate parameter

The updated value of synaptic weight w_{kj} is determined by

$$w_{kj}(n+1) = w_{kj}(n) + \Delta w_{kj}(n) \quad \text{--- (4)}$$

In effect, $w_{kj}(n)$ and $w_{kj}(n+1)$ may be viewed as the old & new values of synaptic weight w_{kj}

$$w_{kj}(n) = Z^{-1} [w_{kj}(n+1)] \quad \text{--- (5)}$$

where Z^{-1} is the unit-delay operator

Z^{-1} represents a storage element

Method of Steepest Descent

In the method of steepest descent, the successive adjustments applied to the weight vector w are in the direction of steepest descent, that is in a direction opposite to the gradient vector $\nabla \mathcal{E}(w)$.

For convenience of presentation we write

$$g = \nabla \mathcal{E}(w) \quad \text{--- (1)}$$

Accordingly, the steepest descent algorithm is

formally described by

$$w(n+1) = w(n) - \eta g(n) \quad \text{--- (2)} \quad \Delta w_{kj}(n) = \eta e_k(n) z_j(n)$$

$$w(n+1) - w(n) = -\eta g(n)$$

where η is a positive constant called the

stepsize or learning-rate parameter and $g(n)$ is gradient vector evaluated at the point $w(n)$.

In going from iteration n to $n+1$ the algorithm applies the correction

$$\Delta w(n) = w(n+1) - w(n) \quad \text{--- (3)}$$

$$\Delta w(n) = -\eta g(n)$$

First-order Taylor series expansion around $w(n)$ to approximate $\mathcal{E}(w(n+1))$ as

$$\mathcal{E}(w(n+1)) = \mathcal{E}(w(n)) + g^T(n) \Delta w(n)$$

the use of which is justified for small η .

Substituting eqn (3) in this approximate relation yields

$$\begin{aligned} \mathcal{E}(w(n+1)) &= \mathcal{E}(w(n)) - \eta g^T(n)g(n) \\ &= \mathcal{E}(w(n)) - \eta \|g(n)\|^2 \end{aligned}$$

Which shows that, for a positive learning-rate parameter η , the cost function is decreased as the algorithm progresses from one iteration to the next.

* when η is small, the transient response of the algorithm is overdamped.

* when η is large, the transient response of the algorithm is underdamped.

Linear Least-Squares ~~Filter~~ Filter :

Linear least-squares filter has two distinctive characteristics. First, the single ~~neuron~~ neuron around which it is built is linear, as shown in the model

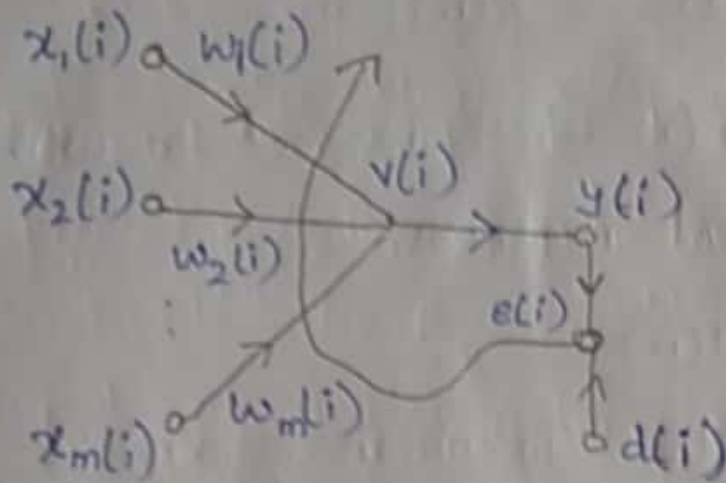


fig: model of neuron

Second the cost function $\mathcal{E}(w)$ used to design the filter consists of the sum of error squares as

$$\text{defined in } \mathcal{E}(w) = \frac{1}{2} \sum_{i=1}^n e^2(i)$$

$$\& y(i) = v(i) = \sum_{k=1}^m w_k(i) x_k(i)$$

$$y(i) = x^T(i) w(i) \quad \text{where } w(i) = [w_1$$

$$\text{where } w(i) = [w_1(i), w_2(i), \dots, w_m(i)]^T$$

$$e(i) = d(i) - y(i)$$

Express the error vector $e(n)$ as follows

$$e(n) = d(n) - [x(1), x(2), \dots, x(n)]^T w(n)$$

$$= d(n) - x(n) w(n) \quad \text{--- (1)}$$

where $d(n)$ is the n -by-1 desired response vector.

$$d(n) = [d(1), d(2) \dots d(n)]^T$$

and $x(n)$ is the n -by- m data matrix

$$x(n) = [x(1), x(2) \dots x(n)]^T$$

Differentiating eqn (1) with respect to $w(n)$ yields the gradient matrix.

$$\nabla e(n) = -x^T(n)$$

Correspondingly, the Jacobian of $e(n)$ is

$$J(n) = -x(n) \quad \text{--- (2)}$$

Substitute eqn (1) & (2) in the below eqn

$$w(n+1) = w(n) - (J^T(n) J(n))^{-1} J^T(n) e(n)$$

$$w(n+1) = w(n) + (x^T(n) x(n))^{-1} x^T(n) (d(n) - x(n) w(n))$$

$$w(n+1) = w(n) + (x^T(n) x(n))^{-1} x^T(n) \cdot d(n) - \cancel{(x^T(n) x(n))^{-1} x^T(n)} \cdot x(n) w(n)$$

$$= w(n) - \frac{x(n) w(n) x^T(n)}{x^T(n) x(n)} + (x^T(n) x(n))^{-1} x^T(n) \cdot d(n)$$

$$w(n+1) = (x^T(n) x(n))^{-1} x^T(n) \cdot d(n) \quad \text{--- (3)}$$

The term $(x^T(n) x(n))^{-1} x^T(n)$ is pseudoinverse of the matrix data matrix $x(n)$.

$$x^+(n) = (x^T(n) x(n))^{-1} x^T(n) \quad \text{--- (4)}$$

Rewrite eqn $w(n+1) = x^+(n) \cdot d(n)$

The weight vector $w(n+1)$ solves the linear least-squares

Wiener Filter: Limiting form of the linear Least-Squares filter for an Ergodic environment

The input vector $x(i)$ and desired response $d(i)$ are drawn from an ergodic environment that is also stationary.

The second-order statistics is described by

* Correlation matrix of the input vector $x(i)$; it is denoted by R_x

* Cross-correlation vector between the input vector $x(i)$ and desired response $d(i)$; it is denoted by r_{xd}

These two quantities are defined as follows, respectively.

$$R_x = E[x(i)x^T(i)]$$

$$= \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^{\infty} x(n)x^T(n)$$

$$= \lim_{n \rightarrow \infty} \frac{1}{n} x^T(n)x(n)$$

$$\& r_{xd} = E[x(i)d(i)]$$

$$= \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n x(i)d(i)$$

$$= \lim_{n \rightarrow \infty} \frac{1}{n} x^T(n) \cdot d(n)$$

Where E denotes the statistical expectation operator.

Accordingly we may reformulate the linear least-squares soln of eqn (1) as follows

$$w(n+1) = (X^T(n)X(n))^{-1} X^T(n) \cdot d(n) \quad \text{--- (1)}$$

$$w_0 = \lim_{n \rightarrow \infty} w(n+1)$$

$$= \lim_{n \rightarrow \infty} (X^T(n)X(n))^{-1} \cdot X^T(n) \cdot d(n)$$

$$= \lim_{n \rightarrow \infty} \frac{1}{n} \underbrace{(X^T(n)X(n))^{-1}}_{R_x} \cdot \lim_{n \rightarrow \infty} \frac{1}{n} \underbrace{X^T(n) \cdot d(n)}_{\gamma_{xd}}$$

$$w_0 = R_x^{-1} \gamma_{xd}$$

Where R_x^{-1} is the inverse of the correlation matrix R_x

The weight vector w_0 is called the Wiener solution to the linear optimum filtering problem.