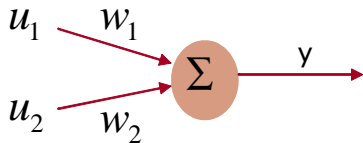


Training of NN - 2

Wen Yu

<https://www.ctrl.cinvestav.mx/~yuw/>
yuw@ctrl.cinvestav.mx



$$y = w_1 u_1 + w_2 u_2$$

The weights are trained as

$$w(k+1) = w(k) + \eta e(k) u_i$$
$$e(k) = y(k) - \hat{y}(k)$$

$$w_1(k+1) = w_1(k) + \eta [y(k) - \hat{y}(k)] u_1$$
$$w_2(k+1) = w_2(k) + \eta [y(k) - \hat{y}(k)] u_2$$

Least square method

But

$$y = w_1 u_1 + w_2 u_2$$

$$y(1) = w_1 u_1(1) + w_2 u_2(1)$$

$$y(2) = w_1 u_1(2) + w_2 u_2(2)$$

Or

$$\begin{bmatrix} y(1) \\ y(2) \end{bmatrix} = \begin{bmatrix} u_1(1) & u_2(1) \\ u_1(2) & u_2(2) \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

So

$$\begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} u_1(1) & u_2(1) \\ u_1(2) & u_2(2) \end{bmatrix}^{-1} \begin{bmatrix} y(1) \\ y(2) \end{bmatrix}$$

But if

$$y(1) = w_1 u_1(1) + w_2 u_2(1)$$

$$y(2) = w_1 u_1(2) + w_2 u_2(2)$$

$$y(3) = w_1 u_1(3) + w_2 u_2(3)$$

Least square method

Gauss (1777-1855). How to design a line

$$y = ax + b$$

by a group of data (x_i, y_i) . For two independent data (x_1, y_1) and (x_2, y_2)

$$a = \frac{y_1 - y_2}{x_1 - x_2}, \quad b = \frac{-x_1 y_2 - x_2 y_1}{x_1 - x_2}$$

Least square method

For n data (x_i, y_i) , how to find a best line? We define

$$e_i = y_i - (ax_i + b)$$

The square error is

$$J = \sum_i^n e_i^2$$

Least square is

$$\min_{(a,b)} J$$

Least square method

The solution for (a, b) is

$$\frac{\partial J}{\partial a} = 0, \quad \frac{\partial J}{\partial b} = 0$$

So

$$a = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2}$$
$$b = \frac{\sum x_i^2 \sum y_i^2 - \sum x_i \sum x_i y_i}{n \sum x_i^2 - (\sum x_i)^2}$$

Least square method

In general case, we use linear regression

$$\hat{y} = w_1 x_1 + w_2 x_2 + \cdots + w_m x_m$$

where $(x_1 \cdots x_m, y)$ is known data, $(w_1 \cdots w_m)$ is parameter

For each data $(x_{i1} \cdots x_{mi}, y_i)$

$$y_1 = w_1 x_{11} + w_2 x_{12} + \cdots + w_m x_{1m} + e_1$$

$$y_2 = w_1 x_{21} + w_2 x_{22} + \cdots + w_m x_{2m} + e_2$$

\vdots

$$y_n = w_1 x_{n1} + w_2 x_{n2} + \cdots + w_m x_{nm} + e_n$$

In matrix form

$$Y = XW + E$$

where $W = [w_1 \cdots w_m]^T$, $Y = [y_1 \cdots y_n]^T$,

$$X = \begin{bmatrix} x_{11} & \cdots & x_{1m} \\ \vdots & \ddots & \vdots \\ x_{n1} & & x_{nm} \end{bmatrix} \in R^{n \times m}$$

Least square method

Because

$$J = E^T E = \sum_i^n e_i^2, \quad \min_W J \rightarrow \text{least square}, \quad \frac{\partial J}{\partial W} = 0$$

$$\begin{aligned} \frac{\partial J}{\partial W} &= \frac{\partial J}{\partial E} \frac{\partial E}{\partial W} = 2E \frac{\partial}{\partial W} (Y - XW) \\ &= 2(Y - XW)^T (-X) \\ &= -2Y^T X + 2W^T X^T X = 0 \end{aligned}$$

So

$$\begin{aligned} Y &= XW + E \\ W &= (X^T X)^{-1} X^T Y \end{aligned}$$

Discrete-time linear system

$$y(k) = -a_1y(k-1) - \dots - a_ny(k-n) \\ + b_1u(k-1) + \dots + b_mu(k-m)$$

It is ARX model

$$Ay(k) = Bu(k)$$

where $A = 1 + a_1q^{-1} + \dots + a_nq^{-n}$, $B = b_1q^{-1} + \dots + b_muq^{-m}$.

With color noise, it is ARMAX model

$$Ay(k) = Bu(k) + C\xi(k)$$

System identification -least square

Parameter-in-linear

$$y = \theta_1 x_1 + \dots + \theta_n x_n$$

where

$$x_i = \phi(u), \quad \text{or } x_i = \phi(Vu)$$

In matrix form

$$Y = X\Theta + E$$

$$\text{where } Y = [y_1 \dots y_m]^T, \quad X = \begin{bmatrix} x_{11} & \dots & x_{1n} \\ \vdots & \ddots & \vdots \\ x_{m1} & & x_{mn} \end{bmatrix} \in R^{m \times n},$$

$$\Theta = [\theta_1 \dots \theta_n]^T, \quad E = [e_1 \dots e_m]^T.$$

$$Y = X\Theta + E, \quad J = E^T E = \sum_i^m e_i^2$$

$\min_{\Theta} J \rightarrow$ least square, $\frac{\partial J}{\partial \Theta} = 0$

$$\begin{aligned} \frac{\partial J}{\partial \Theta} &= \frac{\partial J}{\partial E} \frac{\partial E}{\partial \Theta} = 2E \frac{\partial}{\partial \Theta} (Y - X\Theta) \\ &= 2(Y - X\Theta)^T (-X) \\ &= -2Y^T X + 2\Theta^T X^T X \\ &= 0 \end{aligned}$$

So

$$\Theta = (X^T X)^{-1} X^T Y$$

Linear model

In linear regressive form

$$y(k) = \varphi^T(k) \hat{\theta} + e(k)$$

where

$$\varphi(k) = [-y(k-1) \cdots -y(k-n), u(k-1) \cdots u(k-m)]^T$$
$$\hat{\theta} = [\hat{a}_1 \cdots \hat{a}_n, \hat{b}_1 \cdots \hat{b}_m]^T$$

For l data $y(k)$ and $\varphi(k)$

Least square

$$\hat{\theta} = \left[\sum_{k=1}^l \varphi^T(k) \varphi(k) \right]^{-1} \sum_{k=1}^l y(k) \varphi(k)$$

where $l \gg \max(m, n)$. With this θ

$$\min J = \min \sum_{k=1}^l e^2(k)$$

To find the relation

$$\hat{\theta} = \left[\sum_{k=1}^l \varphi^T(k) \varphi(k) \right]^{-1} \sum_{k=1}^l y(k) \varphi(k)$$

$$\begin{aligned} \theta_{k+1} &= \theta_k + P_{k+1} \varphi(k+1) [y(k+1) - \varphi^T(k+1) \theta_k] \\ P_{k+1} &= P_k - \frac{P_k \varphi^T(k+1) \varphi(k+1) P_k}{1 + \varphi(k+1) P_k \varphi^T(k+1)}, \quad P_k \gg 0 \end{aligned}$$

Extreme learning machines (ELM), ESN

It is feedforward neural network with a single hidden layer, the output active function is linear,

$$\hat{y} = W\phi(Vx) \quad (1)$$

- the weights V in the hidden layer are random
- W will be trained by the least square method

The linear regression is (NN) $\phi_m(Vx_i)$

$$\hat{y} = w_1x_1 + w_2x_2 + \dots + w_mx_m$$

The real data is

$$y_i = w_1\phi_m(Vx_i) + w_2\phi_m(Vx_i) + \dots + w_m\phi_m(Vx_i) + e_i$$

where

$$e_i = y_i - \hat{y}$$

Extreme learning machines (ELM)

$$\begin{aligned} Y &= ZW + E, & Z &= \phi(Vx_i) \\ y_i &= W\phi(Vx_i) + e_i \end{aligned} \quad (2)$$

If we want to

$$J = E^T E = \sum_i^n e_i^2, \quad \min_W J$$

then

$$W = (Z^T Z)^{-1} Z^T Y$$

$$\text{where } Y = [y_1 \cdots y_n], \quad Z = \begin{bmatrix} \phi(Vx_{11}) & \cdots & \phi(Vx_{1m}) \\ \vdots & \ddots & \vdots \\ \phi(Vx_{n1}) & & \phi(Vx_{nm}) \end{bmatrix} \in R^{n \times m}$$

```
for i=1:M "training time"  
    for j=1:N  
        I(j,i)=0; "input to each hidden node"  
        for k=1:L  
            I(j,i)=I(j,i)+W(k,j)*II(k,i);  
        end  
        O(j,i)=(exp(I(j,i))-exp(-I(j,i)))/(exp(I(j,i))+exp(-I(j,i)));  
    end  
    Y(i)=y(i); "target output (teacher)"  
end  
V=Y*pinv(O); "weights in uotput layer, Moore-Penrose pseudoinverse  
training"
```


Variations on learning rate

- Fixed rate much smaller than 1
- Start with large η , gradually decrease its value
- Start with a small η , steadily double it until MSE start to increase
- Find the maximum safe step size at each stage of learning (to avoid overshoot the minimum J when increasing η)
- Adaptive learning rate (Conjugate gradient minimisation)
 - Each weight w_{ij} has its own rate η_{ij}
 - If Δw_{ij} remains in the same direction, increase η_{ij}
 - If Δw_{ij} changes the direction, decrease η_{ij}

Variations on learning rate

- time-varying rate as $\eta(k) = \frac{c}{k}$ can give optimally fast convergence, but it results in slow convergence to bad solutions when c and k are small.
- Modification $\eta(k) = \frac{\eta_0}{1+a_1\frac{c}{k}+a_2\left(\frac{c}{k}\right)^2}$.
- Stable learning $\eta(k) = \frac{\eta_0}{1+\|\varphi(x_k)\|}$
- new learning rate as $\eta(k) = \frac{\eta_0}{1+a_1\frac{c}{k}+\|G(x_k)\|}$

Plant

$$y(k) = \varphi^T(k) \theta^*$$

Model

$$\hat{y}(k) = \varphi^T(k) \theta$$

$$\theta(k+1) = \theta(k) - \eta_k \varphi(k) e(k), \quad e(k) = \hat{y}(k) - y(k) \quad (3)$$

Gradient

$$\eta_k = \eta$$

$$\theta(k+1) = \theta(k) - \eta \varphi(k) e(k)$$

Least square $\eta_k \in R^{n \times n}$

$$\eta_k = \frac{P_k}{1 + \varphi P_k^T \varphi^T}$$

$$\theta_{k+1} = \theta_k + \eta_k \varphi e_k$$

$$P_{k+1} = P_k - \frac{P_k \varphi^T \varphi P_k}{1 + \varphi P_k \varphi^T}, \quad P_1 \gg 0$$

Kalman $\eta_k \in R^{n \times n}$ is matrix

$$\eta_k = \frac{P_k}{R_2 + \varphi P_k \varphi^T}$$

$$\theta_{k+1} = \theta_k + \eta_k \varphi e_k$$

$$P_{k+1} = P_k - \frac{P_k \varphi^T \varphi P_k}{R_2 + \varphi P_k \varphi^T} + R_1$$

where $R_1, R_2 > 0$

Stable learning

Nonlinear system

$$y(k) = f[X(k)]$$

The neural networks is

$$\hat{y}(k) = W_k \phi[X(k)]$$

If the identified nonlinear system can be represented as (matching condition)

$$y(k) = W^* \phi[X(k)]$$

The training law is

$$W_{k+1} = W_k - \eta \phi[X(k)] e(k)$$

where

$$e(k) = \hat{y}(k) - y(k)$$

$$\begin{aligned}\hat{y}(k) &= W_k \phi[X(k)] \\ y(k) &= W^* \phi[X(k)]\end{aligned}$$

Modeling error is

$$\begin{aligned}e(k) &= \hat{y}(k) - y(k) \\ &= [W_k - W^*] \phi[X(k)] \\ &= \tilde{W}_k \phi[X(k)]\end{aligned}$$

where

$$\tilde{W}_k = W_k - W^*$$

We select Lyapunov function as

$$V_k = \|\tilde{W}_k\|^2 = \text{tr} \left\{ \tilde{W}_k^T \tilde{W}_k \right\}$$

Because

$$\begin{aligned}\tilde{W}_{k+1} &= \tilde{W}_k - \eta_k \phi[X(k)] e(k) \\ e(k) &= \tilde{W}_k \phi[X(k)]\end{aligned}$$

From the updating law

$$\begin{aligned}\Delta V_k &= V_{k+1} - V_k = \|\tilde{W}_k - \eta_k \phi[X(k)] e(k)\|^2 - \|\tilde{W}_k\|^2 \\ &= \|\tilde{W}_k\|^2 - 2\eta_k e(k) \phi[X(k)] \tilde{W}_k + \|\eta_k \phi[X(k)] e(k)\|^2 - \|\tilde{W}_k\|^2 \\ &= \eta_k^2 e^2(k) \|\phi[X(k)]\|^2 - 2\eta_k e(k) \phi[X(k)] \tilde{W}_k \\ &= \eta_k^2 e^2(k) \|\phi[X(k)]\|^2 - 2\eta_k e^2(k) \\ &= \left[\eta_k \|\phi[X(k)]\|^2 - 2 \right] \eta_k e^2(k)\end{aligned}$$

We select

$$\eta_k = \frac{2\eta_0}{1 + \|\phi[X(k)]\|^2}, \quad 0 \leq \eta_0 \leq 1$$

then

$$\left[\eta_k \|\phi[X(k)]\|^2 - 2 \right] = 2 \left[\frac{\|\phi[X(k)]\|^2}{1 + \|\phi[X(k)]\|^2} \eta_0 - 1 \right] \leq 2[\eta_0 - 1] \leq 0$$

so

$$\Delta V_k \leq 0$$