

EE5904 Multi-layer Perceptron (MLP)

Last week:

① simple perceptron \leftrightarrow Pattern Recognition

linearly separable $\Leftrightarrow \exists \underline{w}$ s.t. classify correctly!

$\begin{cases} \text{off-line} \\ \text{on-line} \rightarrow \text{Perceptron Convergence Alg.} \end{cases}$

② Regression

input $\begin{cases} x_1(i) \\ x_2(i) \\ \vdots \\ x_m(i) \end{cases} \rightarrow \text{Unknown System} \rightarrow \text{output } y(i)$

$y(i) = v(i) = \underline{w}^T(i) \underline{x}(i)$

learning framework

LLS \leftarrow ① $E(w)$ cost f^2

LMS \leftarrow ② $\begin{cases} \frac{\partial E}{\partial w} = 0 \rightarrow \text{if task is easy } (\frac{\partial E}{\partial w} \approx 0 \text{ has closed form sol}^n) \\ \text{iterated method} \end{cases}$

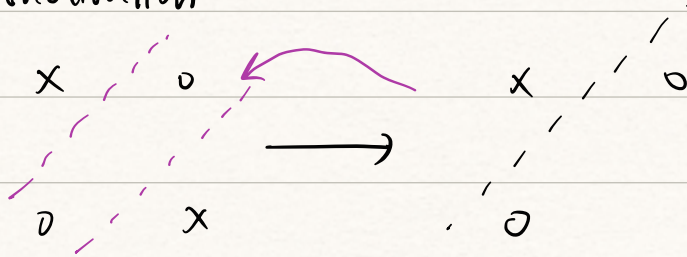
\Downarrow $\begin{cases} \text{GD} \\ \text{Newton's Method} \end{cases}$

Limits

$\begin{cases} \text{① classification} \rightarrow \text{linearly separable} \\ \text{② reg} \rightarrow \text{linear model} \end{cases}$

This week:

MLP \rightarrow motivation \rightarrow XOR



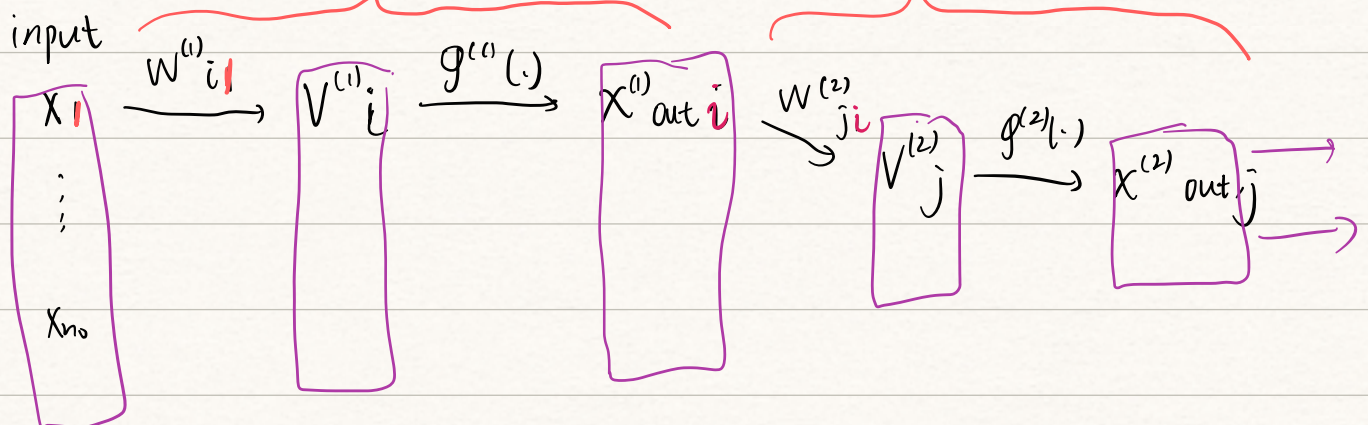
How to design? \Rightarrow in lecture note!

MLP \leftrightarrow David Rumelhart & PDP

① hidden layer \leftrightarrow neuron y_j (not include $V_j \rightarrow$ induced local field)

$$\underbrace{V_j \xrightarrow{\text{sigmoid}} y_j}_{\text{activation } f^-}$$

② BP Algorithm (Learning) \uparrow 1-st hidden neuron



$\underline{d}(n) \rightarrow$ desired output $\underline{y}(n) \rightarrow$ actual

$$\rightarrow \underline{e}(n) = \underline{d}(n) - \underline{y}(n)$$



we want $\underline{w}(n+1) = \underline{w}(n) + \underline{\Delta w}(n)$



$\Delta w(n) := -\eta \nabla E(n)$

Define cost f^2 $E(n) = \frac{1}{2} \underline{e}(n)^T \underline{e}(n)$

$\nabla_w E(n)^T = \frac{\partial E(n)}{\partial w(n)}$

$\Delta w_{ji}^{(3)}(n) = -\eta \frac{\partial E(n)}{\partial w_{ji}^{(3)}(n)}$

$\Delta w_{ji}^{(3)}(n) = -\eta \frac{\partial E(n)}{\partial w_{ji}^{(3)}(n)}$

$E(n) = \frac{1}{2} \sum_{j=1}^{n_3} e_j(n)^2 = \frac{1}{2} \sum_{j=1}^{n_3} (d_j(n) - x_{out,j}^{(3)}(n))^2$

① $\frac{\partial E(n)}{\partial x_{out,j}^{(3)}(n)} = (d_j(n) - x_{out,j}^{(3)}(n)) (-1)$

$= -e_j(n)$

② $y_j(n) = x_{out,j}^{(3)} = \phi^{(3)}(V_j^{(3)}(n)) \frac{\partial x_{out,j}^{(3)}(n)}{\partial V_j^{(3)}(n)} = \phi^{(3)'}(V_j^{(3)}(n))$

③ $\frac{\partial V_j^{(3)}(n)}{\partial w_{ji}^{(3)}(n)} = x_{out,i}^{(2)}(n)$

$V_j^{(3)}(n) = \sum_{i=1}^{n_2} w_{ji}^{(3)}(n) x_{out,i}^{(2)}(n)$

$\frac{\partial E(n)}{\partial w_{ji}^{(3)}(n)} = \frac{\partial E(n)}{\partial x_{out,j}^{(3)}(n)} \cdot \frac{\partial x_{out,j}^{(3)}(n)}{\partial V_j^{(3)}(n)} \cdot \frac{\partial V_j^{(3)}(n)}{\partial w_{ji}^{(3)}(n)}$

$= -e_j(n) \phi^{(3)'}(V_j^{(3)}(n)) x_{out,i}^{(2)}(n)$

$$:= -\delta_j^{(3)}(n) X_{out,i}^{(2)}(n)$$

$$\Rightarrow \Delta W_{ji}^{(3)}(n) = -\eta \frac{\partial E(n)}{\partial W_{ji}^{(3)}(n)} \quad \text{similar to eqn}$$

$$= \eta \underbrace{\delta_j^{(3)}(n)}_{\text{output}} \underbrace{X_{out,i}^{(2)}(n)}_{\text{input}}$$

$W^{(2)} \rightarrow$ 五层网络

$$\frac{\partial E(n)}{\partial W_{ji}^{(2)}(n)} = \sum_{k=1}^{n_3} \left(\frac{\partial E(n)}{\partial X_{out,k}^{(1)}(n)} \frac{\partial X_{out,k}^{(1)}(n)}{\partial V_k^{(3)}(n)} \frac{\partial V_k^{(3)}(n)}{\partial X_{out,j}^{(2)}(n)} \frac{\partial X_{out,j}^{(2)}(n)}{\partial V_j^{(2)}(n)} \frac{\partial V_j^{(2)}(n)}{\partial W_{ji}^{(2)}(n)} \right)$$

$$= \sum_{k=1}^{n_3} \left(-\delta_k^{(3)}(n) W_{kj}^{(3)}(n) \phi^{(2)'}(V_j^{(2)}(n)) X_{out,i}^{(1)}(n) \right)$$

$$= \phi^{(2)'}(V_j^{(2)}(n)) X_{out,i}^{(1)}(n) \sum_{k=1}^{n_3} (-\delta_k^{(3)}(n) W_{kj}^{(3)}(n))$$

$$:= -\delta_j^{(2)}(n) X_{out,i}^{(1)}(n)$$

$$\Rightarrow W_{ji}^{(2)}(n+1) = W_{ji}^{(2)}(n) + \eta \delta_j^{(2)}(n) X_{out,i}^{(1)}(n)$$

output input

$$\delta_j^{(3)} = \underbrace{(d_j - x_{out,i,j})}_{\text{network error}} \underbrace{\phi'(V_j^{(3)})}_{\text{gradient}} \Rightarrow \underbrace{\text{output layer}}_{\text{output error}}$$

$$\delta_j^{(2)} = \left(\sum_{k=1}^{n_3} \underbrace{\delta_k^{(3)}}_{\text{output error}} \underbrace{w_{kj}^{(3)}}_{\text{gradient}} \right) \underbrace{\phi'(V_j^{(2)})}_{\text{error propagate through network}}$$

$$\delta_j^{(1)} = \left(\sum_{k=1}^{n_2} \underbrace{\delta_k^{(2)}}_{\text{output error}} \underbrace{w_{kj}^{(2)}}_{\text{gradient}} \right) \underbrace{\phi'(V_j^{(1)})}_{\text{error propagate through network}}$$

hidden layer → output error propagate backwards

BR

Rate of Learning

Momentum

Origin: $W_{ji}^{(s)}(u+1) = W_{ji}^{(s)}(u) + \eta^{(s)} \delta_j^{(s)} x_{out,i}^{(s-1)}$

↓

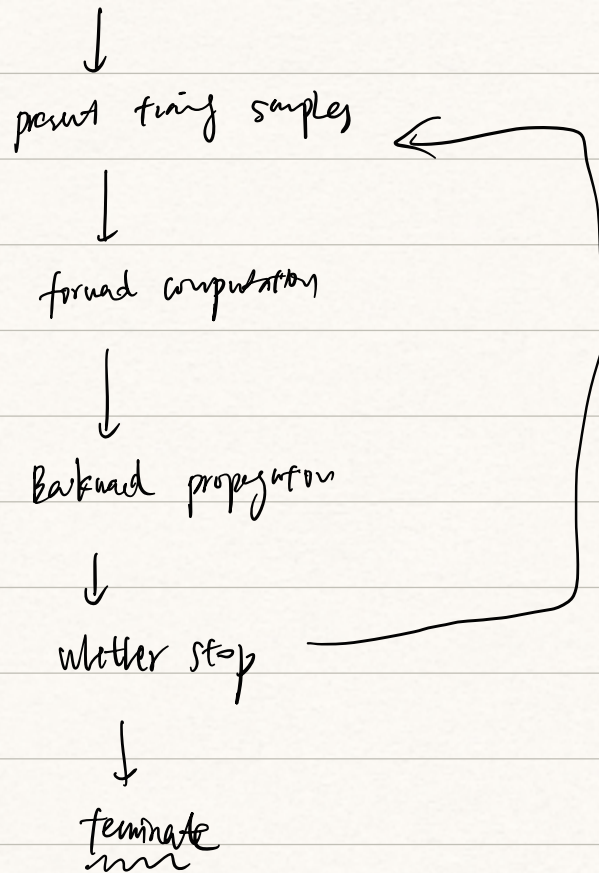
Now: $\Delta W_{ji}^{(s)}(k) = \underbrace{\alpha \Delta W_{ji}^{(s)}(u-1)}_{\text{history}} + \eta^{(s)} \delta_j^{(s)} x_{out,i}^{(s-1)}(u)$

epoch : use whole training samples!

stopping criteria

Summary

Initialization



Limitation: → Black-Box

↓
可解释性差