

DSA5202 LEC4 → Neural Network

1. Shallow NN

→ formulation:
$$f(x) = \sum_{j=1}^M v_j \phi(w_j^T x + b_j)$$

$$= v^T \phi(Wx + b)$$

$$x \in \mathbb{R}^d \quad v \in \mathbb{R}^M$$

$$\begin{cases} W \in \mathbb{R}^{M \times d} \\ b \in \mathbb{R}^M \end{cases}$$

→ universal approximation theorem → on compact set

2. Deep NN

→ formulation:
$$\begin{cases} f(x) = v^T f_T(x) \\ f_{t+1}(x) = \phi(W_t f_t(x) + b_t) \quad t=0, \dots, T-1 \end{cases}$$

$T=1$ → degrade to shallow NN

$$f_0(x) = x$$

3. Training NN

$$\hat{\theta} = \underset{\theta \in \Theta}{\operatorname{argmin}} \quad \frac{1}{N} \sum_{i=1}^N \ell(f_{\theta}(x_i), y_i) \rightarrow \text{ERM}$$

→ using GD / SGD / momentum GD

Recap: our deepest goal:

$$\min_{\theta \in \Theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(f_{\theta}(x), y)]$$

4. Exploration on $\begin{cases} \text{width } M \\ \text{depth } T \end{cases}$ of DNN

→ Define $H_{M,T} = \{ \text{NN with } \underline{\text{width } M} \text{ \& \underline{depth } T} \}$
 $= \{ f : f(x) = V^T f_T(x) \}$.

where : $f_{t+1}(x) = \sigma(W_t f_t(x) + b_t)$

$$\begin{cases} W_t \in \mathbb{R}^{M \times M}, & t=1, 2, \dots, T-1 \\ W_0 \in \mathbb{R}^{M \times d} \end{cases}$$

Richness of Hypothesis Space

generally not true

⇒ we want to show

$$\begin{cases} \textcircled{1} \text{ Width} \Rightarrow H_{M,T} \subseteq H_{M+1,T} \\ \textcircled{2} \text{ Depth} \Rightarrow H_{M,T} \subseteq H_{M,T+1} \end{cases}$$

$\textcircled{1} \forall f \in H_{M,T}, \underline{f(x) = V^T f_T(x)} \text{ \& } \underline{f_{t+1}(x) = \sigma(W_t f_t(x) + b_t)}$

construct \tilde{f} as follows :

$$\begin{cases} \text{a) } \tilde{f}_1(x) = \sigma \left(\begin{bmatrix} W_0 \\ 0 \end{bmatrix} x + \begin{bmatrix} b_0 \\ 0 \end{bmatrix} \right) \\ \text{b) } \tilde{f}_{t+1}(x) = \sigma \left(\begin{bmatrix} W_0 & 1 \\ 0 & 0 \end{bmatrix} \tilde{f}_t(x) + \begin{bmatrix} b_t \\ 0 \end{bmatrix} \right) \\ \text{c) } \tilde{f}(x) = \begin{pmatrix} V \\ 0 \end{pmatrix}^T \tilde{f}_T(x) \end{cases}$$

It is obvious that $f(x) \equiv \tilde{f}(x) \in H_{M+1,T}$

[trick] : utilize the dummy dimension (M+1)

$\textcircled{2} \forall f \in H_{M,T}, \underline{f(x) = V^T f_T(x)} \text{ \& } \underline{f_{t+1}(x) = \sigma(W_t f_t(x) + b_t)}$

suppose $\sigma(\cdot)$ is ReLU

construct \hat{f} as follows :

$$\begin{cases} \text{a) } \hat{f}_T(x) = f_T(x) \end{cases}$$

$$b) \hat{f}_{T+1}(x) = \sigma(\tilde{W}_T f_T(x) + \tilde{b}_T)$$

where $\begin{cases} \tilde{W}_T = I_M \in \mathbb{R}^{M \times M} \\ \tilde{b}_T = \text{sufficiently large} \end{cases}$

$$\begin{aligned} \text{then } \hat{f}(x) &= V^T \hat{f}_{T+1}(x) \\ &= V^T f_T(x) + V^T \tilde{b}_T \end{aligned}$$

$$\Rightarrow \underline{H_{M,T} \subseteq H_{M,T+1} + \text{constant}}$$

only holds for ReLU

5. Gradient Calculation for DNN \rightarrow via Backpropagation

$$\rightarrow \text{objective: } \min_{\theta} \frac{1}{N} \sum_{i=1}^N \ell(f_{\theta}(x_i), y_i)$$

$$\theta = \{ \{W_t\}_{t=0}^{T-1}, \{b_t\}_{t=0}^{T-1}, V \}$$

\rightarrow we only consider gradient derivation for one data (x, y)

\Downarrow Omit bias

$$\begin{cases} f(x) = V^T f_T(x) := x_{T+1} = g_T(x_T, W_T) \\ f_{t+1}(x) = \sigma(W_t f_t(x) + b_t) := x_{t+1} = g_t(x_t, W_t) \end{cases}$$

Our interest is: $\nabla_{W_t} \ell(\underbrace{x_{T+1}}_{\text{prediction of } x}, \underbrace{y}_{\text{label of } x})$

$$= \nabla_{W_t} x_{t+1} \cdot \nabla_{x_{t+1}} \cdot \ell(x_{T+1}, y)$$

$$= \underline{\nabla_{W_t} g_t(W_t, x_t)} \cdot \underline{\nabla_{x_{t+1}} \cdot \ell(x_{T+1}, y)}$$

(Here, we denote $p_t = \nabla_{x_t} \ell(x_{T+1}, y)$)

$$= \underbrace{\nabla_{w_t} g_t(w_t, x_t)}_{\text{easy to compute}} \cdot \underbrace{p_{t+1}}_{\text{manage}}$$

$$p_t = \nabla_{x_t} \ell(x_{T+1}, y)$$

Boundary condition

$$p_{T+1} = \nabla_{x_{T+1}} \ell(x_{T+1}, y)$$

(starting from this)

$$= \nabla_{x_t} x_{t+1} \cdot \nabla_{x_{t+1}} \ell(x_{T+1}, y)$$

$$= \underbrace{\nabla_{x_t} g_t(w_t, x_t)}_{\text{easy to compute}} \cdot p_{t+1}$$

BP Algorithm Summary :

Algorithm 1: Backpropagation for FC-DNN

```

1  $x_0 = x \in \mathcal{R}^d$  for  $t = 0, 1, \dots, T$  do
2    $x_{t+1} = g_t(x_t, W_t) = \sigma(W_t^\top x_t)$ ;
3 end
4 Set  $p_{T+1} = \nabla_{x_{T+1}} \ell(x_{T+1}, y)$ ;
5 for  $t = T, T-1, \dots, 1$  do
6    $\nabla_{W_t} \ell(x_{T+1}, y) = p_{t+1}^\top \nabla_{W_t} g_t(x_t, W_t)$ ;
7    $p_t = [\nabla_{x_t} g_t(x_t, W_t)]^\top p_{t+1}$ ;
8 end
9 return  $\{\nabla_{W_t} \ell(x_{T+1}, y) : t = 0, \dots, T\}$ 

```

} forward pass
(function value)

} backward propagation
(gradient)