## Lec 15 Training of Neural Network

Recall:

· How from (+) to (b)

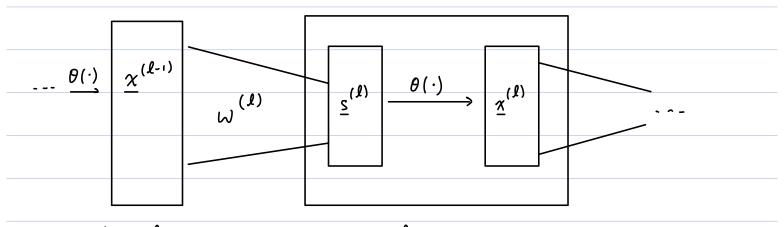
$$S_{j}^{(\ell)} = W_{0,j}^{(\ell-1)} + \sum_{j=1}^{r} W_{i,j} \cdot \chi_{j}^{(\ell-1)}$$

## 1. Vector notation:

$$\underline{S}^{(\ell)} = \begin{bmatrix} S_{\ell}^{(\ell)} \\ \vdots \\ S_{\ell}^{(\ell)} \end{bmatrix} \qquad \theta \left( \underline{S}^{(\ell)} \right) = \begin{bmatrix} \theta \left( S_{\ell}^{(\ell)} \right) \\ \vdots \\ \theta \left( S_{\ell}^{(\ell)} \right) \end{bmatrix} \qquad \underline{\chi}^{(\ell)} = \begin{bmatrix} \chi_{\ell}^{(\ell)} \\ \vdots \\ \chi_{\ell}^{(\ell)} \end{bmatrix} = \begin{bmatrix} \theta \left( \underline{S}^{(\ell)} \right) \end{bmatrix}$$

$$w^{(l)} = \left\{ w_{i,j}^{(l)} \right\} \underset{0 \leq i \leq d}{\circ i \leq d^{(l-1)}} = \left\{ w_{0,1}^{(l)} & w_{0,2}^{(l)} & w_{0,d}^{(l)} \\ \vdots & \vdots & \vdots \\ w_{d^{(l-1)},1}^{(l)} & w_{d^{(l-1)},d^{(l)}} \right\}$$

$$\Rightarrow \underline{S}^{(\ell)} = (\omega^{(\ell)})^{\top} \underline{x}^{(\ell-1)}$$



layer L-1

layer l

2. Forward propagation:

$$\varphi \quad \chi^{(0)} \quad \xrightarrow{\omega^{(1)}} \quad \underline{\leq}^{(1)} \quad \xrightarrow{\theta(\cdot)} \quad \underline{\chi}^{(1)} \quad \xrightarrow{\omega^{(2)}} \quad \underline{\leq}^{(2)} \quad \xrightarrow{\underline{\Rightarrow}} \quad \underline{\leq}^{(2)}$$

$$\cdots \xrightarrow{\chi} (L^{-1}) \xrightarrow{\omega^{(L)}} \xrightarrow{S} (L) \xrightarrow{\vartheta(\cdot)} \chi^{(L)}$$

11 May not be O(·)

1. Input 
$$\underline{\chi}^{(0)} = (1, \chi_1, \chi_2, \dots, \chi_d)$$

for 
$$l = 1, 2, ---$$
 do
$$\underline{S}^{(l)} = (\omega^{(l)})^{T} \underline{\chi}^{(l-1)}$$

$$\underline{\chi}^{(l)} = \begin{bmatrix} 1 \\ \theta(\underline{S}^{(l)}) \end{bmatrix}$$

4 Computation complexity:

. ., ., Edges: 
$$Q = \sum_{l=0}^{l-1} (d^{(l)} + 1) d^{(l+1)}$$

3. Loss function: is let  $\Omega = \{ W^{(1)}, W^{(2)}, --- W^{(L)} \}$  be our model parameters  $en(\Omega)$  be loss function w.r.t example  $(\underline{x}_n, \underline{y}_n)$ 4 e-q. · linear regression for output layer: Square error:  $(x^{(1)} - y_n)^2$ · logistic ? regression for output layer: softmax | log - 1055 : - log P ( yn | xn) 4. <u>SGD</u>: is In iteration K, suppose  $\Omega_{\mathbf{k}}$  is the current set of parameters • Select  $(X_n, y_n)$  in random · compute or en(Dk) i.e.  $\frac{\partial e_{n}(\Sigma_{k})}{\partial \omega_{i,j}^{(l)}}$ ,  $\forall i, j, l$   $\omega_{i,j}^{(l)} \leftarrow \omega_{i,j}^{(l)} - \varepsilon \frac{\partial e_{n}(\Sigma_{k})}{\partial \omega_{i,j}^{(l)}}$ ,  $\forall i, j, l$ 4 How to compute den (2k)? · Idea: use numerical approach — finite difference. · Replace Wij by Wij + SW, while keeping all other weight fixed. Let NK be the resultant set of parameters.  $\frac{\partial e_n(\Sigma_k)}{\partial w_{i,j}^{(l)}} \approx \frac{e_n(w_k^+) - e_n(\Sigma_k)}{\Delta w}$ 

Complexity:
$O(R)$ to compute $en(W_k^t)$ by forward propagation
Repeat for every Wi,j , Q of them
:. Total No. of computation per iteration = $O(0^2)$
Typically $Q \sim 10$ Million $\Rightarrow D(Q^2)$ unacceptable
Ji J