

Lecture - 19

(4) Use CV to evaluate λ

(5) Pick λ^*

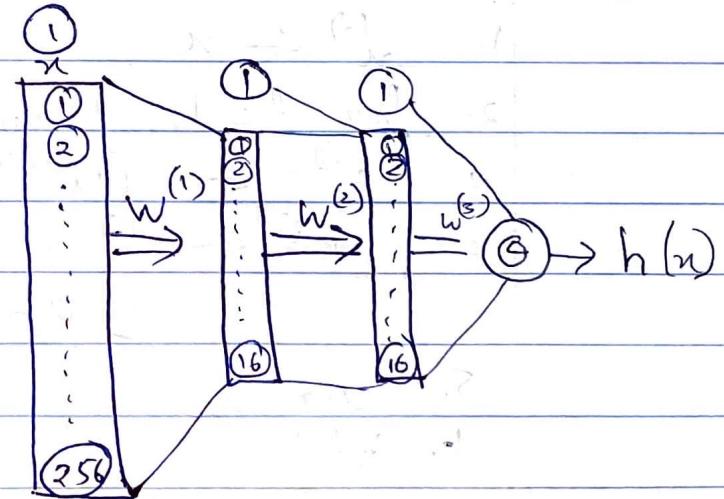
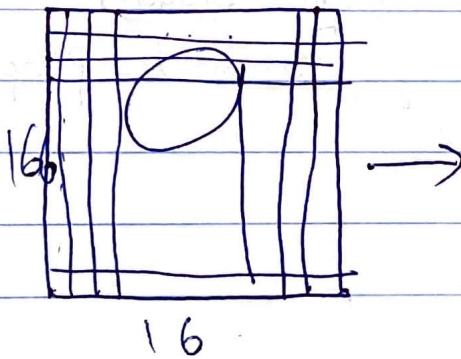
(6) $E_{\text{out}}(g\lambda^*) \approx E_{\text{test}}(g\lambda^*)$

Use 8998 test examples.

Show 99% error bar

$$E_{\text{test}} - \epsilon \leq E_{\text{out}} \leq E_{\text{test}} + \epsilon$$

$$E_{\text{out}} \leq E_{\text{test}} + \sqrt{\frac{1}{2N} \log 2M'}$$



how many nodes?

$$257 + 17 + 17 + 1 \Rightarrow 291$$

How many weights? $\Rightarrow 257 \times 16 + 17 \times 16 + 17 \times 16 \Rightarrow$

$$\text{in } w^{(1)} = 257 \times 16 + 17 \times 16$$

$$\text{in } w^{(2)} = 17 \times 16 + 17 \times 16 \Rightarrow w^{(3)}$$

The linear signal $s^{(l)}$

$s^{(l)}$ is a linear combination (using $W^{(l)}$)
of the outputs of the previous layer $x^{(l-1)}$

$$s^{(l)} = (W^{(l)})^T x^{(l-1)}$$

Forward Propagation : Computing $h(x)$

$$x = x^{(0)} \xrightarrow{w^{(0)}} s^{(1)} \xrightarrow{\theta} n^{(1)} \xrightarrow{w^{(1)}} s^{(2)} \xrightarrow{\theta} n^{(2)} \rightarrow \dots \rightarrow s^{(L)} \xrightarrow{\theta} n^{(L)} = h(x)$$

1. $n^{(0)} \leftarrow n$ [Initialization]
2. for $l = 1 \dots L$ do [Forward propagation]
 3. $s^{(l)} \leftarrow (W^{(l)})^T n^{(l-1)}$
 4. $n^{(l)} \leftarrow \begin{bmatrix} 1 \\ \Theta(s^{(l)}) \end{bmatrix}$
5. end for.
6. $h(x)$ output.

Dimensions of :

$$W^{(l)} = (d^{(l-1)} + 1) \times d^{(l)}$$

$$s^{(l)} = d^{(l)} \times 1$$

$$n^{(l)} = (d^{(l)} + 1) \times 1$$

For classification, regression:

$$E_{in}(h) = E_{in}(w) = \frac{1}{N} \sum_{n=1}^N (h(x_n) - y_n)^2$$

For Logistic regression,

$$E_{in}(h) = \frac{1}{N} \sum_{n=1}^N \ln \left(1 + e^{y_n h(x_n)} \right)$$

Fitting the Data: Minimizing E_{in}

$$E_{in}(h) = E_{in}(w) = \frac{1}{N} \sum_{n=1}^N (h(x_n) - y_n)^2$$

$$\text{where } w = \{w^{(1)}, w^{(2)}, \dots, w^{(L)}\}$$

take gradient descent set it to 0 and solve for variables

$$w(t+1) = w(t) - \eta \nabla E_{in}(w(t))$$

$$w^{(t+1)} \leftarrow w^{(t)} - \eta \frac{\partial E_{in}}{\partial w^{(t)}}$$

Need to compute the gradient of E_{in} at $w(t)$

Gradient of E_{in}

(given weights w),

$$E_{in}(w) = \frac{1}{N} \sum_{n=1}^N e(h(x_n; w), y_n)$$

the gradient of $E_{in}(w)$ is:

$$\frac{\partial E_{in}(w)}{\partial w^{(l)}} = \frac{1}{N} \sum_{n=1}^N \frac{\partial e(h(x_n), y_n)}{\partial w^{(l)}}$$

Try not to use for loops, instead use matrix.

Computing the Gradient:

Numerical, Finite Difference approach.

$$\frac{\partial e(x)}{\partial w_{ij}^{(l)}} \approx \frac{e(x|w_{ij}^{(l)} + \Delta) - e(x|w_{ij}^{(l)} - \Delta)}{2\Delta}$$

Back

$$O((N|w| + N|v|)|w|)$$

↑
just to compute
(e)
tan h]

For every weight we
have to count. So multiply
(w)

Back propagation

$$O((N|w| + N|v|) f(w))$$

A dynamic programming Algorithm.

chain Rule.

Node i only affects the output $h(x)$ through its input $s_i^{(l)}$

$$\frac{\partial e}{\partial w_i^{(l)}} = \frac{\partial s_i^{(l)}}{\partial w_i^{(l)}} \cdot \frac{\partial e}{\partial s_i^{(l)}}$$

how error changes w.r.t. $w_i^{(l)}$

how input signal changes w.r.t. $w_i^{(l)}$

how error changes w.r.t. $s_i^{(l)}$

$$s_i^{(l)} = (w_i^{(l)})^T u^{(l-1)} \Rightarrow \frac{\partial s_i^{(l)}}{\partial w_i^{(l)}} = u^{(l-1)}$$

$$\frac{\partial e}{\partial w_i^{(l)}} = \frac{\partial g_i^{(l)}}{\partial w_i^{(l)}}$$