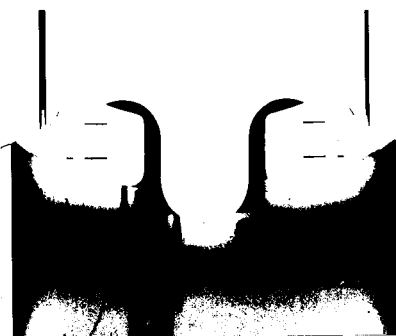


These notes cover the following topics:

1. Brief review of back-propagation in DNN (BP)
2. a) Back-propagation through time (BPTT) in RNN.
b) Vanishing gradient problem in RNN.
3. A short introduction of LSTM.

References:

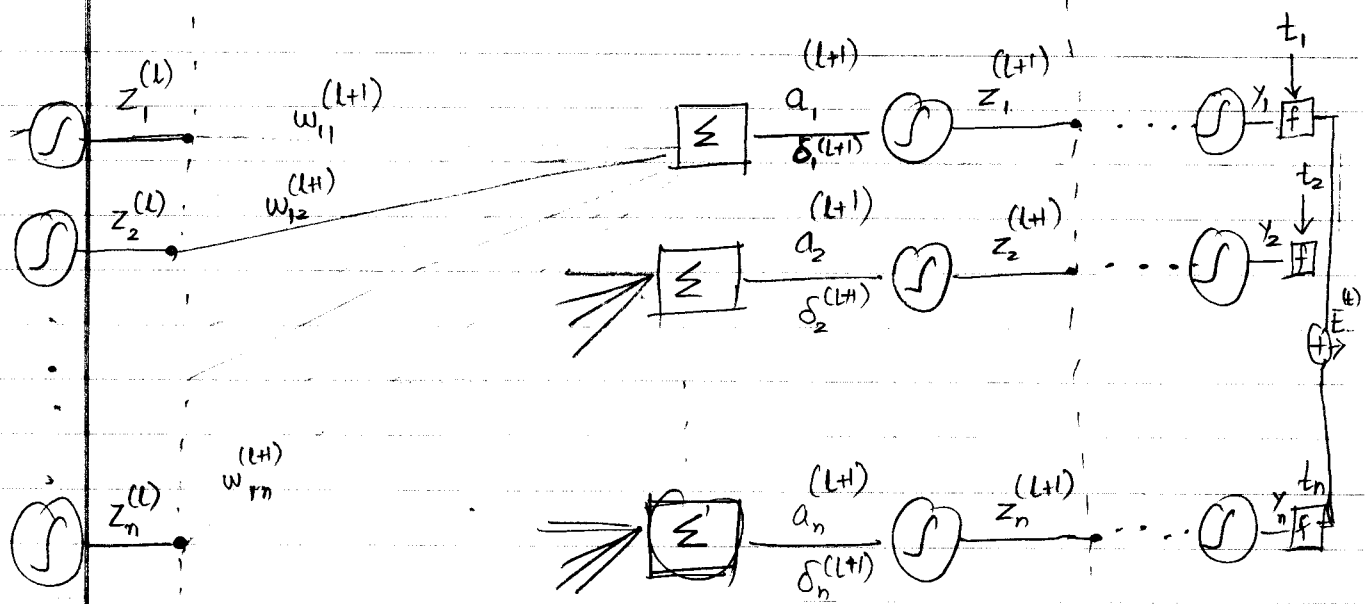
- [1] Bishop, "Neural networks for pattern recognition," 1995, Oxford University Press
- [2] S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory," Neural Computation, vol. 9, no. 8, pp 1735-1780, 1997
- [3] A. Graves, A. R. Mohamed, and G. Hinton, "Speech recognition with deep recurrent neural networks," ICASSP 2013.
- [4] H. Sak, A. Senior, F. Beaufays, "Long short-term memory based recurrent neural network architectures for large vocabulary speech recognition," ArXiv - eprints, Feb 2014.



Bishop's notation
layer

Deep Neural Networks

①



$$a_j^{(l+1)} = \sum_i w_{ji}^{(l+1)} z_i^{(l)} \quad (w_{ji} = j = \text{o/p node}, i = \text{i/p node})$$

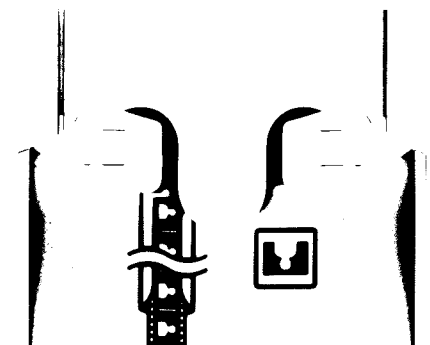
$$z_j^{(l+1)} = g(a_j^{(l+1)}) = g\left(\sum_i w_{ji}^{(l+1)} z_i^{(l)}\right)$$

$$E_{\text{tot}} = \sum_{t=1}^T E^{(t)} \quad (T = \# \text{ training examples})$$

$$\frac{\partial E^{(t)}}{\partial w_{ji}^{(l)}} = \underbrace{\frac{\partial E^{(t)}}{\partial z_j^{(l)}}}_{\delta_j^{(l)}} \underbrace{\frac{\partial z_j^{(l)}}{\partial a_j^{(l)}}}_{z_i^{(l-1)}} \frac{\partial a_j^{(l)}}{\partial w_{ji}^{(l)}}$$

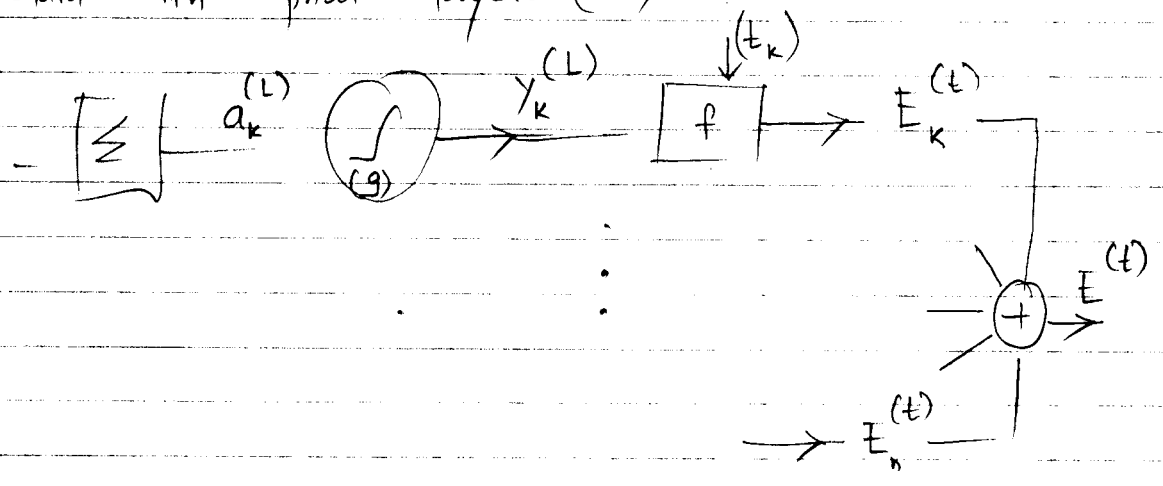
$$\frac{\partial E^{(t)}}{\partial w_{ji}^{(l)}} = \delta_j^{(l)} z_i^{(l-1)} \quad \left[\text{where } \delta_j^{(l)} \triangleq \frac{\partial E^{(t)}}{\partial a_j^{(l)}} \right]$$

$z_i^{(l-1)}$: forward pass calculations, $\delta_j^{(l)}$: backward pass calculations



Calculate δ_j^* :

Case I: Start with final layer (L) : $\therefore l = L$



$$E^{(L)} = \sum_{k=1}^n E_k^{(L)} = \sum_{k=1}^n f(y_k^{(L)}, t_k) \quad (2)$$

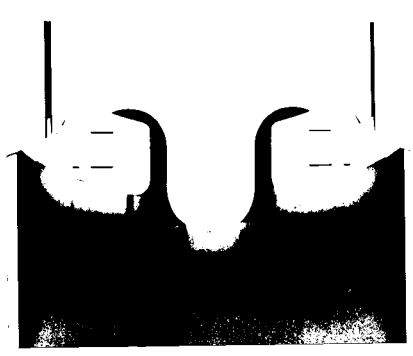
$$\delta_k^{(L)} = \frac{\partial E^{(L)}}{\partial a_k^{(L)}} = \frac{\partial E^{(L)}}{\partial y_k^{(L)}} \frac{\partial y_k^{(L)}}{\partial a_k^{(L)}} \quad (3)$$

$$y_k^{(L)} = g(a_k^{(L)}) \Rightarrow \frac{\partial y_k^{(L)}}{\partial a_k^{(L)}} = g'(a_k^{(L)}) \quad (4)$$

$$\frac{\partial E^{(L)}}{\partial y_k^{(L)}} = \frac{\partial f(y_k^{(L)}, t_k)}{\partial y_k^{(L)}} \quad (5)$$

If $f(y_k^{(L)}, t_k) = \frac{1}{2} (y_k^{(L)} - t_k)^2$

$$\frac{\partial E}{\partial y_k^{(L)}} = (y_k - t_k)$$

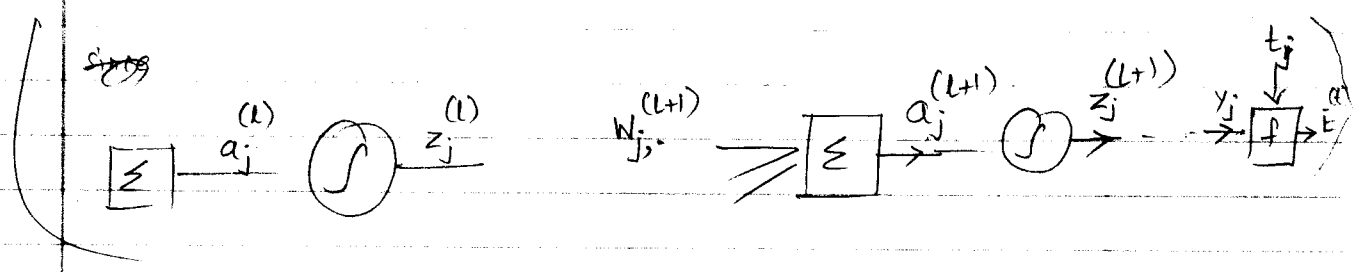


Substituting (4) & (5) in (3)

$$\delta_k^{(L)} = g'(a_k^{(L)}) \frac{\partial E^{(L)}}{\partial y_k^{(L)}} \quad - (6)$$

Case II For the remaining layers (0 < l < L)

$$\delta_j^{(l)} = \frac{\partial E^{(l)}}{\partial a_j^{(l)}} = \sum_{k=1}^n \frac{\partial E^{(l)}}{\partial a_k^{(l+1)}} \frac{\partial a_k^{(l+1)}}{\partial a_j^{(l)}} \quad - (7)$$



$$\frac{\partial E^{(l)}}{\partial a_j^{(l+1)}} = \delta_j^{(l+1)} \quad (\text{by definition}) \quad - (8)$$

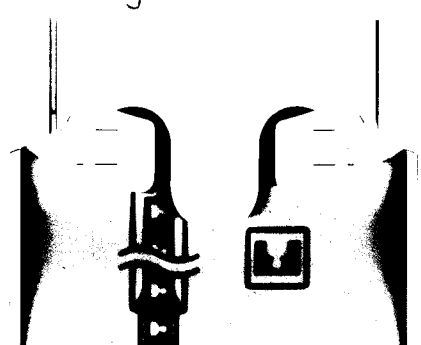
We only need to compute $\frac{\partial a_k^{(l+1)}}{\partial a_j^{(l)}}$

$$a_k^{(l+1)} = \sum_j w_{kj}^{(l+1)} z_j^{(l)}$$

$$= \sum_j w_{kj}^{(l+1)} g(a_j^{(l)})$$

For a particular j ,

$$\frac{\partial a_k^{(l+1)}}{\partial a_j^{(l)}} = w_{kj}^{(l+1)} g'(a_j^{(l)}) \quad (9)$$



(4)

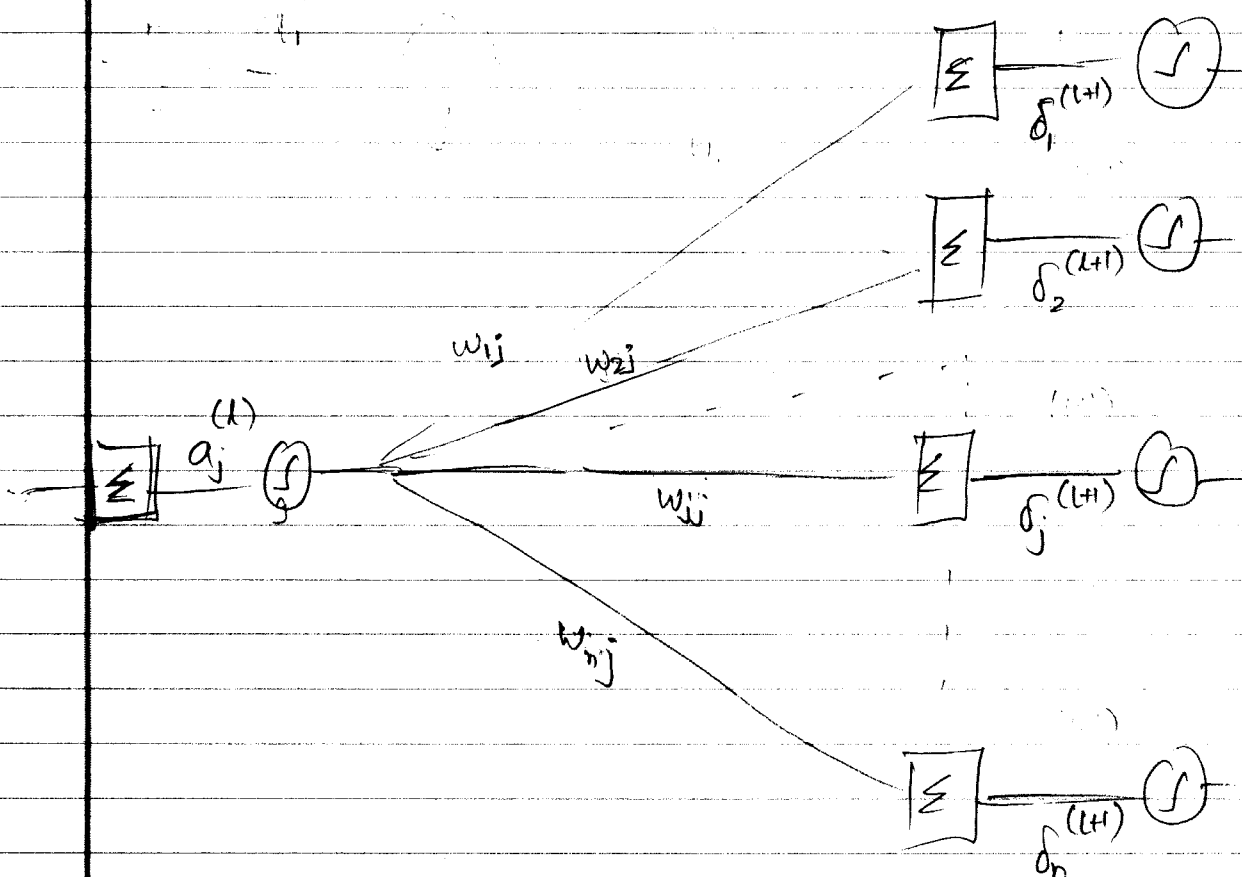
Substitute (8) and (9) in (7)

$$\delta_j^{(l)} = \sum_{k=1}^n \delta_k^{(l+1)} w_{kj}^{(l+1)} g'(a_j^{(l)})$$

$$\delta_j^{(l)} = g'(a_j^{(l)}) \sum_{k=1}^n w_{kj}^{(l+1)} \delta_k^{(l+1)} \quad - (10)$$

(Back-propagation equation)

Back-prop diagram



(5)

From (10) (back prop), we know $\delta_j^{(l)} \forall j, l$,

From forward pass, we know $z_j^{(l-1)}$

$$\frac{\partial E^{(t)}}{\partial w_{ji}^{(l)}} = z_i^{(l-1)} \delta_j^{(l)} \quad (\text{Repeat eqn. (10)})$$

$$\frac{\partial E}{\partial w_{ji}^{(l)}} = \frac{1}{T} \sum_{t=1}^T \frac{\partial E^{(t)}}{\partial w_{ji}^{(l)}} \quad (11)$$

$$w_{ji}^{(l)} \leftarrow w_{ji}^{(l)} - \eta \frac{\partial E}{\partial w_{ji}^{(l)}} \quad (12) \quad (\text{Gradient descent})$$

($\eta > 0$)

$$= w_{ji}^{(l)} - \frac{\eta}{T} \sum_{t=1}^T \frac{\partial E^{(t)}}{\partial w_{ji}^{(l)}}$$

$$= w_{ji}^{(l)} - \frac{\eta}{T} \sum_{t=1}^T \left[z_i^{(l-1)} \delta_j^{(l)} \right]^{(t)}$$

eqn. (1) for a pattern 't'.

Recurrent Neural Networks

(1)

Notations:

DNN (Bishop)

RNN

o/p of node 'i'
last layer:

y_i

$y_i(n)$

$n = \text{at time } n$

i/p of
0th layer

$z_i(0)$

$u_i(n)$

Summed
input

a_i

$s_i(n)$

o/p node 'i'
i/p node 'j'

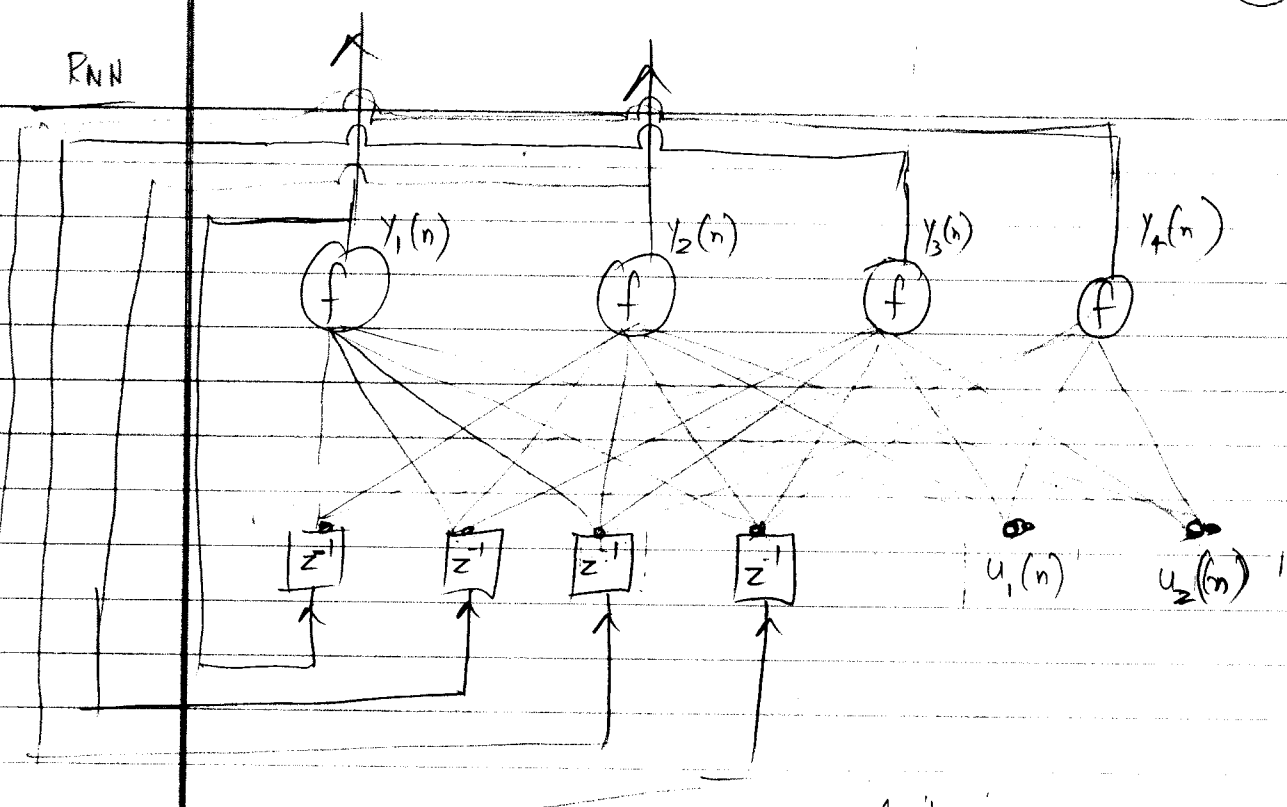
w_{ij}

w_{ij}

Act function

$g(\cdot)$

$f(\cdot)$



$N = 4$ (# neurons in hidden layer)
 $M = 2$ (external i/p s, u_1, u_2)
 $L = 2$ (# of o/p)

Input

$$z_i(n) = \begin{cases} y_i(n-1), & 1 \leq i \leq N \\ u_{i-N}(n), & N+1 \leq i \leq N+M \end{cases} \quad \text{--- (1)}$$

Summed Input

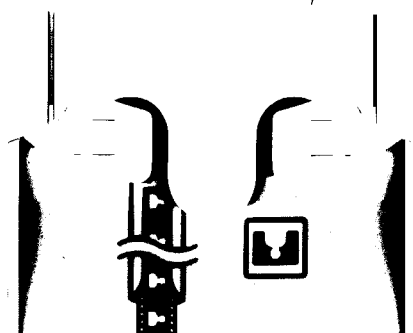
$$s_i(n) = \sum_{j=1}^{N+M} w_{ij} z_j(n) \quad \text{--- (2)}$$

$$= \sum_{j=1}^N w_{ij} y_j(n-1) + \sum_{j=N+1}^{N+M} w_{ij} u_{j-N}(n)$$

Output

$$y_i(n) = f(s_i(n)) \quad \text{--- (3)}$$

(or there could be another non-recurrent type of affine x forms & sigmoid form)



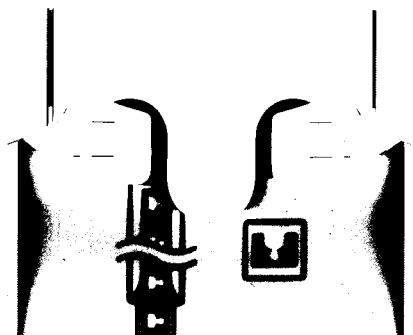
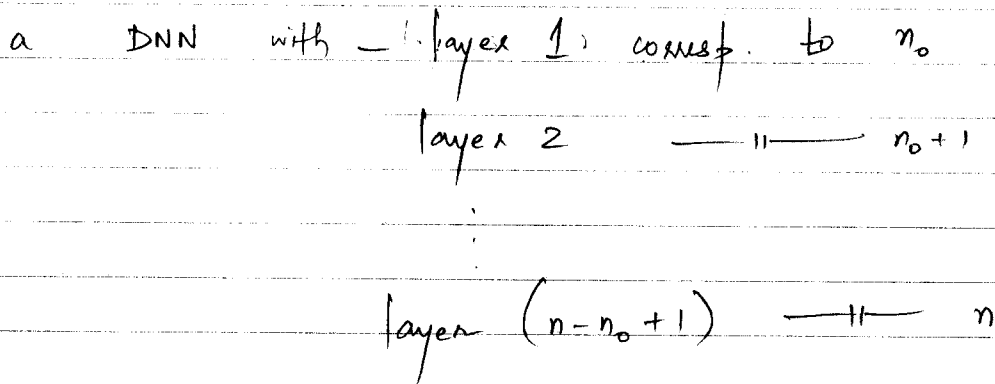
$$E[n_0, n] = \sum_{m=n_0}^n \sum_{k=1}^L e_k^2(m)$$

In DNN, we had $E^{(t)}$ for 1 training pattern which occurs, at time t .

In RNN, 1 training pattern occurs during time $[n_0, n]$.

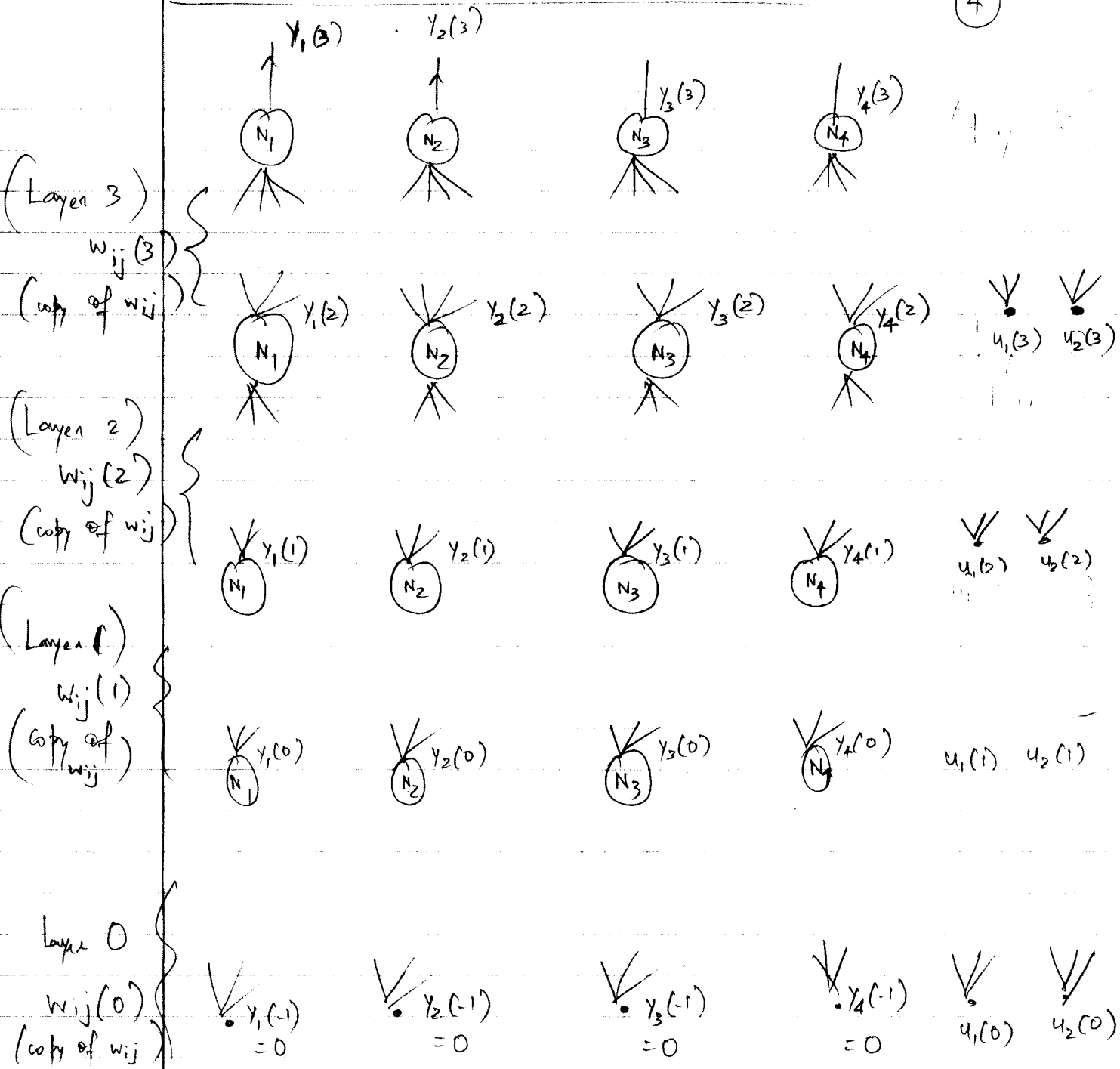
$$E[n_0, n] = \sum_{t=n_0}^n E^{(t)}$$

To calculate the δ^s , we need to unfold the RNN in time. The unfolded RNN becomes



Unfolded RNN in time $[0, 3]$

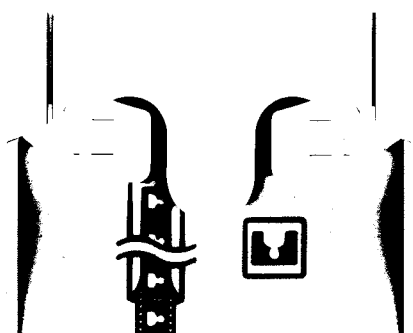
(4)



Input $z_i(n) = \begin{cases} y_i(n-1), & 1 \leq i \leq N \\ u_{i-N}(n), & N+1 \leq i \leq N+M \end{cases} \quad - (4)$

Summed I/P $s_i(n) = \sum_{j=1}^{N+M} w_{ij}(n) z_j(n) \quad - (5)$

Output $y_i(n) = f(s_i(n)) \quad - (6)$
 (on of add another layer affine x-form & sigmoid)



(5)

Q: Why unfolding RNN?

Ans: Because easy to calculate δ 's. Instead of $\delta_j^{(L+1)}$ (in DNN), we have $\delta_j^{(t+1)}$ in unfolded RNN. Then calculating $\delta_j^{(t+1)}$ is easy using BPTT (back-prop through time).

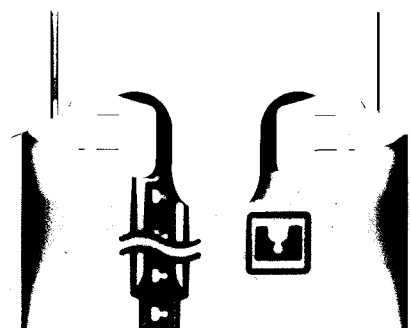
Rewriting the summed i/p equations for RNN & unfolded RNN,

$$\text{(RNN)} \quad s_i(m) = \sum_{j=1}^{N+M} w_{ij} z_j(m) \quad - \quad (3)$$

$$\text{(Unfolded RNN)} \quad s_i(m) = \sum_{j=1}^{N+M} w_{ij}(m) z_j(m) \quad - \quad (5)$$

$$(3) = (5) \quad \text{since} \quad w_{ij}(m) = w_{ij}, \quad m \in [n_0, n_1]$$

During fwd pass, w_{ij} 's are fixed. So, we get the same o/p in RNN & unfolded RNN for a given set of i/p



6

During back-prop, things are different for each layer

$$\frac{\partial E(n_0, n)}{\partial w_{ij}(m)} \neq \frac{\partial E(n_0, n)}{\partial w_{ij}(l)} \quad m \neq l$$

since

$$\frac{\partial E(n_0, n)}{\partial a_i(m)} \neq \frac{\partial E(n_0, n)}{\partial a_i(l)}$$

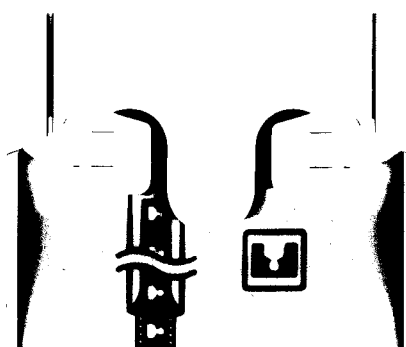
$$\frac{\partial E(n_0, n)}{\partial w_{ij}} = \text{Error gradient of RNN}$$

$$= \sum_{m=n_0}^n \frac{\partial E(n_0, n)}{\partial w_{ij}(m)} \underbrace{\frac{\partial w_{ij}(m)}{\partial w_{ij}}}_{=1} \quad (\text{since } w_{ij}(m) = w_{ij})$$

$$= \sum_{m=n_0}^n \frac{\partial E(n_0, n)}{\partial w_{ij}(m)}$$

$$\frac{\partial E(n_0, n)}{\partial w_{ij}} = \sum_{m=n_0}^n \delta_i^{(m)} z_j^{(m)}, \quad 1 \leq i, j \leq N+M$$

(BPTT equation) (similar to BP but additionally calculated through all layers)



Qualitative
Explanation

(7)

In RNN, $\frac{\partial E(n_0, n)}{\partial w_{ij}}$ = 'change' in $E(n_0)$

" $E(n_0) + \dots + E(n)$ when you perturb w_{ij} .

∴ Need to calculate ΔE for each time step.

In unfolded RNN, when we perturb w_{ij} , we
are perturbing $w_{ij}(m)$, $m = n_0, \dots, n$ by
the same amount. $\frac{\partial w_{ij}(m)}{\partial w_{ij}} = 1$

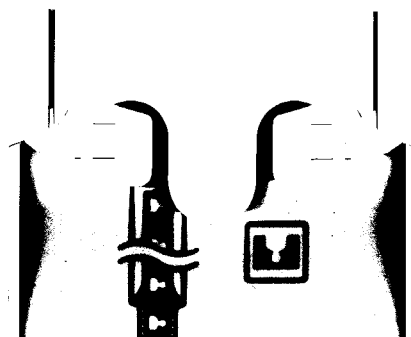
But to calculate change in $E(n)$, we need
to compute $\frac{\partial E(n)}{\partial w_{ij}(n)}$ & $\frac{\partial E(n)}{\partial s_i(n)} = \delta_i(n)$

Similarly, to compute change in $E(n-1)$, we

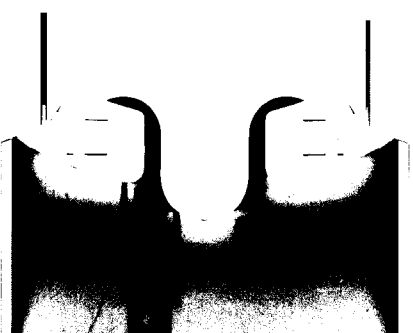
need $\frac{\partial E(n-1)}{\partial w_{ij}(n-1)}$ & $\frac{\partial E(n-1)}{\partial s_i(n-1)} = \delta_i(n-1)$

And, so on.

Finally we add all the $\delta_i(m)$ from
 $m = n_0, \dots, n$.

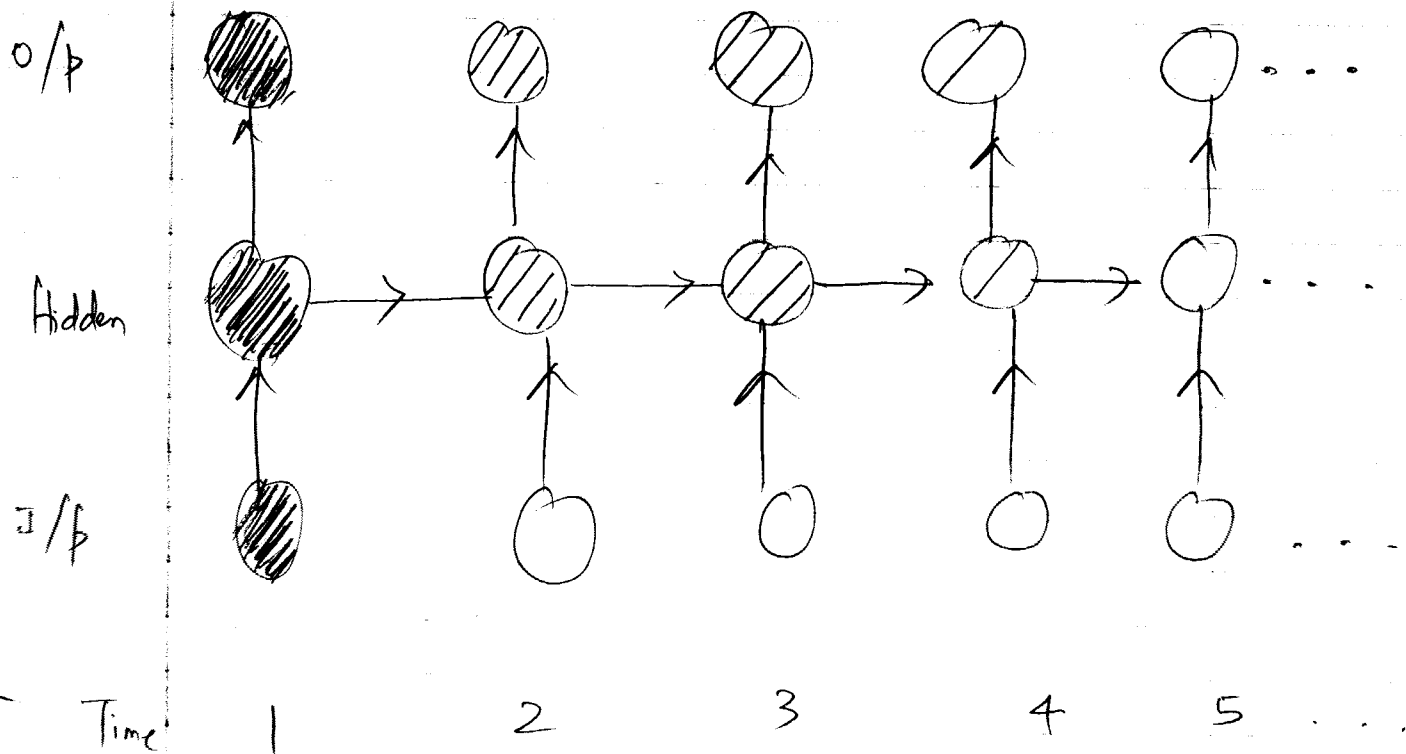
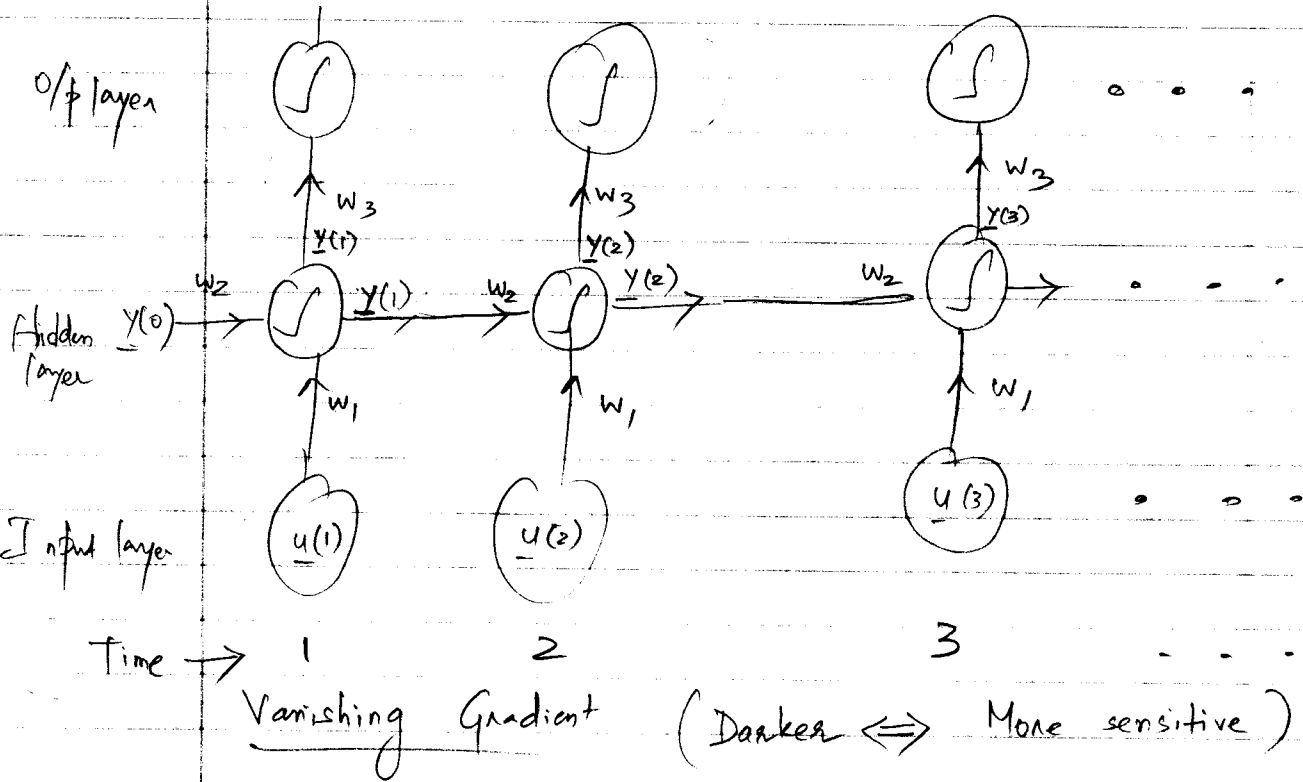


[Blank lined page with a vertical margin line on the left]

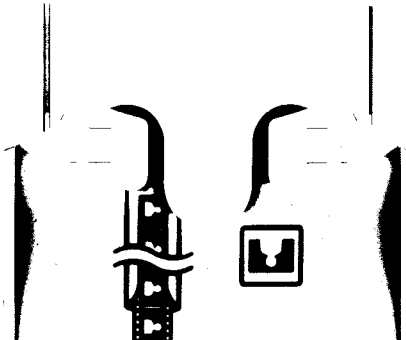


RNN: problem of vanishing (or explosive) gradients 9

Another view of unfolded RNN:



Sensitivity of hidden & o/p layer activations due to i/p at $t=1$
 \downarrow as time \uparrow . N/w forgets i/p at $t=1$.



(10)

Recall, back-prop eqn.,

$$\delta_j(t-1) = f_j'(s_j(t-1)) \sum_i w_{ij} \delta_i(t) \quad - (7)$$

$$\frac{\partial E}{\partial w_{ij}(t-1)} \propto \eta \delta_j(t-1) y_i(t-2)$$

learning rate

∴ Differentiate (7) w.r.t. $\delta_i(t)$, ($t = \text{final time}$)

$$\frac{\partial \delta_j(t-1)}{\partial \delta_i(t)} = f_j'(s_j(t-1)) w_{ij} \quad - (8)$$

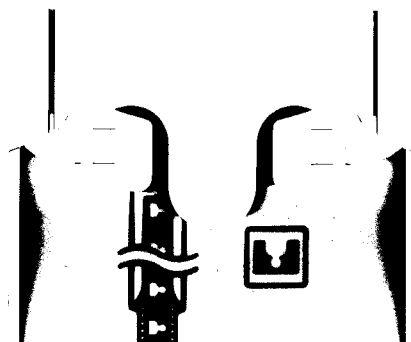
If the back-prop eqn. to compute $\delta_j(t-q)$,

i.e. 'q' time steps behind is,

$$\delta_j(t-q) = f_j'(s_j(t-q)) \sum_i w_{ij} \delta_i(t-q+1) \quad - (9)$$

∴ Differentiate (9) w.r.t. $\delta_i(t)$, ($t = \text{final time}$)

$$\frac{\partial \delta_j(t-q)}{\partial \delta_i(t)} = f_j'(s_j(t-q)) \sum_i w_{ij} \frac{\partial \delta_i(t-q+1)}{\partial \delta_i(t)} \quad - (10)$$



\therefore Gradient of the errors can be computed using the back-prop method (as in (10))

Putting (8) and (10) together, we have (11)

$$\frac{\partial \delta_j(t-q)}{\partial \delta_i(t)} = \begin{cases} f'_j(s_j(t-1)) w_{ij}, & q=1 \\ f'_j(s_j(t-q)) \sum_i w_{ij} \frac{\partial \delta_i(t-q+1)}{\partial \delta_i(t)}, & q > 1 \end{cases}$$

For $q=2$,

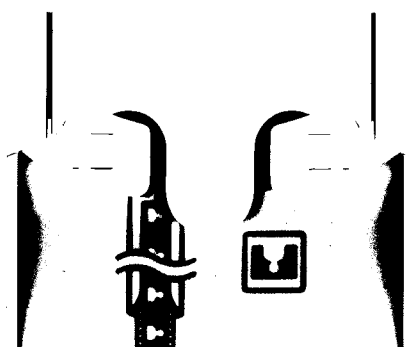
$$\frac{\partial \delta_j(t-2)}{\partial \delta_i(t)} = \sum_k f'_k(s_k(t-1)) w_{ik} f'_j(s_j(t-2)) w_{kj}$$

Replace 'i' by l_0 and 'j' by l_2 , 'k' by l_1

$$\frac{\partial \delta_j(t-2)}{\partial \delta_i(t)} = \sum_{l_1} f'_{l_1}(s_{l_1}(t-1)) w_{l_0 l_1} f'_{l_2}(s_{l_2}(t-2)) w_{l_1 l_2}$$

From this, we can write the back-prop for 'q' time steps behind as,

$$\frac{\partial \delta_j(t-q)}{\partial \delta_i(t)} = \sum_{l_{q-1}}^{\# \text{neurons}} \sum_{l_{q-2}}^{\# \text{neurons}} \dots \sum_{l_1}^{\# \text{neurons}} \prod_{m=1}^q f'_{l_m}(s_{l_m}(t-m)) w_{l_{m-1} l_m} \quad (12)$$



Therefore (12) looks like :

$$\underbrace{(+ \dots +)}_{n \text{ terms}} \underbrace{(+ \dots +)}_{n \text{ terms}} \dots \dots (q-1) \text{ times}$$

There are totally $n \cdot n \cdot \dots = n^{q-1}$ terms
(q-1) times

and each term of those n^{q-1} terms look like this:

$$f'(\dots) w_{ij}$$

$$(f'_j w_{ij} + f'_k w_{ik} + \dots) (f'_{x_1} w_{x_1} + \dots)$$

After x-multiplying, we have

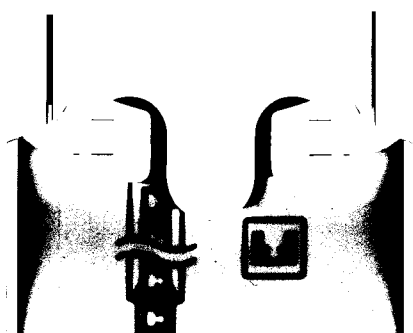
$$(12) = \underbrace{f'_{(1)} w_{(1)}}_1 \underbrace{f'_{(2)} w_{(2)}}_2 \dots \dots \underbrace{f'_{(q-1)} w}_{(q-1)} + \dots$$

Case I

If

$$|f'_{l_m}(s_{l_m}(t-m)) w_{l_{m-1}l_m}| > 1 \quad \forall m$$

and $m^* = \operatorname{argmax}_m |f'_{l_m}(s_{l_m}(t-m)) w_{l_{m-1}l_m}|$



then

$$\left| f'_{L_{m^*}}(s_{L_{m^*}}(t-m^*)) w_{L_{m^*-1}, L_{m^*}} \right| \uparrow \text{exponentially}$$

with q (It has ' $q-2$ ' coefficients each > 1)

This is explosive gradient case.

Case II:

If

$$\left| f'_{L_m}(s_{L_m}(t-m)) w_{L_{m-1}, L_m} \right| < 1 \quad \forall m,$$

then

$$\left| f'_{L_{m^*}}(s_{L_{m^*}}(t-m^*)) w_{L_{m^*-1}, L_{m^*}} \right| \downarrow \text{exponentially}$$

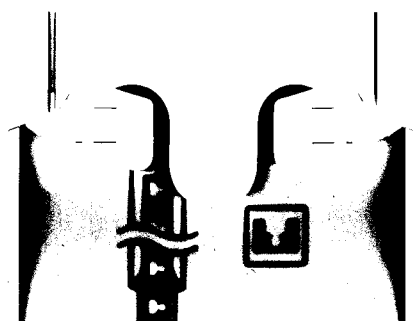
with q (It has ' $q-2$ ' coefficients each < 1)

This is the vanishing gradient case.

Analysis of $f'_c w_c$:

$$\text{Let } g = \left| f'_{L_m}(s_{L_m}(t-m)) w_{L_{m-1}, L_m} \right| = \textcircled{13}$$

If $f = \sigma(\cdot)$, then $f' = f(1-f)$



Therefore, for logistic sigmoid, (13) can be written as

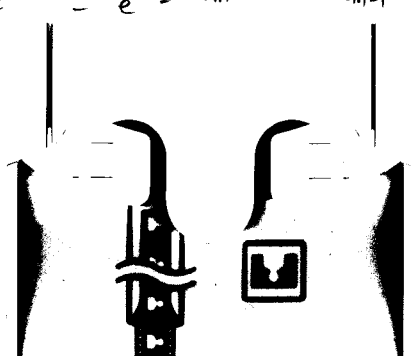
$$g = \left(f_{l_m}(s_{l_m}) - f_{l_m}^2(s_{l_m}) \right) w_{l_{m-1}l_m}$$

We want to find the $w_{l_{m-1}l_m}^*$ for which (13) attains its max.

$$\frac{\partial g}{\partial w_{l_{m-1}l_m}} = \left(f'_{l_m}(s_{l_m}) y_{l_{m-1}} - 2 f_{l_m}(s_{l_m}) f'_{l_m}(s_{l_m}) y_{l_{m-1}} \right) w_{l_{m-1}l_m} + f'_{l_m}(s_{l_m}) = 0$$

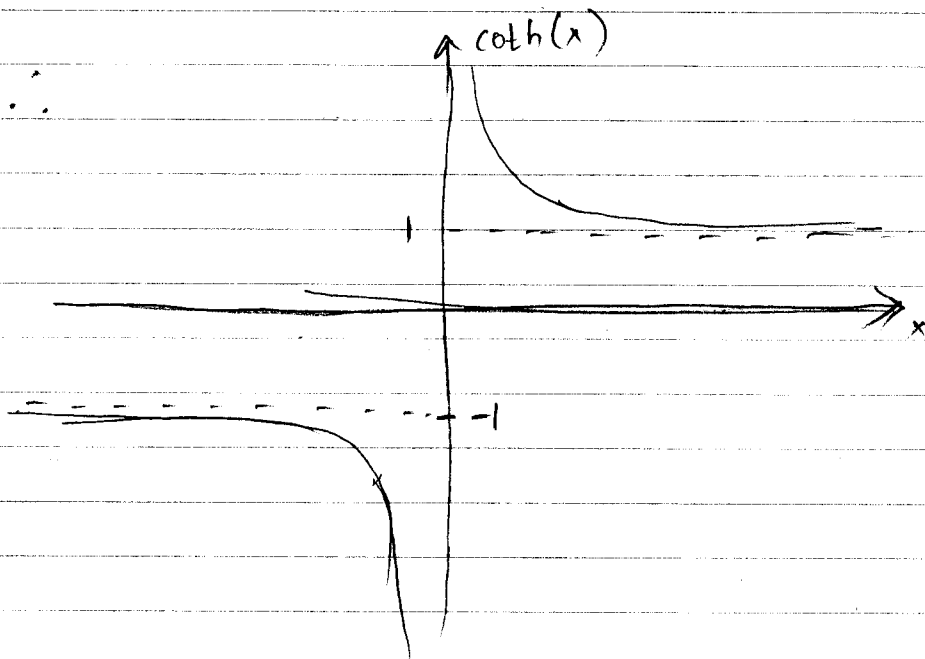
$$\Rightarrow y_{l_{m-1}} \left(1 - 2 f_{l_m}(s_{l_m}) \right) w_{l_{m-1}l_m} + 1 = 0$$

$$\begin{aligned} \Rightarrow w_{l_{m-1}l_m} &= \frac{1}{y_{l_{m-1}} \left(2 f_{l_m}(s_{l_m}) - 1 \right)} \\ &= \frac{1}{y_{l_{m-1}} \frac{1 + e^{-s_{l_m}}}{1 - e^{-s_{l_m}}}} = \frac{1}{y_{l_{m-1}}} \frac{1 + e^{-s_{l_m}}}{1 - e^{-s_{l_m}}} \\ &= \frac{1}{y_{l_{m-1}}} \frac{e^{\frac{1}{2}s_{l_m}} + e^{-\frac{1}{2}s_{l_m}}}{e^{\frac{1}{2}s_{l_m}} - e^{-\frac{1}{2}s_{l_m}}} = \frac{1}{y_{l_{m-1}}} \coth\left(\frac{1}{2}s_{l_m}\right) \end{aligned}$$



(5)

$$w_{l_{m-1} l_m}^* = \frac{1}{\gamma_{l_{m-1}}} \coth\left(\frac{1}{2} s_{l_m}\right)$$



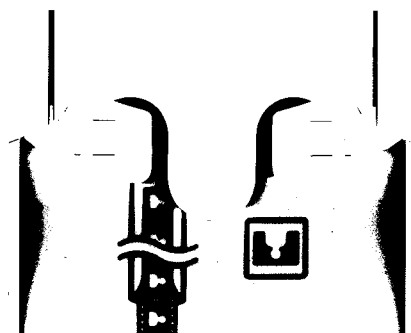
$$\therefore \text{If } |w_{l_{m-1} l_m}| \rightarrow \infty \Leftrightarrow s_{l_m} \rightarrow 0.$$

$$g = f'_{l_m}(s_{l_m}) w_{l_{m-1} l_m}$$

$$= f_{l_m}(s_{l_m}) (1 - f_{l_m}(s_{l_m})) w_{l_{m-1} l_m}$$

$$= \frac{1}{1 + e^{-s_{l_m}}} \frac{e^{-s_{l_m}}}{1 + e^{-s_{l_m}}} w_{l_{m-1} l_m}$$

$$= \frac{1}{(1 + e^{-s_{l_m}})^2} \frac{w_{l_{m-1} l_m}}{e^{s_{l_m}}} \approx \frac{1}{(1+0)^2} \frac{\infty}{\text{faster } \infty} \rightarrow 0$$



(6)

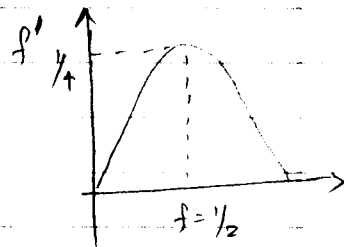
Now, what if $|w_{l_{m-1}, l_m}^x|$ is small?

Note $f'_{l_m} = f(1-f) = h$

$$\frac{\partial h}{\partial f} = 1 - 2f = 0 \Rightarrow f = \frac{1}{2}$$

$$\frac{\partial^2 h}{\partial f^2} = -2 < 0 \therefore f = \frac{1}{2} \text{ at max}$$

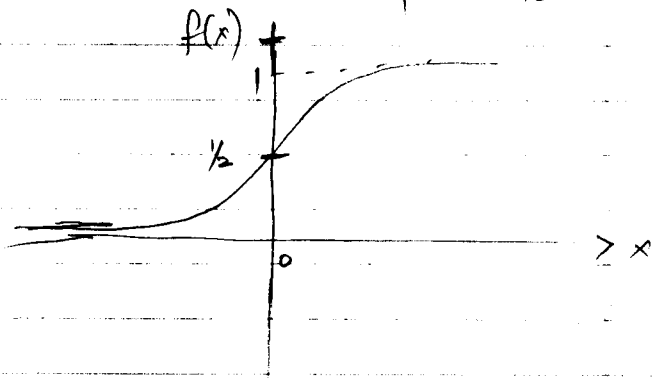
$$\left(f'_{l_m}\right)_{\max} = \frac{1}{2} \left(1 - \frac{1}{2}\right) = \frac{1}{4}$$



$$f = \frac{1}{2} \Rightarrow \frac{1}{2} = \frac{1}{1 + e^{-x}}$$

$$\Rightarrow e^{-x} = 1$$

$$\Rightarrow x = 0$$

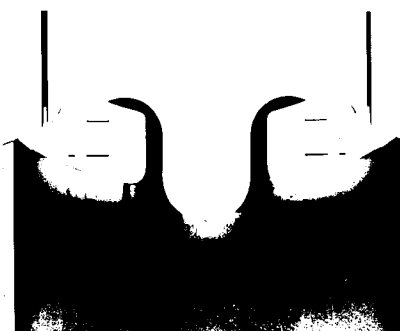


f' achieves steepest slope at $x = \frac{1}{2}$.

The steepest slope is $\left(\frac{f'}{f}\right)_{\max} = \frac{1}{4}$.

$$g = \left| f'_{l_m}(s_{l_m}) w_{l_{m-1}, l_m} \right| < 1$$

$$= \left| \frac{1}{4} w_{l_{m-1}, l_m} \right| < 1 \Rightarrow |w_{l_{m-1}, l_m}| < 4.0$$



(17)

∴ When (a) $|w_{l_{m-1}l_m}| \rightarrow \infty \Rightarrow g \rightarrow 0$

(b) $|w_{l_{m-1}l_m}| < 4 \Rightarrow g \rightarrow \text{value less than 1}$

Very Large weights, error vanishes

Small weights, error vanishes

To avoid this problem, we need

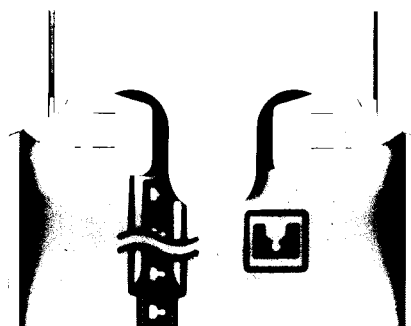
$$f'_{l_m}(s_{l_m}) w_{l_{m-1}l_m} = 1 \quad \left(\begin{array}{l} \text{assume we are interested} \\ \text{in the weight connecting} \\ \text{'l_{m-1}' to 'l_m'} \end{array} \right)$$

$$\Rightarrow \int f'_{l_m}(s_{l_m}) w_{l_{m-1}l_m} d(s_{l_m}) = \int d(s_{l_m})$$

$$\Rightarrow f_{l_m}(s_{l_m}) w_{l_{m-1}l_m} = s_{l_m}$$

$$\Rightarrow f_{l_m}(s_{l_m}) = \frac{s_{l_m}}{w_{l_{m-1}l_m}}$$

∴ $f_{l_m}(s_{l_m}) \propto s_{l_m}$ ∴ $f(\cdot)$ has to be linear



(18)

In addition, unit j 's activation has to be constant :

$$\begin{aligned} y_j(t+1) &= f_j(s_j(t+1)) = f_j(w_{jj} y_j(t)) \\ &= f_j\left(w_{jj} \frac{s_j(t)}{w_{jj}}\right) \\ &= y_j(t) \end{aligned}$$

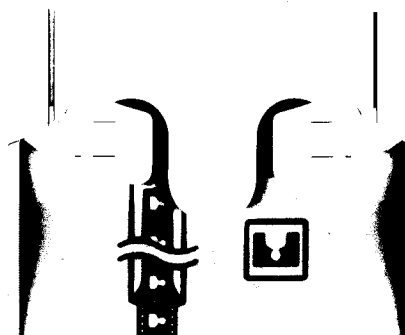
$$\therefore y_j(t+1) = y_j(t)$$

\therefore Constant Error Carousel we need

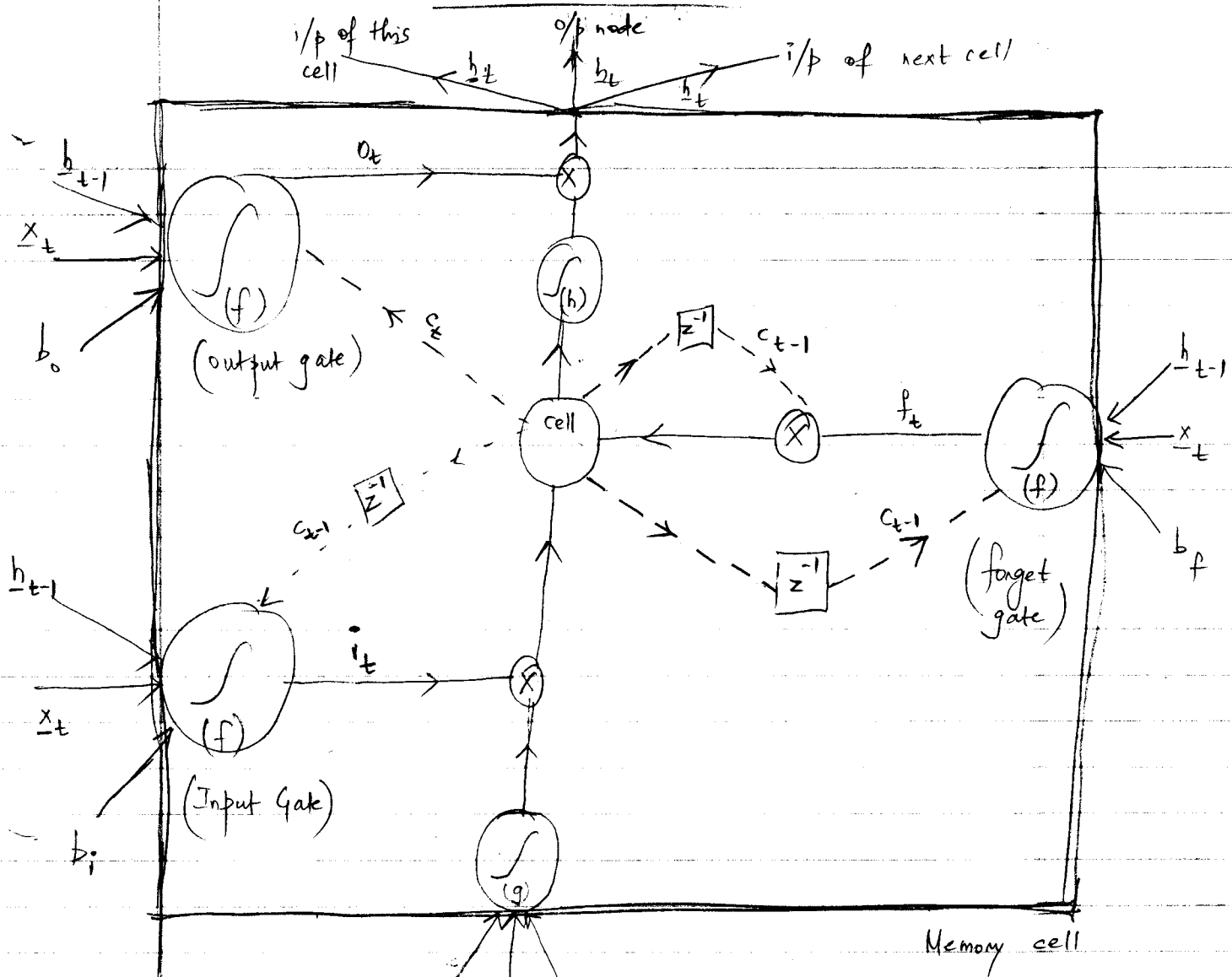
$$(a) \quad y_j(t) = f_j(s_j(t)) \quad \& \quad s_j(t) \quad (\text{linear})$$

$$(b) \quad y_j(t+1) = y_j(t) \quad (\text{constant})$$

This is implemented in LSTM cell.



LSTM block



Gates: i/p, o/p, forget ($f(\cdot)$)
 Cell i/p activation: $g(\cdot)$
 Cell o/p activation: $h(\cdot)$

$$i_t = \sigma(W_{ix} \underline{x}_t + W_{ih} \underline{h}_{t-1} + W_{ic} c_{t-1} + b_i)$$

$$f_t = \sigma(W_{fx} \underline{x}_t + W_{fh} \underline{h}_{t-1} + W_{fc} c_{t-1} + b_f)$$

$$c_t = f_t c_{t-1} + i_t g(W_{cx} \underline{x}_t + W_{ch} \underline{h}_{t-1} + b_c)$$

$$o_t = \sigma(W_{ox} \underline{x}_t + W_{oh} \underline{h}_{t-1} + W_{oc} c_t + b_o)$$

$$\underline{h}_t = o_t \cdot h(c_t)$$

$$\underline{y}_t = W_{yh} \underline{h}_t + b_y$$

