

3) LSTM

13 December 2025 02:37 PM

for GRU
z gate. update gate and reset gate

$$y_t = w_y \cdot x_t + w_{yh} \cdot h_{t-1} + b_y$$

$$z_t = w_{zx} \cdot x_t + w_{zh} \cdot h_{t-1} + b_z$$

candidate memory $\rightarrow \tilde{h}_t = \tanh(w \cdot x_t + (h_{t-1} \oplus y_t) + b)$

$$\boxed{\tilde{h}_t = \tilde{h}_t z_t + h_{t-1} (1-z_t)} \quad \rightarrow \text{reset}$$

\downarrow
next state memory.

$$h_t = z_t \tilde{h}_t + (1-z_t) h_{t-1}$$

$$h_t = z_t \tilde{h}_t + h_{t-1} - z_t h_{t-1}$$

$$h_t = z_t (\tilde{h}_t - h_{t-1}) + h_{t-1}$$

$$h_t = z_t (\tilde{h}_t - h_{t-1}) + h_{t-1}$$

$$h_t = 0.001 (\tilde{h}_t - h_{t-1}) + h_{t-1}$$

$$h_t = 0.001 \tilde{h}_t + 0.999 h_{t-1}$$

$$h_t \approx 0.999 h_{t-1} \rightarrow \text{only } 0.001 \text{ of the memory is retained.}$$

forget gate.

$$f_t = \sigma(w_f [x_t, h_{t-1}] + b_f)$$

$$i_t = \sigma(w_i [x_t, h_{t-1}] + b_i)$$

$f_t \rightarrow$ forget of current input.

$$c_t = c_{t-1}$$

$$c_t = c_{t-1} + \Delta t$$

$c_t = f_t \odot c_{t-1} + i_t \odot g_t$

\uparrow forget gate
 \uparrow input gate
 \uparrow candidate content.

$$\tilde{h}_t = \tanh(w_x \cdot x_t + w_h (h_{t-1} \oplus y_t) + b_h)$$

equations for
LSTM.

$$\left\{ \begin{array}{l} f_t = \sigma(w_{fx} \cdot x_t + w_{fh} \cdot h_{t-1} + b_f) \\ i_t = \sigma(w_{ix} \cdot x_t + w_{ih} \cdot h_{t-1} + b_i) \end{array} \right\}$$

$$h_t = z_t \tilde{h}_t + (1-z_t) h_{t-1}$$

$$\begin{aligned}
 & \left(\begin{array}{l} i_t = \sigma(w_{ix} \cdot x_t + w_{ih} \cdot h_{t-1} + b_i) \\ g_t = \tanh(w_{gx} \cdot x_t + w_{gh} \cdot h_{t-1} + b_g) \end{array} \right) \\
 & c_t = f_t \odot c_{t-1} + i_t \odot g_t \quad \xrightarrow{\text{express it as long term memory.}} \\
 & d_t = \sigma(w_{dx} \cdot x_t + w_{dh} \cdot h_{t-1} + b_d) \\
 & h_t = d_t \odot \tanh(c_t) \\
 & \text{Return long term Memory, short term Memory}
 \end{aligned}$$

$$c_t = c_{t-1} + \Delta c$$

$$\begin{aligned}
 & f_t = \sigma(w_{fx} \cdot x_t + w_{fh} \cdot h_{t-1} + b_f) \\
 & i_t = \sigma(w_{ix} \cdot x_t + w_{ih} \cdot h_{t-1} + b_i) \\
 & g_t = \tanh(w_{gx} \cdot x_t + w_{gh} \cdot h_{t-1} + b_g) \\
 & c_t = f_t \odot c_{t-1} + i_t g_t \quad \xrightarrow{\text{long term memory.}}
 \end{aligned}$$

forget gate
 input gate
 candidate gate.
 output gate.

$$\begin{aligned}
 & o_t = \sigma(w_{ox} \cdot x_t + w_{oh} \cdot h_{t-1} + b_o) \\
 & h_t = o_t \odot \tanh(c_t) \quad \xrightarrow{\text{short term memory.}}
 \end{aligned}$$

GRU

$$r_t = \sigma(x_t \cdot w_{rx} + h_{t-1} \cdot w_{rh} + b_r) \rightarrow \text{Reset Gate}$$

$$z_t = \sigma(x_t \cdot w_{zx} + h_{t-1} \cdot w_{zh} + b_z) \rightarrow \text{Update Gate}$$

h_{t-1} = prev time step's hidden state \rightarrow compressed memory.

$$r_{h,t} = \frac{\sinh}{\cosh}(x_t \cdot w_x + (h_{t-1} * r_t) \cdot w_h + b)$$

$$h_t \leftarrow h_{t-1} + (1 - z_t) h_{t-1} \rightarrow \text{what happened in this time step}$$

new information \rightarrow information to ?
erase
to write

so, Updated memory = $h_t = z_t * h_{t-1} + (1 - z_t) h_{t-1}$

$$\text{cp problem, } h_t = z_t (h - h_{t-1}) + h_{t-1}$$

$$\Rightarrow h_t = z_t (h - h_{t-1}) + h_{t-1}$$

for $z_t = 0.00001 \rightarrow$ do not update old memory, keep new memory
as old memory, $h_t = 0.9999 h_{t-1}$

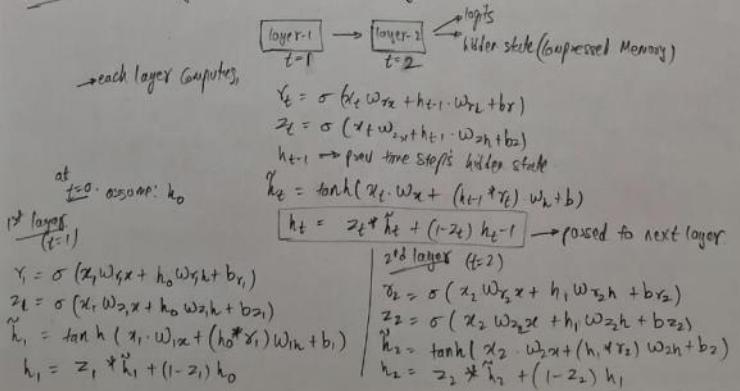
$$\text{for } k^{\text{th}} \text{ time step, } h_t = (0.9999)^{t-k} h_{t-1} \rightarrow \text{only 99.99% memory is retained.}$$

Memory decay \Rightarrow Memory loss.

clipped:
 -> No memory decay \rightarrow change the design so that no
 -> 100% old memory \rightarrow Memory loss.
 retention must be possible \Rightarrow freeze the memory.
 -> Gradient won't vanish while BPTT

why gradient vanishes or explodes?

for Example, Consider 2-layered GAN



Trainable Parameters:

Layer 1: $\rightarrow w_{1x}, w_{1h}, b_1$ $\rightarrow w_{2x}, w_{2h}, b_2$ $\rightarrow w_x, w_h, b$	Layer 2: $\rightarrow w_{0x}, w_{0h}, b_0$ $\rightarrow w_{1x}, w_{1h}, b_1$ $\rightarrow w_{2x}, w_{2h}, b_2$	<u>18 trainable parameters</u> $L \rightarrow \text{loss}$ $\left\{ \frac{\partial L}{\partial W} \right\} \text{ for grad of loss w.r.t. trainable parameters.}$
---	---	---

#BPTT: Let's see $\frac{\partial L}{\partial w_{2x}}$ use chain rule.

Weight of next gate for layer 2 projecting from RERE

$$\frac{\partial L}{\partial w_{2x}} = \frac{\partial L}{\partial z_2} \cdot \frac{\partial z_2}{\partial w_{2x}} = \prod \left(\frac{\partial L}{\partial z_i} \cdot \frac{\partial z_i}{\partial w_{2x}} \right)$$

We know that, ignore bias for now.

$$y = w_{2x}z_2 + b_2$$

$$y = w_{2x}(z_2w_{1x} + h_1) + b_2$$

$$y = w_{2x}[z_2w_{1x} + (1 - z_2)(z_1w_{1h} + h_0) + b_2]$$

$$\text{and } \frac{\partial L}{\partial z_1} = \text{grad}_h(z_1, w_{1x}x + (h_0 * z_1)w_{1h} + b_1)$$

$$\text{again } z_1 = \sigma(w_{1x}x + h_0w_{1h} + b_0)$$

So eventually $y = f(w_{2x}x)$

so, $\frac{\partial L}{\partial w_{2x}} = \frac{\partial y}{\partial z_2} \times \frac{\partial z_2}{\partial h_1} \times \frac{\partial h_1}{\partial z_1} \times \frac{\partial z_1}{\partial w_{1x}}$ Multiplicative Decay.

anyways, $\frac{\partial L}{\partial w} = \prod \frac{\partial L}{\partial z_i} \frac{\partial z_i}{\partial w_{2x}}$ {Something > 1 → gradient explode. Something < 1 → gradient vanish.}

This multiplication is the reason of Gradient exploding or Gradient vanishing.

Note: This won't happen because then, the model's weights and biases will not get updated, we can't afford this.

#Change-2:

- 1) In new design Gradient vanishing/explode doesn't happen
- 2) This can be achieved by introducing linear function's memory update so that if gradient passes then it will be unaffected.
 \Rightarrow Gradient free memory update, $f(a) = a \Rightarrow [f'(a) = 1]$
 Linear gradient

#LSTMs

By implementing those changes discussed in GRU's Architecture, we can implement new architecture such that,

1) No memory decay will happen.

2) 100% old memory retention.

3) Gradient noise Varify/Explode

In GRU, we had one gate that handled both the renewal and Updation of old memory and new memory respectively, we now need separate gates for that,

so,

$$h_t = z_t * h_t + (1 - z_t) * h_{t-1}$$

\Rightarrow no Decay $\Rightarrow h_t = h_{t-1}$ (must be possible)

$\Rightarrow h_t = h_{t-1} + \Delta t \rightarrow$ new information.

\Rightarrow must be linear (No Gradient decay/explode)

$\Rightarrow h_t = f(h_{t-1} + \Delta t)$; where $f_t \rightarrow$ linear funcn.

for control over memory introduce gates,

$$c_t = f_t * c_{t-1} + i_t * g_t$$

and express,

$\bullet o_t \rightarrow$ output gate $\bullet i_t \rightarrow$ input gate $\bullet f_t \rightarrow$ Candidate gate $\bullet c_{t-1} \rightarrow$ previous long term memory $\bullet f_t \rightarrow$ forget gate

$$\text{Updated memory} \leftarrow o_t * \tanh(c_t)$$

- $\bullet f_t = \sigma(x_t \cdot W_{fx} + h_{t-1} \cdot W_{fh} + b_f)$ \rightarrow how much to forget from c_{t-1}
- $\bullet i_t = \sigma(x_t \cdot W_{xi} + h_{t-1} \cdot W_{hi} + b_i)$ \rightarrow How strongly we need it to update old memory
- $\bullet g_t = \tanh(x_t \cdot W_{gx} + h_{t-1} \cdot W_{gh} + b_g)$ \rightarrow what to write in old memory, as a new information
- $\bullet c_t = f_t * c_{t-1} + i_t * g_t$ \rightarrow new information $\left\{ \begin{array}{l} \text{Long Term Memory} \\ \text{Short Term Memory} \end{array} \right.$
- $\bullet o_t = \sigma(x_t \cdot W_{xo} + h_{t-1} \cdot W_{ho} + b_o)$
- $\bullet h_t = o_t * \tanh(c_t) \rightarrow$ "Short Term Memory".

$y = W \cdot h_t + b$.

Note: It is necessary to keep passing short term memory as well as long term memory to next time step in Training as well as Inference mode.

$c_t \rightarrow$ What happened before 100th time step | $h_t \rightarrow$ what happened before 2nd time step

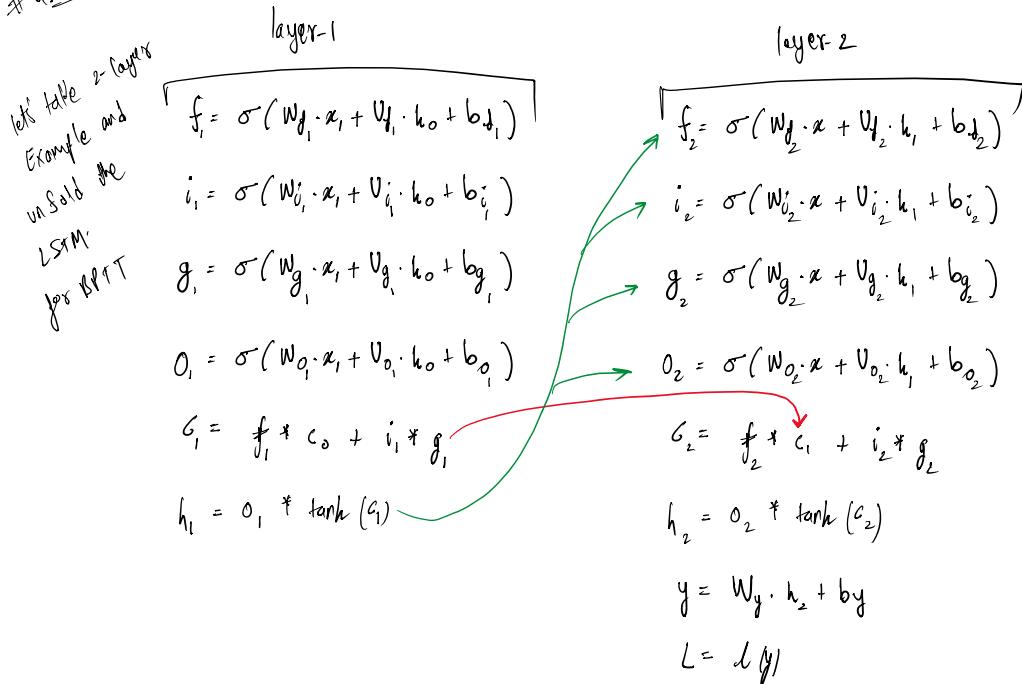
\rightarrow GRU_L, GRU_F, GRU_Y

\rightarrow Astro_L, Astro_F, Astro_Y

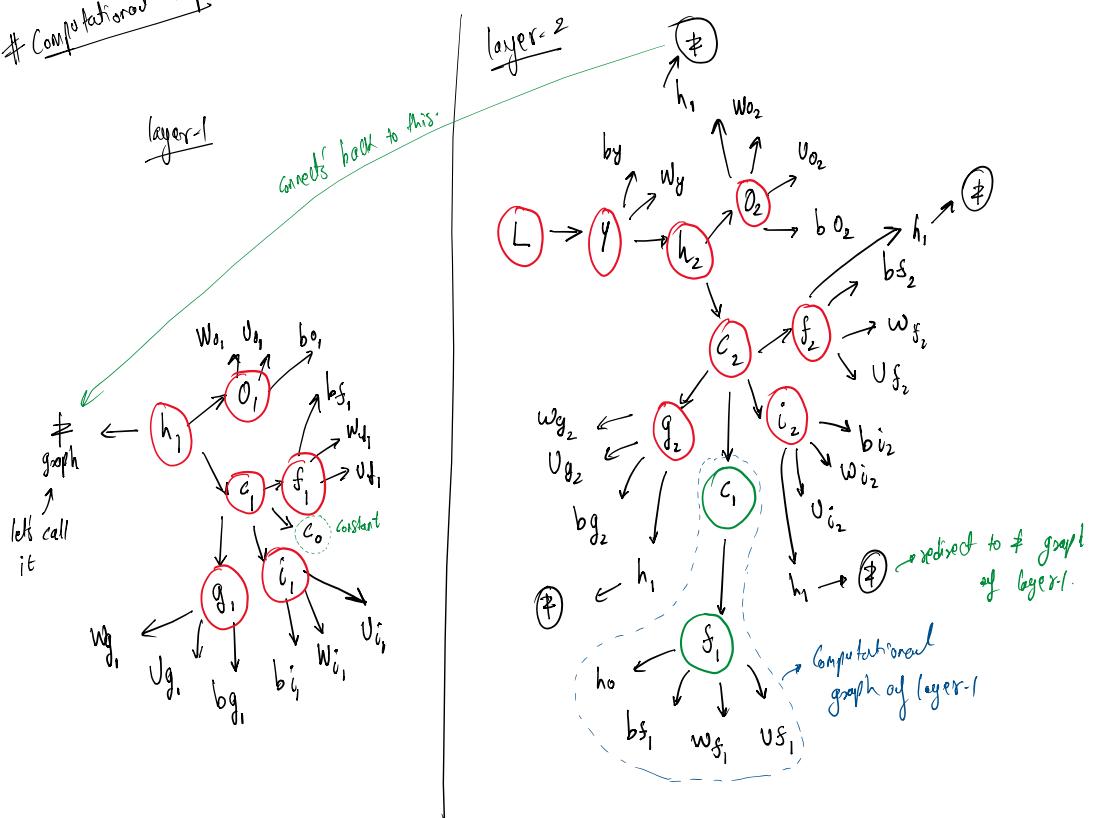
\rightarrow Scribe_L, Scribe_F, Scribe_Y

\rightarrow Lenovithia

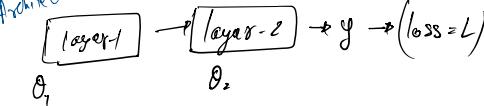
#Gradient Analysis



Computational Graph



Architecture:



$$\left. \begin{array}{l} \theta_1 : \{w_{s1}, u_{s1}, b_{s1}, w_{g1}, u_{g1}, b_{g1}, w_{o1}, b_{o1}, u_{o1}, w_i, b_i, u_i\} \\ \theta_2 : \{w_{s2}, u_{s2}, b_{s2}, w_{g2}, u_{g2}, b_{g2}, w_{o2}, b_{o2}, u_{o2}, w_i, b_i, u_i\} \end{array} \right\}$$

to optimize w_{s1}

$$\frac{\partial L}{\partial w_{s1}} = \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial h_2} \cdot \frac{\partial h_2}{\partial c_2} \cdot \frac{\partial c_2}{\partial f_2} \cdot \frac{\partial f_2}{\partial g_2} \cdot \frac{\partial g_2}{\partial i_2} \cdot \frac{\partial i_2}{\partial s_1} \cdot \frac{\partial s_1}{\partial w_{s1}}$$

$$\text{but, } c_2 = f_2 * g_2 + i_2 * g_2 \quad \text{for we get this by Design.}$$

$$\frac{\partial c_2}{\partial f_2} = f_2 \Rightarrow \frac{\partial c_2}{\partial g_2} \times f_2$$

Now $f_2 \rightarrow \text{Sigmoid function} \rightarrow [0 < f_2 < 1] \rightarrow \text{Stabilized gradient } \frac{\partial}{\partial w}$.

so model can now learn to stabilize this gradient flow.

Generally, # Constant Error Carousel

$$\frac{\partial C_t}{\partial C_{t-K}} = \prod_{j=t-K+1}^t f_j ; \quad 0 < f_j < 1$$

further analysis

$$C_t = f_t * C_{t-1} + i_t * g_t$$

$$\Rightarrow \frac{\partial C_t}{\partial C_{t-1}} = f_t$$

$$\frac{\partial C_t}{\partial C_{t-2}} = f_t * f_{t-1}$$

$$\frac{\partial C_t}{\partial C_{t-1}} = f_t * f_{t-1} * C_{t-2} + i_t * g_t$$

$$\frac{\partial C_t}{\partial C_{t-2}} = f_t * f_{t-1} * f_{t-2}$$

$$\frac{\partial C_t}{\partial C_{t-1}} = f_t * f_{t-1} * f_{t-2} * C_{t-3} + i_t * g_t$$

$$\frac{\partial C_t}{\partial C_{t-3}} = f_t * f_{t-1} * f_{t-2} * f_{t-3}$$

$$\vdots$$

$$\frac{\partial C_t}{\partial C_{t-K}} = \prod_{j=t-K+1}^t f_j$$

error at time t is sent backward.

$\overbrace{C_{t-1}}^{\text{old}}$ $\overbrace{C_t}^{\text{new}}$ $\overbrace{\delta C_{t-k}}^{\text{old}} \quad j = t-k+1$

if $f_t = 0.999$
 sequence length > 100 } $\frac{\delta C_t}{\delta C_{t-k}} = (0.999)^{100} \approx 0.9047 \times 1$

No Gradient explosion
 but Gradient will eventually shrink
 and it's delayed.
 if any forget gated output is small
 then it will effect and vanishes the gradient.

Compression of Information in LSTMs

$$C_t = f_t * C_{t-1} + i_t * g_t$$

↗ How much old information to keep
 ↘ How much new info to add.

Imagine: $C_{t-1} = \begin{bmatrix} 0.9 & 0.586 \\ -0.43 & 0.12 \end{bmatrix}$ → old memory

we need to update
this old memory
by adding new memory

$$C_t = f_t * \begin{bmatrix} 0.9 & 0.58 \\ -0.43 & 0.12 \end{bmatrix} + i_t * g_t$$

↗ then add new information
 ↘ Compress old memory

so basically, at every layer of LSTM.

new Memory = (Compress old memory) + new information

Issue 2:- Each time step sees the compressed version of earlier time step. } $\Rightarrow C_{t+k}$ is the compressed version of C_{t+k-1}

\Rightarrow Every information is compressed.

Temporal Locality

Decision about memory are made when the information arrives, not when it is needed later.

LSTM must be able to decide the importance of current information before seeing the future context

model must decide now, whether some information will be useful far in future

$$\begin{aligned} f_t &= \sigma(W_f x_t + U_f h_{t-1} + b_f) \\ i_t &= \sigma(W_i x_t + U_i h_{t-1} + b_i) \\ g_t &= \tanh(W_g x_t + U_g h_{t-1} + b_g) \\ c_t &= f_t \odot c_{t-1} + i_t \odot g_t \\ o_t &= \sigma(W_o x_t + U_o h_{t-1} + b_o) \\ h_t &= o_t \odot \tanh(c_t) \end{aligned}$$

This is where the model decides whether a piece of information is useful or not, but it is only able to see current input



There is no later revisions.

$$c_0 \rightarrow \text{initial memory}$$

$$\dots \quad c_1 \quad c_2 \quad c_3 \quad c_4 \quad c_5 \quad c_6 \quad c_7 \quad c_8 \quad c_9 \quad c_{10}$$

$c_0 \rightarrow$ initial memory:

The diagram shows a red arrow originating from c_0 and pointing to a black rectangular box containing the text "There is no later Revision.".

at $t=2$ $c_2 = (c_1 * s_1 + i_1 * g_1) \rightarrow$ most decide How strongly
 we need to store information
 related to 'Keys' assume model decides
 to allocate s_1 as 0.01 \rightarrow (very small) \rightsquigarrow (Not important)
 but, later at $t=9$
 we need a subject associated to verb are.
 so can we go back in time and change the
 importance of 'Keys'? \Rightarrow No.

issue-3 :- Right Now, $c_t \rightarrow c_{t+1} \rightarrow c_{t+2} \rightarrow c_{t+3} \dots c_{t+k-1}$

- if something happened at time step $\underline{t+3}$
 - then it cannot be changed \Rightarrow • No mechanism to reopen any CTR
 - No access to earlier time step
 - each time step sees the compressed version of earlier timesteps

\Rightarrow [Memory update is Local in time, but requirements are Global in time]