

* Forget Gate - $(\sigma, 1)$

$$f_t = \sigma(W_f h_{t-1} + U_f x_t + b_f)$$

The forget gate decides how much of previous hidden states to forget & how much of new information to retain.

$\rightarrow U_f x_t$

* Output gate - [regulating flow of information in & out of cell.]

$$o_t = \sigma(W_o h_{t-1} + U_o x_t + b_o)$$

The output state determines which information from cell state should be output at current time step

* New Candidate values for cell state - $(W h_{t-1} + U x_t + b)$

$g_t = \tanh[W h_{t-1} + U x_t + b]$

Input modulation gate [current input to LSTM cell & previous hidden state of cell]

$$i_t = \sigma(W_i h_{t-1} + U_i x_t + b_i)$$

* Memory of network at time t

$$s_t = \underbrace{f_t}_{\text{forget gate}} \odot \underbrace{s_{t-1}}_{\text{previous cell state}} + \underbrace{i_t}_{\text{input gate}} \odot \underbrace{g_t}_{\text{New candidate cell state}}$$

$0.7 s_{t-1} + 0.6 g_t$

... + t. (4)

Two internal states involved - Cell State (C) & Hidden state (H)
 Cell state at time $t = z_t$ [responsible for carrying information through time.]

↳ updation done on basis of
 input gate (i_t)
 forget gate (f_t)

[forget gate decides how much of previous hidden state to forget & how much of new input to retain.]

[Input state decides how much of previous hidden state to keep & how much of new input to incorporate.]

The hidden state (h_t) is filtered version of cell state & is computed based on the current input & the cell state using output gate.

Hidden state (layman terms) → Compressed representation of relevant information from the input sequence that is useful for current prediction task.

* Input gate at time t

$$i_t = \sigma \left(\underbrace{w_i h_{t-1}}_{\substack{\uparrow \\ \text{Input state decides how much of previous hidden state to keep}}} + u_i x_t + \underbrace{b_i}_{\substack{\rightarrow \text{Bias for input gate}}} \right)$$

Input state decides how much of previous hidden state to keep & how much of new input to incorporate

$$a = 1/a \quad b = 1/b$$

* Hidden state -

Output of network at time t

$$h_t = \sigma_t \odot \tanh(p_t)$$

$$a = \begin{pmatrix} a_1 \\ a_2 \\ a_3 \end{pmatrix} \quad b = \begin{pmatrix} b_1 \\ b_2 \\ b_3 \end{pmatrix}$$

$$\begin{pmatrix} a_1 \times b_1 \\ a_2 \times b_2 \\ a_3 \times b_3 \end{pmatrix}$$

* Output - $\hat{y}_t = g(\underline{Vp_t} + \underline{c})$

$$\checkmark i_t = \sigma(W_i h_{t-1} + u_i x_t + b_i)$$

$$\checkmark f_t = \sigma(W_f h_{t-1} + u_f x_t + b_f)$$

$$\checkmark \sigma_t = \sigma(W_\sigma h_{t-1} + u_\sigma x_t + b_\sigma)$$

$$\checkmark g_t = \tanh(W g_{t-1} + u g_t + b)$$

$$\underline{p_t} = f_t \odot p_{t-1} + i_t \odot g_t \checkmark$$

$$h_t = \sigma_t \odot \tanh(p_t)$$

$$\hat{y}_t = g(\underline{Vp_t} + \underline{c})$$

Question: Consider An LSTM Network With One Hidden Layer Of 20 Nodes Used For Predicting The Next Word In A Text Corpus. No Bias Is...

Consider an LSTM network with one hidden layer of 20 nodes used for predicting the next word in a text corpus. No bias is used in any of the nodes. The corpus is of length 1000 words and there are 100 unique words. Assume a 10 dimensional word embedding module outside of the LSTM network, whose output is fed to the word predictor LSTM network. What will be the total number of trainable weights in the LSTM network?

We have x_t of shape 10×1

x_t is 10 dimensional word embedding.

$$x_t = \begin{pmatrix} 0.5 \\ 0.1 \\ 0.7 \\ 0.11 \\ \vdots \\ 0.9 \end{pmatrix}_{10 \times 1}$$

Following the LSTM architecture, we first need to evaluate parameter sizes for i_t , o_t , f_t , g_t .

Equation	Parameter size
$i_t = \sigma(\underbrace{w_i}_{20 \times 20} \underbrace{h_{t-1}}_{20 \times 1} + \underbrace{u_i}_{10 \times 1} \underbrace{x_t}_{20 \times 1})$ <p>(20×1)</p>	u_i is of size 20×10 . [1 hidden layer with 20 nodes] w_i is of size 20×20 $\therefore h_{t-1}$ will be of size 20×1
$o_t = \sigma(\underbrace{w_o}_{20 \times 20} \underbrace{h_{t-1}}_{20 \times 1} + \underbrace{u_o}_{10 \times 1} \underbrace{x_t}_{20 \times 1})$	u_o is of size 20×10 . w_o is of size 20×20 .
$f_t = \sigma(\underbrace{w_f}_{20 \times 20} \underbrace{h_{t-1}}_{20 \times 1} + \underbrace{u_f}_{10 \times 1} \underbrace{x_t}_{20 \times 1})$	w_f is of size 20×20 u_f is of size 20×10 .
$g_t = \sigma(\underbrace{w}_{20 \times 20} \underbrace{h_{t-1}}_{20 \times 1} + \underbrace{u}_{10 \times 1} \underbrace{x_t}_{20 \times 1})$	u is of size 20×10 w is of size 20×20 .
Total parameters	$20 \times 10 \times 4 + 20 \times 20 \times 4$ $= 2400$

we will eventually have

$$\hat{y}_t = g(Vx_t + c)$$

$$\Rightarrow \boxed{\hat{y}_t = g(Vx_t)}$$

$$V \quad x_t$$
$$100 \times 20 \quad 20 \times 1$$

Here V will be a matrix of order 100×20

[$\because x_t$ is a vector of order 20×1 hence.

$$x_t = f_t \odot x_{t-1} + i_t \odot g_t]$$

$$\text{Total trainable parameters} = \frac{2400}{=} + \underbrace{100 \times 20}_{=}$$
$$= 4400$$

Two gates for GRU \rightarrow reset gate & update gate.

* Reset gate

$$r_t = \sigma(W_r \cdot h_{t-1} + U_r x_t + b_r)$$

Reset gate determines how much of past knowledge to forget

* Update gate

$$u_t = \sigma(W_u \cdot h_{t-1} + U_u x_t + b_u)$$

update gate determines how much of past knowledge to be passed on to future.

* New Candidate value for all state

$$g_t = \tanh(W_c (\underbrace{r_t}_{\text{reset}} * h_{t-1}) + U_c x_t + b_c)$$

$$\Rightarrow h_t = u_t * g_t + (1 - u_t) h_{t-1} \rightarrow \text{Memory unit at time } t$$

$$h_t = 0.8 g_t + 0.2 h_{t-1}$$

$$z_t = h_t$$

$$\hat{y}_t = g(V z_t + c) \rightarrow [\text{output layer}]$$