

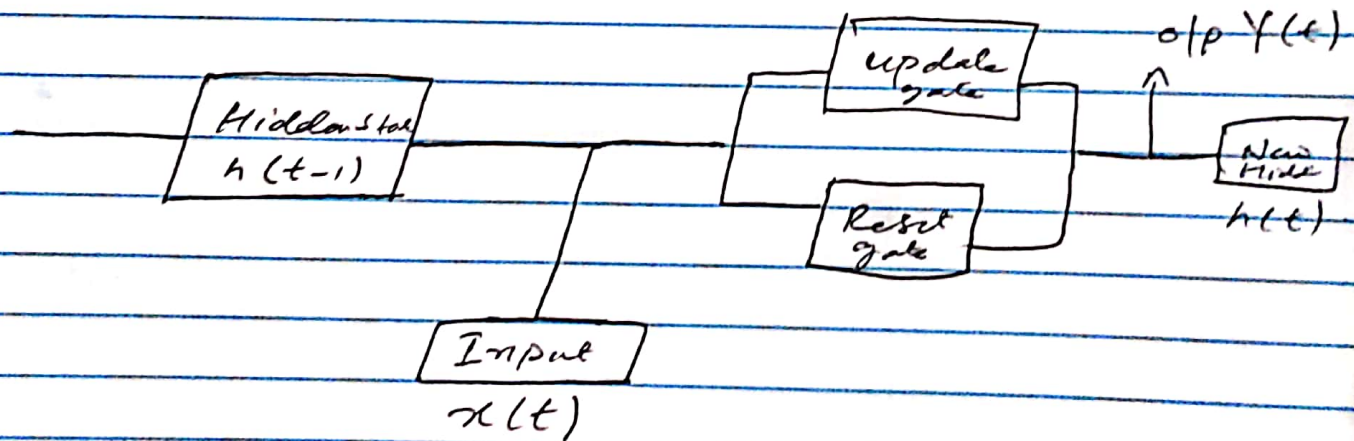
Gated Relu units, useful for solving the vanishing gradient problem.

* Architecture Similar to LSTM, while there are some differences.

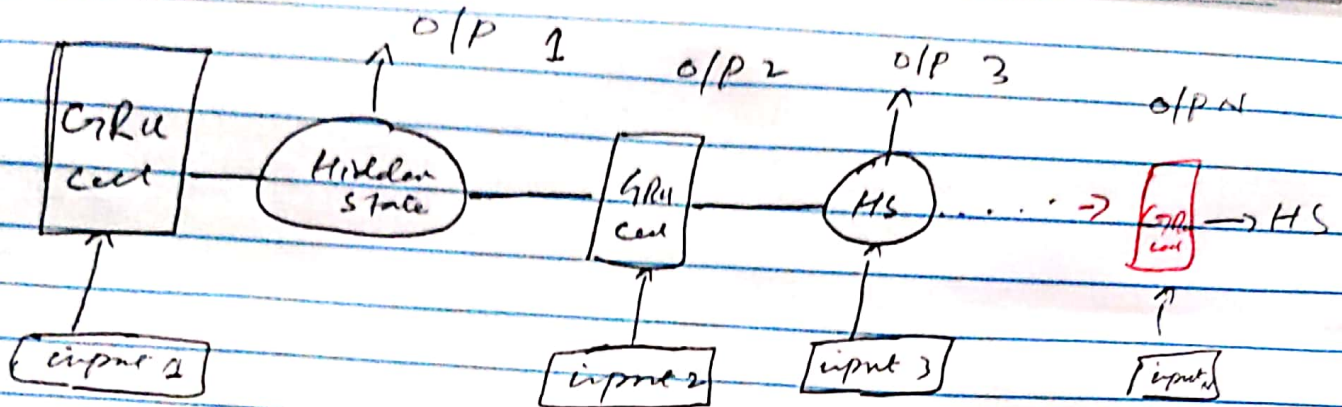
- 1- LSTM - 3 gates (forget, input and output gates)
- 2- GRU - 2 gates (update, reset)

* GRU is less complex, well suited for small datasets

* Gates Allow you to control the flow of information. How much information to be kept and how much information to forget. This helps in tracking the Long term dependency problem.



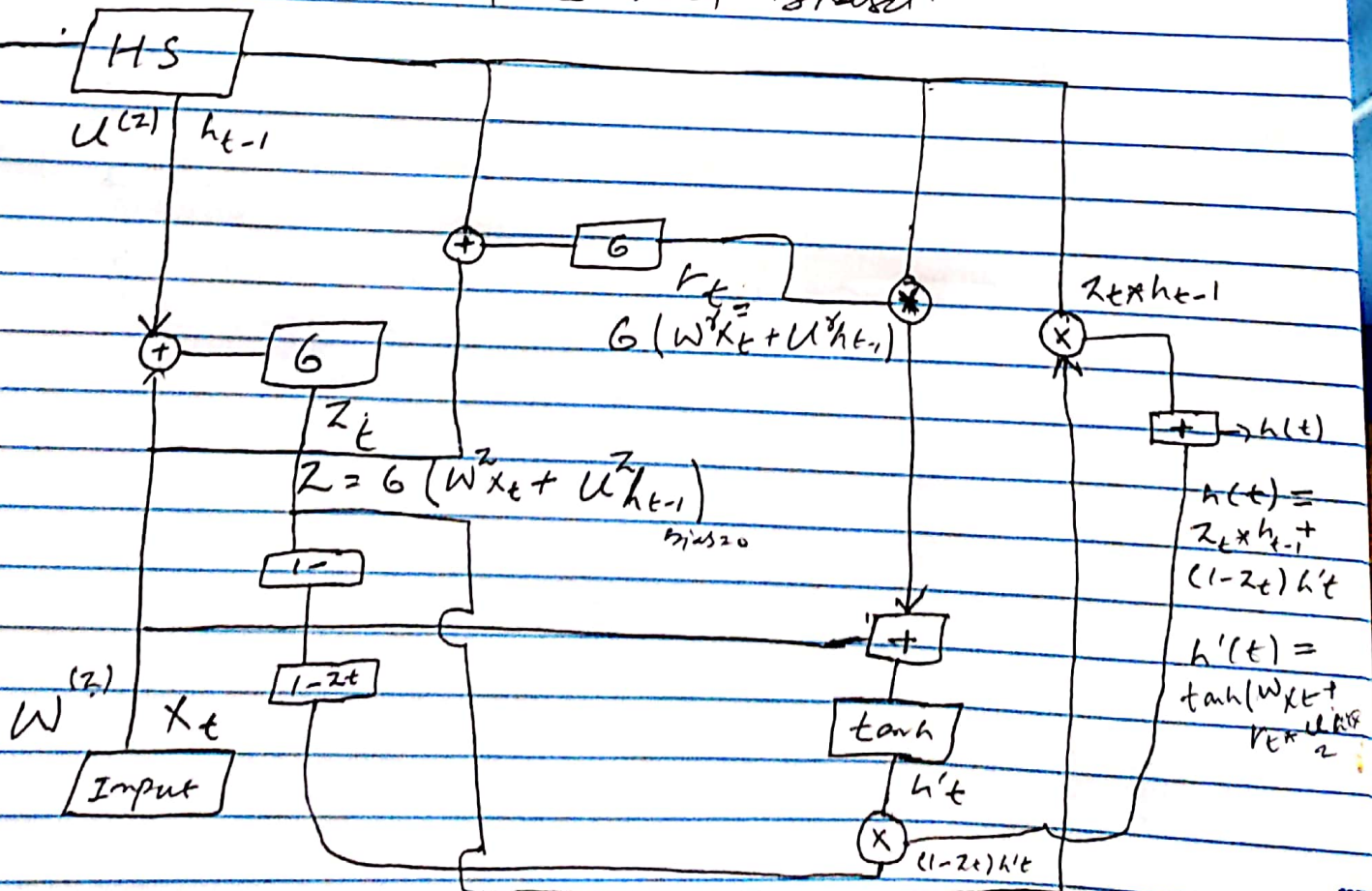
GRU Architecture.



Update Gate:- { How much to keep the information }
 Reset Gate:- { How much to forget }.

* weights are associated at every stage.

* Update Gate :- z_t works on i/p $x(t)$, $h(t-1)$ hidden state and biases.



$$z_t = \sigma(W^z x_t + U^z h_{t-1})$$

(update)

$$bias = 0$$

$$r_t = \sigma(W^r x_t + U^r h_{t-1})$$

(reset)

$$h'_t = \tanh(W^h x_t + r_t * U^h h_{t-1})$$

(Intermediate)

$$h_t = (1 - z_t) h'_t + z_t * h_{t-1}$$

(current unit)

A worked example

'text = MathMathMathMathMath'

Preprocessing steps :-

S-1 Convert text into numeric values and form a dictionary.

Dictionary: {'h': 0, 'a': 1, 't': 2, 'm': 3}

Our encoded o/p

MathMath = [3, 1, 2, 0, 3, 1, 2, 0].

S-2 Create batches of data

Let's Put the following Settings

Batch size (B) = 2

Sequence size (S) = 3

Vocabulary (V) = 4

Output (O) = 4

$$\boxed{[3, 1, 2, 0, 3, 1]}, \boxed{[2, 0, 3, 1, 2, 0]}, \dots$$

Batch 1

Batch 2

Create mini-batches of 3

$$\left[\begin{bmatrix} [3, 1, 2] \\ [0, 3, 1] \end{bmatrix} \begin{bmatrix} [2, 0, 3] \\ [1, 2, 0] \end{bmatrix} \begin{bmatrix} [3, 1, 2] \\ [0, 1, 3] \end{bmatrix} \right]$$

Transpose Sequences

$$\left[\begin{bmatrix} [3, 1, 2] \\ [0, 1, 3] \end{bmatrix} \begin{bmatrix} [2, 0, 3] \\ [1, 2, 0] \end{bmatrix} \begin{bmatrix} [3, 1, 2] \\ [0, 1, 3] \end{bmatrix} \right] = \begin{bmatrix} \begin{bmatrix} 3 \\ 0 \end{bmatrix} \begin{bmatrix} 1 \\ 3 \end{bmatrix} \begin{bmatrix} 2 \\ 1 \end{bmatrix} \\ \begin{bmatrix} 2 \\ 1 \end{bmatrix} \begin{bmatrix} 0 \\ 2 \end{bmatrix} \begin{bmatrix} 3 \\ 0 \end{bmatrix} \\ \begin{bmatrix} 3 \\ 0 \end{bmatrix} \begin{bmatrix} 1 \\ 3 \end{bmatrix} \begin{bmatrix} 2 \\ 1 \end{bmatrix} \end{bmatrix}$$

5 - One hot encoding

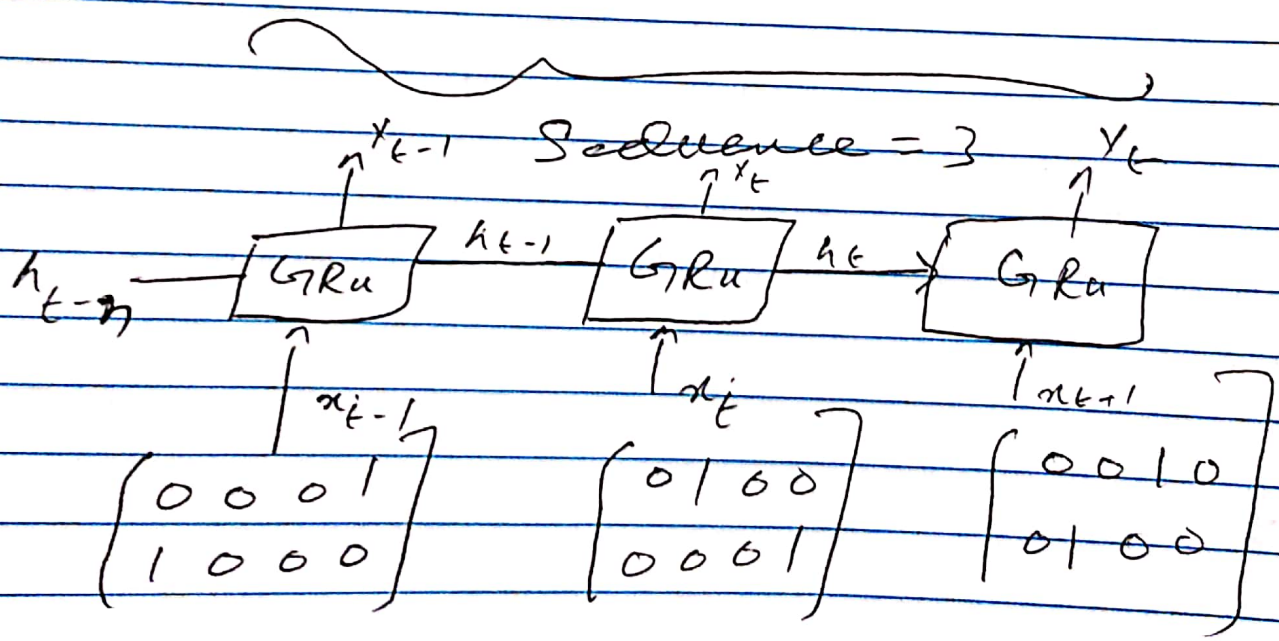
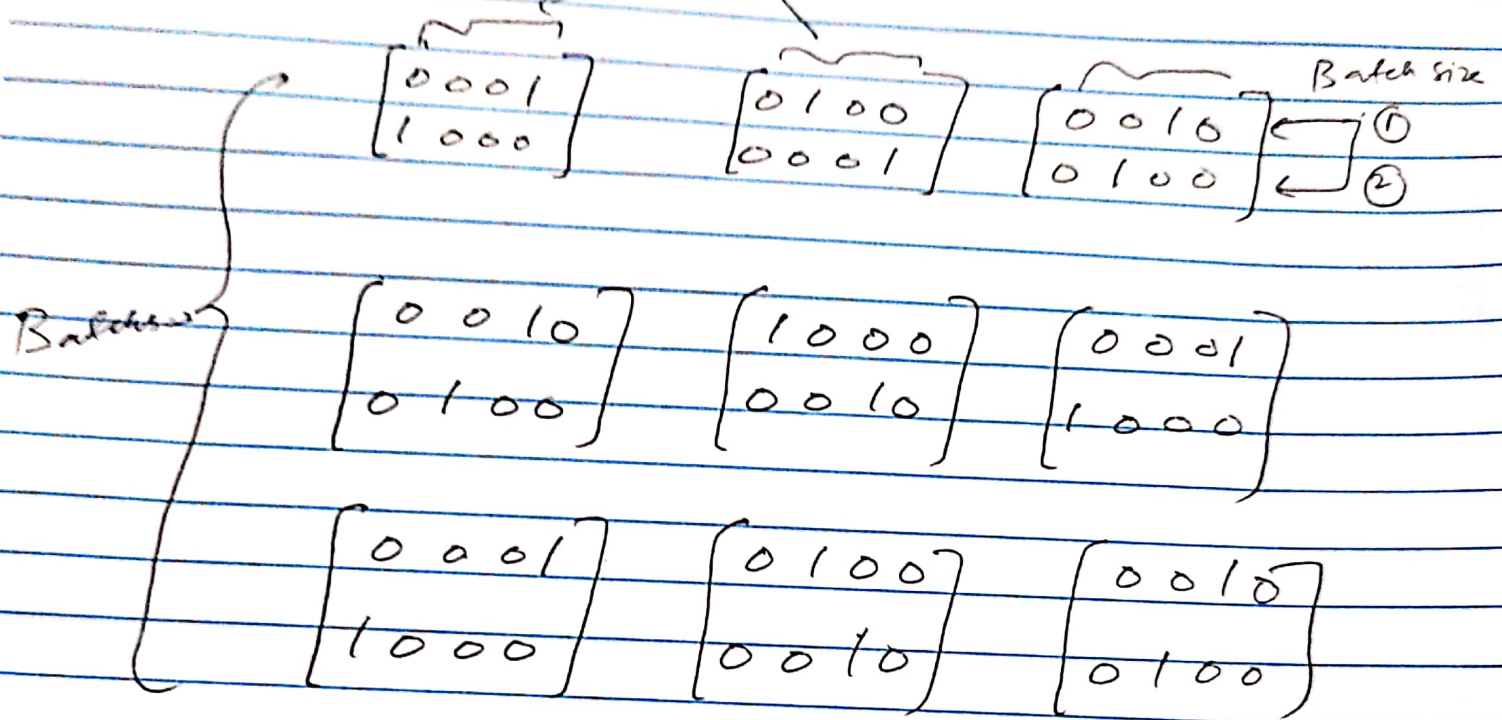
$$\begin{bmatrix} \begin{bmatrix} 3 \\ 0 \end{bmatrix} \begin{bmatrix} 1 \\ 3 \end{bmatrix} \begin{bmatrix} 2 \\ 1 \end{bmatrix} \\ \begin{bmatrix} 2 \\ 1 \end{bmatrix} \begin{bmatrix} 0 \\ 2 \end{bmatrix} \begin{bmatrix} 3 \\ 0 \end{bmatrix} \\ \begin{bmatrix} 3 \\ 0 \end{bmatrix} \begin{bmatrix} 1 \\ 3 \end{bmatrix} \begin{bmatrix} 2 \\ 1 \end{bmatrix} \end{bmatrix} = \begin{bmatrix} \begin{bmatrix} 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix} \\ \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \end{bmatrix} \\ \begin{bmatrix} 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix} \end{bmatrix}$$

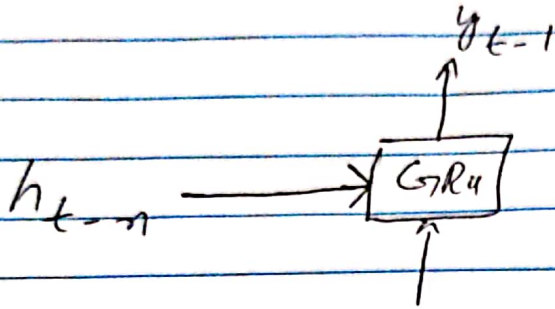
$$x(t-1) \quad x(t) \quad x(t+1)$$

Dimensions of our dataset

Number of batches \times Number of letter in each sequence
 \times Size of batch \times Vocabulary

$\{ 3 \times 3 \times 2 \times 4 \} \Rightarrow \text{Shape}$
Vocabulary = 4





Batch
Sequence 1

$$\begin{bmatrix} 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \end{bmatrix}$$

S-2 Define weight Matrices: 1. W_{hx}

Hidden Size = 2, $\Rightarrow 4 \times 2$

Random initialization

Vocabulary = 4

$$\begin{bmatrix} 0.6614 & 0.2669 \\ 0.0617 & 0.6213 \\ 0.4519 & -0.1661 \\ -1.5228 & 0.3817 \end{bmatrix}$$

W_2