

Gated RNN: The Gated Recurrent Unit (GRU) RNN

Background:-

- Have shown success in several applications involving sequential / temporal data.

For ex: In speech recognition, music synthesis, NLP, machine translation etc.

- Just like LSTM, GRUs perform well in long sequence applications.
- Introduced by Kyunghyun Cho in 2014 and offers an improvement over LSTM and Vanilla RNN.

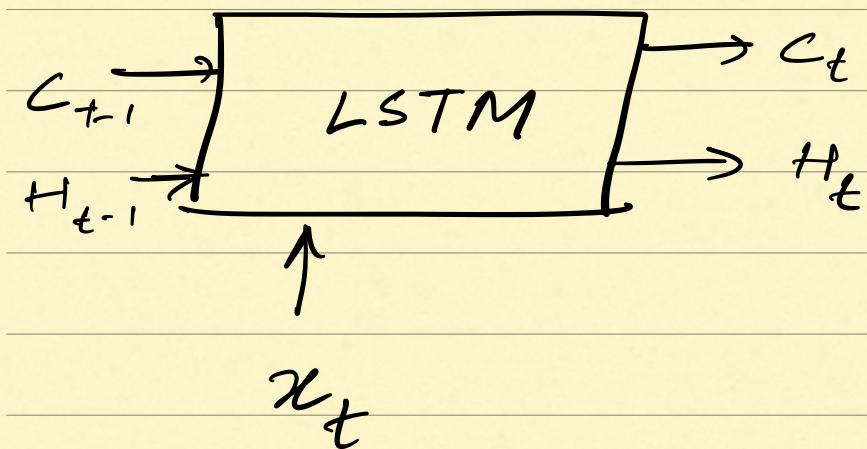
- Its' success can be attributed to the "gate signaling" that controls how the present input-block and previous memory are fused to update the present activation and produce present state.
- Only flipside: Increase in parameterization through the added gate-networks.
- Unlike LSTM, GRU employs 2 gate-networks (couples two LSTM gates).



How?

- GRU couples the input gate and the forget gate by a direct sum
- GRU also removes one non-linearity in the network "signal path" by feeding back the memory cell state.

The Standard GRU RNN :-





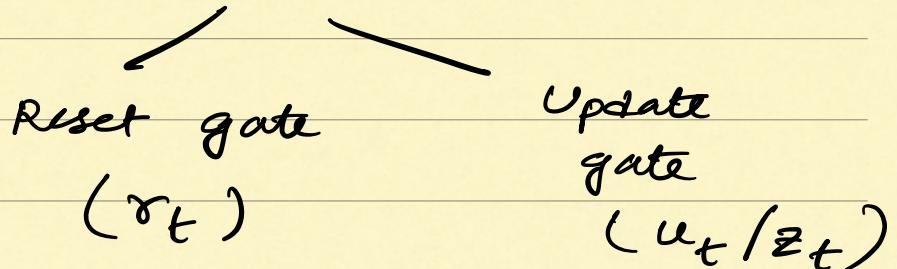
Doesn't have
a separate
cell state (C_t).

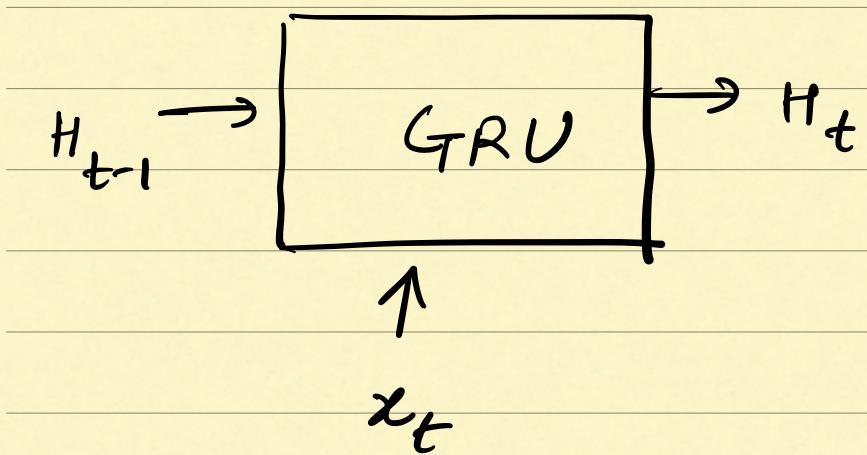
Due to simpler
architecture, GRUs
are faster to train.

Architecture of GRU:-

GRU reduces the gating signals to two from three in the LSTM architecture.

The two gates are





- At each timestamp t , it takes an input x_t and the prev. hidden state H_{t-1} , and subsequently outputs the new hidden state H_t .

Reset gate (Short Term Memory) :-

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

weight matrices (for simplicity)

Range: $0 < r_t < 1$

Update gate (Long Term Memory) :-

$$u_t / z_t = \sigma(x_t \circ v_u +$$

$$H_{t-1} \cdot w_u$$

$$+ b_u)$$

D

* How GRU works ?

— Two step process :-

Step 1: Generate Candidate hidden state

$$\hat{h}_t = \tanh(x_t \cdot v_g + (r_t \cdot h_{t-1}) \cdot w_g)$$

↑
↓

Tells us how much influence the previous hidden state can have on the candidate state (by adjusting r_t).

'6 $r_t \rightarrow 0 \Rightarrow$ Info. from the previous hidden state is ignored.

'6 $r_t \rightarrow 1 \Rightarrow$ entire info. from the previous hidden state is

considered.

Step 2: calculate the current hidden state (H_t).

— Instead of using a separate gate (like in LSTM), in GRU we use a single update gate to control both the historical info.

(H_{t-1}) as well as the new info (\hat{H}_{t-1}).

— Eqn: $H_t = u_t \cdot H_{t-1} +$

$$(1-u_t) \cdot \hat{H}_t$$

"Hence, the current hidden memory

state is a linear interpolation b/w previous hidden memory states and the candidate new state."

16 $u_t \rightarrow 0$: New hidden state won't have much info from prev. hidden state

$$H_t \sim \hat{H}_t$$

16 $u_t \rightarrow 1$: Previous step states for units are copied over to the current step.

$$H_t \sim H_{t-1}$$

Extra :-

The GRU model as presented
in the original paper :-
(Chung et al. 2014)

$$\tilde{h}_t = g \cdot (v_h \cdot (r_t \odot h_{t-1}) + w_h \cdot s_t + b_h)$$

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

* The two gates were presented as :-

$$z_t = \sigma (w_z h_{t-1} + w_z s_t + b_z)$$

$$\sigma_t = \sigma(u_r h_{t-1} + w_s s_t + b_r)$$

$g(\cdot) \Rightarrow$ Activation non-linearity

Typically tanh (can also be implemented as a ReLU).

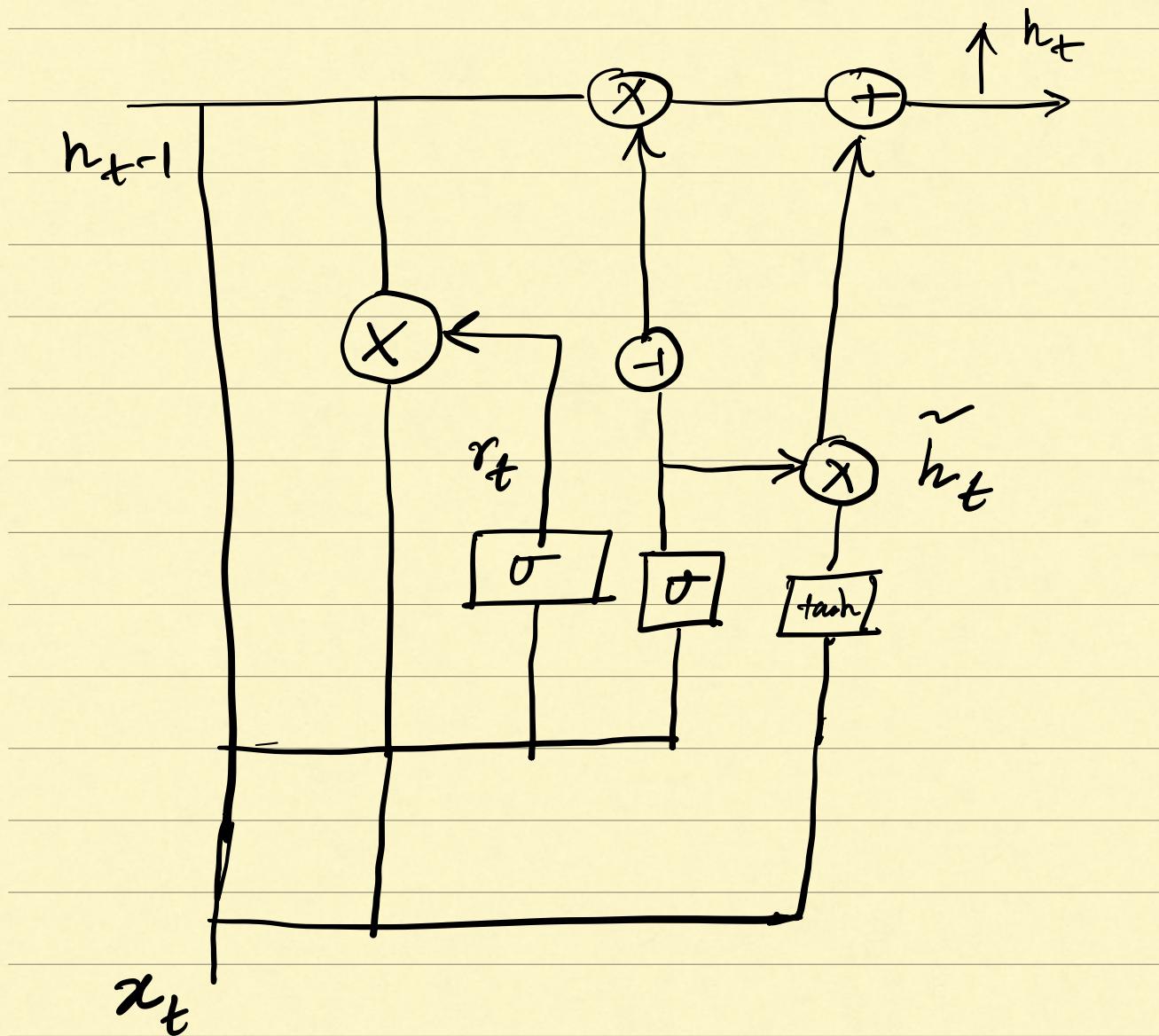
• \Rightarrow Hadamard multiplication

Weighted sum is implemented via element wise H.M. to gating signals.

Gates intuitively resemble "analog switches".

- GRU is similar to LSTM, however with one less external gating signal in the coupling interpolation.

- This eliminates one gating signal and the associated parameters.



* How GRU overcomes the problem of vanishing / exploding gradients during backpropagation ?

⇒ In RNN BP , the component that may lead to vanishing or exploding gradients in the component $\frac{\partial h_t^{(i)}}{\partial h_k^{(i)}}$, which backpropagation the error at sequence step t to sequence step k so that the model learns distant dependencies or correlations.

$$\frac{\partial h_t^{(i)}}{\partial h_k^{(i)}} = (u_{ii})^{t-k} \prod_{g=k+1}^t \sigma'(z_g^{(i)})$$

When $(t-k)$ is large, a vanishing gradient condition will arise when the gradients of the activations in the hidden state and / or the weights are less than 1 since the product of $(t-k)$ in them would force the overall product to near zero.

Sigmoid and tanh gradients are often less than 1 and saturate fast where they have near-zero gradients, thus making the problem of v.g. more sense.

Similarly, exploding gradients can happen when the weight-connection $w_{ii}^{(i)}$ b/w the i^{th} hidden to the i^{th} hidden unit

is greater than 1 since the product of $(t-k)$ in then would make the term $(u_{ik})^{t-k}$ exponentially large for large values of $(t-k)$.

$$\text{When } u_t \rightarrow 1 \Rightarrow h_t^{(i)} \approx h_{t-1}^{(i)}$$

$$t \in k$$

k : set of all hidden units for which $z_t^{(i)} \geq 0$

$$\text{now, } \frac{\partial h_t^{(i)}}{\partial h_{t-1}^{(i)}} \approx 1$$

This will ensure that the notorious term $\frac{\partial h_t^{(i)}}{\partial h_k^{(i)}}$ is also 1

Why?

$$\Rightarrow \frac{\partial h_t^{(i)}}{\partial h_k^{(i)}} = \prod_{g=t+1}^t \frac{\partial h_g^{(i)}}{\partial h_{g-1}^{(i)}}$$

$$= 1 \cdot 1 \cdots 1$$

($t-k$) times

$$= 1$$

This allows the hidden states to be copied over many sequence steps w/o activation and hence the chances of a v.g. diminish and the model is able to learn temporally long-distance correlation b/w words.

GRU VIS LSTM :-

Correspondence among two
gated RNNs :-

$$i_t \leftrightarrow z_t$$

$$f_t \leftrightarrow (1 - z_t)$$

$$c_t \leftrightarrow h_t$$

$$\tilde{z}_t \rightarrow \tilde{h}_t$$

Using LSTM notations, we can
exprn GRU eqns :-

$$\tilde{c}_t = g (v_c h_{t-1} + w_c s_t + b_c)$$

$$c_t = (1 - i_t) \odot c_{t-1} + i_t \odot \tilde{c}_t$$

$$h_{t-1} = o_t \odot c_{t-1}$$

✳ Summary :-

Hochreiter & Schmidhuber (1997)

Bengio (1994)

\Rightarrow It is difficult
to capture long-term
dependencies using simple
RNNs because the gradients
tend to vanish or explode
with long sequences.

- Two particular architectures :-

① LSTM (1997)

② GRU (2014)