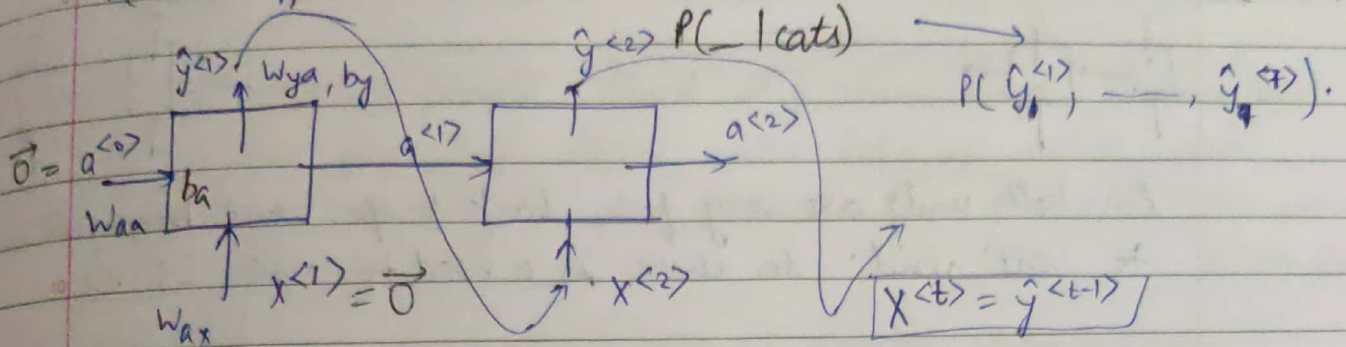


date
14/11/19

* Language Modelling - What is the probab. associated with the sentence or word?

$$P(\text{The apple \& the pair salad}) \approx 3.3 \times 10^{-13}$$

$$P(\text{The apple \& the pear salad}) \approx 5.7 \times 10^{-10} \quad \checkmark \quad (\text{more common})$$



$a^{<0>}$ is initialised to $\vec{0}$ (No input)

$x^{<1>}$ = most frequently occurring words is not a good idea as all generated sent. by RNNs will start with 'The'.

$$a^{<t>} = \tanh(g_1(w_{ax} x^{<t>} + w_{aa} a^{<t-1>} + b_a))$$

$$\hat{y}^{<t>} = \underset{\text{softmax}}{g_2(w_{ya} a^{<t>} + b_y)}$$

Idea:- We select the 1st word, based on the probab of words in the corpus. eg- $P(\text{The}) = 0.7$, there are 70% chances of choosing 'The' out of all possible words/choices

$$P(\text{cats}) = 0.0004$$

$$P(\hat{y}^{<1>} | \text{cats})$$

$$P(\hat{y}^{<3>} | \text{cats average})$$

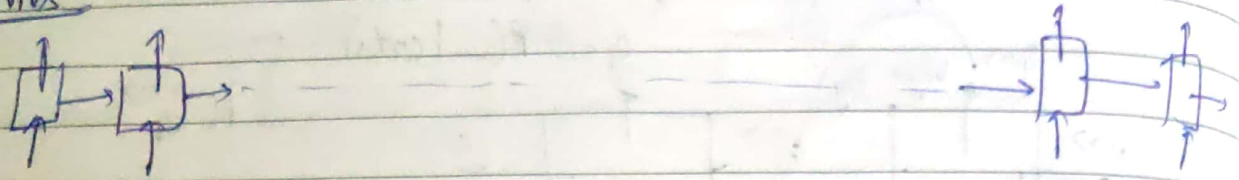
- $\langle \text{EOS} \rangle$: End of sent. marker
- To deal with $\langle \text{UNK} \rangle$, we can even have a model ~~with~~ at character level as no chr. (ASCII) can be unknown.

LSTMs

RNNs suffer from Vanishing Gr. problem much more than DNNs or CNNs.

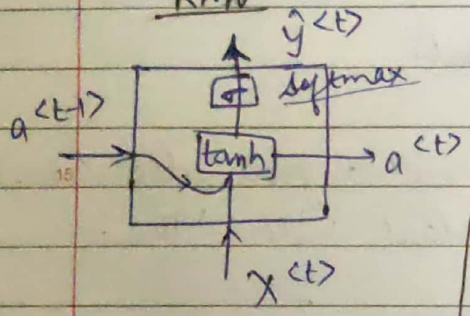
The cat, which already ate fruits like apple, pear, etc., was full.

RNNs



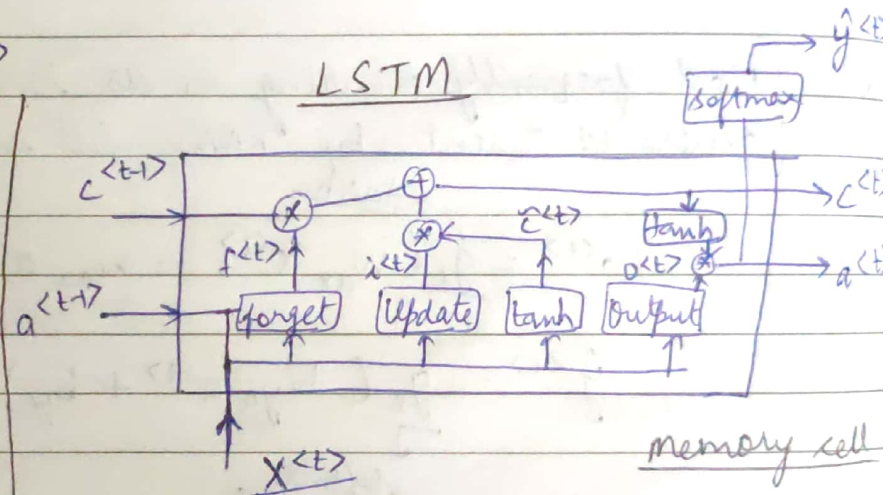
As both units are very far, back prop. won't propagate error to 'cat' unit. So, there is a need to use LSTMs.

RNN

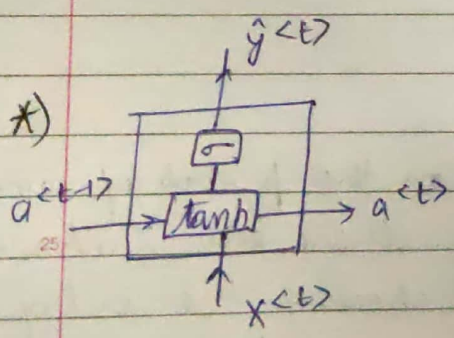


Doesn't capture long term dependencies.

LSTM



memory cell



(forget gate)

$$\Gamma_f = \sigma(W_f [a^{(t-1)}, x^{(t)}] + b_f)$$

$$\Gamma_u = \sigma(W_u [a^{(t-1)}, x^{(t)}] + b_u)$$

$$a^{(t)} = g_1(W_{ax} x^{(t)} + W_{aa} a^{(t-1)} + b_a)$$

$$\Gamma_o = \sigma(W_o [a^{(t-1)}, x^{(t)}] + b_o)$$

$$\hat{y}^{(t)} = g_2(W_{ya} a^{(t)} + b_y)$$

$$c^{(t)} = \Gamma_u * \tilde{c}^{(t)} + \Gamma_f * c^{(t-1)}$$

update $c^{(t-1)}$
& change it to $c^{(t)}$.

$$\tilde{c}^{(t)} = \tanh(W_c [a^{(t-1)}, x^{(t)}] + b_c)$$

$$a^{(t)} = \Gamma_o * c^{(t)}$$

In our terminology, $W_f = [W_{fa}, W_{fx}]$. Composed of 2 wts.
(Andrew Ng combined both in 1 slide).

$\tilde{c}^{(t)}$: cell state (for memorizing)

$a^{(t)}$: Activation of previous unit

Before LSTMs, people tried skip connections to maintain long term dep. But there were lots of extra conn. resulting in very high # parameters. Hence, LSTMs were invented.

GRUs (much simpler than LSTMs).

$$\tilde{c}^{(t)} = \tanh(W_c [\Gamma_r * c^{(t-1)}, x^{(t)}] + b_c).$$

$$\Gamma_u = \sigma(W_u [\tilde{c}^{(t-1)}, x^{(t)}] + b_u).$$

(Relevance gate)

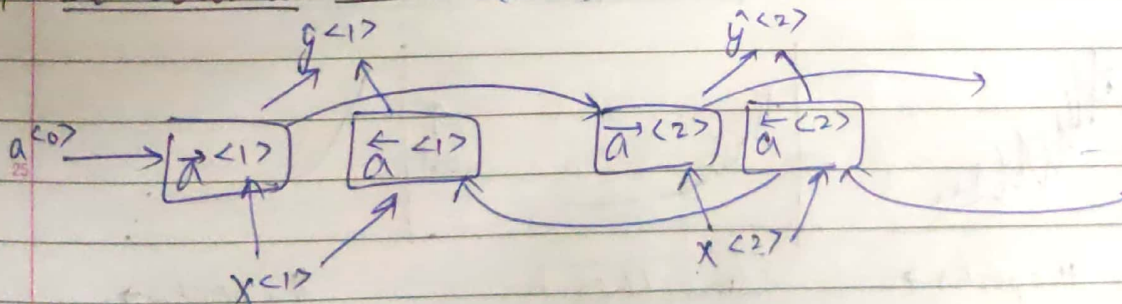
$$\Gamma_r = \sigma(W_r [\tilde{c}^{(t-1)}, x^{(t)}] + b_r).$$

$$\tilde{c}^{(t)} = \Gamma_u * \tilde{c}^{(t)} + (1 - \Gamma_u) * \tilde{c}^{(t-1)}.$$

$$a^{(t)} = \tilde{c}^{(t)}.$$

[has only 2 gates - update and Relevance].

Bidirectional RNNs (BRNN)

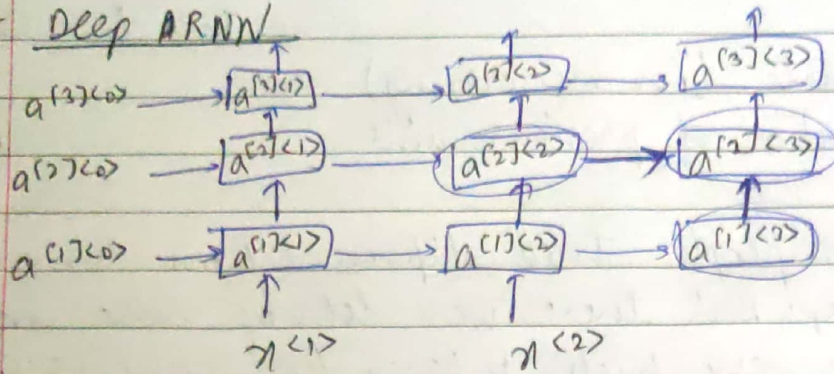


$$\hat{y}^{(t)} = g(W_{ya} a^{(t)} + b_y) \quad \text{--- one dir}$$

$$= g(\underbrace{W_{y\vec{a}} \vec{a}^{(t)}}_{\text{fwd}} + \underbrace{W_{y\bar{a}} \bar{a}^{(t)}}_{\text{Bwd}} + b_y)$$

$a^{(1)} \langle 2 \rangle$ means 1st layer, 2nd unit / timestamp.

* Deep ARNN



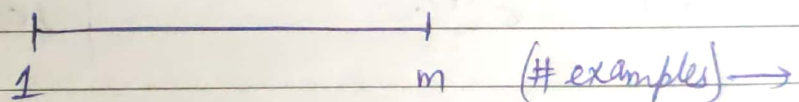
$$a^{(2)} \langle 3 \rangle = g(W_a^{(2)} [a^{(2)} \langle 2 \rangle, a^{(1)} \langle 3 \rangle] + b_a^{(2)}).$$

Turing complete archi

date
21/11/19

SGD (Mini-batch Gr. Descent)

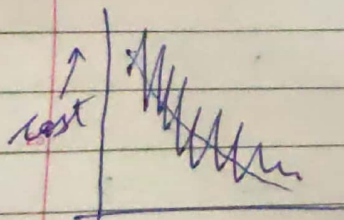
We can consider following # examples in 1 Gr desc update -



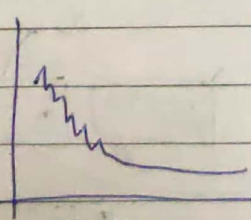
for 1 ex.

$1 \leq i \leq m$

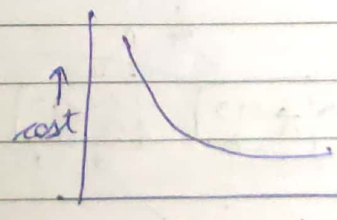
for m exs.



epochs \rightarrow
(Stochastic GD)



Smoother
(Mini-Batch GD)



epochs \rightarrow
(Batch GD)

$x^{(1)}$ \rightarrow Mini-batch 1.
 $x^{(1)}$ \rightarrow 1st th. ex

$x^{(2)}$ \rightarrow Mini-batch 2.
 $x^{(2)}$ \rightarrow 2nd th. ex

Stochastic Gr. Descent → When batch size = 1. We train on just 1 ex. before updating the parameters (ex. chosen randomly) epoch-wise

With Batch-GrD, learning is very slow

With SGD, learning is very fast

While Mini-Batch GrD takes good of both the things.

We require:

- 1) Learning should be moderately fast
- 2) Should reach to min as fast as possible.

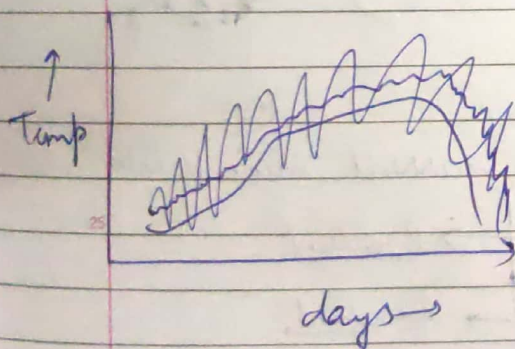
Moving Avg. : Taking avg of obs. within a window.
Assign higher wts. to nearer obs. & so on

$$\begin{cases} V_0 = 0 \\ V_1 = 0.9 V_0 + 0.1 \theta_1 \\ V_2 = 0.9 V_1 + 0.1 \theta_2 \\ \vdots \\ V_t = 0.9 V_{t-1} + 0.1 \theta_t \end{cases}$$

$$V_t = \beta V_{t-1} + (1-\beta) \theta_t \quad \text{exponentially weighted average}$$

$\beta=1$, just considering past avg

$\beta=0$, just considering current value



	$\beta = 0.98$
	$\beta = 0.9$
	$\beta = 0.5$

Momentum

$$V_0 = 0$$

$$V_0 = \beta V_0 + (1-\beta) \theta_1$$

$$V_0 = \beta V_0 + (1-\beta) \theta_2$$

$$V_{100} = 0.9 V_{99} + 0.1 \theta_{100}$$

$$V_{100} = \underbrace{0.9}_{[\beta]} (0.9 V_{98} + 0.1 \theta_{99}) + 0.1 \theta_{100}$$

Date
26/11/19

Exp. Weighted Averages

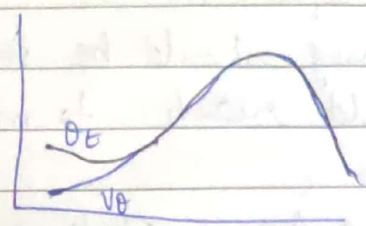
$$V_0 = 0$$

Repeat ϵ

Let new θ_t

$$V_0 := \beta V_0 + (1-\beta) \theta_t$$

Initially, V_0 will be very less & then it slowly tries to match up with moving weighted averages.

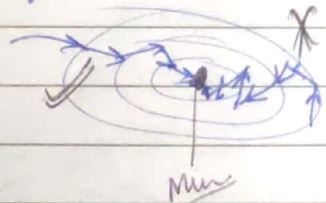


GD with Momentum

Try to make update step more efficient

$$w := w - \alpha \frac{\partial J}{\partial w}$$

Momentum helps in moving down gradients to reach min quickly. It disallows sideways movements when we move away from minima.



Momentum

On iter t :

compute $\partial J / \partial w$, $\partial J / \partial b$ on current mini-batch

$$V_{dw} = \beta V_{dw} + (1-\beta) \frac{\partial J}{\partial w} \rightarrow \text{gradient}$$

$$V_{db} = \beta V_{db} + (1-\beta) \frac{\partial J}{\partial b} \rightarrow \text{gradient}$$

$$w := w - \alpha V_{dw} ; \quad b := b - \alpha V_{db}$$

$$V_{dw} = \beta V_{dw} + (1-\beta) \frac{\partial J}{\partial w}$$

α = learning rate
 β = momentum

$\frac{V_{dw}}{1-\beta_1^t}$ (Bias correction to upgrade the V_{dw} value).

RMS prop
 $S_{dw} = \beta S_{dw} + (1-\beta) dw^2$
 $S_{db} = \beta S_{db} + (1-\beta) db^2$
 $\left(\frac{\partial J}{\partial w}\right)^2$ is element-wise multiplication

$$w := w - \alpha \frac{dw}{\sqrt{S_{dw} + \epsilon}} ; \quad b := b - \alpha \frac{db}{\sqrt{S_{db} + \epsilon}}$$

Adam opti.

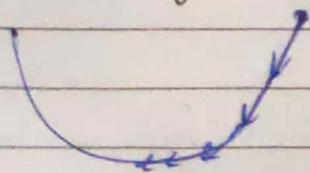
V_{dw}, V_{db} (by momentum)
 S_{dw}, S_{db} (by RMS prop)

$$V_{dw}^{corrected} = V_{dw} / (1 - \beta_1^t) ; \quad S_{dw}^{corr} = S_{dw} / (1 - \beta_2^t)$$

$$w := w - \alpha \frac{V_{dw}^{corr}}{\sqrt{S_{dw}^{corr} + \epsilon}} ; \quad b := b - \alpha \frac{V_{db}^{corr}}{\sqrt{S_{db}^{corr} + \epsilon}}$$

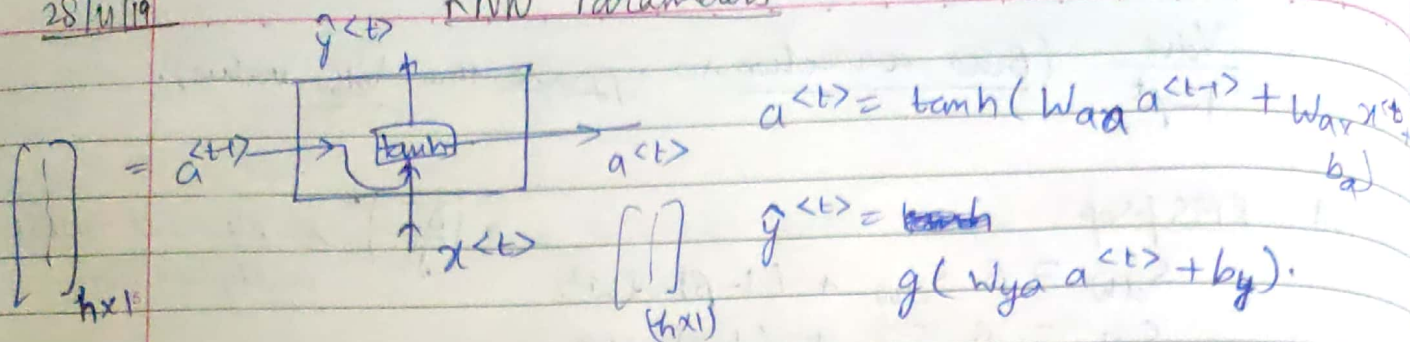
$$[\beta = 0.9, \quad \beta_2 = 0.999, \quad \epsilon = 10^{-8}]$$

lr decay



As # epochs \uparrow , $lr \downarrow$ so that we don't miss minima

28/11/19

RNN Parameters

$$a^{<t>} = \tanh(w_{aa} a^{<t-1>} + w_{ax} x^{<t>} + b_a)$$

$$g^{<t>} = \tanh(g(w_{ga} a^{<t>} + b_g))$$

$$x^{<t>} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}_{n \times 1}$$

{input can be video matrix, word embeddings, or one-hot vector}

$$\textcircled{1} \quad \begin{bmatrix} \quad \end{bmatrix}_{h \times n} \cdot \begin{bmatrix} \quad \end{bmatrix}_{n \times 1} + \begin{bmatrix} \quad \end{bmatrix}_{h \times h} \cdot \begin{bmatrix} \quad \end{bmatrix}_{h \times 1} + \begin{bmatrix} \quad \end{bmatrix}_{h \times 1} = \begin{bmatrix} \quad \end{bmatrix}_{h \times 1}$$

$(w_{ax}) \quad (x^{<t>}) \quad (w_{aa}) \quad (a^{<t-1>}) \quad (b_a) \quad (a^{<t>})$

$$\begin{aligned} \# \text{parameters} &= (h \times n) + (h \times h) + (h \times 1) \\ &= \boxed{h(n+h+1)} \end{aligned}$$

Combining w_{aa} & w_{ax} .

$$\begin{bmatrix} \quad \end{bmatrix}_{h \times (n+h)} \cdot \begin{bmatrix} \quad \end{bmatrix}_{(n+h) \times 1} + \begin{bmatrix} \quad \end{bmatrix}_{h \times 1} = \begin{bmatrix} \quad \end{bmatrix}_{h \times 1}$$

$(w_a) \quad (x^{<t>}, a^{<t-1>}) \quad (b_a) \quad (a^{<t>})$

$$\Rightarrow a^{<t>} = \tanh(w_a [x^{<t>}, a^{<t-1>}] + b_a)$$

$$\begin{matrix}
 \left[\right]_{w \times h} & \left[\right]_{h \times 1} & + & \left[\right]_{w \times 1} & = & \left[\right]_{w \times 1} \\
 \text{Wya} & \text{a}^{<t>} & & \text{by} & & \text{y}^{<t>}
 \end{matrix}$$

LSTM parameters

$$\Gamma_f^a = \sigma \left(W_f \left[\underbrace{x^{<t>}_{(h \times 1)} \underbrace{a^{<t-1>}_{(h \times 1)}}_{(n+h) \times 1} \right] + b_f \right)$$

$$\# \text{ parameters} = 4 \times [h \times (n+h) + h].$$

↳ 4 gates (forget, update, O/P, tanh).

$$\# \text{ parameters in GRVs} = 3 \times [h \times (n+h) + h].$$