# Dealing with Vanishing Gradients in RNN

Tirtharaj Dash

Dept. of CS & IS and APPCAIR
BITS Pilani, Goa Campus

November 6, 2021

Where we are:

→ forward computation in RNN

→ Backprop equations ( BPTT )

Issues:

① Exploding gradients → Easy.

② Vanishing gradients

Truncated. Gradient
BPTT clipping.

by value by Norm

Vanishing gradient problem: (Hard)

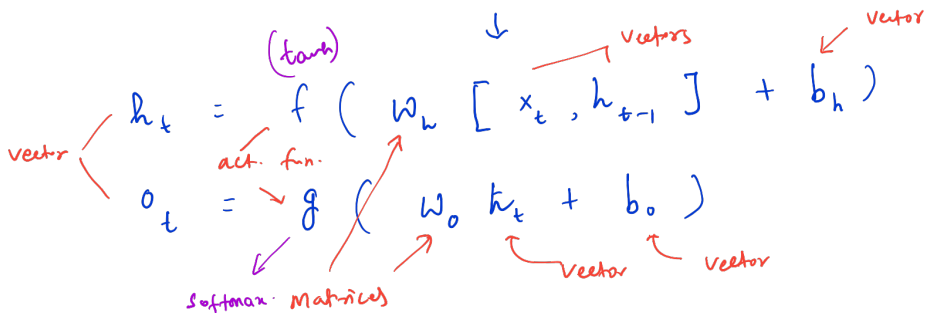$\longrightarrow$ Deal with. it using

( Architectural change to RNN).

GRU          LSTM

RNN Unit: or RNN cell:

$$h_t = f(x_t, h_{t-1} \; ; \; W_h)$$

"$\hat{y}_t$"

$$o_t = g(h_t; W_o)$$

(tanh)

$$h_t = f(W_h [x_t, h_{t-1}] + b_h)$$

vector — vectors — vector

act. fun.

$$o_t = g(W_o h_t + b_o)$$

softmax. Matrices

vector — vector

$o_t$

$g$

$f$

$h_t$

$h_{t-1}$

$x_t$

RNN unit/cell

Varding grad.

Long-term dependencies ( Long-range dependency )

"That [guy] , who we met at the airport ---- ,
                                    [was] one of my students "

Singular.

. guy — was

The `model has to "remember"
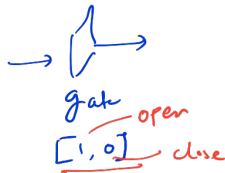                                    ↓
                                Memory.

In RNN:

hidden state is updated.

hidden state overwritten.

"guy" ———— "wao"

$t = 2$          $t = 35$

( GRU: Gated Recurrent Unit. )

<u>GRU</u> :   uses   a   memory   cell.  (addition   variable)

(Basic)

       relevance  ( reset )   gates

       update   gates.



gate — open

$[1, 0]$ — close

$\rightarrow$   candidate   memory

$$\tilde{c}_t = \tanh\left( W_c \left[ c_{t-1}, x_t \right] + b_c \right)$$

$\rightarrow$   update   gate

$$z_t = \sigma\left( W_u \left[ c_{t-1}, x_t \right] + b_u \right)$$

$\rightarrow$   Memory   cell

hidden state

$$c_t = z_t * \tilde{c}_t + (1 - z_t)\, c_{t-1}$$

$$h_t = c_t$$

$$C_t = \underline{\underline{Z_t}} * \tilde{C}_t + \underline{(1 - Z_t)} \ \underline{C_{t-1}}$$

push    new       forget    old memory

info       relevance

$\underline{\underline{0.8}}$   0.99         $\underline{\underline{0.2}}$   0.01

$Z_t = 0.8$   $Z_w = 0$

$C_t = 1$   $C_t = \cdots\cdots\cdots\cdots$       $C_t =$

That   | guy | ,     -------   ---    ,   | was | one of - - -

GRU cell:



GRU cell:

$O_t$

softmax  $\sigma$

$T$

$C_{t-1}$

$= h_{t-1}$

$C_t$

$\tilde{C_t}$     $Z_t$

tanh     $\sigma$

New
Enformation

Weightings.( relevance/
how much open
the gate

$x_t$

The original GRU:

$$\tilde{C}_t = \tanh\left( W_c \left[ Z_r * C_{t-1}, x_t \right] + b_c \right)$$

How relevant is $C_{t-1}$ in computing $c_t$

Update $\quad Z_t = \sigma\left( W_u \left[ x_t, h_{t-1} \right] + b_u \right)$

relevance (reset) $\quad Z_r = \sigma\left( W_r \left[ x_t, h_{t-1} \right] + b_r \right)$

$$c_t = Z_t \, \tilde{C}_t + (1 - Z_t) \, c_{t-1}$$

$$h_t = c_t$$

LSTM: Long Short-Term Memory. $(1991)$

candidate memory $\quad \tilde{C}_t = \tanh\left(W_c\left[x_t, h_{t-1}\right] + b_c\right) \Leftarrow$

update $\quad Z_u = \sigma\left(W_u\left[x_t, h_{t-1}, c_{t-1}\right) + b_u\right)$

forget $\quad Z_f = \sigma\left(W_f\left[x_t, h_{t-1}, c_{t-1}\right] + b_f\right]$

output $\quad Z_o = \sigma\left(W_o\left[x_t, h_{t-1}, c_{t-1}\right] + b_o\right)$

new memory $\quad C_t = Z_u * \tilde{C}_t + Z_f * C_{t-1}$

hidden state $\quad h_t = Z_o * \tanh\left(C_t\right)$

Peephole.

* GRU
* RNN

Read about:
peephole connection.

Bidirectional RNNs

$$O_t = g\left(\omega_0 \left[\overrightarrow{h_t}, \overleftarrow{h_t}\right] + b_0\right)$$



* RNN
* GRU
* LSTM.

I think, "Deep learning is fun."

I think, "Deep sea is dangerous"

$\longrightarrow$ L to R.

$\longleftarrow$ R t

# Deep RNNs:

Grammar of the language.



$(L=L)$
$h_0$

$(L=2)$
$h_0$

$(L=1)$
$h_0$

$x_1$   $x_2$   $x_3$   $x_T$

Deep
Stack