Recap: RNN → tanh

$$a^{\langle t \rangle} = g\left(W_a\left[a^{\langle t-1 \rangle}, x^{\langle t \rangle}\right] + b_a\right)$$



Softmax

$\hat{y}^{\langle t \rangle}$

$h^{\langle t-1 \rangle}$   tanh   $h^{\langle t \rangle}$

$x^{\langle t \rangle}$

1) Reset gate: Influence of the previous hidden state

2) Update gate: Influence of a newly computed update

3) Proposing an update hidden state

4) Computing updated hidden state

## Sample Sentence:

The [cat,] which already ate,...., [was] full.

GRU unit will have new memory unit called $c$ = memory cell.

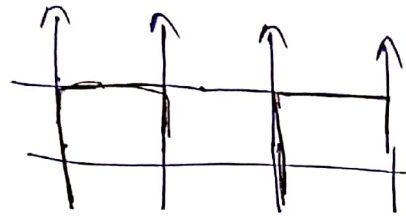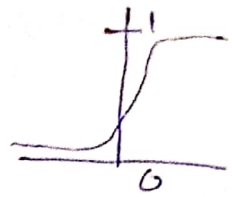will be used to remember if cat is singular or plural

$$c^{\langle t \rangle} = h^{\langle t \rangle} = a^{\langle t \rangle}$$

Candidate for $c^{\langle t \rangle}$
for $c$ for replacement

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

GRU → gamma u $\Gamma_u$ = 0 → 1
gate
assume always 0 or 1        update

$$\Gamma_u = \sigma \left( w_u [ c^{<t-1>}, x^{<t>}] + b_u \right) \quad \text{eqn (2)}$$



Gate

G for gamma

gate will decide wether we update it

$\tilde{c}^{<t>}$ will be the candidate value

Assume singular = 1, plural = 0
GRU will memorize the value antill was
jobs of gate is to decide when to update
the value.

When we see The __cat__ → subject

forget when we see __was__

gate → Candidate value          → old value

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \left(1 - \Gamma_u\right) * c^{<t-1>} \quad \text{eqn (3)}$$

E.g. if $\Gamma_u = 1$ then set the candidate value.
    and update
    for all values in the middle $\Gamma_u = 0$ should

The __cat__, which already ate...., __was__ full.

if $1 - \Gamma_u = 1$ ; $\Gamma_u = 0$ ; $c^{<t>} = c^{<t-1>}$ $\lg$ ③

$c^{<t-1>}$
$c^{<t-1>}$
$= a$

Softmax $\rightarrow$ $g^{<t>}$

$c^{<t>}$ new value for memory cell

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

element wise multiplication

$\tilde{c}^{<t>}$   $\Gamma_u$

tanh   σ

$x^{<t>}$

$\Gamma_u = 0.00001$

does not suffer from vanishy gradient problem
b/c if $\Gamma_u \approx 0$ then $1 - \Gamma_u = 1$
and $c^{<t>} = c^{<t-1>}$

allowed NN to train over long range
dependencies

$c^{<t>}$ can be a vector. eg. 100-dimensn
vector. then $\tilde{c}^{<t>}$ will also be
same dimension & $\Gamma_u$ will also be
same dimension.

if gate is a 100-dimensional then
it we will decide which bits we would like
to update. $\Gamma_u$ can have middle values
Element-wise multiplication. allows
which bits of the memory cell needs to be updated

it can help in keeping some bits Page 4
constant ~~as~~ while updating other bits.

E.g

The cat, which already ate .... was full

↓ Some bits for singular plural

↓ other bits to remember food ate

use only a subset of bits

Eq ① ② and ③

## Full GRU

Introduce $\Gamma_r \to$ how relevant is $c^{<t-1>}$ to compute the next candidate $\tilde{c}^{<t>}$

Reset gate... Influence of previous hidden state

$$v \to \Gamma_r = \sigma\left(w_r\left[c^{<t-1>}, x^{<t>}\right] + b_r\right) \quad \text{Eq} \, ④$$

Researchers have evaluated different combinations to assess long range affects

$$\tilde{h} \to \tilde{c}^{<t>} = \tanh\left(W_c\left[\Gamma_r * c^{<t-1>}, x^{<t>}\right] + b_c\right) — ①$$

Update Grate: Determines influence of an update proposal on the new hidden state

$$u \to \Gamma_u = \sigma\left(W_u\left[c^{<t-1>}, x^{<t>}\right] + b_u\right) — ②$$

$$h \to c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>} — ③$$

common version.