

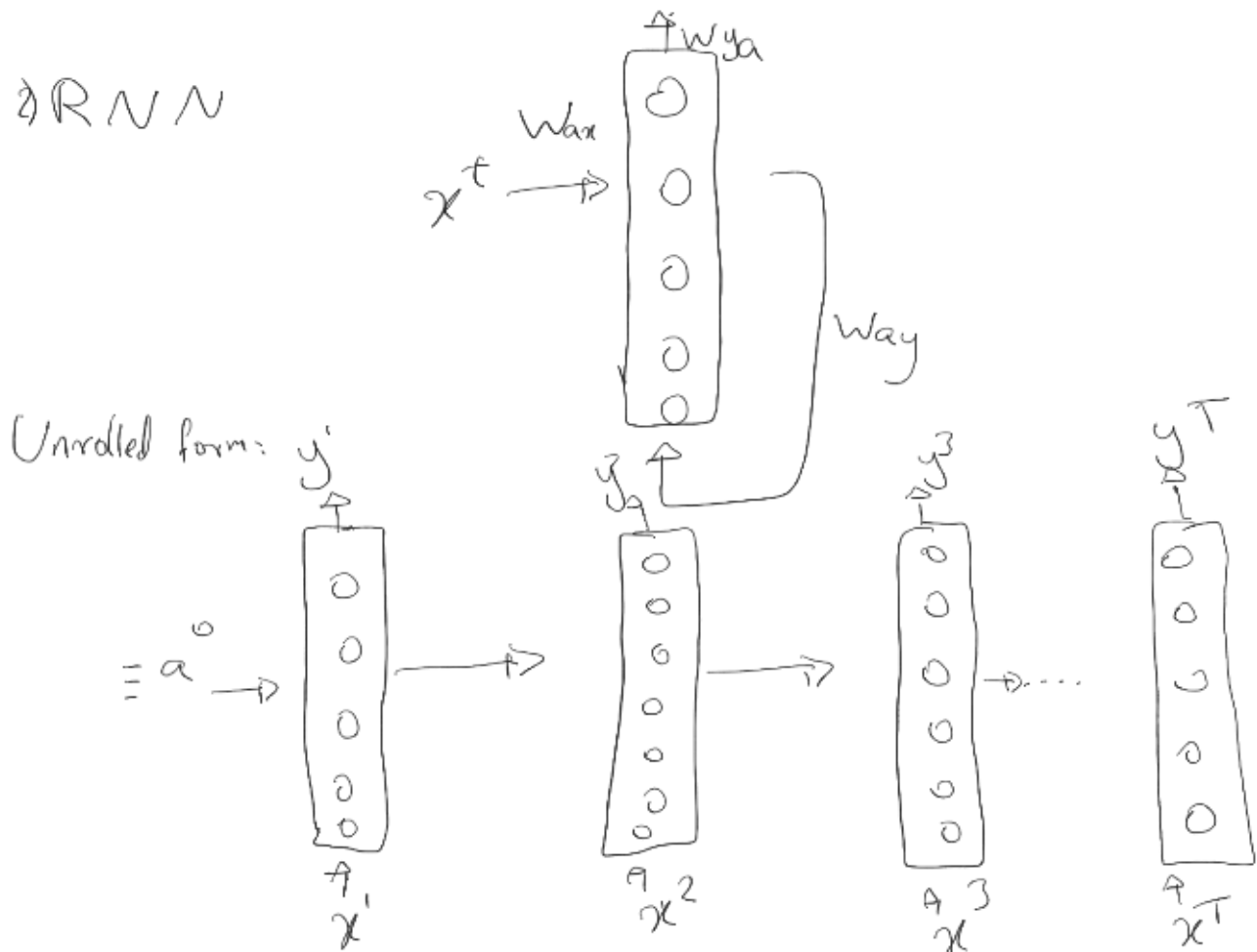
1) what a standard NN does.

for sequences.

1- input length

2- feature sharing.

2 RNN



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

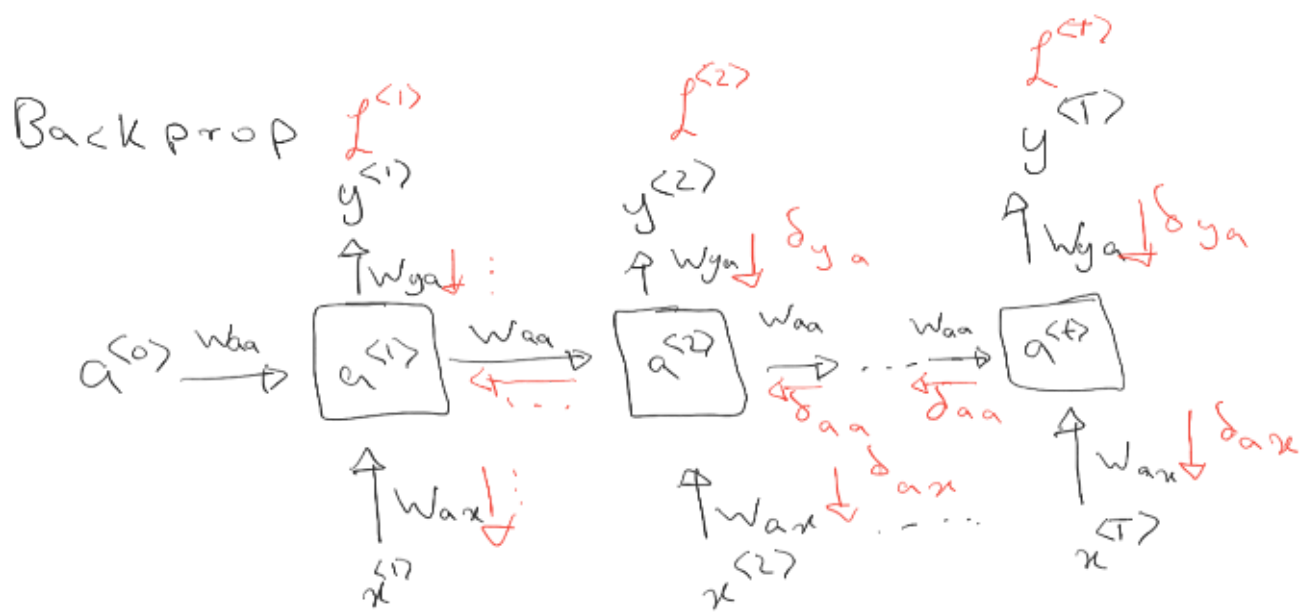
activation func

$$y^{(t)} = g(w_{ya} a^{(t)} + b_y)$$

activation fun

$$a^{(t)} = f\left(w_a \begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix} + b_a\right)$$

where  $W_a = [W_{a1} \dots W_{an}]$

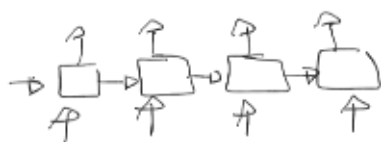


$$L^{(t)} = -y^{(t)} \log \hat{y}^{(t)} - (1 - y^{(t)}) \log (1 - \hat{y}^{(t)})$$

$$L = \sum_{t=1}^T L^{(t)}$$

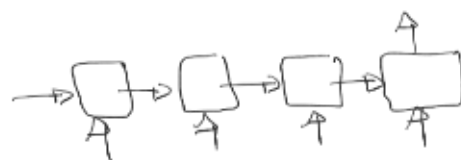
RNN types from input/output lengths

$T_y = T_x = T$   
many-to-many  
"NER, translation?"

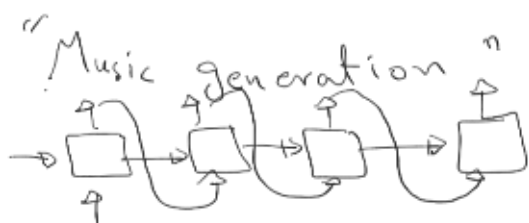


$$T_y = T, T_x = 1$$

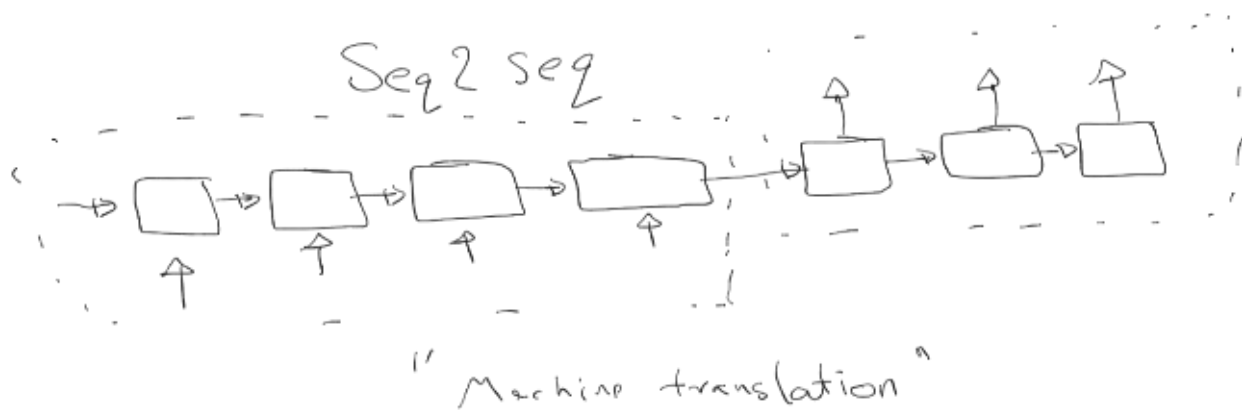
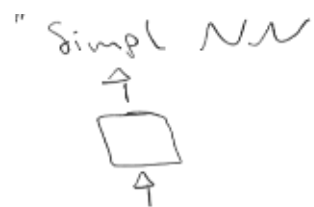
$T_y = 1, T_x = T$   
many-to-one  
"Text classification"



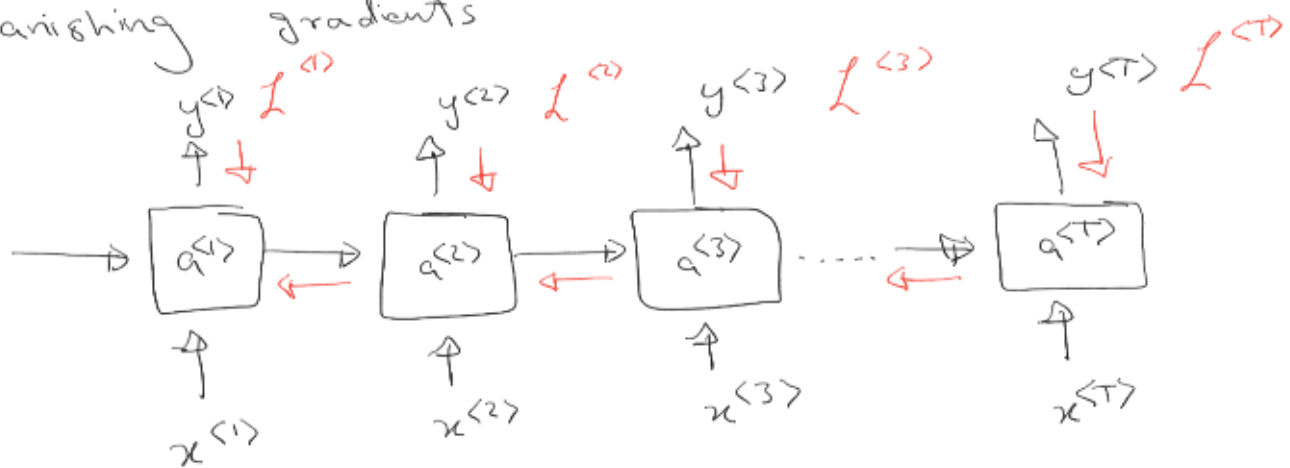
$$T_x = T, T_y = 1$$



$$T_y = T_1, T_x = T_2$$



Vanishing gradients



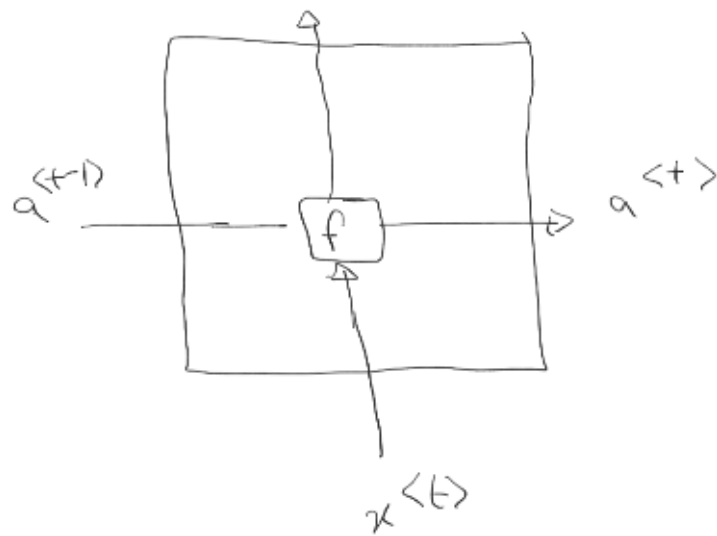
gradient computed at later time steps will have a problem to affect computations in earlier time steps.

GRU:

for RNN:

$$a^{(t)} = f(w_a \begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix} + b_a)$$

$$\hat{y}^{(t)} = g(w_y a^{(t)} + b_y)$$

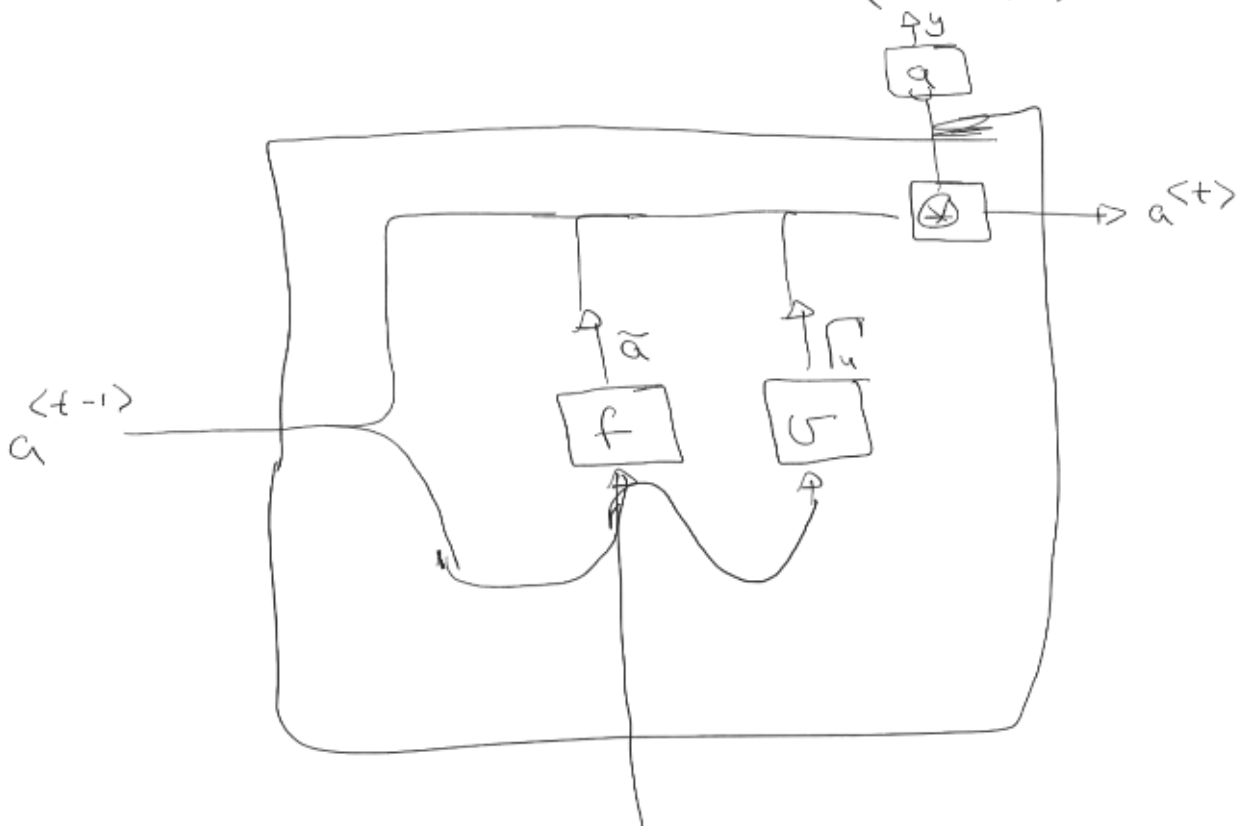


for GRU we have

$$\tilde{a}^{(t)} = f \left( W_a \begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix} + b_a \right)$$

$$\Gamma_u = \overset{\text{sigmoid}}{\sigma} \left( W_u \begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix} + b_u \right)$$

$$a^{(t)} = \Gamma_u \tilde{a}^{(t)} + (1 - \Gamma_u) a^{(t-1)} \quad \oplus$$



$$\frac{1}{x^{(t)}}$$

Standard GRU:

relevance gate.  $\rightarrow$

$$\Gamma_r = \sigma \left( W_r \begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix} + b_r \right)$$

$$\tilde{a}^{(t)} = f \left( W_a \begin{bmatrix} \Gamma_r * a^{(t-1)} \\ x^{(t)} \end{bmatrix} + b_a \right)$$

$$\Gamma_u = \sigma \left( W_u \begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix} + b_u \right)$$

$$a^{(t)} = \Gamma_u \tilde{a}^{(t)} + (1 - \Gamma_u) a^{(t-1)}$$

LSTM:

$$\tilde{c}^{(t)} = f \left( W_c \begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix} + b_c \right)$$

$$\Gamma_u = \sigma \left( W_u \begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix} + b_u \right)$$

$$\Gamma_f = \sigma \left( W_f \begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix} + b_f \right)$$

$$\Gamma_o = \sigma \left( W_o \begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix} + b_o \right)$$

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

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

PC

$$- \quad = \quad 1_0 \quad / \quad \backslash \quad -$$

