

Sequence Model.

2. 10

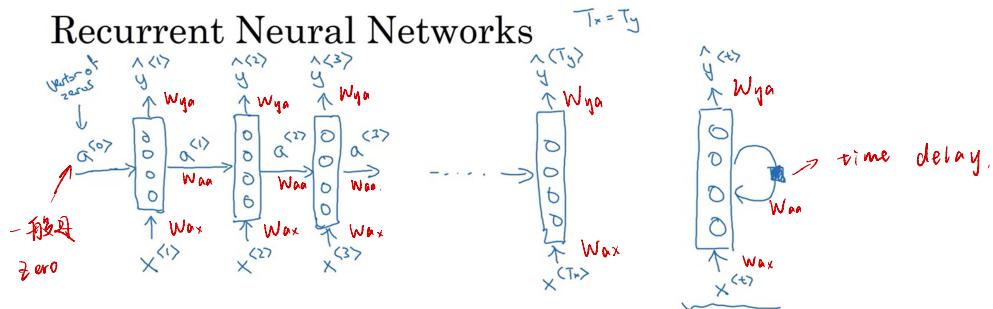
- sequence data: speech recognition
- DNA sequence analysis
- video activity recognition.

Recurrent Neural Network.

why not just use standard Network?

Doesn't share feature learned across different position of text.

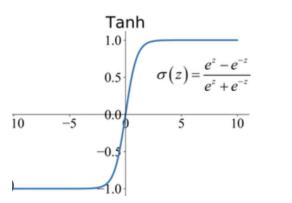
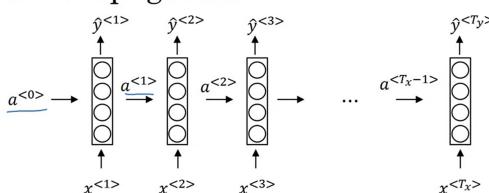
Recurrent Neural Networks



RNN 缺点是仅使用序列中较早信息来进行预测。

这可以用 Bi-directional RNN 解决。

Forward Propagation



activation function
↙ tanh / Relu.
常用

$$a^{<0>} = \vec{0}$$

$$a^{<1>} = g(W_{aa} a^{<0>} + W_{ax} x^{<1>} + b_a)$$

$$\hat{y}^{<1>} = g(W_{ya} a^{<1>} + b_y) \leftarrow \text{= sigmoid, } K\%.$$

$$a^{<2>} = g(W_{aa} a^{<1>} + W_{ax} x^{<2>} + b_a)$$

$$\hat{y}^{<2>} = g(W_{ya} a^{<2>} + b_y)$$

$$a^{(t)} = g \left(W_{aa} \overset{a^{(t-1)}}{\underset{(100, 100)}{\overbrace{a}}} + W_{ax} \overset{x^{(t)}}{\underset{(100, 10000)}{\overbrace{x}}} + b_a \right)$$

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

$$\overset{100}{\downarrow} \left[\overset{W_{aa}}{\overbrace{a^{(t-1)}}} ; \overset{W_{ax}}{\overbrace{x^{(t)}}} \right] = W_a \cdot \underset{(100, 10000)}{\overbrace{(100, 10000)}}$$

$$\left[a^{(t-1)}, x^{(t)} \right] = \left[\begin{array}{c} a^{(t-1)} \\ \hline x^{(t)} \end{array} \right] \underset{10,000}{\uparrow} \underset{10,000}{\downarrow} \underset{10,100}{\uparrow}$$

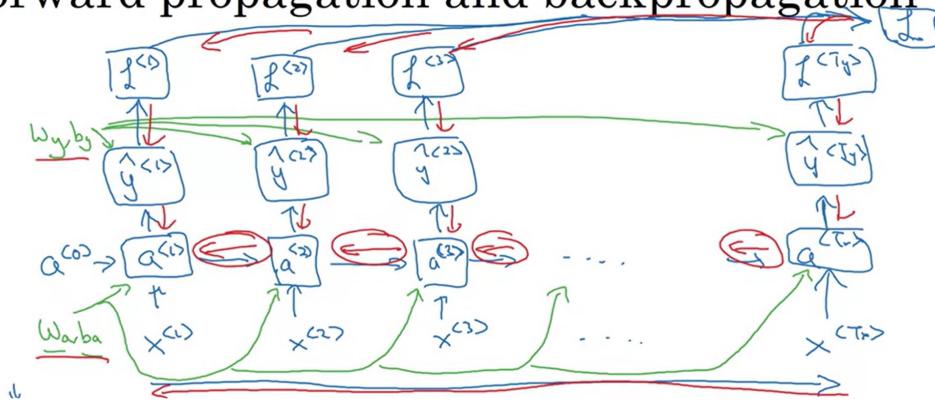
$$\left[W_{aa} \mid W_{ax} \right] \left[\begin{array}{c} a^{(t-1)} \\ \hline x^{(t)} \end{array} \right] = W_{aa} a^{(t-1)} + W_{ax} x^{(t)}$$

△ Back propagation.

$$\sum_{t=1}^{T_y} (\hat{y}^{(t)}, y^{(t)}) = -\hat{y}^{(t)} \log \hat{y}^{(t)} - (1 - \hat{y}^{(t)}) \log (1 - \hat{y}^{(t)})$$

$$J(\hat{y}, y) = \sum_{t=1}^{T_y} J^{(t)}(\hat{y}^{(t)}, y^{(t)})$$

Forward propagation and backpropagation



可以看作从最后一层逐步调整权重参数。

3.1 - Basic RNN Backward Pass

Begin by computing the backward pass for the basic RNN cell. Then, in the following sections, iterate through the cells.

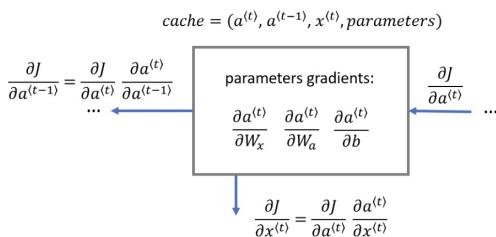


Figure 6: The RNN cell's backward pass. Just like in a fully-connected neural network, the derivative of the cost function J backpropagates through the time steps of the RNN by following the chain rule from calculus. Internal to the cell, the chain rule is also used to calculate $(\frac{\partial J}{\partial W_{ax}}, \frac{\partial J}{\partial W_{aa}}, \frac{\partial J}{\partial b})$ to update the parameters (W_{ax}, W_{aa}, b_a) . The operation can utilize the cached results from the forward path.

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

$$\frac{\partial \tanh(x)}{\partial x} = 1 - \tanh^2(x)$$

element-wise multiplication.

$$d\tanh = d\alpha_{next} \times (1 - \tanh^2(W_{ax}x^{(t)} + W_{aa}a^{(t-1)} + b_a))$$

$$dW_{ax} = d\tanh \cdot x^{(t)T}$$

$$f(x) = \tanh(x)$$

$$f'(x) = 1 - \tanh^2(x)$$

$$dW_{aa} = d\tanh \cdot a^{(t-1)T}$$

指定最后-1

$$db_a = \sum_{\text{batch}} d\tanh$$

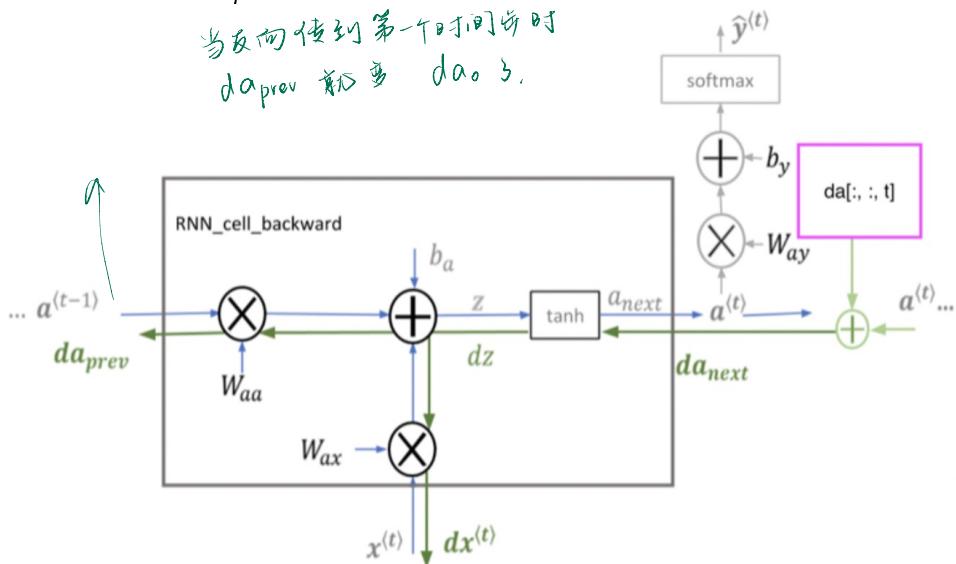
np.sum(dtanh, axis=-1, keepdim=True)

$$dx^{(t)} = W_{ax}^T \cdot d\tanh$$

保留在 dba 和 dtanh 相同
维度数

$$da_{prev} = W_{aa}^T \cdot d\tanh$$

当反向传播到第一个时间步时
da_{prev} 需要 da₀。



* Vector Derivatives,

Ex1: \vec{y} : length C
 W : C row, D column. $\vec{y} = W \vec{x}$
 \vec{x} : length D.

我们要计算.

$$\frac{\partial \vec{y}_3}{\partial \vec{x}_7}$$

$$\vec{y}_3 = \sum_{j=1}^D W_{3,j} \vec{x}_j$$

$$= W_{3,1} \vec{x}_1 + W_{3,2} \vec{x}_2 + \dots + W_{3,D} \vec{x}_D$$

$$\frac{\partial \vec{y}_3}{\partial \vec{x}_7} = \frac{\partial}{\partial \vec{x}_7} (-\dots + W_{3,7} \vec{x}_7 + \dots)$$

$$= 0 + 0 + \dots + \frac{\partial}{\partial \vec{x}_7} (W_{3,7} \vec{x}_7) + \dots + 0$$

$$= W_{3,7}$$

也就是说我们得 $\frac{\partial \vec{y}_i}{\partial \vec{x}_j} = W_{ij}$

所以 $\frac{\partial \vec{y}}{\partial \vec{x}} = W$

Ex2: $\vec{y} = \vec{x} W$ (反过来)

$$\vec{y}_3 = \sum_{i=1}^D \vec{x}_i W_{i,3}$$

$$\frac{\partial \vec{y}_3}{\partial \vec{x}_7} = W_{7,3}$$

$$\frac{\partial \vec{y}}{\partial \vec{x}} = W$$

$$Ex 3: \frac{d\vec{y}}{dW} \quad ?$$

$$\vec{y}_3 = \vec{x}_1 W_{1,3} + \dots + \vec{x}_D W_{D,3}$$

$$\frac{\partial \vec{y}_3}{\partial W_{7,8}} = 0 \quad \text{not in } \vec{w}_3.$$

$$\frac{\partial \vec{y}_3}{\partial W_{2,3}} = \vec{x}_2,$$

$$\therefore \frac{\partial \vec{y}_j}{\partial W_{i,j}} = \vec{x}_i$$

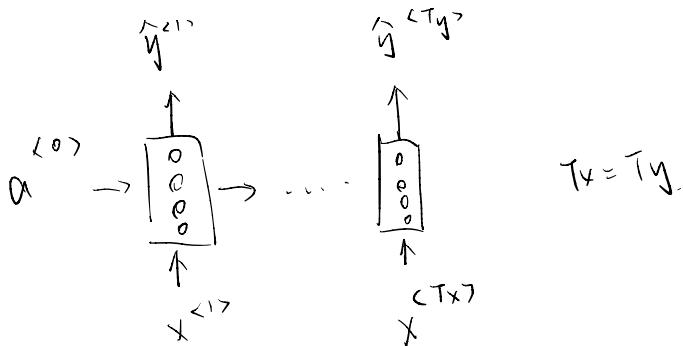
$$\text{if it is 3-D}, \quad F_{i,j,k} = \frac{\partial \vec{y}_i}{\partial W_{j,k}}$$

$$\text{then, } F_{i,j,k} = \vec{x}_j$$

for more detail it's in **vecDerivs.pdf**.

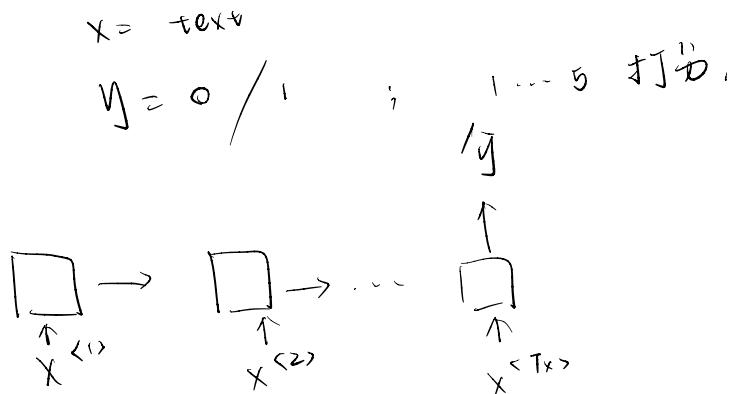
△ Architecture

- Many to Many



- Many to one

Ex: Sentiment classification

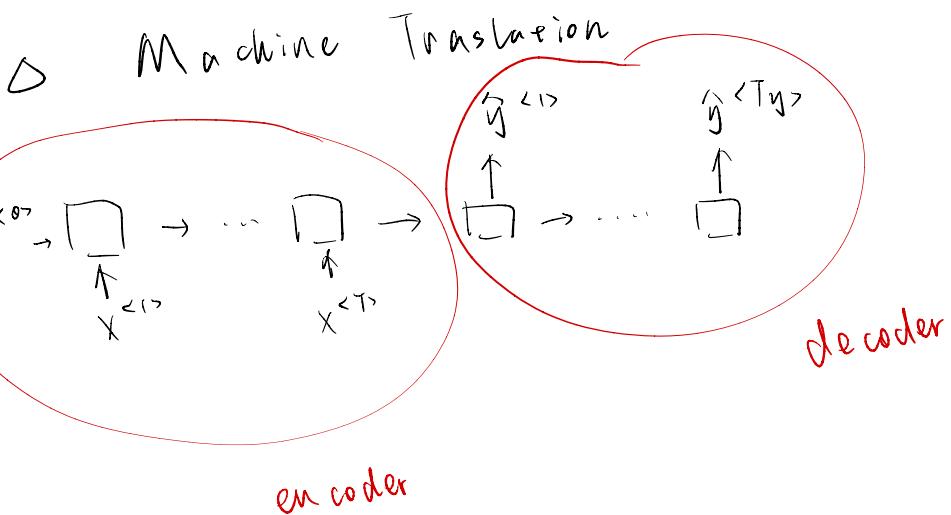
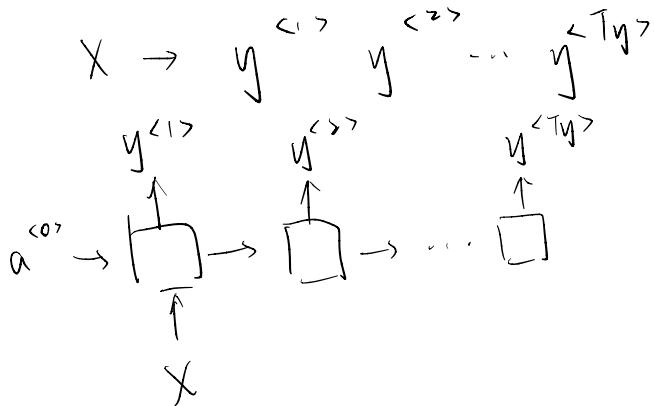


- One to One



- One to Many

Music Generation.



• Language modelling with RNN.

Training set: large corpus of english text.

Tokenize

Cats average 15 hours of sleep a day. $\downarrow <\text{EOS}>$

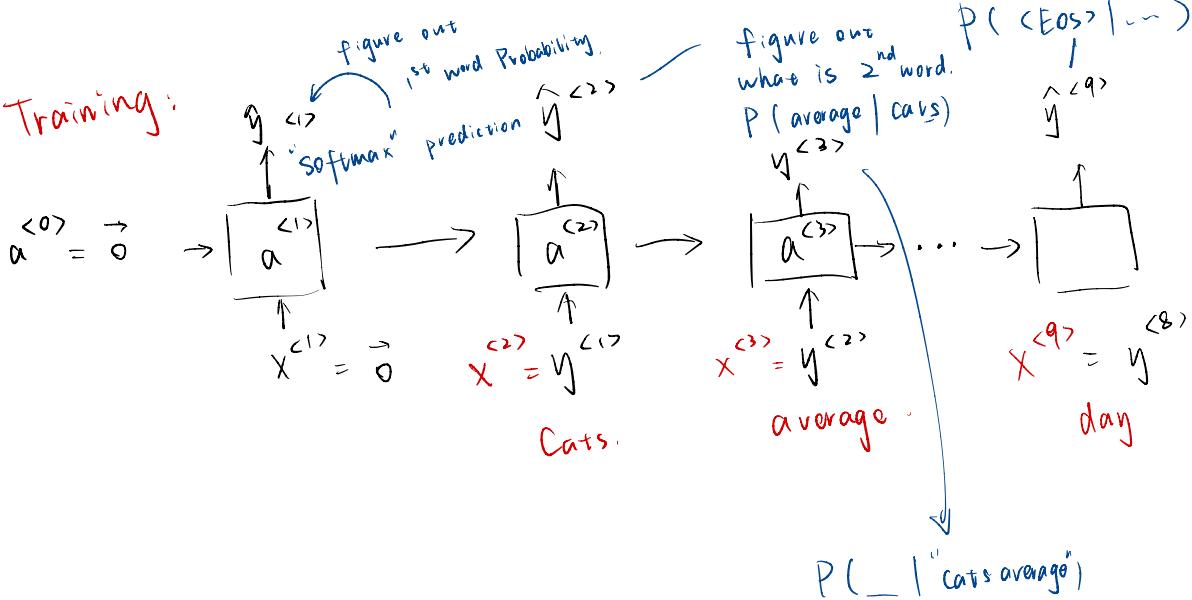
$$\begin{matrix} y^{<1>} & y^{<2>} & y^{<3>} & \dots & y^{<n>} & y^{<n+1>} \\ x^{<1>} = y^{<1>} \end{matrix}$$

The Egyptian ~~Mau~~ is a bread of cat. $<\text{EOS}>$

10,000

Ex: Cats average 15 hours of sleep a day. $<\text{EOS}>$.

Training:

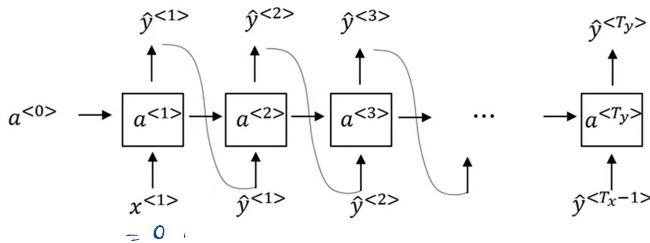


$$P(y^{<1>}, y^{<2>}, y^{<3>})$$

$$= P(y^{<1>}) \cdot P(y^{<2>} | y^{<1>}) \cdot P(y^{<3>} | y^{<1>}, y^{<2>})$$

Sampling:

<EOS>



First word randomly sample according to softmax distribution:
n.p. random choice to sample according to distribution defined by vector probability.

- **Step 1:** Input the "dummy" vector of zeros $x^{(1)} = \vec{0}$.
 - This is the default input before you've generated any characters. You also set $a^{(0)} = \vec{0}$

- **Step 2:** Run one step of forward propagation to get $a^{(1)}$ and $\hat{y}^{(1)}$. Here are the equations:

hidden state:

$$a^{(t+1)} = \tanh(W_{ax}x^{(t+1)} + W_{aa}a^{(t)} + b) \quad (1)$$

activation:

$$z^{(t+1)} = W_{ya}a^{(t+1)} + b_y \quad (2)$$

prediction:

$$\hat{y}^{(t+1)} = \text{softmax}(z^{(t+1)}) \quad (3)$$

- Details about $\hat{y}^{(t+1)}$:

- Note that $\hat{y}^{(t+1)}$ is a (softmax) probability vector (its entries are between 0 and 1 and sum to 1).
- $\hat{y}_i^{(t+1)}$ represents the probability that the character indexed by " i " is the next character.
- A `softmax()` function is provided for you to use.

Vanishing gradients with RNNs

The **cat**, which ate fish, ... , **was** full
 The **Cats**, which ate fish, ... , **were** full

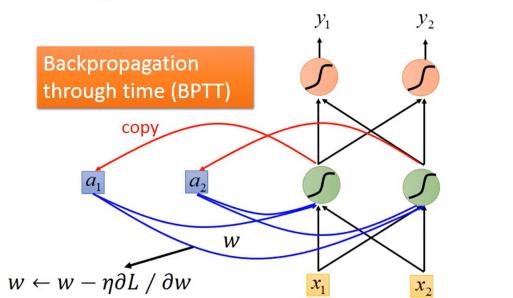
language 非常长的依赖性，较早的措辞会影响后面的内容。

RNN 不善于 long, 因为训练 very deep neural network

会出现 vanishing gradients problem (梯度消失)
 背景中的梯度将很难传播回来影响早期层的权重。

$\Rightarrow \hat{y}^{(t)} \text{ 只受最近 } n \text{ 个 } x^{(t)}$

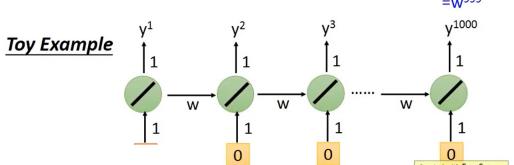
Learning



Exploding gradients. NaN.

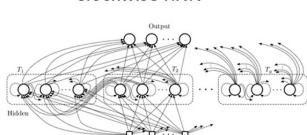
Why?

$w = 1$	$y^{1000} = 1$	Large $\frac{\partial L}{\partial w}$	Small Learning rate?
$w = 1.01$	$y^{1000} \approx 20000$	small $\frac{\partial L}{\partial w}$	Large Learning rate?
$w = 0.99$	$y^{1000} \approx 0$		
$w = 0.01$	$y^{1000} \approx 0$		



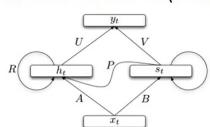
Helpful Techniques

Clockwise RNN



[Jan Koutnik, JMLR'14]

Structurally Constrained Recurrent Network (SCRN)



[Tomas Mikolov, ICLR'15]

Vanilla RNN Initialized with Identity matrix + ReLU activation function [Quoc V. Le, arXiv'15]

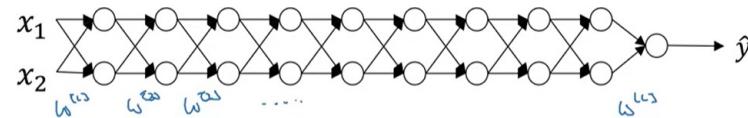
➤ Outperform or be comparable with LSTM in 4 different tasks

[Created with EverCam.]

六:

Vanishing / Exploding Gradients

Normal Network 为什么:



假设我们使用 $g(z) = z$ $b^{[l]} = 0$.

$$y = W^{[1]} W^{[2]} \dots W^{[L-1]} X$$

$\underbrace{W^{[1]} W^{[2]} \dots W^{[L-1]}}_{Z^{[l]}}$ X

$$Z^{[l]} = W^{[l]} X$$

\exp vs \tanh

$$\text{假设 } W^{[0]} = \begin{bmatrix} 0.9 & 0 \\ 0 & 0.5 \end{bmatrix}$$

$$\hat{y} = W^{[L]} \begin{bmatrix} 0.5 & 0 \\ 0 & 0.6 \end{bmatrix}^{-1} X \cdot \begin{bmatrix} 0.5 \\ 0.6 \end{bmatrix}$$

$a^{[l]}$

$$a^{[l]} = g(Z^{[l]}) = Z^{[l]}$$

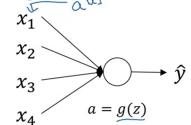
$$a^{[l+1]} = g(W^{[l+1]} a^{[l]})$$

$W^{[l]} >]$ explode big

$W^{[l]} <]$ small.

Weight Initialization for Deep Network

Single neuron example



$$z = w_1 x_1 + w_2 x_2 + \dots + w_n x_n$$

large $n \rightarrow$ smaller w_i

$$\text{Var}(w_i) = \frac{1}{n}$$

$$W^{[l]} = \text{np.random.randn(shape)} \times \text{np.sqrt}\left(\frac{1}{n^{[l-1]}}\right)$$

$$\text{if use ReLU} \quad \text{Var} = \frac{2}{n}$$

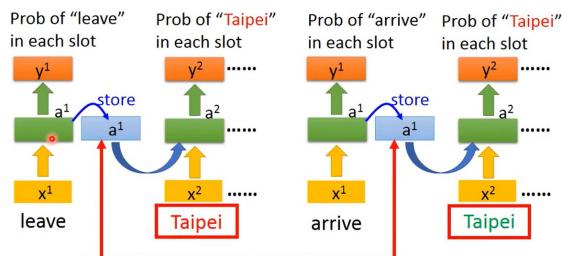
Xavier Initialization

$$\sqrt{\frac{1}{n^{[l-1]}}} < \sqrt{\frac{2}{n^{[l-1]} \times n^{[l]}}}$$

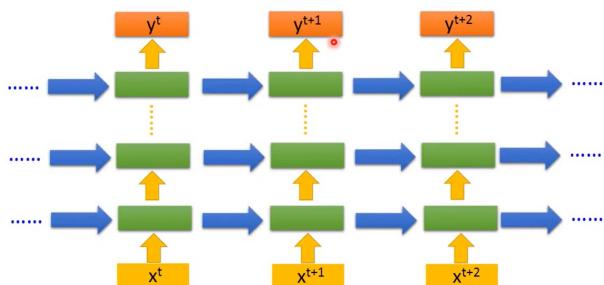
Yoshua Bengio.

常見基底 RNN .

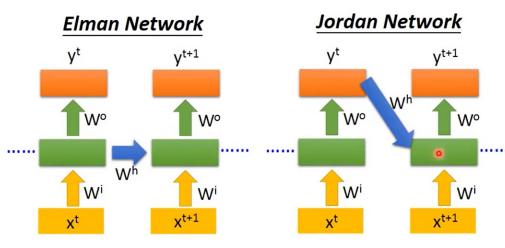
RNN



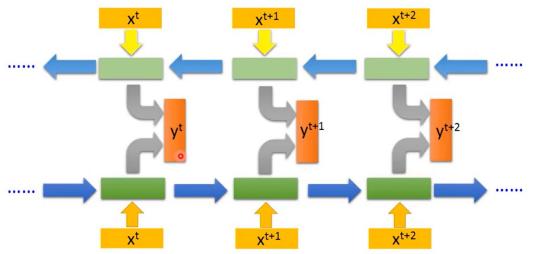
Of course it can be deep ...



Elman Network & Jordan Network

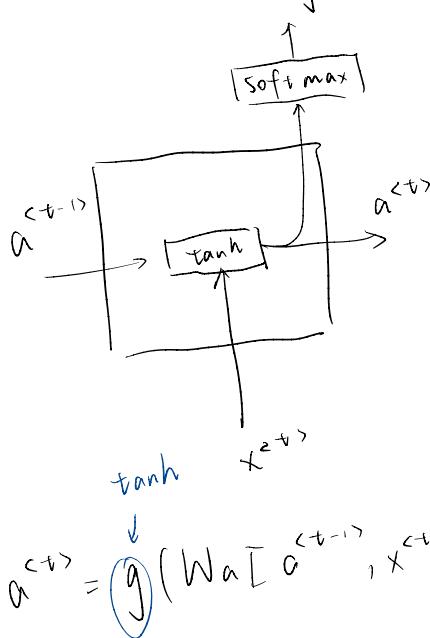


Bidirectional RNN

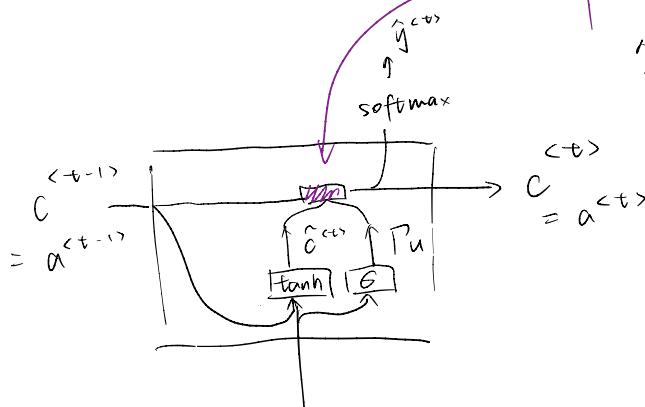


• GRU (Gated Recurrent Unit).

RNN Unit



$$a^{(t)} = g (W_a [a^{(t-1)}, x^{(t)}] + b_a)$$



GRU (Simplify)

• 一个新变量 "c" = Memory cell
用来记忆且 can ~~还~~ cats.

$$c^{(t)} = a^{(t)}$$

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

$$\text{Gate: } \tilde{\Gamma}_u = \sigma (W_u [c^{(t-1)}, x^{(t)}] + b_u)$$

$\tilde{\Gamma}_u$ \downarrow
sigmoid function.
 \downarrow

$\tilde{\Gamma}_u$ 0 或 1.

$$c^{(t)} = \tilde{\Gamma}_u * \tilde{c}^{(t)} + (1 - \tilde{\Gamma}_u) * c^{(t-1)}$$

等于更新, 等于保留旧值.

element wise

Full GRU

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

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

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

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

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

$$L \leq T M.$$

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

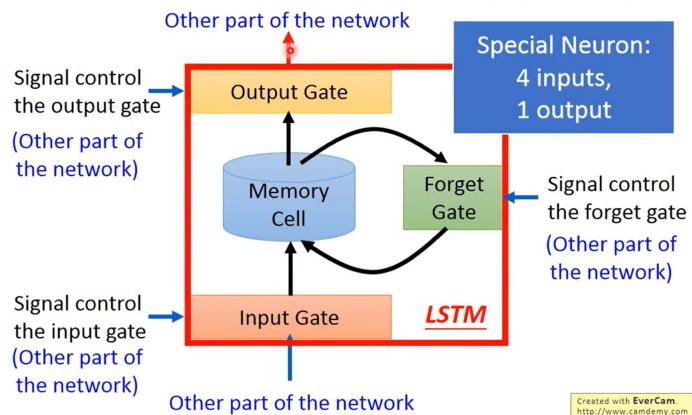
update $\tilde{r}^{(t)} = \sigma(W_r [a^{(t-1)}, x^{(t)}] + b_r)$

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

output $\tilde{c}^{(t)} = \tilde{r}^{(t)} \times \tilde{c}^{(t-1)} + \tilde{f}^{(t)} \times c^{(t-1)}$

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

Long Short-term Memory (LSTM)



但两者各有优缺点

LSTM Back Prop

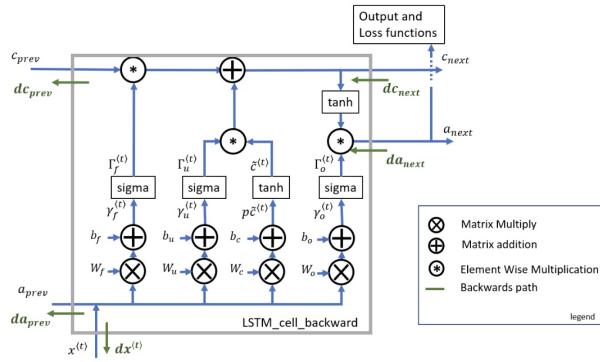


Figure 8: LSTM Cell Backward. Note the output functions, while part of the 'lstm_cell_forward', are not included in 'lstm_cell_backward'

The equations for the LSTM backward pass are provided below. (If you enjoy calculus exercises feel free to try deriving these from scratch yourself.)

2. Gate Derivatives

Note the location of the gate derivatives ($\gamma_{..}$) between the dense layer and the activation function (see graphic above). This is convenient for computing parameter derivatives in the next step.

$$d\gamma_o^{(t)} = da_{next} * \tanh(c_{next}) * \Gamma_o^{(t)} * (1 - \Gamma_o^{(t)}) \quad (7)$$

$$d\tilde{c}^{(t)} = (dc_{next} * \Gamma_u^{(t)} + \Gamma_o^{(t)} * (1 - \tanh^2(c_{next})) * \Gamma_u^{(t)} * da_{next}) * (1 - (\tilde{c}^{(t)})^2) \quad (8)$$

$$d\gamma_u^{(t)} = (dc_{next} * \tilde{c}^{(t)} + \Gamma_o^{(t)} * (1 - \tanh^2(c_{next})) * \tilde{c}^{(t)} * da_{next}) * \Gamma_u^{(t)} * (1 - \Gamma_u^{(t)}) \quad (9)$$

$$d\gamma_f^{(t)} = (dc_{next} * c_{prev} + \Gamma_o^{(t)} * (1 - \tanh^2(c_{next})) * c_{prev} * da_{next}) * \Gamma_f^{(t)} * (1 - \Gamma_f^{(t)}) \quad (10)$$

3. Parameter Derivatives.

$$dW_f = d\gamma_f^{(t)} \begin{bmatrix} a_{prev} \\ x_t \end{bmatrix}^\top \quad db_f = \sum_{batch} d\gamma_f^{(t)}$$

$$dW_u = d\gamma_u^{(t)} \begin{bmatrix} a_{prev} \\ x_t \end{bmatrix}^\top \quad db_u = \sum_{batch} d\gamma_u^{(t)}$$

$$dW_c = d\tilde{c}^{(t)} \begin{bmatrix} a_{prev} \\ x_t \end{bmatrix}^\top \quad db_c = \sum_{batch} d\tilde{c}^{(t)}$$

$$dW_o = d\gamma_o^{(t)} \begin{bmatrix} a_{prev} \\ x_t \end{bmatrix}^\top \quad db_o = \sum_{batch} d\gamma_o^{(t)}$$

Finally, you will compute the derivative with respect to the previous hidden state, previous memory state, and input.

$$da_{prev} = W_f^T d\gamma_f^{(t)} + W_u^T d\gamma_u^{(t)} + W_c^T d\tilde{pc}^{(t)} + W_o^T d\gamma_o^{(t)} \quad (19)$$

Here, to account for concatenation, the weights for equations 19 are the first n_a , (i.e. $W_f = W_f[:, : n_a]$ etc...)

$$dc_{prev} = dc_{next} * \Gamma_f^{(t)} + \Gamma_o^{(t)} * (1 - \tanh^2(c_{next})) * \Gamma_f^{(t)} * da_{next} \quad (20)$$

$$dx^{(t)} = W_f^T d\gamma_f^{(t)} + W_u^T d\gamma_u^{(t)} + W_c^T d\tilde{pc}^{(t)} + W_o^T d\gamma_o^{(t)} \quad (21)$$

where the weights for equation 21 are from n_a to the end, (i.e. $W_f = W_f[:, n_a :]$ etc...)

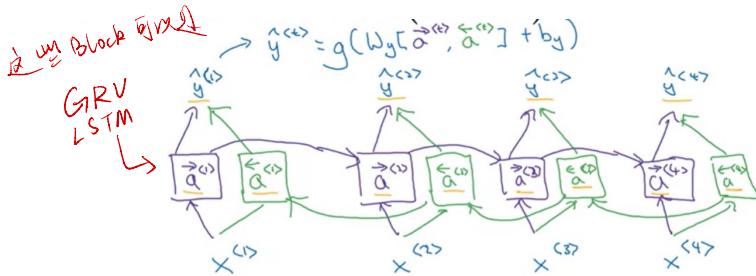
取除 n_a 以外

的 β

只考虑与 a_{prev}
相关的列
取前 $n-a$ 列

问题
 $W_u \quad W_o \quad W_o$

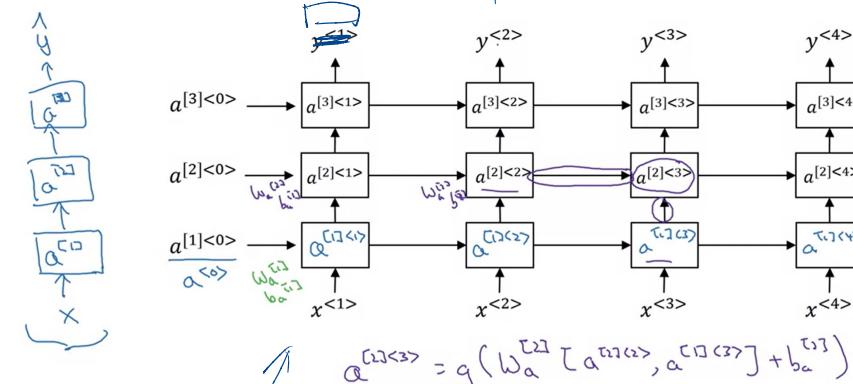
• Bidirectional RNN (BRNN)



Acyclic graph

• Deep RNN

Deep RNN example 可以这样级
深有水平.



Andrew Ng

RNN
GRU
LSTM

ConvLSTM

ConvLSTM is a type of recurrent neural network for spatio-temporal prediction that has convolutional structures in both the input-to-state and state-to-state transitions. The ConvLSTM determines the future state of a certain cell in the grid by the inputs and past states of its local neighbors. This can easily be achieved by using a convolution operator in the state-to-state and input-to-state transitions (see Figure). The key equations of ConvLSTM are shown below, where $*$ denotes the convolution operator and \odot the Hadamard product:

$$i_t = \sigma(W_{xi} * X_t + W_{hi} * H_{t-1} + W_{ci} \odot C_{t-1} + b_i)$$

$$f_t = \sigma(W_{xf} * X_t + W_{hf} * H_{t-1} + W_{cf} \odot C_{t-1} + b_f)$$

$$C_t = f_t \odot C_{t-1} + i_t \odot \tanh(W_{xc} * X_t + W_{hc} * H_{t-1} + b_c)$$

$$o_t = \sigma(W_{xo} * X_t + W_{ho} * H_{t-1} + W_{co} \odot C_t + b_o)$$

$$H_t = o_t \odot \tanh(C_t)$$

$$\begin{aligned} A \odot B \\ \begin{bmatrix} a & b \\ c & d \end{bmatrix} \odot \begin{bmatrix} h & i \\ j & k \end{bmatrix} \\ = \begin{bmatrix} a \times h & b \times i \\ c \times j & d \times k \end{bmatrix} \end{aligned}$$

Next-Frame Video Prediction
with convLSTM.
[https://keras.io/examples/vision/
conv_lstm/](https://keras.io/examples/vision/conv_lstm/)

A \times B convolution operator
(fft + convolve)

$$(1, 2, 3) * (4, 5, 6) = (4, 13, 28, 27, 18)$$

$$(1 + 2x + 3x^2)(4 + 5x + 6x^2) = 4 + 13x + 28x^2 + 27x^3 + 18x^4$$

$$5 + 8 = 13.$$

	1	$2x$	$3x^2$
4	4	8	12
5	5	10	15
6	6	12	18