

深度学习

2019年3月15日 18:12

查准率 (precision) 和查全率

F_1 分数的定义是这个公式: $\frac{2}{\frac{1}{P} + \frac{1}{R}}$

端到端的深度学习

end-to-end deep learning

Train/dev/test distributions

让你的开发集和测试集来自同一分布

不是所有程序都可以用CUDA加速的, 一般的python程序就不可以, 像tf、pt的CUDA加速修改版可以, 专门的CUDA编程可以, 能快很多。

CUDA_VISIBLE_DEVICES=1 python xxx.py

指定GPU进行CUDA加速

注意: 要大写CUDA_VISIBLE_DEVICES

动量梯度下降法

Momentum

RMSprop

$$S_{dW} = \beta S_{dW} + (1 - \beta) dW^2$$

$$S_{db} = \beta S_{db} + (1 - \beta) db^2$$

$$W := W - \alpha \frac{dW}{\sqrt{S_{dW}}}, \quad b := b - \alpha \frac{db}{\sqrt{S_{db}}}$$

On iteration t:

Compute dW, db on the current mini-batch

$$v_{dW} = \beta v_{dW} + (1 - \beta) dW$$

$$v_{db} = \beta v_{db} + (1 - \beta) db$$

$$W = W - \alpha v_{dW}, \quad b = b - \alpha v_{db}$$

Adam optimization algorithm

$$V_{dW} = 0, S_{dW} = 0, V_{db} = 0, S_{db} = 0$$

On iteration t:

Compute dW, db using current mini-batch

$$V_{dW} = \beta_1 V_{dW} + (1 - \beta_1) dW, \quad V_{db} = \beta_1 V_{db} + (1 - \beta_1) db \leftarrow \text{"momentum"} \beta_1$$

$$S_{dW} = \beta_2 S_{dW} + (1 - \beta_2) dW^2, \quad S_{db} = \beta_2 S_{db} + (1 - \beta_2) db^2 \leftarrow \text{"RMSprop"} \beta_2$$

$$V_{dW}^{\text{corrected}} = V_{dW} / (1 - \beta_1^t), \quad V_{db}^{\text{corrected}} = V_{db} / (1 - \beta_1^t)$$

$$S_{dW}^{\text{corrected}} = S_{dW} / (1 - \beta_2^t), \quad S_{db}^{\text{corrected}} = S_{db} / (1 - \beta_2^t)$$

$$W := W - \alpha \frac{V_{dW}^{\text{corrected}}}{\sqrt{S_{dW}^{\text{corrected}} + \epsilon}}, \quad b := b - \alpha \frac{V_{db}^{\text{corrected}}}{\sqrt{S_{db}^{\text{corrected}} + \epsilon}}$$

Adam:

Momentum + RMSprop

使用 Adam 算法, 首先你要初始化, $v_{dW} = 0, S_{dW} = 0, v_{db} = 0, S_{db} = 0$, 在第 t 次迭

Hyperparameters choice:

$\rightarrow \alpha$: needs to be tune

$\rightarrow \beta_1$: 0.9 $\rightarrow (dW)$

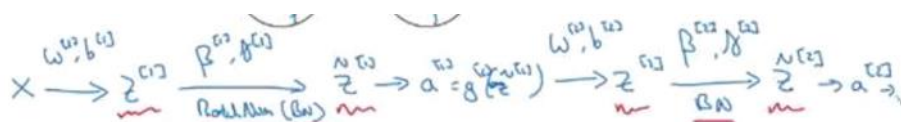
$\rightarrow \beta_2$: 0.999 $\rightarrow (dW^2)$

$\rightarrow \epsilon$: 10^{-8}

Adam: Adaptive moment estimation

Batch Norm是将 z 归一化

$$z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$



When transfer learning makes sense 云课堂

Task from A \rightarrow B

迁移学习 (transfer learning)

- : 1.预训练
- 2.微调

- Task A and B have the same input x .
- You have a lot more data for Task A than Task B.
- Low level features from A could be helpful for learning B.

多任务学习

端到端学习

卷积运算: 边缘检测

Padding (图像四周填充)

Same 卷积, 那意味你填充后, 你的输出大小和输入大小是一样的

Valid 卷积意味着不填充

卷积步长 (stride): 卷积框移动的长度

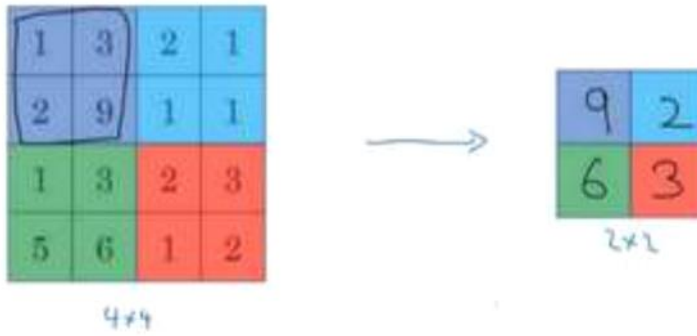
如果你用一个 $f \times f$ 的过滤器卷积一个 $n \times n$ 的图像,

你的 **padding** 为 p , 步幅为 s ,

输出于是变为 $\frac{n+2p-f}{s} + 1$

池化层 (Pooling layers)

Max pooling
超参数 $f=2$ $s=2$

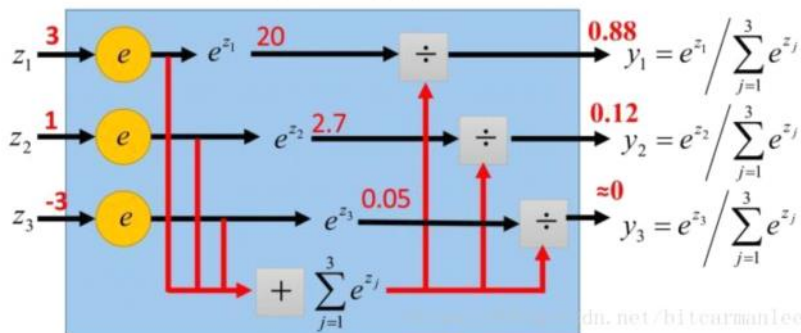


- Softmax layer as the output layer

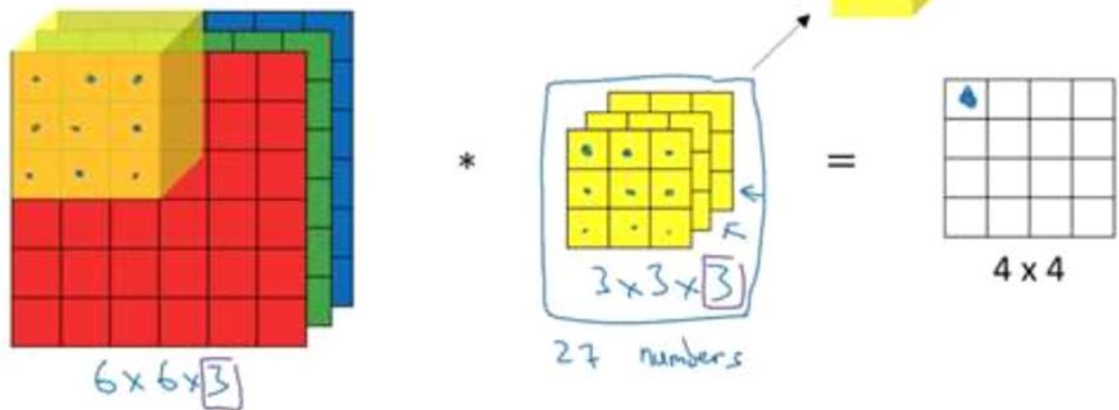
Probability:

- $1 > y_i > 0$
- $\sum_i y_i = 1$

Softmax Layer



Convolutions on RGB image



Summary of notation

If layer l is a convolution layer:

$f^{[l]}$ = filter size

$p^{[l]}$ = padding

$s^{[l]}$ = stride

Input: $n_H^{[l-1]} \times n_W^{[l-1]} \times n_C^{[l-1]} \leftarrow$

Output: $n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$

$$n_{HW}^{[l]} = \left\lfloor \frac{n_{HW}^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1 \right\rfloor$$

Outline

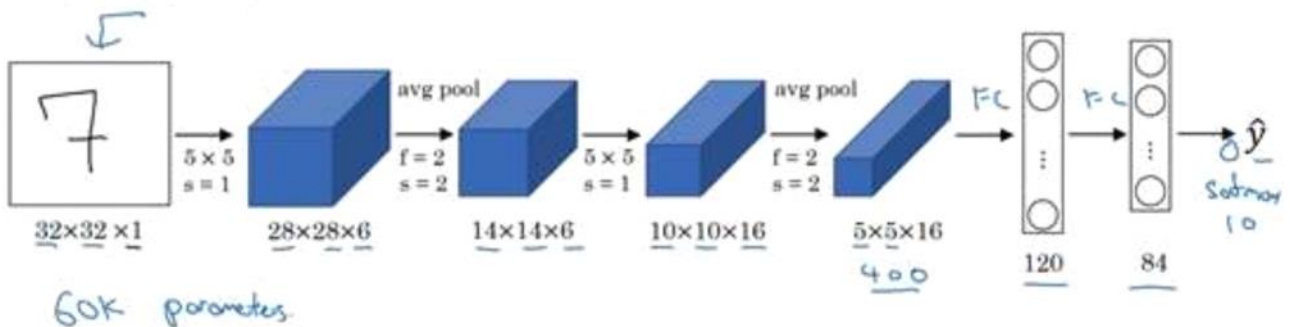
Classic networks:

- LeNet-5 \leftarrow
- AlexNet \leftarrow
- VGG \leftarrow

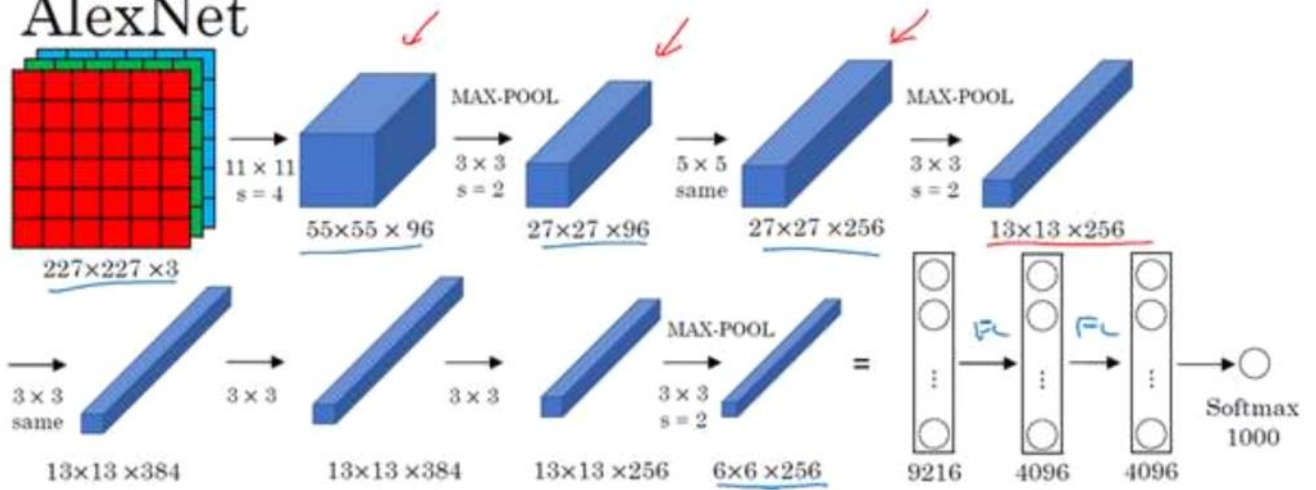
ResNet (152)

Inception

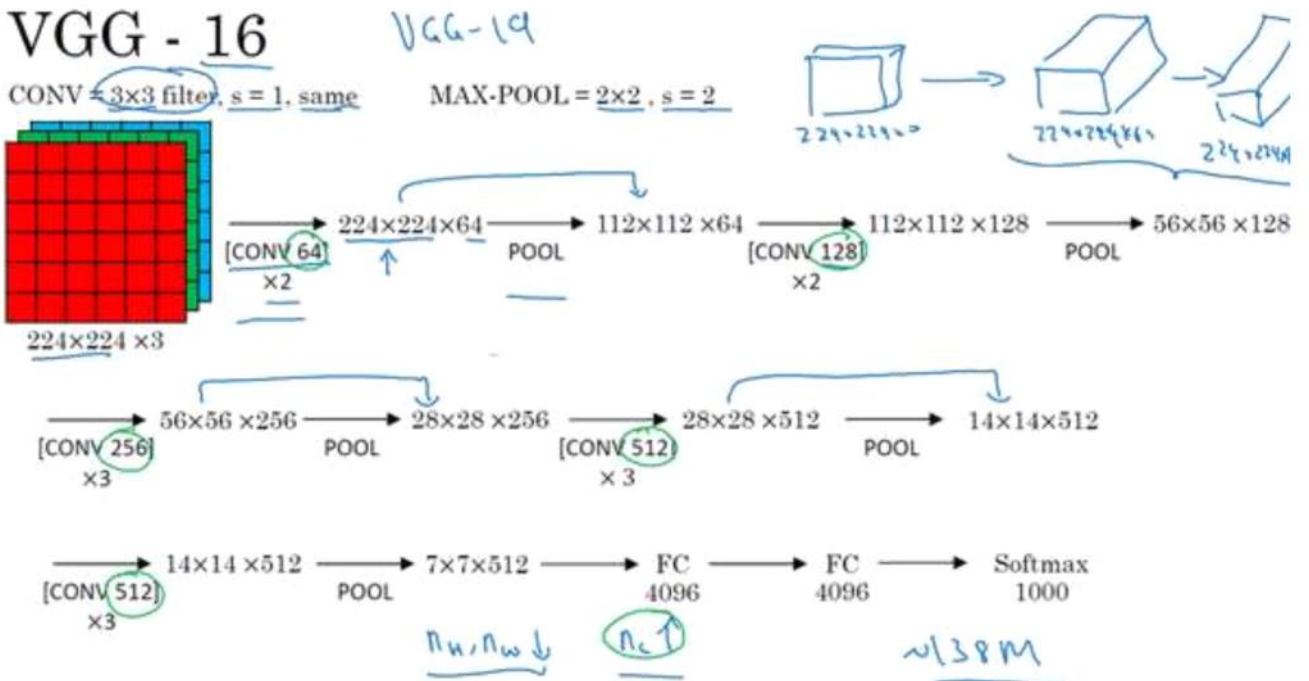
LeNet - 5



AlexNet

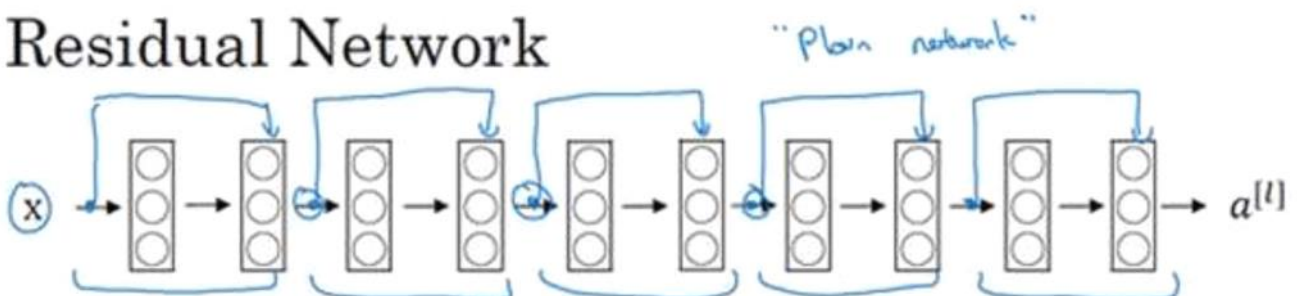


VGG - 16

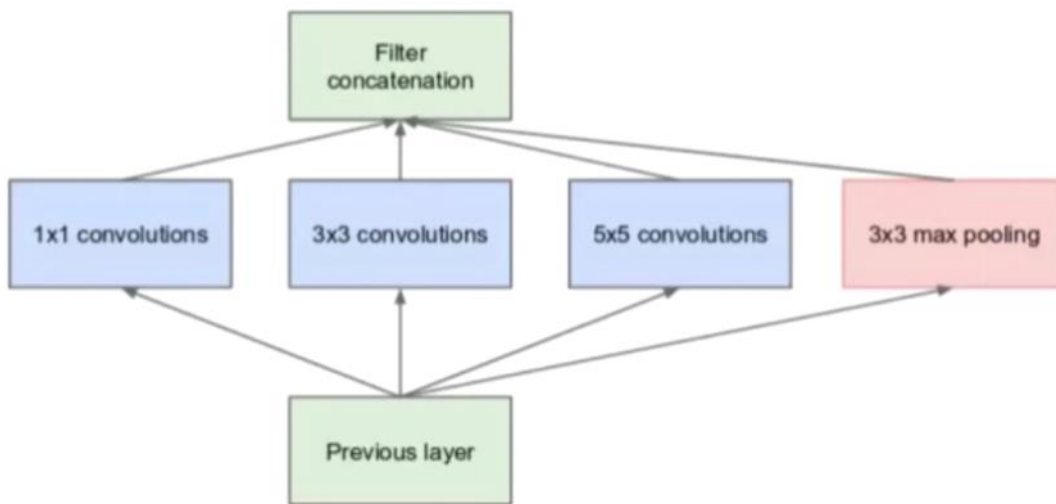


残差网络 (ResNets) : 跳跃连接 (Skip connection)

Residual Network

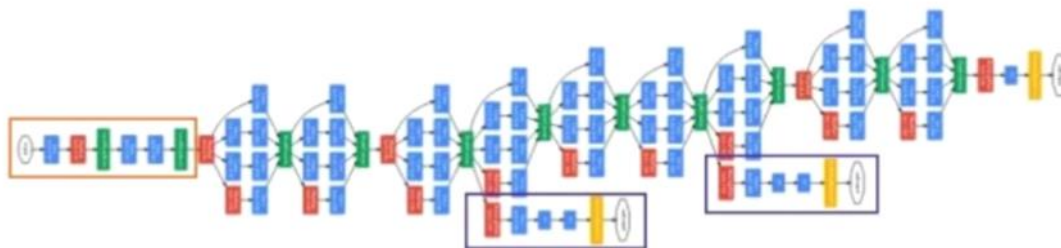


Inception 网络



(a) Inception module, naïve version

googLeNet是由Inception模块堆积而成:



使用开源的实现方案

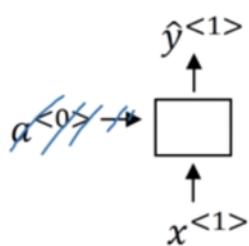
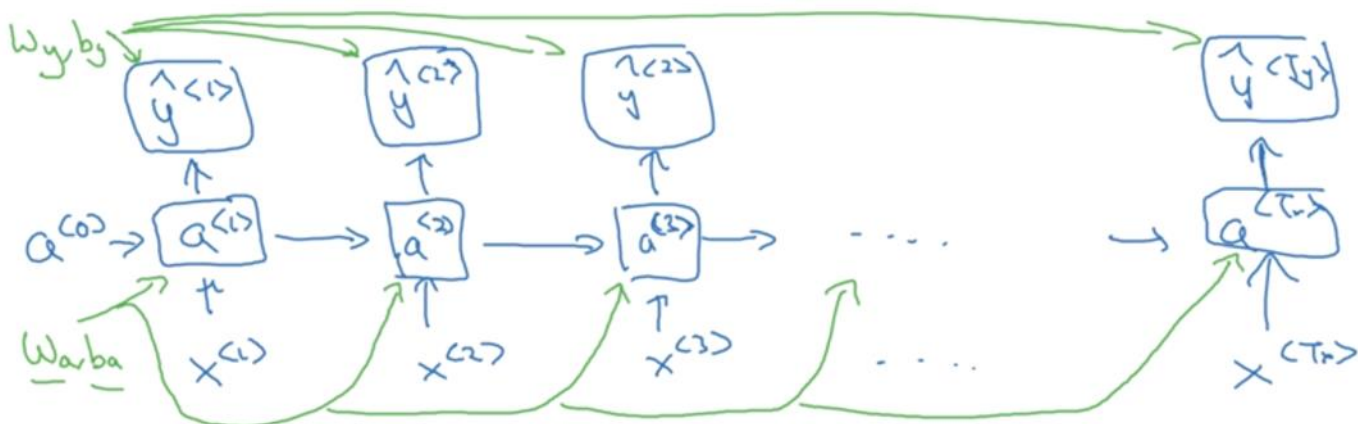
迁移学习:

如果你有越来越多的数据, 你需要冻结的层数越少, 你能够训练的层数就越多。

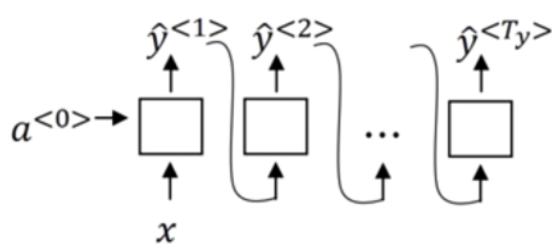
序列模型 (Sequence Models)

RNN

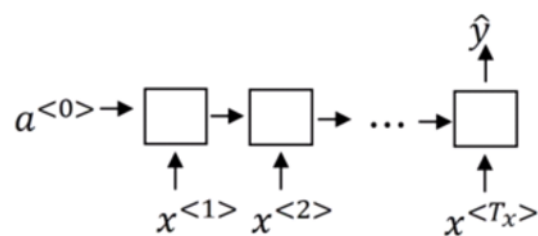




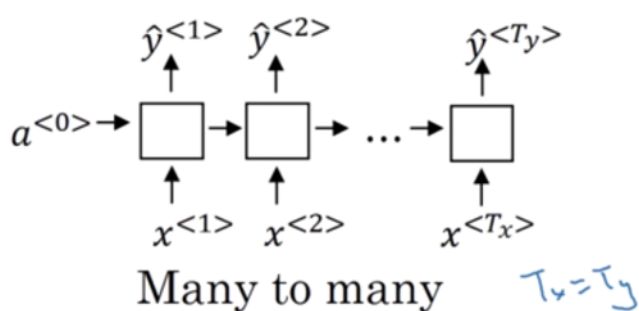
One to one



One to many

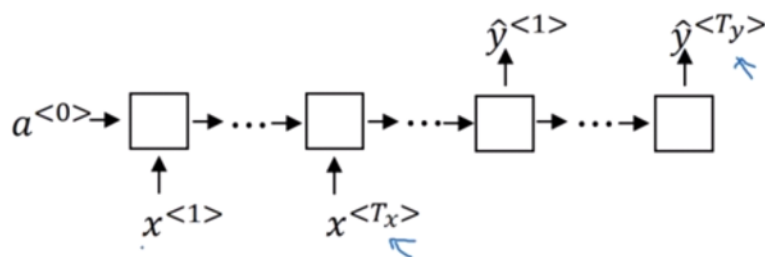


Many to one



Many to many

$T_x = T_y$



Many to many

GRU and LSTM

1 GRU

$$\tilde{c}^{<t>} = \tanh(W_c[\Gamma_r * \underline{c}^{<t-1>}, x^{<t>}] + b_c) \quad \textcircled{3} \quad \tilde{c}^{<t>} = \tanh(W_c[a^{<t-1>}, x^{<t>}] + b_c)$$

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

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

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

$$\text{(forget)} \quad \Gamma_r = \sigma(W_r[a^{<t-1>}, x^{<t>}] + b_r)$$

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

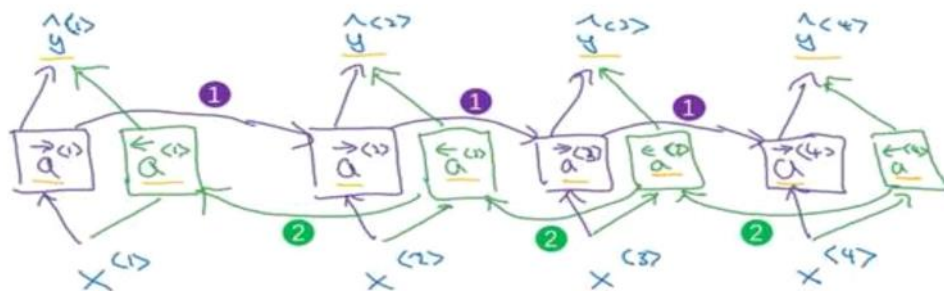
$$\text{(output)} \quad \Gamma_o = \sigma(W_o[a^{<t-1>}, x^{<t>}] + b_o)$$

$$\underline{a}^{<t>} = \underline{c}^{<t>}$$

$$\textcircled{10} \quad \underline{c}^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \Gamma_r * \underline{c}^{<t-1>}$$

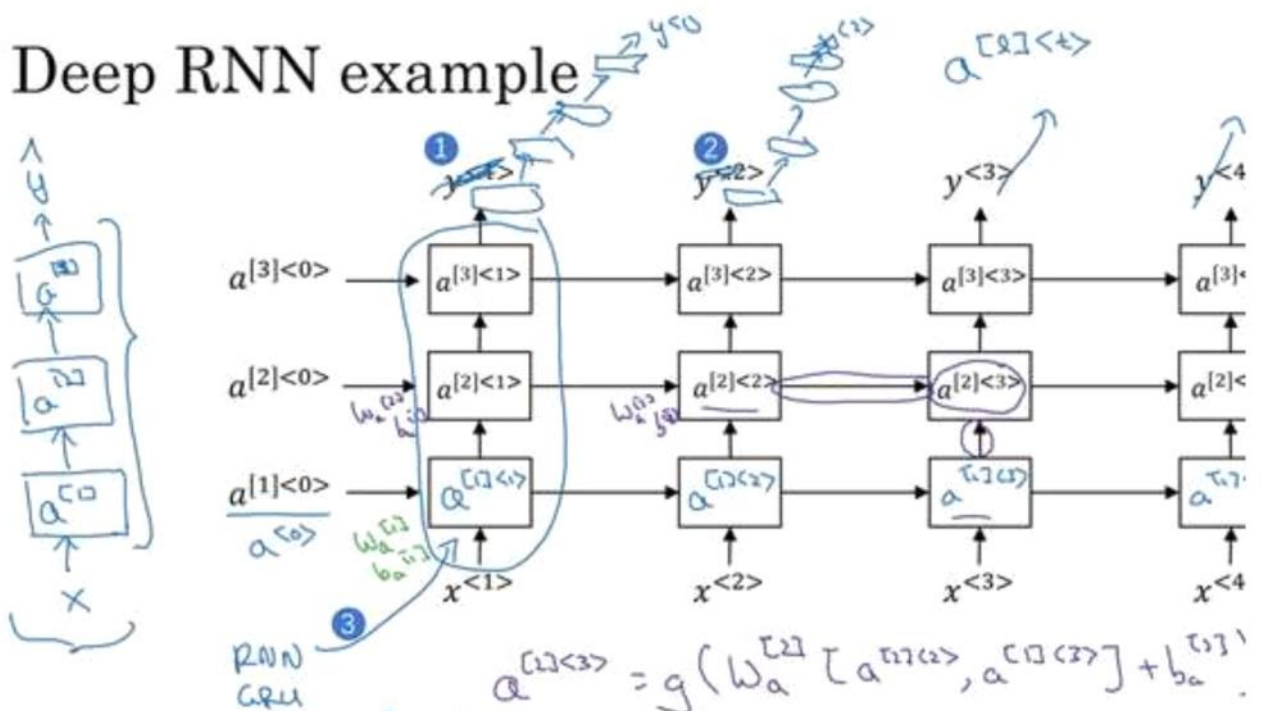
$$\underline{a}^{<t>} = \Gamma_o * \underline{c}^{<t>}$$

Bidirectional RNN (BRNN)

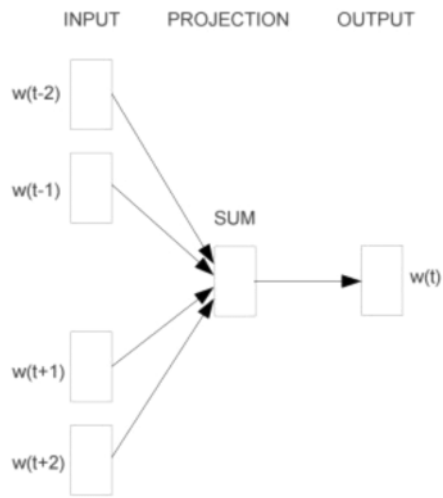


Acyclic graph

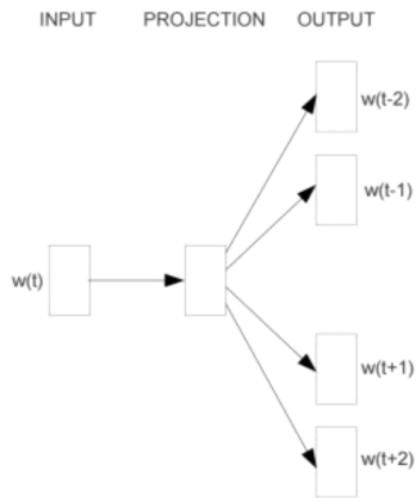
Deep RNN example



Word2Vec



CBOW



Skip-gram

<https://blog.csdn.net/pnnngchg>