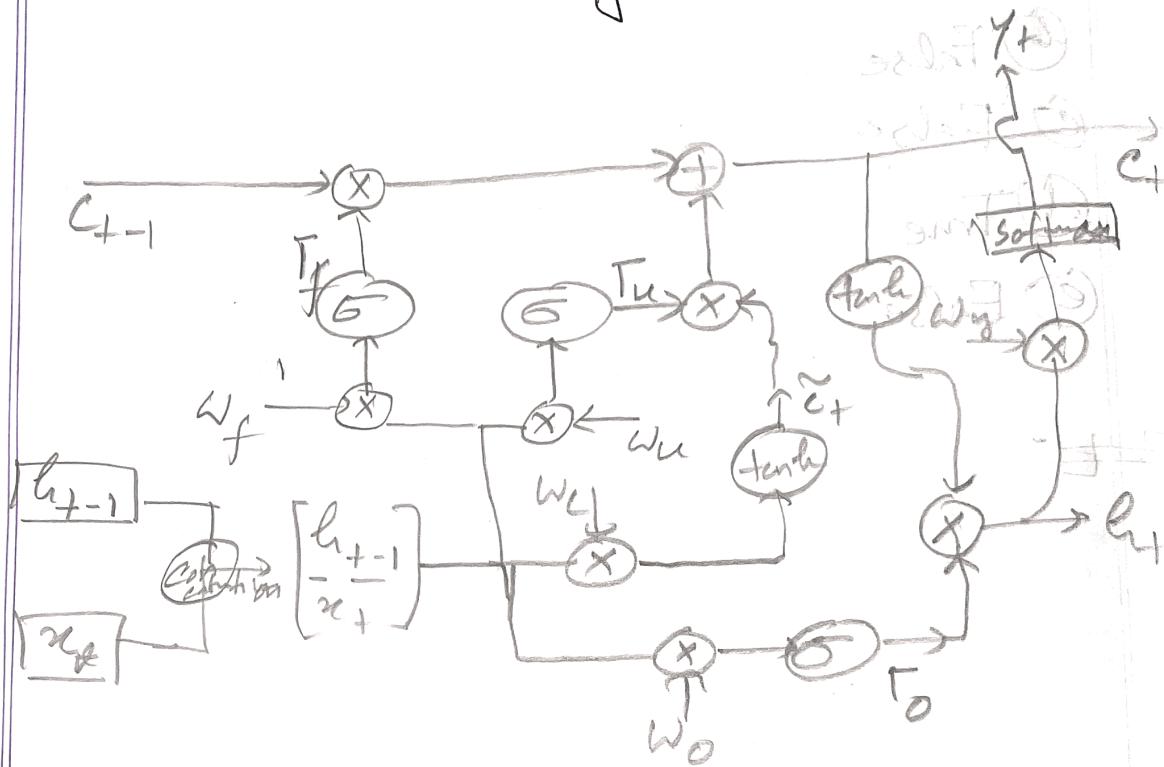


1. (a) True
(b) False
(c) False
(d) True
(e) False

~~Set~~

5. LSTM Network Diagram



In the above diagram c_{t-1} , c_t are cell states. ~~h_{t-1}~~ & h_t are hidden states. $(t-1)$ & t are previous and next time steps. x_t is the input.

For easier multiplication, h_{t-1} & x_t are concatenated together as $\begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix}$.

for updating and removing cell state, two other parameters are used.

$$f_f = \sigma(w_f \cdot \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix})$$

$$f_u = \sigma(w_u \cdot \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix})$$

These values are used to update with current and forget previous cell states.

$$\tilde{c}_t = \tanh(w_c \cdot \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix})$$

$$c_t = f_u \cdot \tilde{c}_t + f_f \cdot c_{t-1}$$

For new hidden state value,

$$h_t = f_o \tanh(c_t) \text{ where } f_o = \sigma(w_o \cdot \begin{bmatrix} h_{t-1} \\ x_t \end{bmatrix})$$

$$y = \text{softmax}(w_y \cdot h_t)$$

This is the output using new hidden state.

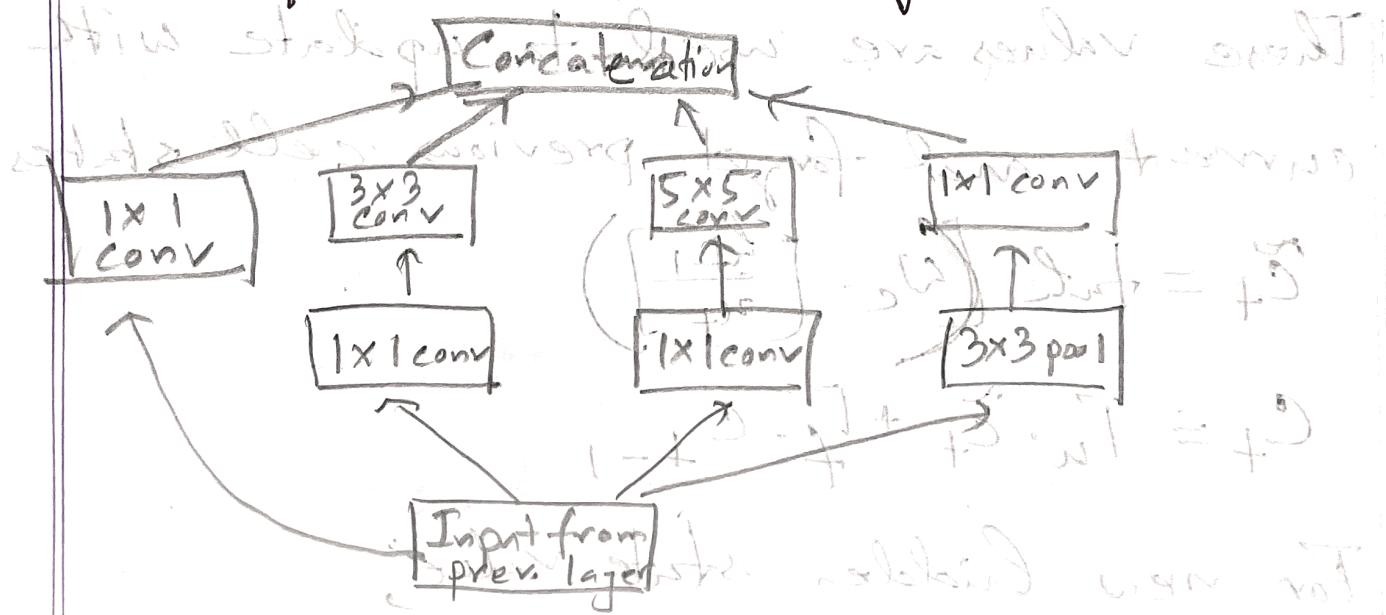
w_o, w_u, w_f, w_c are two parameters
of h_{t-1} and u_t added together
for easiness matrix multiplication

- known size of h_t and u_t

$$\begin{pmatrix} h_{t-1} \\ u_t \end{pmatrix} \cdot \begin{pmatrix} w_o \\ w_u \end{pmatrix} = g_t$$

$$\begin{pmatrix} h_{t-1} \\ u_t \end{pmatrix} \cdot \begin{pmatrix} w_f \\ w_c \end{pmatrix} = f_t$$

6. (a) Inception module in GoogleNet:



⑥ GoogLeNet uses 1×1 convolutional layer before any other convolution which reduces dimension of the input and eventually number of operations and ~~parameters~~ overall parameters. Other than this, ~~VGG-Net uses~~ GoogLeNet uses parallel convolution maintaining the same output spatial dimension and finally concatenates them. VGG-Net does not perform such parallel convolution or 1×1 convolution. ~~This~~ VGG-Net uses normal convolution. Other than this, VGG-Net ~~does not~~ uses fully connected network before final layer whereas GoogLeNet does not use FC layers.

③ Bottleneck layers are those layers which is used to reduce the dimension of the input. This is done using 1×1 convolution before convolving with higher number of filters. This reduces computation numbers, number of parameters and time needed overall. With the Bottleneck layer, highly convolutional network formation and estimation is performed with less computation.

7. a)

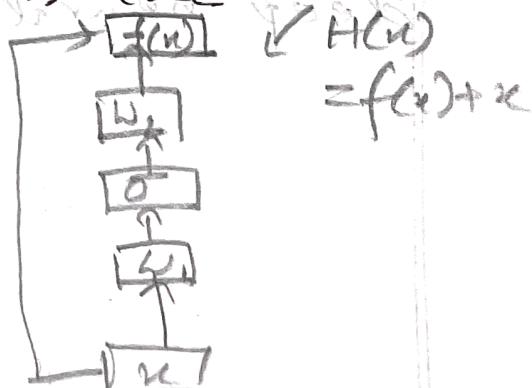
Stacking many layers of VGG-net or GoogLeNet does not improve accuracy. This is because of the degradation problem.

In the Res-Net paper, authors mentioned that even though they used identity mappings for the added layers, the accuracy decreased.

b) Residual Network solves this problem using shortcut connection. The $H(n)$ previously,

$$H(n) = \omega_2 \circ \omega_1(n)$$

but for ResNet,



$$H(x) = W_2 \circ W_1(x) + \epsilon$$

which improves the gradient
at earlier stage and helps
the network to approximate
values.

In the paper, the authors showed
even with 1202 layers, the
training error reduced. Before
this paper, 22-layer was the max
number of layers used for
deep neural networks.

Learning(ϵ)

(ϵ , μ) \rightarrow (0, 0)

initial and final

Q.

$$\Gamma_u = \sigma(w_u \cdot \left[\frac{c_{t-1}}{x_t} \right])$$

$$\tilde{c}_t = \tanh(w_c \cdot \left[\frac{c_{t-1}}{x_t} \right])$$

$$c_t = \Gamma_u \cdot \tilde{c}_t + (1 - \Gamma_u) \cdot c_{t-1}$$

$$\hat{y}^+ = \text{softmax}(w_y \cdot c_t)$$

$$\frac{\partial c_t}{\partial \Gamma_u} = \tilde{c}_t - c_{t-1}$$

$$\Gamma_u = \sigma(m) \quad [\text{Let } m = w_u \cdot \left[\frac{c_{t-1}}{x_t} \right]]$$

$$\frac{\partial \Gamma_u}{\partial m} = \sigma(m) [1 - \sigma(m)]$$

$$\frac{\partial c_t}{\partial m} = \frac{\partial c_t}{\partial \Gamma_u} \cdot \frac{\partial \Gamma_u}{\partial m}$$

$$= (\tilde{c}_t - c_{t-1}) [\sigma(m) [1 - \sigma(m)]]$$

$$\frac{\partial m}{\partial w_u} = \left[\frac{c_{t-1}}{n_t} \right] \cdot \left[1 - \sigma(w_u \cdot \frac{c_{t-1}}{n_t}) \right]$$

$$\begin{aligned} \frac{\partial c_t}{\partial w_u} &= \frac{\partial c_t}{\partial m} \cdot \frac{\partial m}{\partial w_u} \\ &= (c_t - c_{t-1}) \cdot \left[\sigma(w_u \cdot \frac{c_{t-1}}{n_t}) \right] \left[1 - \sigma(w_u \cdot \frac{c_{t-1}}{n_t}) \right] \\ &\quad \times \left[\frac{c_{t-1}}{n_t} \right] \end{aligned}$$

$$\frac{\partial c_t}{\partial \tilde{c}_t} = \text{[Put } c_{t-1} = n \text{]} \quad \tilde{c}_t = \tanh(n) \quad [\text{Answer}]$$

$$\frac{\partial \tilde{c}_t}{\partial n} = 1 - \tanh^2(n) \quad [\text{Let } n = w_c \cdot \frac{c_{t-1}}{n_t}]$$

$$\frac{\partial n}{\partial w_c} = \left[\frac{c_{t-1}}{n_t} \right]$$

$$\frac{\partial C_t}{\partial w_C} = \frac{\partial C_t}{\partial \hat{C}_t} \cdot \frac{\partial \hat{C}_t}{\partial n} \cdot \frac{\partial n}{\partial w_C}$$

$$= T_u \cdot \left[1 - \tanh^2 \left(\omega_C \cdot \left[\frac{C_t - 1}{n} \right] \right) \right] \cdot \left[\frac{C_t - 1}{n} \right]$$

[Answer]

2. (a) Let $x = \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{bmatrix}$ is the input.

and $f = \begin{bmatrix} f_{11} & f_{12} \\ f_{21} & f_{22} \end{bmatrix}$ is the filter. (for 3×3 output)

Transpose conv. of x and $f =$

$$\begin{array}{c}
 \begin{array}{|c|c|} \hline
 x_{11}f_{11} & x_{11}f_{12} \\ \hline
 x_{11}f_{21} & x_{11}f_{22} \\ \hline
 \end{array} +
 \begin{array}{|c|c|} \hline
 f_{11}x_{12} & x_{12}f_{12} \\ \hline
 f_{12}x_{21} & x_{21}f_{22} \\ \hline
 \end{array} +
 \begin{array}{|c|c|} \hline
 f_{21}x_{22} & x_{22}f_{12} \\ \hline
 f_{22}x_{11} & x_{11}f_{22} \\ \hline
 \end{array} +
 \begin{array}{|c|c|} \hline
 f_{11}x_{21} & x_{21}f_{12} \\ \hline
 f_{12}x_{22} & x_{22}f_{22} \\ \hline
 \end{array} \\
 \begin{array}{|c|c|} \hline
 x_{11}f_{11} & x_{11}f_{12} \\ \hline
 x_{11}f_{21} & x_{11}f_{22} \\ \hline
 \end{array} +
 \begin{array}{|c|c|} \hline
 x_{11}f_{11} & x_{11}f_{12} \\ \hline
 x_{12}f_{11} & x_{12}f_{12} \\ \hline
 \end{array} +
 \begin{array}{|c|c|} \hline
 x_{12}f_{12} & x_{12}f_{22} \\ \hline
 x_{21}f_{11} & x_{21}f_{12} \\ \hline
 \end{array} +
 \begin{array}{|c|c|} \hline
 x_{21}f_{11} & x_{21}f_{12} \\ \hline
 x_{22}f_{11} & x_{22}f_{22} \\ \hline
 \end{array} \\
 \begin{array}{|c|c|} \hline
 x_{11}f_{21} & x_{11}f_{22} \\ \hline
 x_{21}f_{11} & x_{21}f_{12} \\ \hline
 \end{array} +
 \begin{array}{|c|c|} \hline
 x_{12}f_{21} & x_{12}f_{22} \\ \hline
 x_{21}f_{12} & x_{21}f_{22} \\ \hline
 \end{array} +
 \begin{array}{|c|c|} \hline
 x_{22}f_{11} & x_{22}f_{12} \\ \hline
 x_{22}f_{21} & x_{22}f_{22} \\ \hline
 \end{array} +
 \begin{array}{|c|c|} \hline
 x_{22}f_{12} & x_{22}f_{22} \\ \hline
 x_{22}f_{21} & x_{22}f_{22} \\ \hline
 \end{array}
 \end{array}
 \end{array}$$

each input is multiplied/convolved with filter/kernel and finally all are added together based on stride.

$$⑥ \quad x = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \quad M = \begin{bmatrix} -1 & 0 & 1 \\ -1 & 2 & 1 \\ -1 & 0 & 1 \end{bmatrix}$$

Unsampling =

$$\begin{array}{c|ccc} -1 & | & 0 & 1 \\ \hline -1 & | & 2 & 1 \\ -1 & | & 0 & 1 \end{array} + \begin{array}{c|ccc} -1 & | & -2 & 0 \\ \hline -1 & | & 4 & 2 \\ -2 & | & 0 & 2 \end{array} + \dots$$

+ Down X for every unsampling

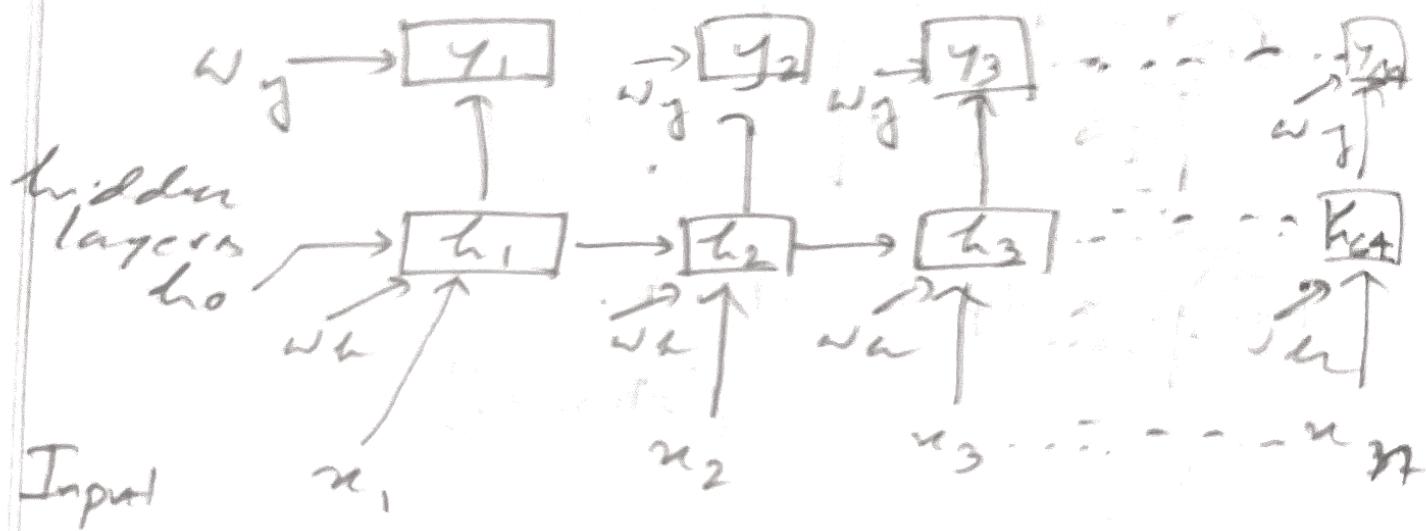
$$\begin{array}{c|ccc} -3 & | & 0 & 3 \\ \hline -3 & | & 6 & 3 \\ -3 & | & 0 & 3 \end{array} + \begin{array}{c|ccc} 1 & | & 0 & 4 \\ \hline -4 & | & 8 & 4 \\ -4 & | & 0 & 4 \end{array}$$

\rightarrow	-1	-2	1	1	2
\rightarrow	-4	-4	8	6	
\rightarrow	-4	0	12	6	
\rightarrow	-3	-4	3	4	
\rightarrow					

Final output 4×4 .

[Answer]

4. @ 6



Since there are 64 nodes, ~~and~~ there
 h_{64} in the hidden layers.

The inputs and outputs are
 $x_1, x_2, x_3, \dots, x_n$ respectively.
 y_1, y_2, \dots, y_n

(e) hidden layer ~~and~~ at $t+1$ time step:

$$h_{t+1} = \tanh(W_h [u_t | h_t])$$

$$y_{t+1} = \text{softmax}(W_y h_{t+1})$$

further loss is RSM derivative of

(f) input dimension 32×1

hidden state dimension 64×1

output dimension 32×1

(g) Here, 2 parameters are available
 $\rightarrow W_h$ & W_y .

dimension of $W_h \rightarrow 64 \times 64$

$$W_h \rightarrow 64 \times 32$$

So dimension of $W_h \rightarrow 64 \times 96$

Dimension of $W_y \rightarrow 32 \times 64$,

9. ② Deep convolutional GAN uses ReLU activation for all layers in generator except for the output which uses tanh and LeakyReLU for discrimination. The original GAN used maxout activation function. Also, DCGAN removed FC layers which is used in GAN. DCGAN also uses batchnorms for both generator and discriminator.

③ DC of 1 problem layer receive

be

⑩ DCGAN should fail when number of layers are more for degradation problem. And since \oplus with lower layers, they only learn from the receptive fields, the outputs should be somewhat similar.

Layer	Input Size	Kernel Size	Output Size	Pad	No. of Params (Neglecting bias)	No. of Ops
L-1-Conv-1	64x64x3	7x7x3	64x64x32	3x3	4704	19267584
L-1-Conv-2	64x64x32	3x3x32	64x64x64	1x1	18432	75497472
L-1-Pool	64x64x64	3x3	32x32x64	0.5x0.5(s = 2) So 1 padding in either side	0	2359296
L2-Conv-1	32x32x64	3x3x64	32x32x128	1x1	73728	75497472
L2-Conv-2	32x32x128	3x3x128	32x32x64	1x1	73728	75497472
L2-Conv-3	32x32x64	3x3x64	32x32x16	1x1	9216	9437184
L-2-Pool	32x32x16	3x3	16x16x16	0.5x0.5 (s =2) 1 padding either side	0	147456
L3-FC-1	4096	NA	1000	N/A	4096000	4096000
L3-FC-2	1000	NA	1000	N/A	1000000	1000000
L3-Softmax	1000	NA	10	N/A	10000	10000

Prithul Sarker