

1. a) Zero Initialisation (All weights same so grade will always be same)

b) L1 Regularisation

c) Instead of training large number of networks for ensemble, we can use DROPOUT which is efficient as we don't need to train all the different networks from scratch. Instead, at each training instance, with some prob  $p$ , drop some neuron in the network. So, we get a diff network everytime but also weights are shared across them.

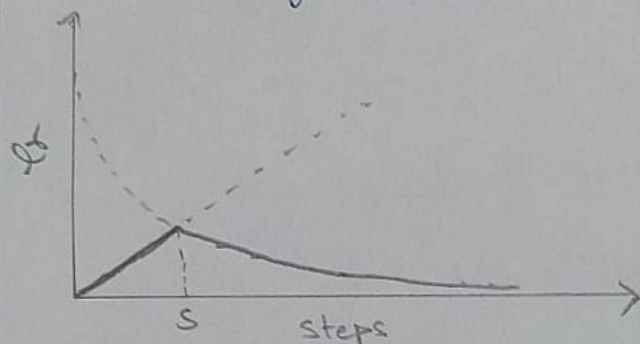
d) The Read, Write and Forget gates

$$c) \quad W' = \frac{W - F + 2P}{S} + 1$$

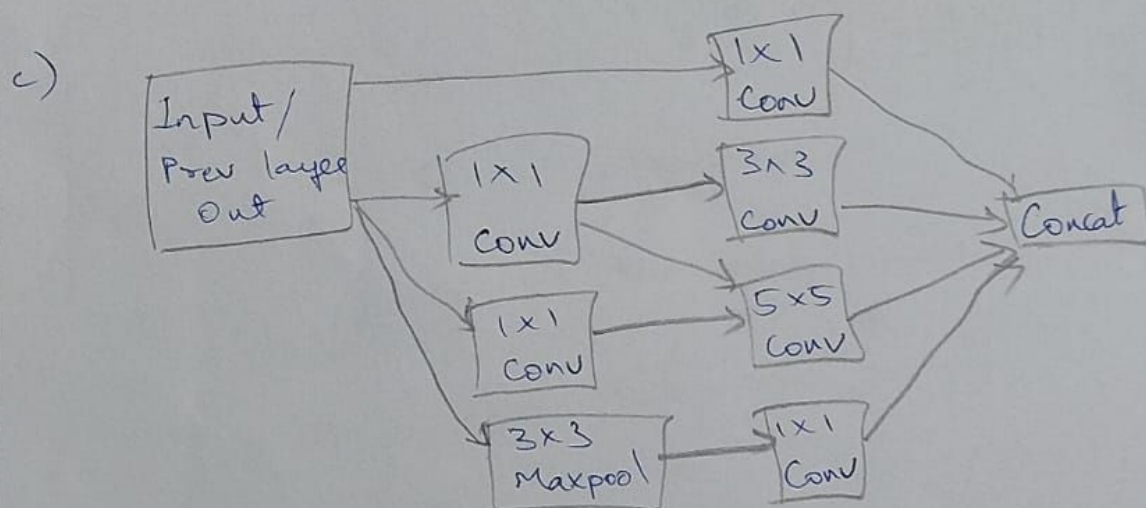
$$H' = \frac{H - F + 2P}{S} + 1$$

2. a)  $\text{Gelu}(x) = x \phi(x)$  ↗ Std Gaussian CDF

b) Learning Rate with warmup



increase linearly up to step  $S$  (warmup)  
then decrease  
eg.  $lr = \begin{cases} t \cdot k & \text{if } t \leq S \\ t^{-1/2} & \text{if } t > S \end{cases}$  ↗ const



d) GRU  $\Rightarrow$  Gates  $\Rightarrow$

$$o_t = \sigma(W_o s_{t-1} + U_o x_t + b_o)$$

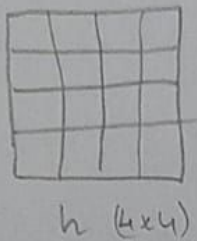
$$i_t = \sigma(W_i s_{t-1} + U_i x_t + b_i)$$

States  $\Rightarrow$

$$\tilde{s}_t = \sigma(W(s_t \odot s_{t-1}) + U x_t + b)$$

$$s_t = (1 - i_t) \odot s_{t-1} + i_t \odot \tilde{s}_t$$

e)  $f_{\text{att}}(h_i, c_j) = V_{\text{att}}^T \tanh(W_{\text{att}}[h_i; c_j])$



$2 \times 2$  maxpool



$m$  (2x2)

$\frac{\partial L}{\partial m}$  is known. In back prop, due to max pooling, the gradient propagates to only the 'max' cells in  $h$ . All other cells get gradient 0.

$$\frac{\partial L}{\partial h_{ij}} = \frac{\partial L}{\partial m} \quad \text{if} \quad h_{ij} = m_{\lfloor \frac{i-1}{2} \rfloor + 1, \lfloor \frac{j-1}{2} \rfloor + 1}$$

0 otherwise.

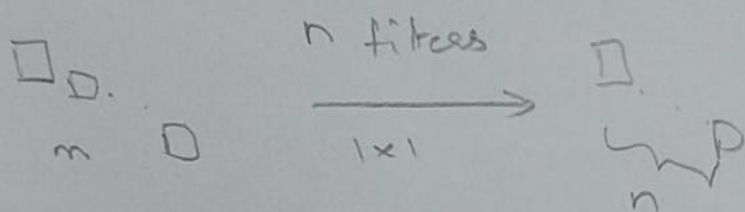
$m_{\lfloor \frac{i-1}{2} \rfloor + 1, \lfloor \frac{j-1}{2} \rfloor + 1}$  is the max value of the cells around  $h_{ij}$  when doing max pooling.



5. The  $n$ -dim input can be represented as  $m \times (1 \times 1)$ .

We can get  $n$ -dim output by applying ' $n$ ' filters that reduce the  $m \times (1 \times 1)$  to  $1 \times 1$ .

We use filter of size  $(1 \times 1)$  and do  $1 \times 1$  convolution to get  $1 \times 1$  from the  $m \times (1 \times 1)$  and we do using  $n$  such filters.



6. Output size after  $7 \times 7$  filter conv will be smaller than for  $3 \times 3$  assuming  $S=1, P=0$ .

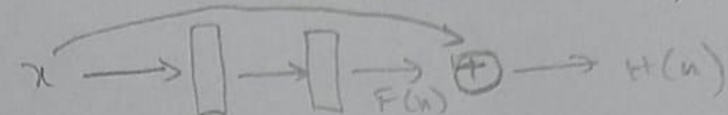
Max pool layer during backprop only propagates gradients to the cells which were maximum in their neighbourhood.

For  $7 \times 7$ , we get 1 cell where grad is propagated for every 49 neighbourhood whereas for  $3 \times 3$  it is 1 in every 9.

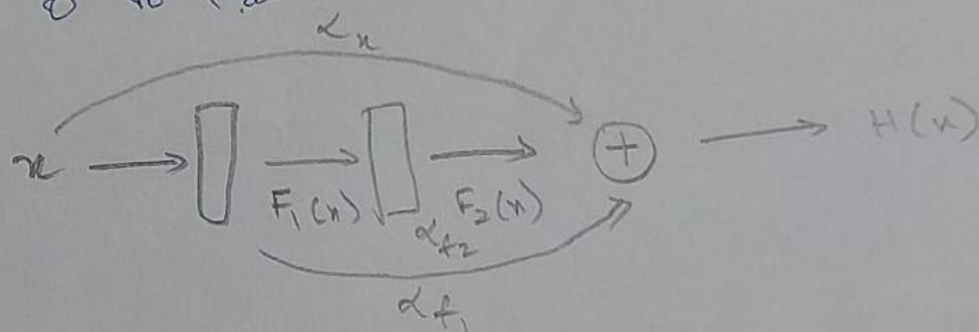
Hence gradients are propagated in more cells for  $3 \times 3$  than for  $7 \times 7$  and hence it can get trained faster.

So,  $3 \times 3$  will be better.

7. In Residual Network, we hardcode the connection as,

$$H(x) = F(x) + x$$


To prevent this hardcoding we can make it adaptive by simply making  $H(x)$  as a linear combination of each layer's outputs. Weights are fixed to be between 0 to 1 and are learnt.



$$H(x) = \alpha_x \odot x + \alpha_{f_1} \odot F_1(x) + \alpha_{f_2} \odot F_2(x)$$

Hence by learning the  $\alpha$ s, we can control the flow of information from all layers to  $H(x)$  and not just  $x$  and  $F(x)$  like before.

If some layer output does not contribute, its  $\alpha$  will go to 0 while training.

This is more adaptive.



$$\begin{aligned}
 9. \quad E((y - \hat{f}(n))^2) &= E((f(n) + e - \hat{f}(n))^2) \\
 &= E((f(n) - E(\hat{f}(n)) + e + E(\hat{f}(n)) - \hat{f}(n))^2) \\
 &= E((f(n) - E(\hat{f}(n)))^2) + E(e^2) + E((E(\hat{f}(n)) - \hat{f}(n))^2) \\
 &\quad + 2E(e(f(n) - E(\hat{f}(n)))) + 2E(e(E(\hat{f}(n)) - \hat{f}(n))) \\
 &\quad + 2E((f(n) - E(\hat{f}(n)))(E(\hat{f}(n)) - \hat{f}(n)))
 \end{aligned}$$

The  $2E$  terms go to 0 as  $e, \hat{f}(n)$  are all independent.  
 Covariance = 0. Also,  $E(e^2) = \text{Var}(e)$ ,  $E((E(\hat{f}(n)) - \hat{f}(n))^2) = \text{Var}(\hat{f}(n))$

$$\Rightarrow (f(n) - E(\hat{f}(n)))^2 + \text{Var}(e) + \text{Var}(\hat{f}(n))$$

First term is bias of  $\hat{f}(n)$ ,  $\text{Var}(e)$  is irreducible error

$$\Rightarrow (\text{Bias})^2 + \sigma^2 (\text{irreducible}) + \text{Variance} //$$

10. Feed forward layer is used in encoder in transformers
- Add more model complexity
  - Convert attention layer output to size as expected by next layer / decoder
  - Add non-linearity

If we don't use feed forward layer, we will have to feed the output of attention layer to decoder. This can model only linear functions as attention is linear and hence the model cannot perform well for real life data and tasks with lot of non-linearity.



11. To ensure that only attentions <sup>of words</sup> inside the window are aggregated for a time step 'P' in decoder, we can use a mask vector specific to time step P which is point wise multiplied with the attention outputs from encoder and then aggregated.

mask<sub>P</sub> is a vector of size M and has all 0s except at the indices P-k to P+k.

$$\text{mask}_P = \begin{bmatrix} 0 & 0 & \dots & 0 & 1 & \dots & 1 & 0 & \dots & 0 \end{bmatrix}$$

Index    1    2        ...    P-k-1   P-k                    P+k   P+k+1    M

$$\therefore C_P = \text{mask}_P \odot \alpha_P \quad \text{mask}_P \propto_P$$

12. Identifying human-animal interactions on college campuses. To detect if someone is feeding the animals, etc.

Input: Image from CCTV

Output: Detect event.

Neural Net: CNN arch.