

Mohamed Ehab Fathy

2002597

PAGE  
DATE

## Assignment 2

1)  $w. shape = (500 \times 500 \times 3, 100)$

$b. shape = 100$

2)  $5 \times 5 \times 10 = 250$  parameters

3) derivative filters for edge detection (Sobel)

\* left images<sup>m</sup> Filter  
(vertical edges)

|   |   |    |
|---|---|----|
| 1 | 0 | -1 |
| 2 | 0 | -2 |
| 1 | 0 | -1 |

\* right images Filter  
(horizontal edges)

|    |    |    |
|----|----|----|
| 1  | 2  | 1  |
| 0  | 0  | 0  |
| -1 | -2 | -1 |

4) exp. moving avg.  $\rightarrow$  momentum term in Adam  
( $\beta_1$ ) ( $V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) dw$ )

\* as  $\beta_1 \uparrow$ ,  $V_{dw}$  becomes more affected

by the previous values (has longer memory)

so it becomes smoother and more sluggish  
if the gradient is noisy (changing sign) and will

become faster if gradients are in one direction Hi . Star

$$5) \mu_{B_m} = \frac{1}{m} \sum_{i=1}^m z^i$$

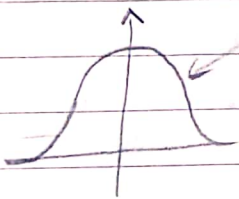
$$\sigma_{B_m}^2 = \frac{1}{m} \sum_{i=1}^m (z^i - \mu_{B_m})^2$$

$$\hat{z}_i = \frac{(z_i - \mu_{B_m})}{\sqrt{\sigma_{B_m}^2 + \epsilon}}$$

$\epsilon$ , small value for numerical stability

\* improves gradients in the network, since

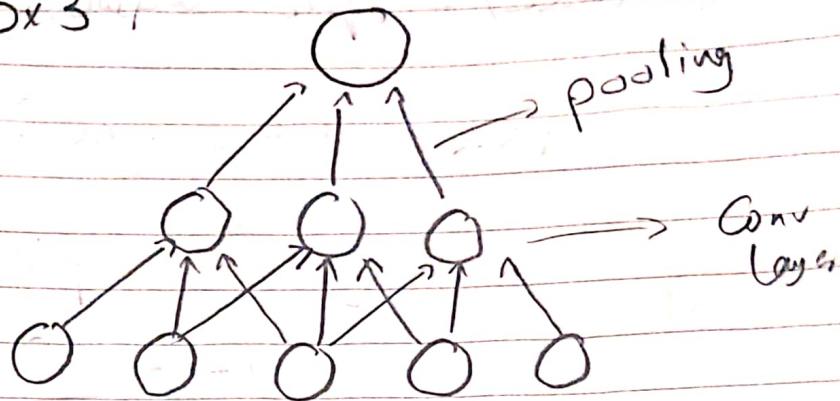
activation gradients becomes zero in high values



, hence avoids layers death  
(stopped learning)

\* Allows higher learning rate

6)  $5 \times 5 \times 3$



7)

$$\text{output\_shape} = (128, 1, x, x)$$

$$x = \frac{128 - 7 + 2 \times 3}{2} + 1 = 64.5$$

Floor  $\rightarrow$  64

8) inverted drop out scales  $\times$  the output of activations by inverse of keep probability  $q = 1 - p$ , so, it's not required to scale network's output at testing or evaluation <sup>during training</sup>

9) because it's not shift invariant, accordingly during training, if the required classes changes it's location in the image, network will produce



different output, preventing the learning algorithm from finding the weights that detects the features of the specific class, regardless to the requirement of large weights size

$$10) a[0] = -2 \times 0 + 1 \times 4 = 4$$

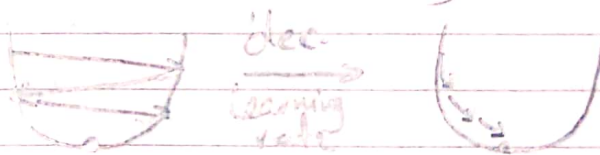
$$a[1] = -2 \times 4 + 1 \times 1 = -7$$

$$a[2] = -2 \times 1 + 1 \times -1 = -3$$

$$a[3] = -2 \times -2 + 1 \times 3 = 7$$

$$a[4] = -2 \times 3 + 1 \times 0 = -6$$

11) learning rate is decreased, since constant error is probably due to overshooting around local minima.



12. Since they are shift invariant, so the same filter will search for a feature across the whole image regardless to the feature's location, this also decreases network's size, since the whole image's pixels share the same filters (weights).

13. Since during each training iteration, some of network's weights are disabled by probability  $(1-p)$ , but during prediction we can't do the same since we want to get benefit of the whole network, also we can't activate all weights without any change, since this will produce output value different from training's, so we will scale the weights by " $p$ ".

14)

\* Momentum

$$\theta_{t+1} = \theta_t + V_{t+1}$$

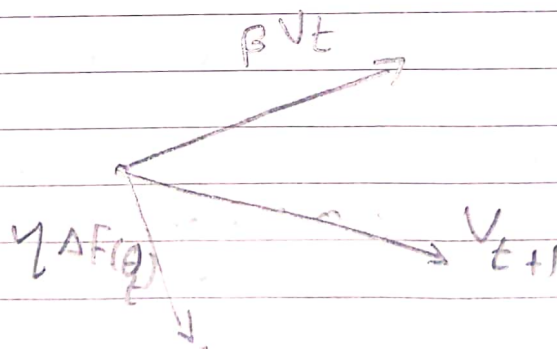
$$V_{t+1} = \beta V_t - \eta \Delta F(\theta_t)$$

\* NAG

$$\theta_{t+1} = \theta_t + V_{t+1}$$

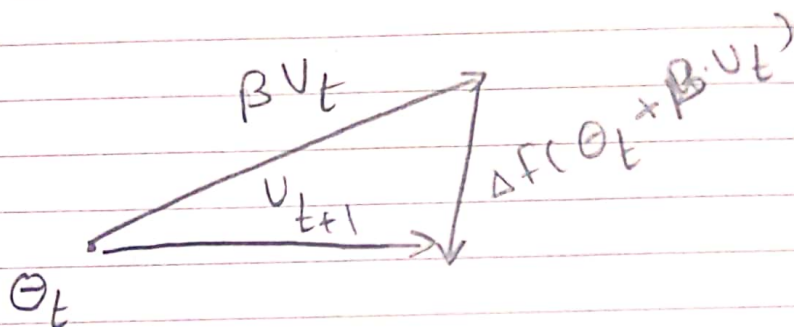
$$V_{t+1} = \beta V_t - \eta \Delta F(\theta_t + \beta V_t)$$

\* in momentum gradient is compute at the point before update





\* in NAG the gradient is computed after applying the velocity term to the previous parameter, so if this velocity term is bad (driving model into higher loss), we will be aware of and correct it



$$15) G_t = \sum_{i=0}^n \Delta F(t) \Delta F(t)^T$$

$$w_{t+1} = w_t - \frac{\eta}{\sqrt{\text{Diag}(G_t) + \epsilon I}} \Delta F$$

\* the learning rate decreases over time due to accumulation of squared gradients summation