

Assignment #2

1) $500 \times 500 \times 3 = 750000$

- weight matrix $\rightarrow 750000 \times 100$
- bias $\rightarrow 1 \times 100$

2) 10 filters, kernel size 5×5

no of parameters = $5 \times 5 \times 10 = 250 + 10 \text{ bias} = 260$

3) • first image \rightarrow vertical edge detection $\begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix}$

• second image \rightarrow horizontal edge detection $\begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{bmatrix}$

4) Adam optimizer makes use of the RMSProp optimization idea, which has the formula

$$\bullet w_{t+1} = w_t - \frac{\eta_t}{(v_t + \epsilon)^{1/2}} \left(\frac{\partial L}{\partial w_t} \right)$$

$$\bullet v_t = \beta v_{t-1} + (1 - \beta) \left(\frac{\partial L}{\partial w_t} \right)^2$$

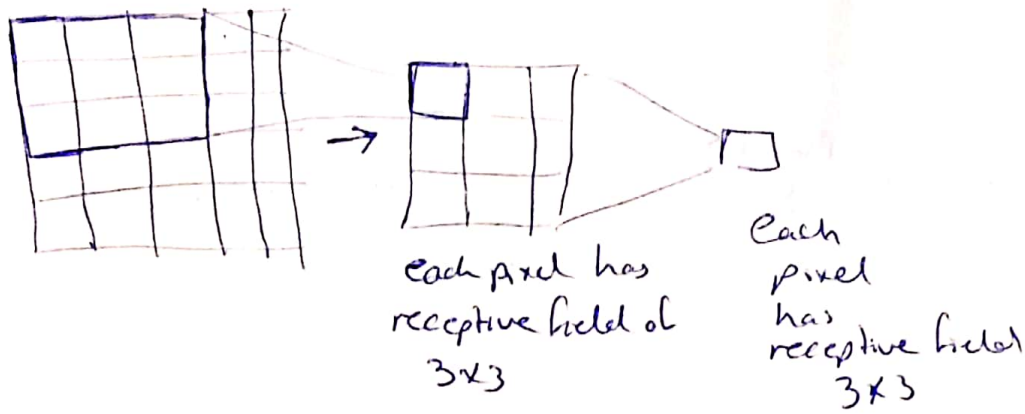
So, we divide at each step by the total of previous gradients. So, when the number of steps increase, this means that the weight updates become smaller & the exponential moving average depends less on recent updates.

5) $\mu_z = \frac{1}{n} \sum_{i=1}^n z_i$, $\sigma_z^2 = \frac{1}{n} \sum_{i=1}^n (z_i - \mu_z)^2$

$$z_i = \frac{z_i - \mu_z}{\sqrt{\sigma_z^2}}$$

- It is used for faster & more robust training, as well as to reduce overfitting effect.

6)



→ In Total, The receptive field of every unit of The pooling layer is 5×5

7) • input $\rightarrow C \times H \times W = 96 \times 128 \times 128$

• filters $\rightarrow D \times C \times H_F \times W_F = 128 \times 96 \times 7 \times 7$

$S=2$, $P=3$

$$\begin{aligned} \text{• output dimensions} &= \frac{\text{input dim} - \text{kernel size} + 2 \times P}{S} + 1 \\ &= \frac{128 - 7 + 2 \times 3}{2} + 1 = 64 \end{aligned}$$

• dimension of output $= 128 \times 64 \times 64$

8) With dropout, we turn off some of the neurons with a probability (p), so to invert this effect, we have to scale the activations at test time with $(1/p)$ to cancel the dropout effect.

With inverted dropout, This scaling is done at training time, This leads to faster inference & Test Time

9) Because features in images are location invariant, this means that if the network learns a feature in some position of image it should be able to generalize the feature. Also, flattening an image to feed into fully connected network may mess up positional data of the image, and runs the nature of image data.

10) •
$$\text{OIP size} = \frac{\text{IIP size} - \text{kernel size} + 2P}{S} + 1$$
$$= \frac{4 - 2 + 2P}{1} + 1 = 5 \rightarrow \therefore \text{Padding} = 1$$

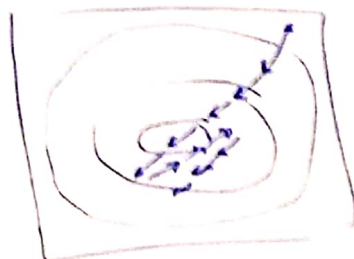
• Convolution = $(0, 4, 1, -1, 3, 0) * (-2, 1)$
 $= (4, -7, -3, 5, -6)$

11) The learning rate is decreased (Learning Rate Decay). This leads to the model taking smaller weight update steps & being able to reach the local minimum which couldn't be reached with the higher learning rate which can be shown from the flat Loss Curve.

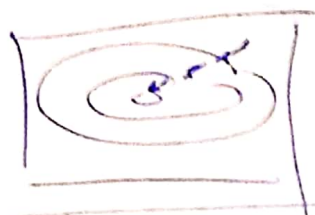
12) Convolutional filters are able to learn features available at any location of the image (location invariant). So they are more suitable for the nature of image data. Also with convolutional layer, less number of learnable parameters is needed to extract features from the image than what is needed with FC layers.

- 3) at train $\rightarrow p$ percentage of neurons are turned off
• at test \rightarrow All neurons are on, so activations need to be scaled by factor of $(1/p)$

- 14) Standard momentum only takes in consideration recent updates when updating the parameters, this may lead to overshoots or noisy paths on the loss curve



With Nesterov accelerated momentum, a new term is taken into consideration which is a look ahead term, this helps for smoother optimization



- 15) Adagrad $\rightarrow w_t = w_{t-1} - \eta_t \frac{\partial L}{\partial w_{t-1}}$

$$\eta_t = \frac{\eta}{\sqrt{\delta_t + \epsilon}}$$

$$\delta_t = \sum_{i=1}^t \left(\frac{\partial L}{\partial w_{t-1}} \right)^2$$