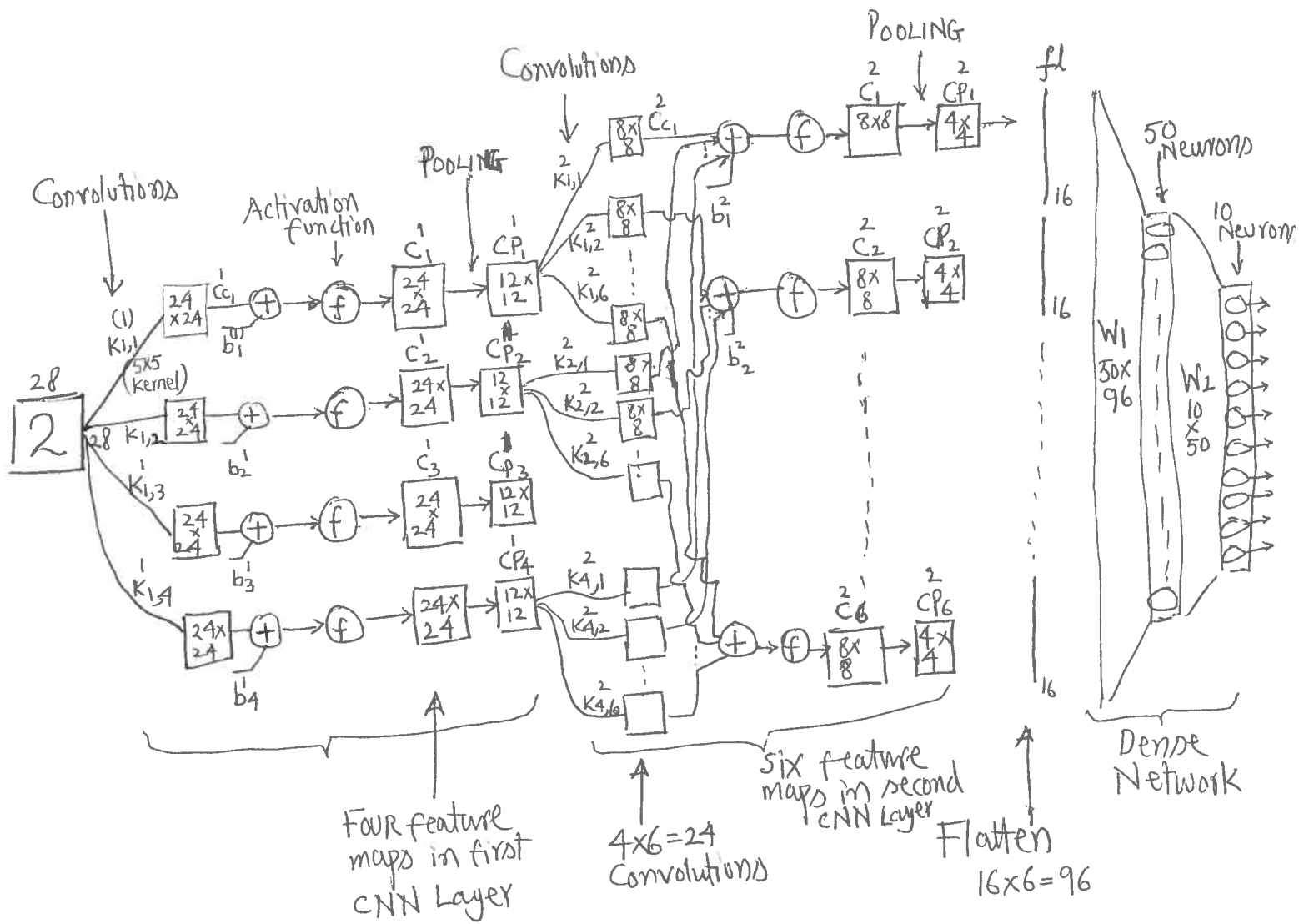


# ARCHITECTURE OF A CNN WITH 4 Feature maps in First Layer and 6 Feature maps in second Layer ①



Number of Parameters to be Learned:

$$\begin{aligned}
 & 4 \times \underbrace{5 \times 5}_{\text{kernel}} + \underbrace{4 \times 1}_{\text{bias}} + 24 \times \underbrace{5 \times 5}_{\text{kernel}} + \underbrace{6 \times 1}_{\text{bias}} + \underbrace{50 \times 96}_{W_1} + \underbrace{50 \times 1}_{\text{bias}} + \underbrace{10 \times 50}_{W_2} + \underbrace{10 \times 1}_{\text{bias}} \\
 & = 6170
 \end{aligned}$$

## Kernel, Bias and Weights Initialization

(2)

$$K_{l,m}^1 = \text{Uniform} \left( \pm \sqrt{\frac{4}{(1+4) \times 5^2}} \right)$$
$$m = 1..4$$

$$K_{m,n}^2 = \text{Uniform} \left( \pm \sqrt{\frac{6}{(4+6) \times 5^2}} \right)$$
$$n = 1..6$$

$$W_1 = \text{Uniform} \pm \left( \sqrt{\frac{6}{96+50}} \right)$$

$$W_2 = \text{Uniform} \pm \left( \sqrt{\frac{6}{10+50}} \right)$$

All biases initialized to zeros.

## FORWARD PASS

Suppose input image is called  $I$   
then  $C_m^1 = f(I \otimes K_{l,m}^1 + b_m^1)$

$$C_m^1(i,j) = f \left( \sum_{u=-2}^2 \sum_{v=-2}^2 I(i-u, j-v) \cdot K_{l,m}^1(u,v) + b_m^1 \right)$$

$$m = 1..4$$

$\otimes$  = convolution

Assuming four feature maps in first CNN Layer and performing  $5 \times 5$  convolutions.

If we are using average pooling, then

$$C_P^1(i,j) = \frac{1}{4} \sum_{u=0}^1 \sum_{v=0}^1 C_m^1(2i-u, 2j-v)$$
$$i,j = 1..12$$

(3)

Similarly in the second layer,

$$C_n^2(i,j) = f\left(\sum_{m=1}^4 \sum_{u=-2}^2 \sum_{v=-2}^2 C_{P_m}^1(\delta-u, j-v), K_{m,n}^2(u,v) + b_n^2\right)$$

$$n = 1..6$$

Size of  $C_n^2$  is  $8 \times 8$ , so  $i, j$  will vary  $1..8$

The average pooling in second CNN layer can be described as

$$C_{P_n}^2(i,j) = \frac{1}{4} \sum_{u=0}^1 \sum_{j=0}^1 C_n^2(2i-u, 2j-v)$$

$$i, j = 1..4$$

The output of 6 pooling layers in the second CNN layer is flattened into  $16 \times 6 = 96$  values. These 96 values are fed to a regular dense Neural Network with 50 hidden layer neurons and 10 output layer neurons.

$$a_1 = f(\text{flatten} \times W_1 + b_1) \quad \text{— output of hidden layer}$$

$$a_2 = f(a_1 \times W_2 + b_2) \quad \text{— output}$$

For the last layer, the activation function for a multi-class problem is usually Softmax, and the loss function as Cross Entropy.

# Backpropagation for CNN:

(4)

If last layer is Softmax, and loss function is Cross Entropy, then  $\delta$  in last layer is given by

$$\delta_i = s_i - y_i \quad (\text{see derivation of Softmax handout})$$

$$\text{or } \delta = S - Y$$

$$\Delta W_2 = \frac{\partial L}{\partial W_2} = \frac{\partial L}{\partial s_i} \times \frac{\partial s_i}{\partial W_2} = \delta \times a_1^T$$

$$\Delta b_2 = \delta$$

$$\delta_1 = \text{delta in hidden layer} = W_{2 \times 8}^T \delta_{\text{last layer}} \times f' \quad \text{L derivative of activation function}$$

$$\Delta W_1 = \frac{\partial L}{\partial a_2} \times \frac{\partial a_2}{\partial W_1} = \delta_1 \times \text{flatten}^T$$

$$\Delta b_1 = \delta_1 \quad \text{96x1 in our example}$$

$$\delta_{\text{flatten}} = \text{delta in flattening layer}$$

$$\delta_{\text{flatten}} = W_1^T \times \delta_1$$

$$\delta C_n^2(i,j) = \frac{1}{4} \delta_{\text{flatten}}([i/2], [j/2]) \cdot f' \quad \text{L } \delta_{\text{flatten}} \text{ is copied to corresponding places in } \delta C_n^2 \text{ and divided by 4}$$

$i, j = 1, 2, \dots, 8$

derivative of activation function

Computing  $\Delta K_{m,n}^2$  is tricky. Note that the equation for forward pass for  $C_n^2$  is (from page 3)

$$C_n^2(i,j) = f\left(\sum_{m=1}^4 \sum_{u=-2}^2 \sum_{v=-2}^2 C_{P_m}^1(i-u, j-v) \cdot K_{m,n}^2(u,v) + b_n^2\right)$$

$$\Delta K_{m,n}^2(u,v) = \sum_{i=1}^8 \sum_{j=1}^8 \delta C_n^2(i,j) \cdot \frac{\partial C_n^2(i,j)}{\partial K_{m,n}^2(u,v)}$$

This is because each kernel element is involved in computing the  $8 \times 8$  values in  $C_n^2$

$$\Delta K_{m,n}^2(u,v) = \sum_{i=1}^8 \sum_{j=1}^8 \delta C_n^2(i,j) \cdot \frac{\partial}{\partial K_{m,n}^2(u,v)} \left( \sum_{m=1}^4 \sum_{u=-2}^2 \sum_{v=-2}^2 C_{P_m}^1(i-u, j-v) \cdot K_{m,n}^2(u,v) + b_n^2 \right)$$

$$\Delta K_{m,n}^2(u,v) = \sum_{i=1}^8 \sum_{j=1}^8 \delta C_n^2(i,j) C_{P_m}^1(i-u, j-v)$$

$$\Delta K_{m,n}^2(u,v) = \delta C_n^2(i,j) \otimes (C_{P_m}^1 \text{ rotated by } 180^\circ)$$

L convolution

derivative is for a particular value of  $u, v$  and an  $i, j$

for convolution, it needs to be  $u-i, j-v$  so rotate it by  $180^\circ$

To determine  $\Delta b_n^2$ , note that it is of size  $1 \times 1$

(5)

$$\Delta b_n^2 = \sum_{i=1}^8 \sum_{j=1}^8 \delta c_n^2(i,j) \times \frac{\partial}{\partial b_n^2} \left( \sum_{m=1}^4 \sum_{u=-2}^2 \sum_{v=-2}^2 C_{P_m}^{\dagger}(i-u, j-v) \cdot K_{m,n}^2(u,v) + b_n^2 \right)$$

$$\boxed{\Delta b_n^2 = \sum_{i=1}^8 \sum_{j=1}^8 \delta c_n^2(i,j)}$$

To determine  $\Delta K_{l,m}^1$ , first we need to determine  $\delta C_{P_m}^1$  and  $\delta C_m^1$

$$\begin{aligned} \delta C_{P_m}^1(i,j) &= \sum_{n=1}^6 \sum_{u=-2}^2 \sum_{v=-2}^2 \delta c_n^2(i+u, j+v) \cdot \frac{\partial}{\partial C_{P_m}^1(i,j)} \left( \sum_{m=1}^4 \sum_{u=-2}^2 \sum_{v=-2}^2 C_{P_m}^1(i,j) \cdot K_{m,n}^2(u,v) + b_n^2 \right) \\ &= \sum_{n=1}^6 \sum_{u=-2}^2 \sum_{v=-2}^2 \delta c_n^2(i+u, j+v) \cdot K_{m,n}^2(u,v) \end{aligned}$$

$$= \sum_{n=1}^6 \sum_{u=-2}^2 \sum_{v=-2}^2 \delta c_n^2(i-(-u), j-(-v)) \cdot K_{m,n}^2 \text{rot. } 180^\circ(-u, -v)$$

$$\boxed{\delta C_{P_m}^1 = \sum_{n=1}^6 \delta c_n^2 \otimes K_{m,n}^2 \text{rot. } 180^\circ} \rightarrow \text{Note that this is full convolution.}$$

$\begin{matrix} 12 \times 12 \\ \text{C} \end{matrix} \quad \begin{matrix} 8 \times 8 \\ \text{C} \end{matrix} \quad \begin{matrix} 5 \times 5 \\ \text{C} \end{matrix}$

$$\delta C_m^1(i,j) = \frac{1}{4} \delta C_{P_m}^1(\lceil i/2 \rceil, \lceil j/2 \rceil) \cdot f'$$

$$i, j = 1 \dots 24$$

derivative of activation function

$$\begin{aligned} \Delta K_{l,m}^1(u,v) &= \sum_{i=1}^{24} \sum_{j=1}^{24} \delta C_m^1(i,j) \cdot \frac{\partial}{\partial K_{l,m}^1(u,v)} \left( \sum_{u=-2}^2 \sum_{v=-2}^2 I(i-u, j-v) \cdot K_{l,m}^1(u,v) + b_m^1 \right) \\ &= \sum_{i=1}^{24} \sum_{j=1}^{24} \delta C_m^1(i,j) \cdot I(i-u, j-v) \end{aligned}$$

$$\boxed{\Delta K_{l,m}^1 = \delta C_m^1 \otimes I \text{rot. } 180^\circ}$$

$$\Delta b_m^1 = \sum_{i=1}^{24} \sum_{j=1}^{24} \delta C_m^1(i,j) \cdot \frac{\partial}{\partial b_m^1} \left( \sum_{u=-2}^2 \sum_{v=-2}^2 I(i-u, j-v) \cdot K_{l,m}^1(u,v) + b_m^1 \right)$$

$$\boxed{\Delta b_m^1 = \sum_{i=1}^{24} \sum_{j=1}^{24} \delta C_m^1(i,j)}$$

⑥

The model parameters will be updated as:

$$K'_{i,m} = K_{i,m} - \lambda \cdot \Delta K'_{i,m}$$

$$b'_m = b_m - \lambda \cdot \Delta b'_m$$

$$K^2_{m,n} = K^2_{m,n} - \lambda \cdot \Delta K^2_{m,n}$$

$$b^2_n = b^2_n - \lambda \Delta b^2_n$$

$$W_1 = W_1 - \lambda \Delta W_1$$

$$b_1 = b_1 - \lambda \Delta b_1$$

$$W_2 = W_2 - \lambda \Delta W_2$$

$$b_2 = b_2 - \lambda \Delta b_2$$