

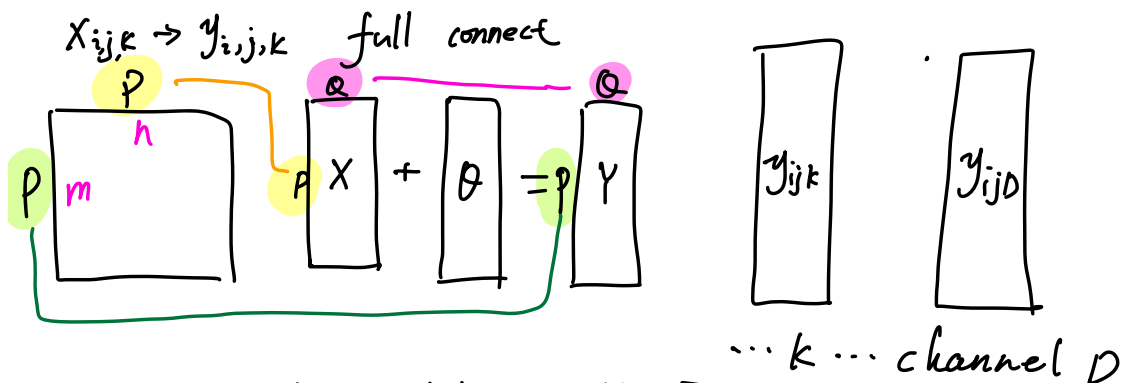
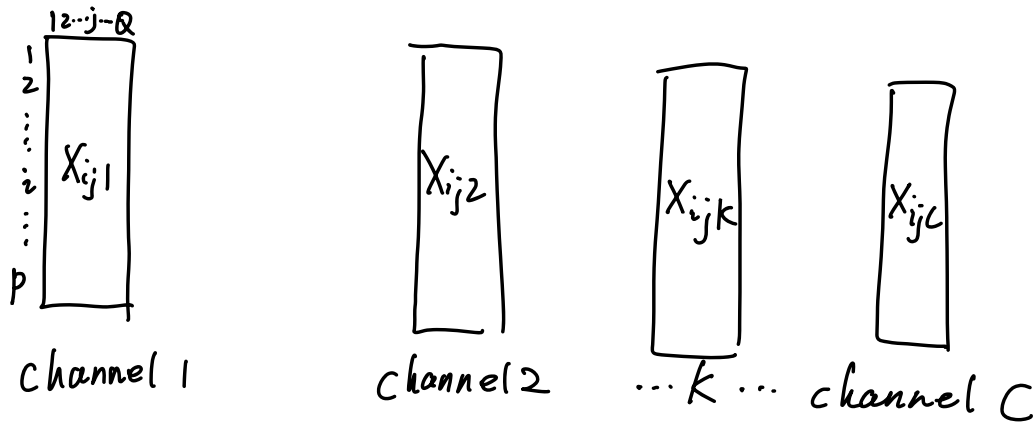
4. In an image application, the input feature maps of spatial size  $P \times Q$  with  $C$  channels and the output feature maps of spatial size  $P \times Q$  with  $D$  channels of a layer of neural network are expressed as  $x_{i,j,k}, 1 \leq i \leq P, 1 \leq j \leq Q, 1 \leq k \leq C$ , and  $y_{i,j,k}, 1 \leq i \leq P, 1 \leq j \leq Q, 1 \leq k \leq D$ , respectively.
- Express the output feature maps  $y_{i,j,k}$  in terms of the input feature maps  $x_{i,j,k}$  via the scalar network parameters of a layer of fully connected neural network. What is the number of learnable parameters? (5 Marks)
  - Express the output feature maps in terms of the input feature maps via the scalar network parameters of a layer of spatial convolutional neural network of filter size  $3 \times 3$ . What is the number of learnable parameters? (5 Marks)
  - Express the output feature maps in terms of the input feature maps via the scalar network parameters of a layer of spatial convolutional neural network of filter size  $1 \times 1$ . What is the number of learnable parameters? (5 Marks)
  - If we use a single index for the 2D spatial position to express the input and output feature maps by  $x_{i,k}, 1 \leq i \leq PQ, 1 \leq k \leq C$ , and  $y_{i,k}, 1 \leq i \leq PQ, 1 \leq k \leq D$ , respectively, re-express the answer to part (c) without using the bias. What is the number of learnable parameters? (5 Marks)
  - Arrange all  $x_{i,k}$  in part (d) into a  $PQ \times C$  matrix,  $\mathbf{X}$ , and arrange all  $y_{i,k}$  in part (d) into a  $PQ \times D$  matrix,  $\mathbf{Y}$ . Re-express the answer to part (d) in the matrix format  $\mathbf{X}$  and  $\mathbf{Y}$ . (Note that  $PQ \times C$  matrix has  $PQ$  rows and  $C$  columns.) What can you conclude about the relation between a convolutional neural network and a Transformer? (5 Marks)

**[Hint for Question 4:** If inputs and outputs of a layer of network are expressed as  $x_i, 1 \leq i \leq P$  and  $y_i, 1 \leq i \leq Q$ , respectively, the outputs of a layer of fully connected neural network can be expressed in terms of the inputs as  $y_i = \sum_{l=1}^P w_{l,i} x_l + b_i, 1 \leq i \leq Q$ , where  $w_{l,i}, b_i$  are the scalar parameters called the weights and biases of the network, respectively.]

22-52-Q4

(a) Q:  $X_{ijk} \rightarrow y_{ijk}$  ? parameters?

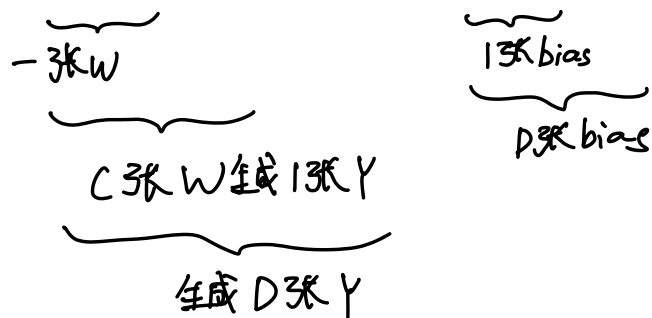
Solution ① understand



$$\sum \begin{bmatrix} W_{11} & W_{12} & \dots & W_{1n} & \dots & W_{1p} \\ W_{21} & W_{22} & \dots & W_{2n} & \dots & W_{2p} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ W_{m1} & W_{m2} & \dots & W_{mn} & \dots & W_{mp} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ W_{p1} & W_{p2} & \dots & W_{pn} & \dots & W_{pp} \end{bmatrix} \begin{bmatrix} X_{11} & \dots & X_{1Q} \\ X_{21} & \dots & X_{2Q} \\ \vdots & \ddots & \vdots \\ X_{i1} & \dots & X_{iQ} \\ \vdots & \ddots & \vdots \\ X_{p1} & \dots & X_{pQ} \end{bmatrix} + \begin{bmatrix} b_{11} & \dots & b_{1Q} \\ b_{21} & \dots & b_{2Q} \\ \vdots & \ddots & \vdots \\ b_{i1} & \dots & b_{iQ} \\ \vdots & \ddots & \vdots \\ b_{p1} & \dots & b_{pQ} \end{bmatrix}$$

$$y_{ijk} = \sum_{l=1}^C \sum_{m=1}^P \sum_{n=1}^P W_{mnlk} X_{ijl} + b_{ijk} (1 \leq k \leq D)$$

parameters  $P \times P \times C \times D + P \times Q \times D$



(b) CNN?



$$y_{i,j,k} = \sum_{u=-1}^1 \sum_{v=-1}^1 \sum_{l=1}^C w_{u,v,l,k} x_{i-u,j-v,l} + b_k$$

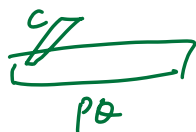
$$3 \times 3 \times C \times D + D = 9 \times C \times D + D$$

$$(c) y_{i,j,k} = \sum_{l=1}^C w_{l,k} x_{i,j,l} + b_k$$

$C \times D + D$  bias  
ReLU D个特征图

$$(d) y_{i,k} = \sum_{l=1}^C w_{l,k} x_{i,l}$$

$C \times D$



$$(e) Y = XW$$



① Both architectures apply matrix multiplications to process input

②  $Y = XW$  shows the CNN perform a linear transformation on the input features at each spatial position, identical across all position.

It is similar to the linear layers used in Transformer, where input are transformed via weight matrices