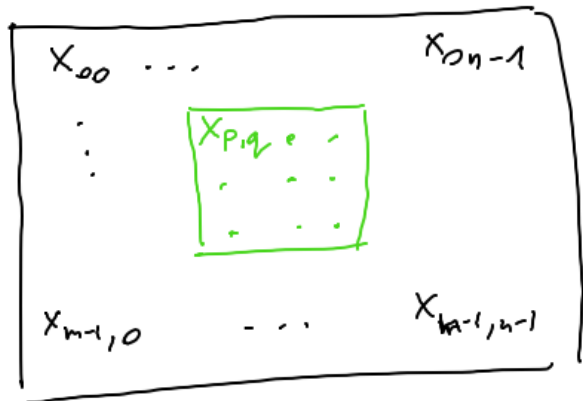


Convolutional Neural Networks (CNN)

LeCun 180 (1989)

input - a matrix $m \times n$

$$(X_{ij})_{i=0, j=0}^{m-1, n-1}$$



kernel of size $k \times k \rightarrow$ weights $(a_{ij})_{i,j=0}^{k-1}$, b - bias

for $p=0, 1, \dots, m-k$, $q=0, 1, \dots, n-k$

$$\text{net}_{p,q} = \sum_{i=0}^{k-1} \sum_{j=0}^{k-1} a_{ij} X_{i+p, j+q} + b$$

activation function, usually ReLU

$$\text{out}_{p,q} = \varphi(\text{net}_{p,q})$$

e.g

input

1.0	0.9	0.3	0	0
1.0	0.8	0.5	0.2	0
1.0	0.9	0.7	0.2	0
1.0	0.8	0.7	0.3	0.1

kernel

1	0	-1
1	0	-1
1	0	-1

bias = 0

net output ($b_{bias} = 0$)

1.5	2.2	1.5
1.1	1.8	1.8

bias = -1.5

0	0.7	0
-0.4	0.3	0.3

↓ $\phi(\cdot)$ ReLU

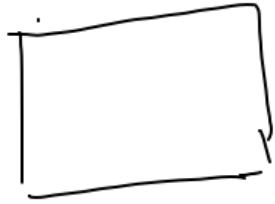
0	0.7	0
0	0.3	0.3



lighter region

light

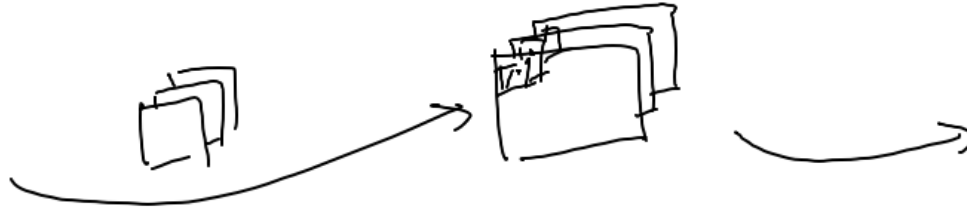
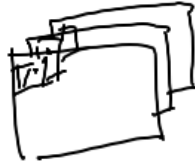
input



several kernels

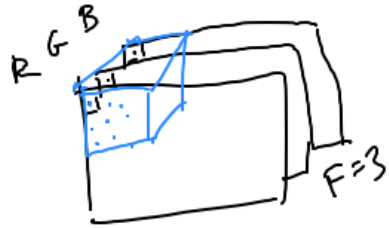


several output matrices



"3x3 kernel" means in fact
"3x3x F" kernel

colour image input



kernel



+ 1 bias



If we have

$$(x_{ij}^{(f)})_{i=0, \dots, m-1, j=0, \dots, m-1, f=1, \dots, F}$$

as input

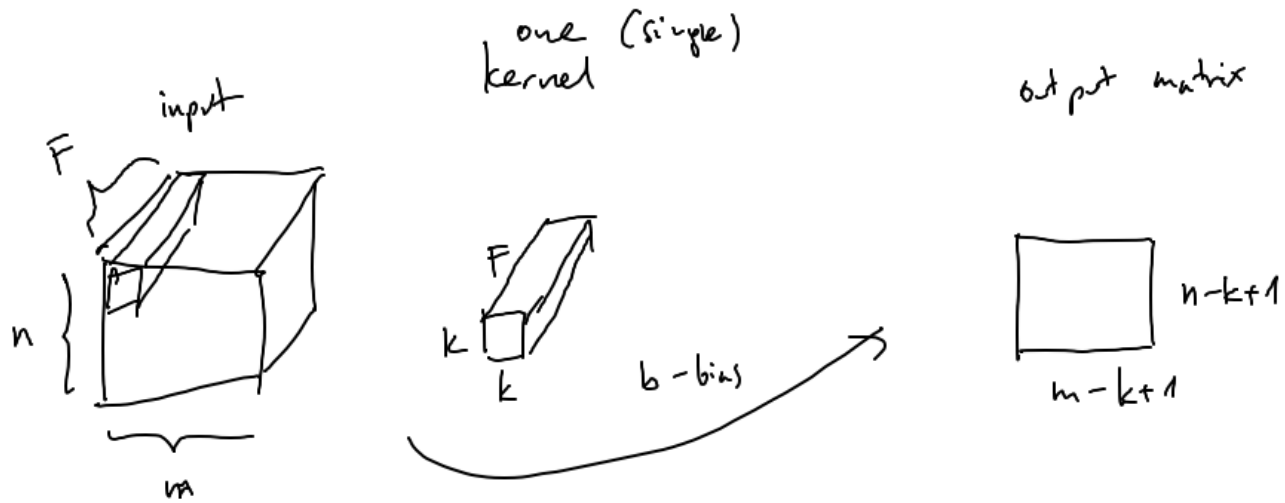
(i.e., F channels of size $m \times n$),
matrices

then the kernel should have $k \times k \times F + 1$ weights:

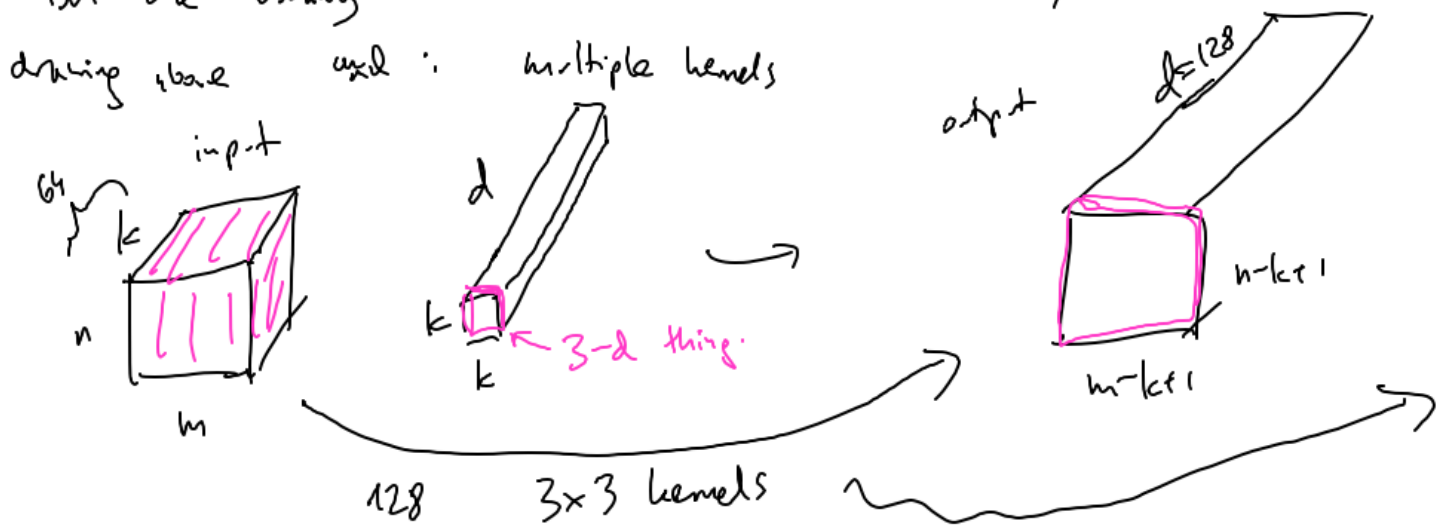
$$(a_{ij}^{(f)})_{i=0, \dots, k-1, j=0, \dots, k-1, f=1, \dots, F}$$

for $p=0, \dots, m-k, q=0, \dots, n-k$:

$$\text{out}_{pq} = \left(\sum_{i=0}^{k-1} \sum_{j=0}^{k-1} \sum_{f=1}^F a_{ij}^{(f)} x_{i+p, j+q}^{(f)} \right) + b$$



But one usually takes more than 1 kernel, ~~that~~ so we skip the dimension F in the drawing above and use: multiple kernels



each kernel: $(3 \times 3 \times 64 + 1)$

128 kernels: $128 \cdot (3 \cdot 3 \cdot 64 + 1)$

= weights

Padding - it allows one to keep the size of the original (image) after convolution: matrix

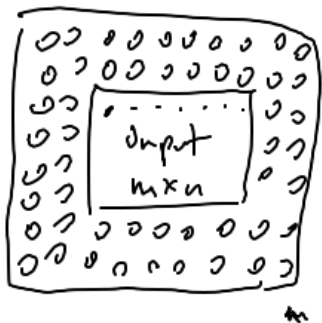
k is usually odd



"Same padding"

pad the input $\frac{k-1}{2}$ - zeros in each direction

e.g.
 $k=5$



convolution with $k \times k$ kernels



after padding:

$$m \times n \rightarrow (m+k-1) \times (n+k-1)$$

$$m \times n$$

the same size

Another advantage of same padding: each input pixel influences ~~the same~~ number of outputs
a more similar

Pooling

transforms an input matrix to a smaller matrix
no weights

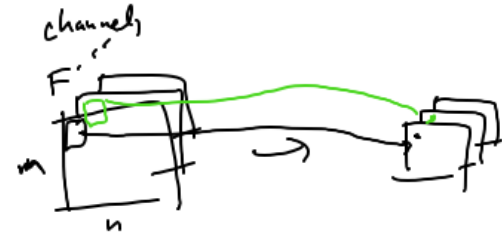
Example:

10x10 matrix

1.0	0			
0.309				
-	0.142			
-	0.403			
-			0.8	
-				
-				
-				

5x5 output

1.0				
	0.4			
			0	



2x2 max pooling: Split the input matrix into 2x2 squares, take maximum of each square as the output
(without overlapping)

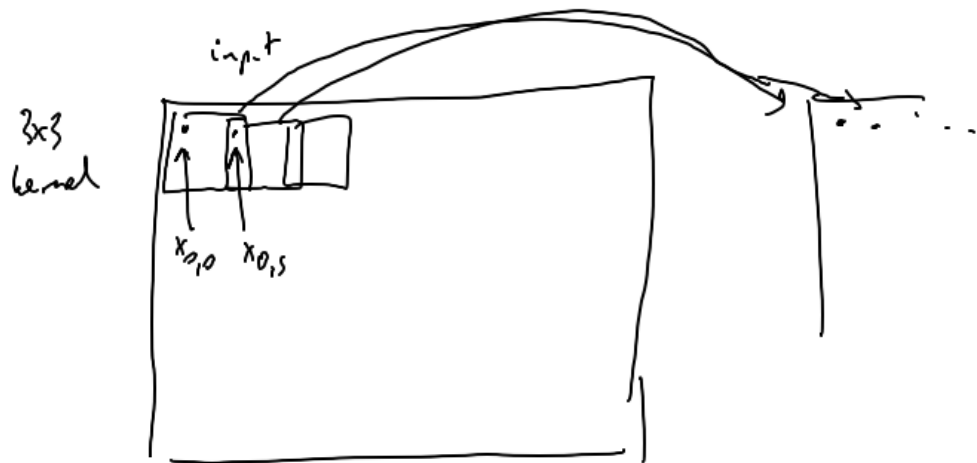
{ average pooling
was used
also

3x3 with overlapping



Stride)

This means that we discard some outputs of the convolutional layer



stride = s means that we take for $p = 0, 1, \dots, \lfloor \frac{m-k}{s} \rfloor$, $q = 0, 1, \dots, \lfloor \frac{n-k}{s} \rfloor$:

$$\text{net}_{p,q} = \sum_{i=0}^{k-1} \sum_{j=0}^{k-1} \sum_{f=1}^F a_{i,j}^{(f)} \cdot x_{ps+i, qs+j}^{(f)}$$

previously: $s = 1$

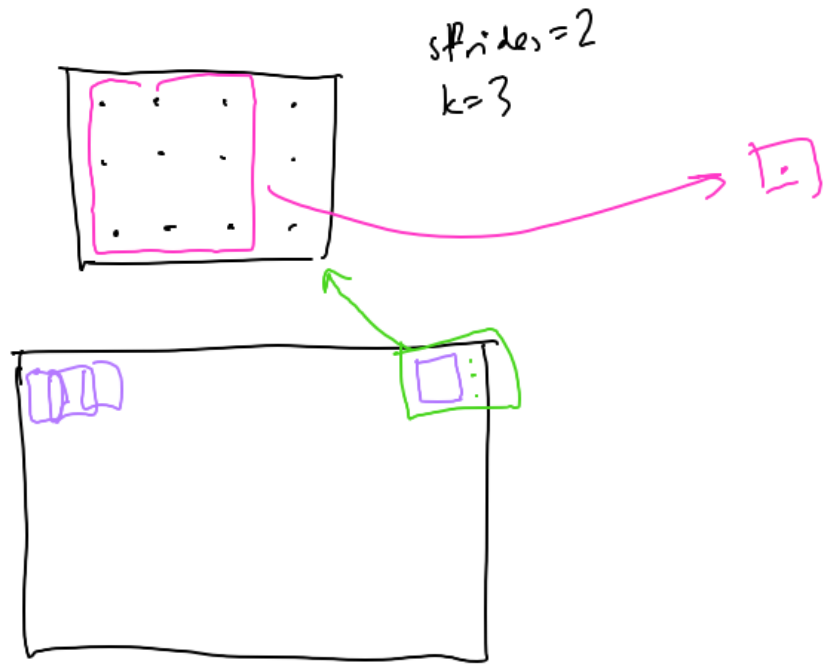
$$ps+i \leq m-1$$

$$ps+k-1 \leq m-1$$

$$ps \leq m-k$$

$$p \leq \frac{m-k}{s}$$

Strides and padding should be chosen carefully to avoid loss of information:
(not necessarily a problem)



CNN:



at the end: usually fully connected layers