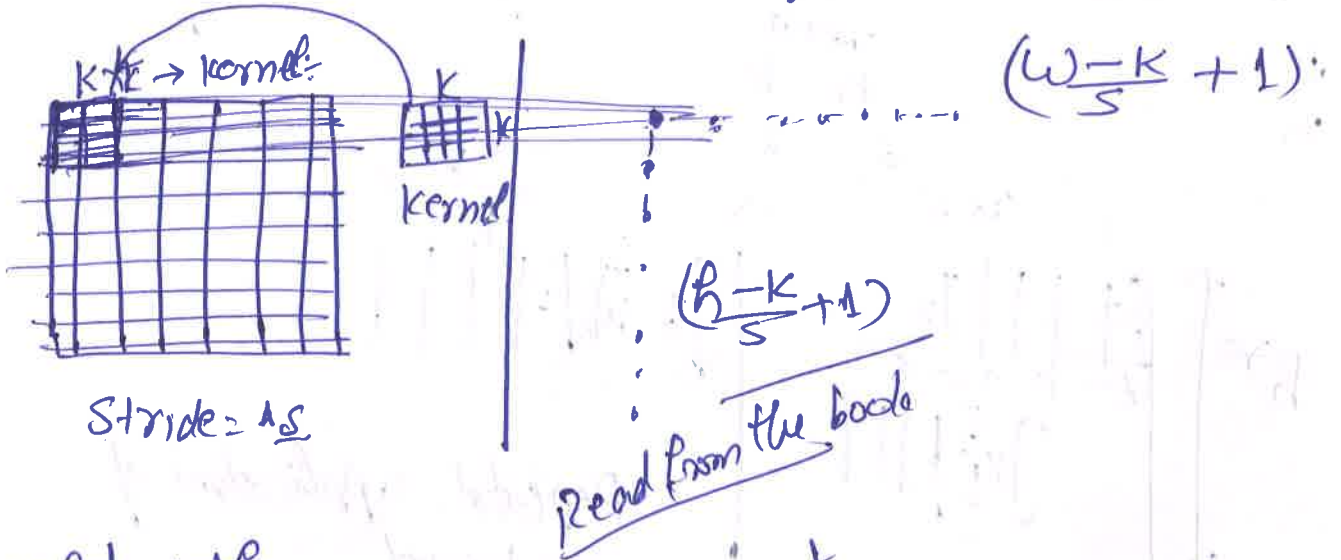


①

COL865
Deep Learning
Convolutional Neural Networks
(CNNs)

First Real Architecture-

At the heart, - convolutional layer



What is

What is a kernel?:-

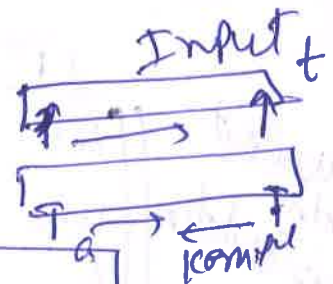
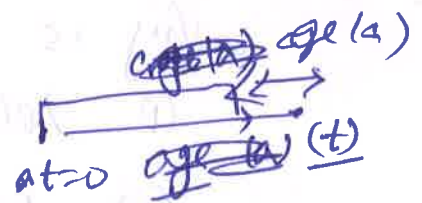
Inspired by the convolution operation:-

$$s(t) = \int_{-\infty}^{\infty} x(a) w(t-a) da$$

:- Two Real valued functions:-

$$x: \mathbb{R} \rightarrow \mathbb{R}$$

$$w: \mathbb{R} \rightarrow \mathbb{R}$$



$$s(t) = (x * w)(t) = (w * x)(t)$$

For discrete case:- Range:- all possible values where x is defined

$$s(t) = \sum_{a=-\infty}^{\infty} x(a) w(t-a)$$

$$x: \mathbb{I} \rightarrow \mathbb{R}$$

$$w: \mathbb{I} \rightarrow \mathbb{R}$$

$$x'(t) = \sum_a x(i-a) k(a)$$

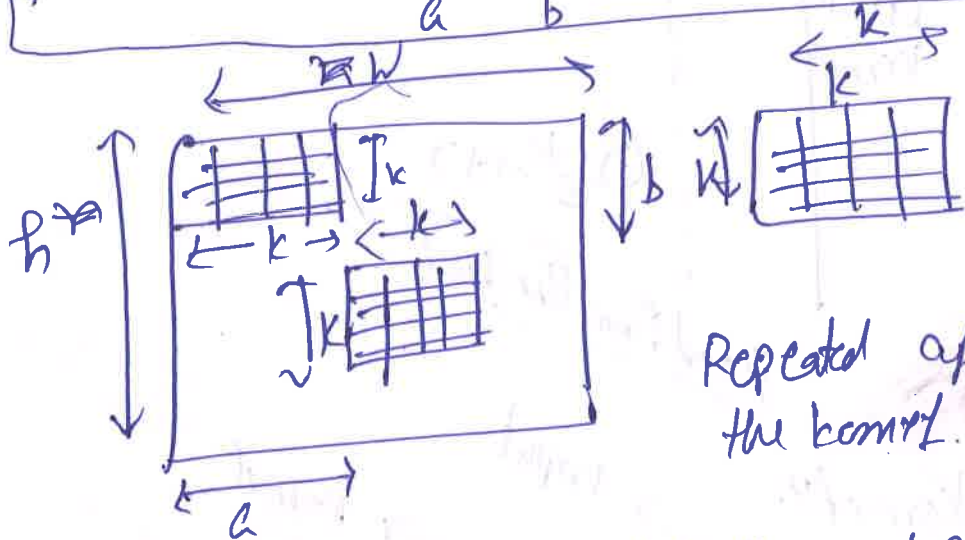
Flipped kernel

⇒

$$X(i) = \sum_a X(i+a) w(a)$$

Two-dimensional:-

$$X^{(d+1)}(i, j) = \sum_{a=1}^a \sum_{b=1}^b X^{(d)}(i+a, j+b) w(a, b)$$



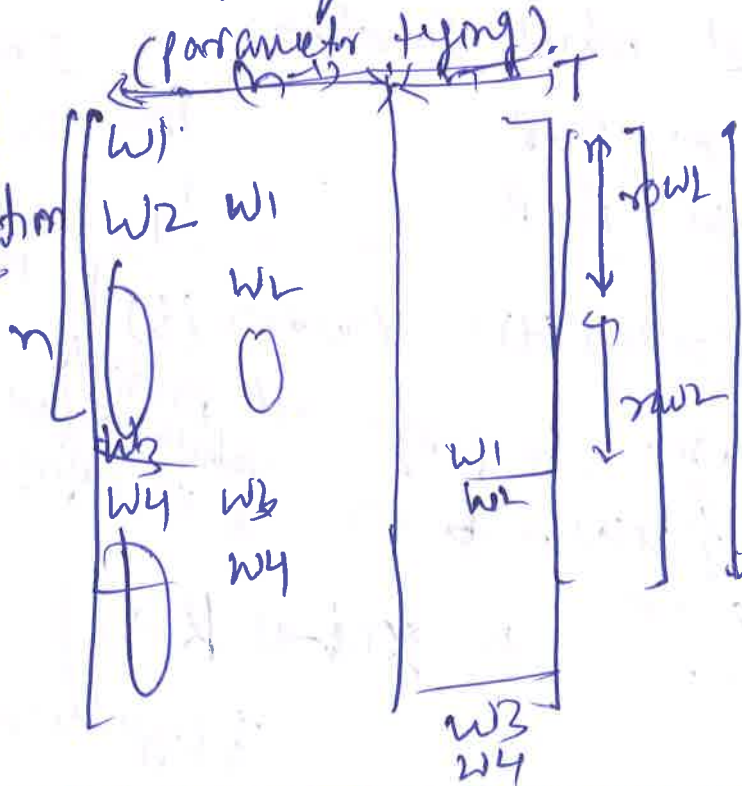
Repeated application of the kernel.

Why is this a useful thing to do?

- ① Localized operation
- ② Significantly less number of parameters.

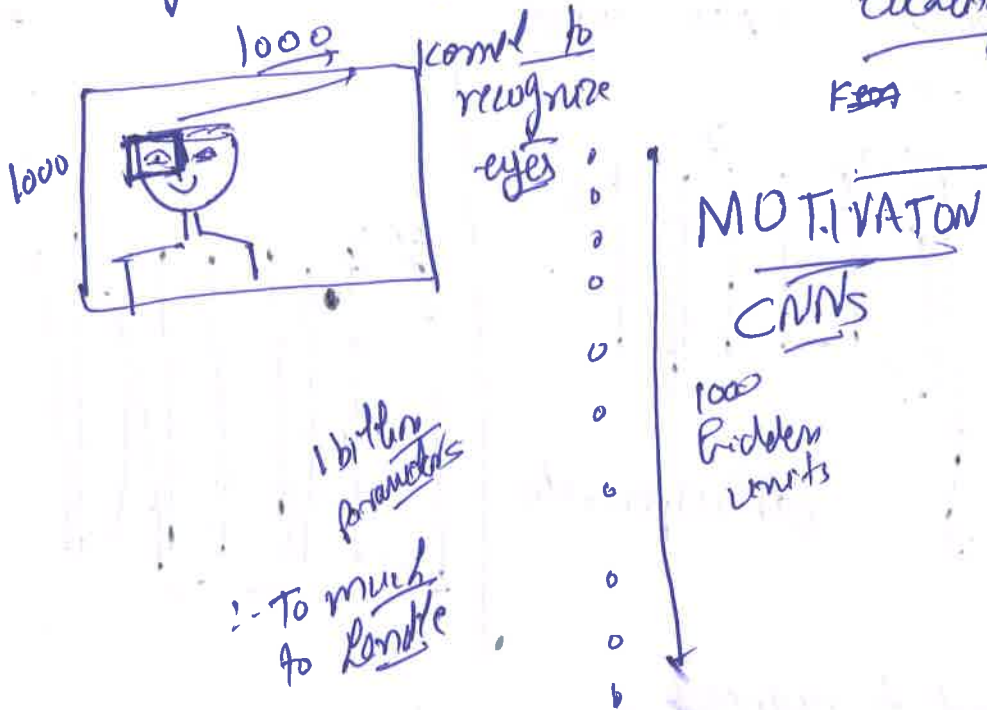
$$\begin{bmatrix} w_1 & w_2 \\ w_3 & w_4 \end{bmatrix}$$

Equivalent matrix representation

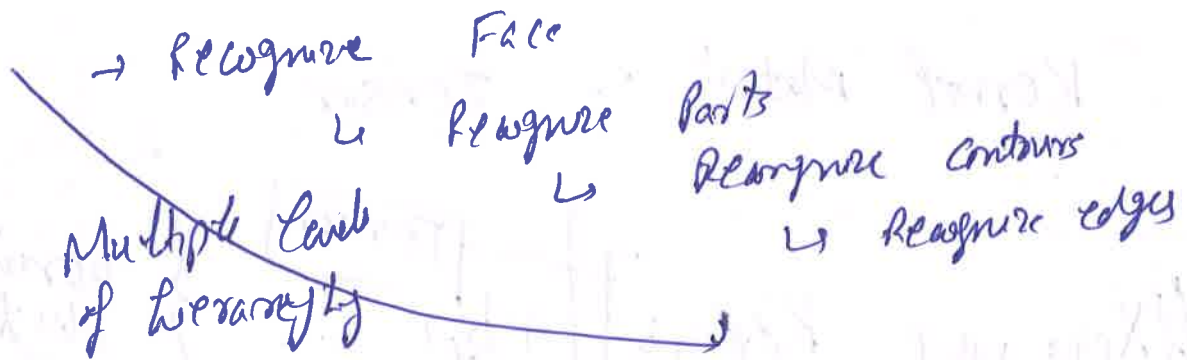


$$k^2 (n-k+1)^2$$

② Why this operation? ~~Locate~~ Exploit spatial locality

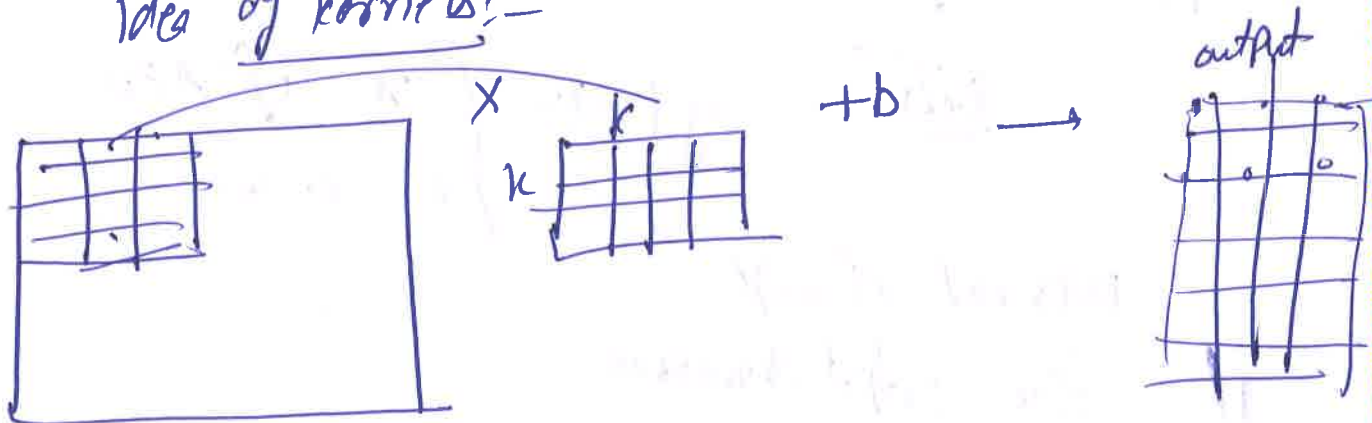


Exploit:- ① Spatial arrangement (Locality) (locality)
② Compositionality

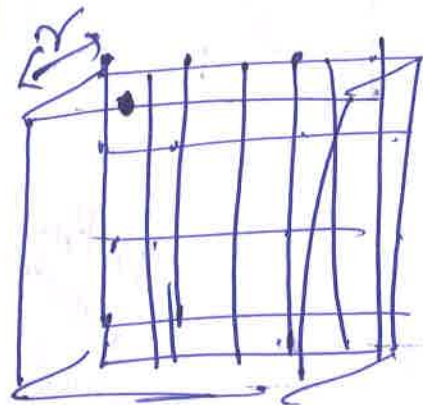
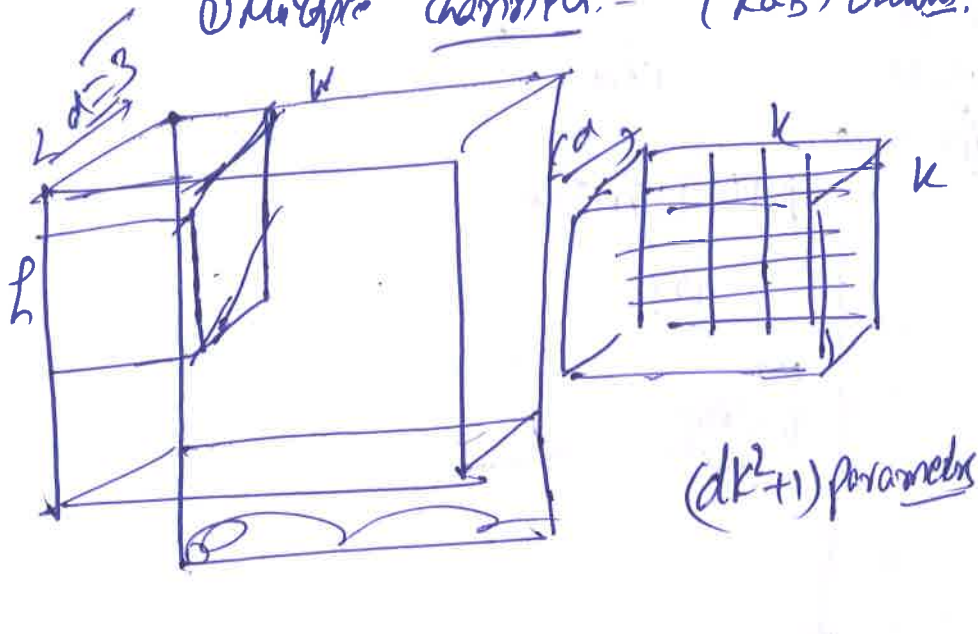


Instead of full matrix multiplication, we ~~space~~ sparse representation

Idea of kernels:-



④ * 1. Multiple kernels: - (RGB) values.



But - we need to recognize multiple such "objects"

\Rightarrow multiple kernels: δ

Kernel "Matrix" :- Tensor

$$\left[\begin{matrix} x'_{i,j,e} \\ x_{i+a,j+b,e} \\ K_{e,a,b} \end{matrix} \right] + b_e \quad \left. \begin{matrix} \text{Tensor} \end{matrix} \right\} \text{Derive the formula when wrong stride } s \neq 1$$

Followed by non-linearity.

ReLU

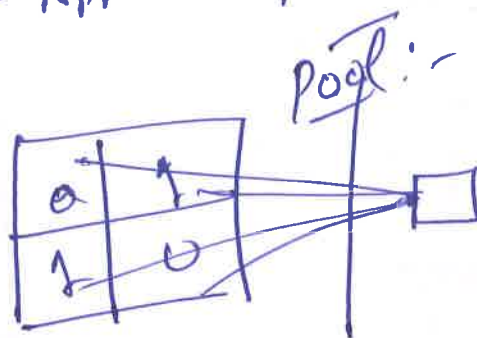
$$g(x) = \begin{cases} x & \text{if } x \geq 0 \\ 0 & \text{o.w} \end{cases}$$

Does not change the output dimension

③

(Max)-Pooling Operation

:- Applied after convolution:-



} Helps to reduce the size of the representation

$$x'(i,j) =$$

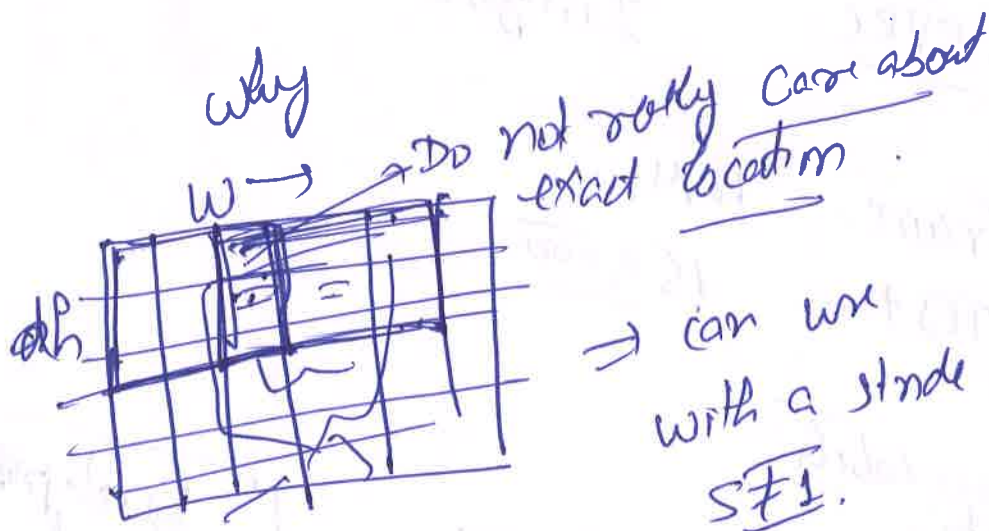
$$\begin{aligned} & \text{Max}_{a,b} x(i+a, j+b) \\ & a \in \{0,1\} \\ & b \in \{0,1\} \end{aligned}$$

Three steps:-

① Linear kernel application

② Non-linear transformation

③ pooling



Another important idea:-

zero-padding

Output layer (shrinks)

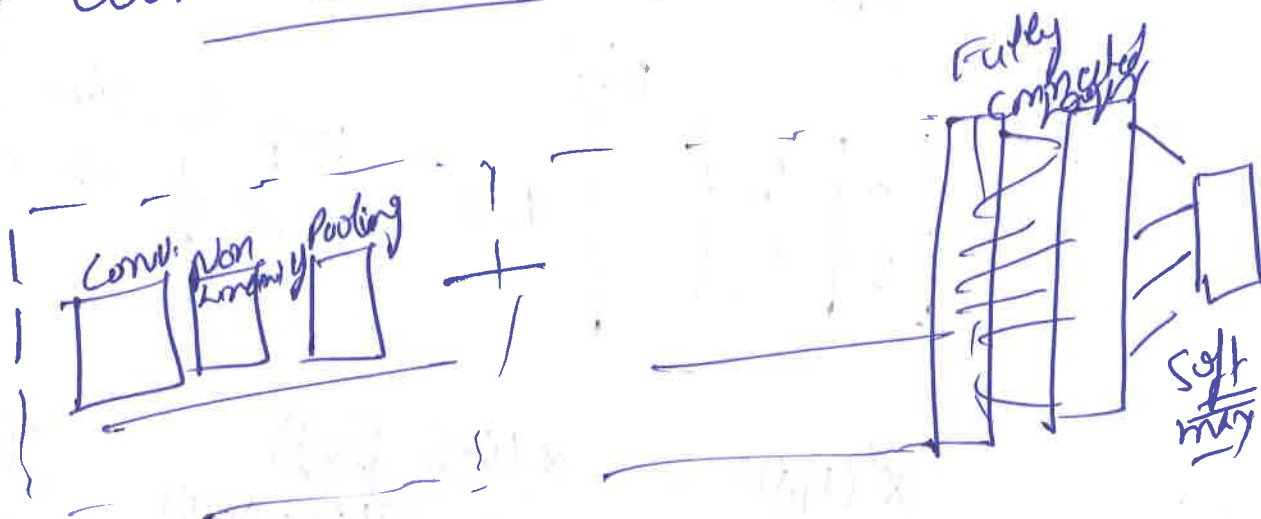
$$(w-k+1) \times (h-k+1)$$

⇒ Pad zeroes:- $(K-1)$ zeros on each side. 339

Add more zeroes than $(K-1)$ - 'Full' convolution in the book

~~Homework~~

Overall Architecture -



Alexnet. (2012)

ILSVRC

Imagenet Challenge

Train: - 12M

Test: - 150,000

Top 5% error rate: - 1000 labels
26.2% to 15.4%

Backpropagation
read from the book

Notion of equivariance.

But not to scaling & rotation

f is equivariant w.r.t g

$$f(g(u)) = g(f(u))$$

convolution (f) is equivariant w.r.t translation (g)