# DS 4440    Conv Nets

So far: Simple feed-forward Networks



$$\begin{bmatrix} \square & & & \\ & & & \\ & & & \\ & & & \square \end{bmatrix} \rightarrow [\; \square \ldots \square \ldots \square \ldots \square \;]$$

$(1,1) \quad (k-1, \ell) \quad (d, d)$

$(k, \ell)$

$X \in \mathbb{R}^{d^2}$

$d^2$

**Problem:** This discards Structure

Consider

$$\begin{bmatrix} \boxed{6} & & \\ & & \\ & & \end{bmatrix} \quad \text{vs.} \quad \begin{bmatrix} & & \\ & & \boxed{6} \\ & & \end{bmatrix}$$

$$\downarrow \qquad\qquad\qquad\qquad \downarrow$$
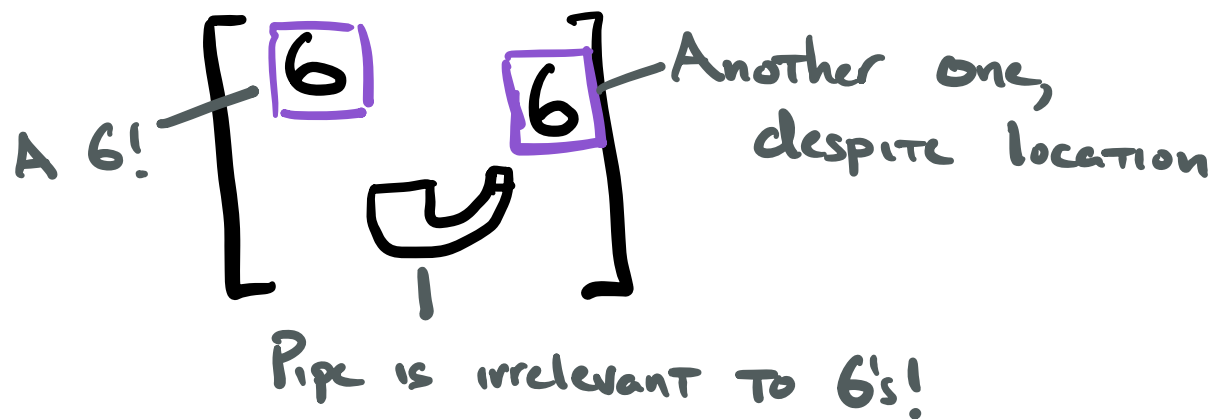
$$[\; \boxed{\;} \; ] \qquad\qquad [\qquad \boxed{\;} \;]$$

The **parameters** associated with these 6's are distinct!

**Upshot** Model must learn to spot 6 in all locations.

But: a 6 is a 6! **Invariance**

A 6! [ 6    6 ] Another one, despite location

Pipe is irrelevant to 6's!

Enter **Convolutional Neural Networks** (aka ConvNets aka CNNs)

**Basic Idea:** Slide **Windows** ☐ over inputs to yield **local features**

$W \in \mathbb{R}^{10 \times 10}$

Dense    Out

 = local features or activations

Formally assume region A starting at $(A_x, A_y)$.

$$\boxed{\ } = \sigma\left(\sum_{i,j} W_{i+j} \cdot A_{ij}\right)$$

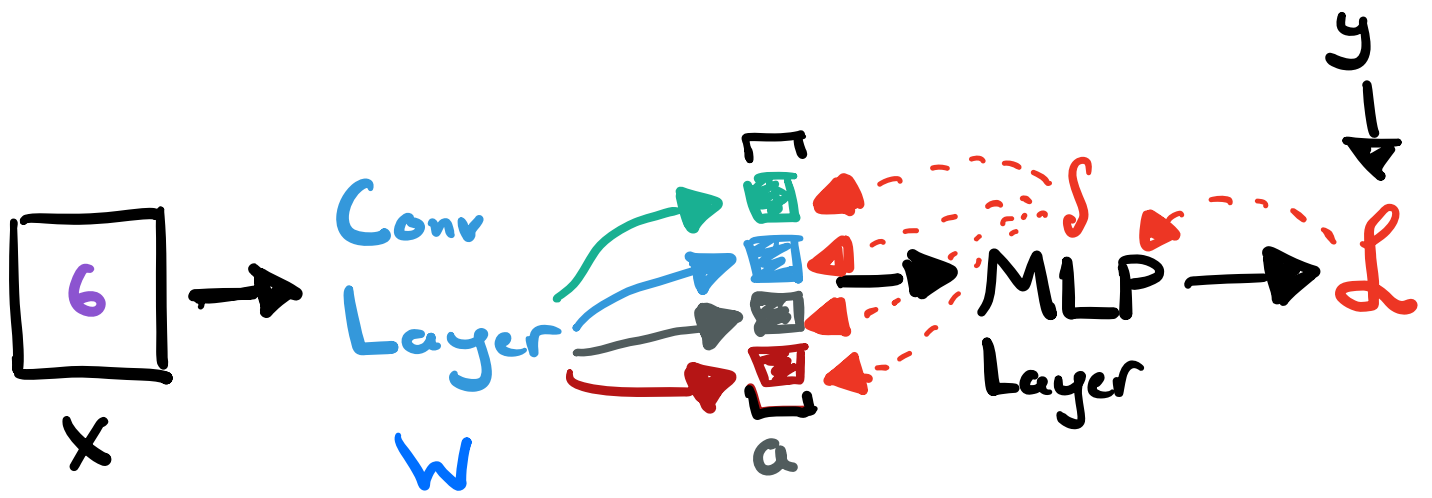Sum over window size!

$$= \sigma\left(\sum_{ij} W_{i+j} X_{(A_x + j)(A_y + i)}\right)$$

$$\boxed{\ } = \sigma\left(\sum_{i,j} W_{i+j} \cdot B_{ij}\right)$$

The key is that weights $W$ are Shared across regions. $W$ is called the Kernel or filter.

Q What about $\nabla_w \mathcal{L}$ here?

$$\vec{\nabla}_W \mathcal{L} = \sum_{a_k \in a} \left( \vec{\nabla}_W a_k \cdot \frac{\partial \mathcal{L}}{\partial a_k} \right)$$

A concrete example

\* Adapted from "Deep Learning" by Goodfellow

$$3 \begin{bmatrix} a & b & c & d \\ e & f & g & h \\ i & j & k & l \end{bmatrix} \overset{4}{\phantom{.}} \otimes 2 \begin{bmatrix} w & x \\ y & z \end{bmatrix} \overset{2}{\phantom{.}}$$

Input                    Kernel

$$\longrightarrow 2 \begin{bmatrix} w \cdot a + x \cdot b \\ + y \cdot e + z \cdot f & w \cdot b + x \cdot c \\ + y \cdot f + z \cdot g & w \cdot c + x \cdot d \\ + y \cdot g + z \cdot h \\ w \cdot e + x \cdot f \\ + y \cdot i + z \cdot j & w \cdot f + x \cdot g \\ + y \cdot j + z \cdot k & w \cdot g + x \cdot h \\ + y \cdot k + z \cdot l \end{bmatrix} \overset{3}{\phantom{.}}$$
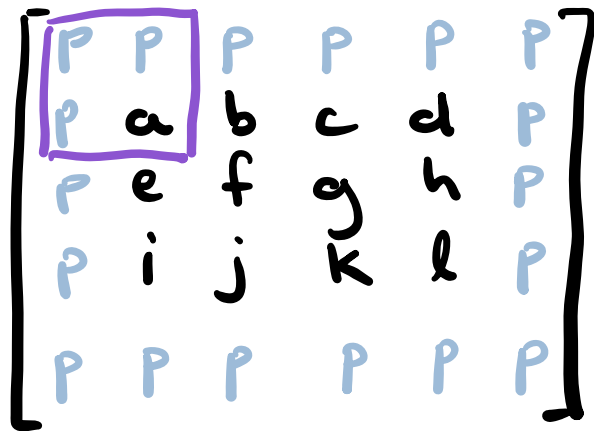
Activation

In the Simplist case — as here —

$$(X_h - W_h + 1, X_w - W_w + 1)$$

Input height    Kernel height    Input Width    Kernel Width

# Padding

$$\begin{bmatrix} P & P & P & P & P & P \\ P & a & b & c & d & P \\ P & e & f & g & h & P \\ P & i & j & K & \ell & P \\ P & P & P & P & P & P \end{bmatrix}$$
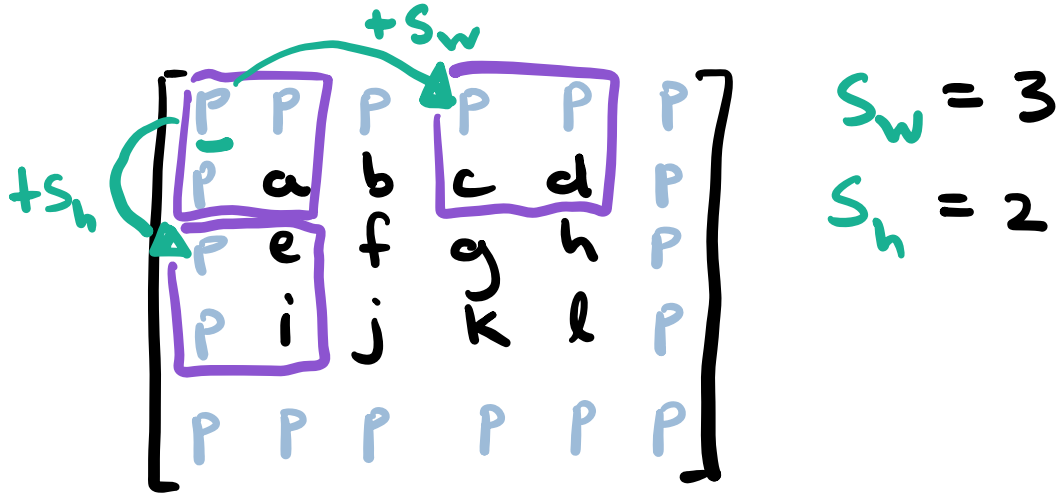
Here we have Padded height & Width by 1.

→ Output Shape

$$(X_h + P_h - W_h + 1, \\ X_w + P_w - W_w + 1)$$

Stride size has to do with how we pass kernels over inputs — have assumed size of 1 above.
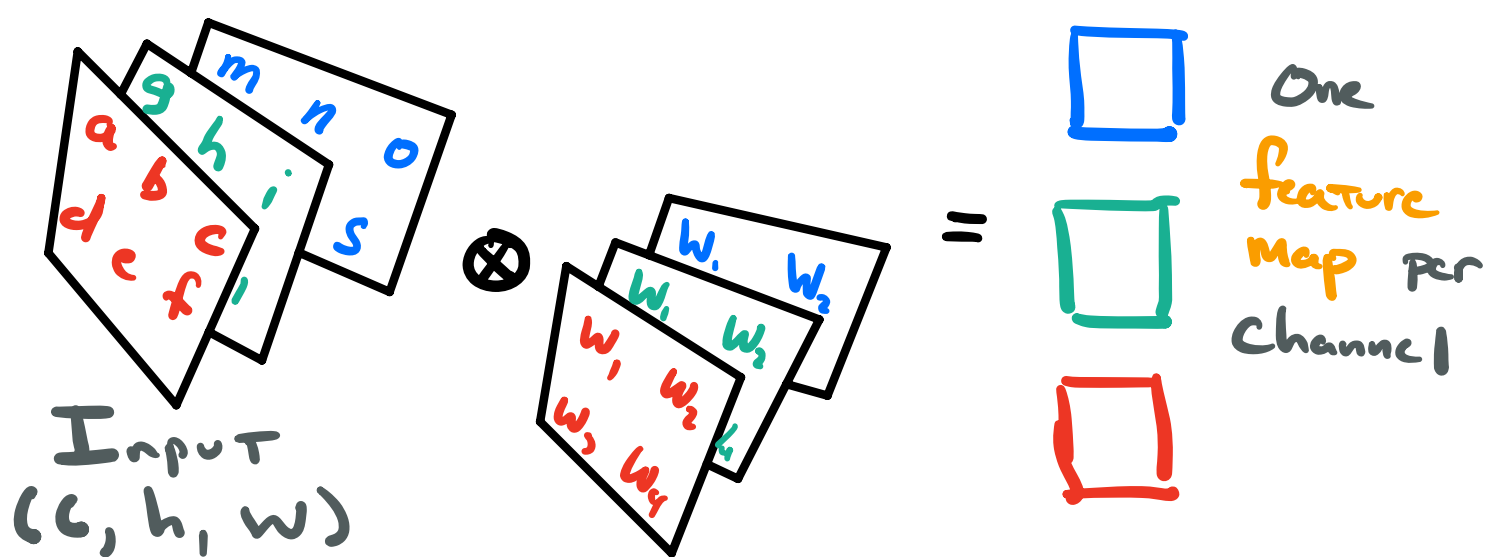
$S_w = 3$

$S_h = 2$

Also affects **Output Size**

$$\left( \frac{(X_h + P_h - W_h + S_h)}{S_h}, \frac{(X_w + P_w - W_w + S_w)}{S_w} \right)$$
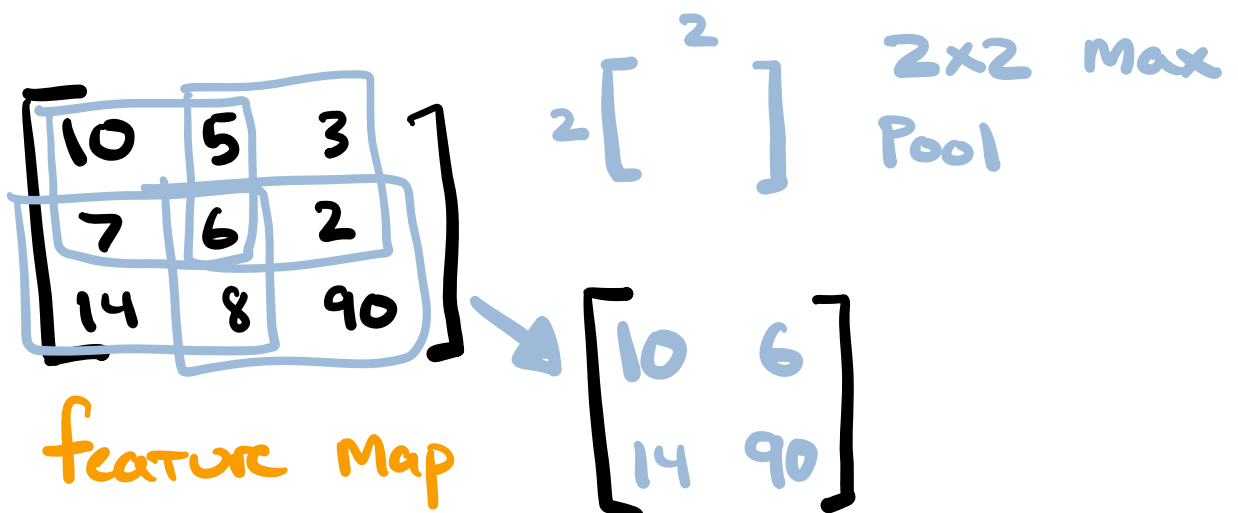
## Higher Dimensions

Above we assumed 2d inputs but we can extend to arbitrary Tensors. Useful, e.g. for RGB images.



Input $(C, h, w)$

One feature map per channel

Can <u>S</u>um for final map, or use
another Pooling strategy.

feature → Pool → Smaller (often Scalar)
Map              Output $f$

<u>M</u>ax <u>pooling</u> extract Max Value

$$2\begin{bmatrix} & 2 \\ & \end{bmatrix}$$ 2×2 Max Pool

$$\begin{bmatrix} 10 & 5 & 3 \\ 7 & 6 & 2 \\ 14 & 8 & 90 \end{bmatrix}$$   $$\begin{bmatrix} 10 & 6 \\ 14 & 90 \end{bmatrix}$$
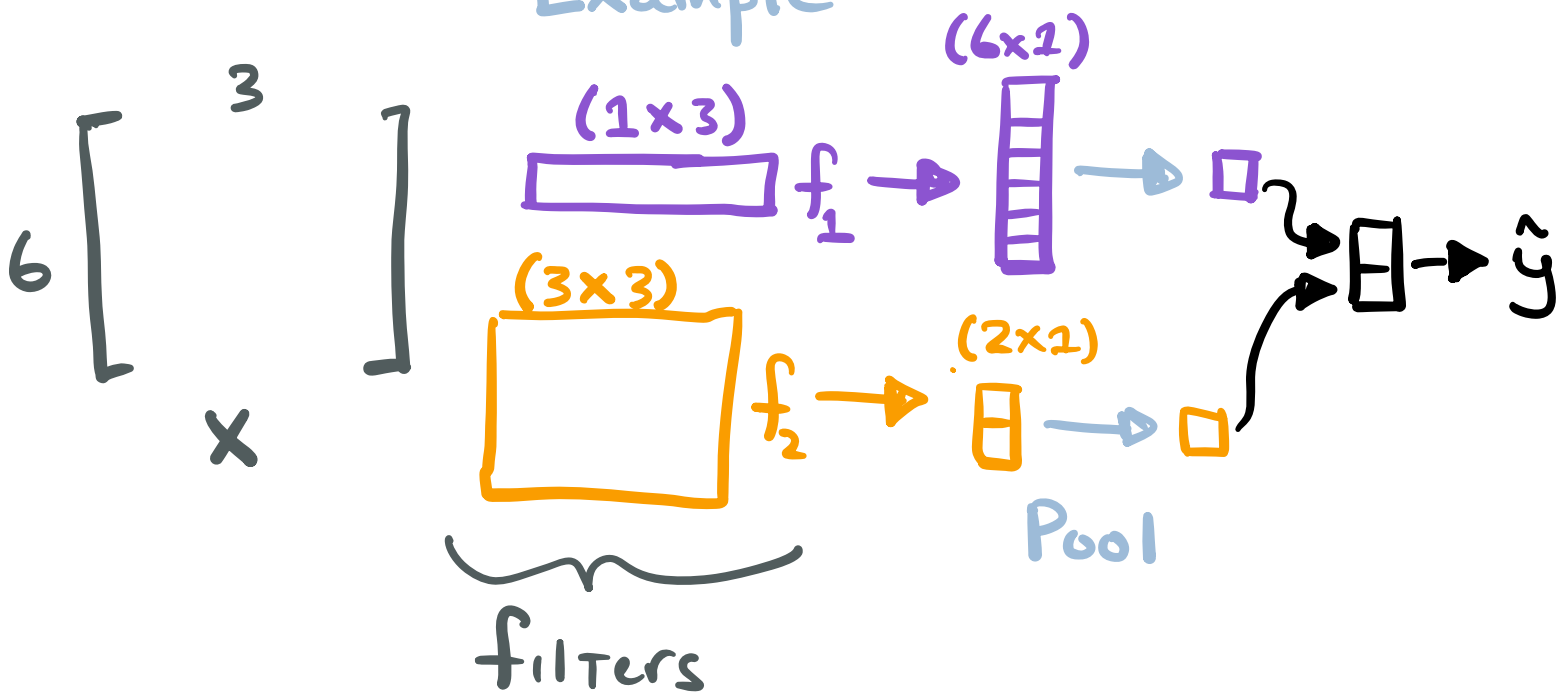
feature Map

Q What happens to ▽ Through
Pooling (Max) operation?

<u>I</u>NTUITION Individual filters may
<u>Specialize</u> at recognizing Something
(Say, Cats) → Max pooling is like
asking "does This window contain a Cat?"
repeatedly (More on HW4!)

Often we define multiple independent filters, with same or varying size.

Example

$$6 \begin{bmatrix} & 3 & \\ & & \\ & & \\ & & \end{bmatrix}$$
X

(1×3)
$f_1 \rightarrow$

(6×1)

$\rightarrow$ □

(3×3)
$f_2 \rightarrow$

(2×1)
$\rightarrow$ □

$\rightarrow \hat{y}$

Pool

filters

Concretely, consider:

$$\begin{bmatrix} 5 & 3 & 1 \\ 6 & 8 & 1 \\ 2 & 3 & 4 \end{bmatrix}$$
X

• No padding
• Two (1×2) filters
• Max pooling (1×2)

[6  8] $f_1$    [-2  3] $f_2$

$$\begin{bmatrix} 6\cdot5 + 8\cdot3 & 6\cdot3 + 8\cdot1 \\ 6\cdot6 + 8\cdot8 & 6\cdot8 + 8\cdot2 \\ 6\cdot2 + 8\cdot3 & 6\cdot3 + 8\cdot4 \end{bmatrix} = \begin{bmatrix} 54 & 26 \\ 100 & 56 \\ 36 & 50 \end{bmatrix} \rightarrow \begin{bmatrix} 54 \\ 100 \\ 50 \end{bmatrix}$$

feature map 1

Pooling (1x2)

$$\begin{bmatrix} -2\cdot5 + 3\cdot3 & -2\cdot3 + 3\cdot1 \\ -2\cdot6 + 3\cdot8 & -2\cdot8 + 3\cdot1 \\ -2\cdot2 + 3\cdot3 & -2\cdot3 + 3\cdot4 \end{bmatrix} = \begin{bmatrix} -1 & -3 \\ 12 & -13 \\ 5 & 6 \end{bmatrix} \rightarrow \begin{bmatrix} -1 \\ 12 \\ 6 \end{bmatrix}$$

$\oplus \rightarrow \begin{bmatrix} 54 & 100 & 50 & -1 & 12 & 6 \end{bmatrix}$