

# Zhang - Convolutional Neural Networks

Images are grid shape data

MLPs flatten this grid into vectors

MLP vectors are invariant to the order  
of the "pixels"

Same result regardless of order

However, nearby pixels tend to be related // important info held in these relationships

CNNs (Le Cun, 1995)

Image net (Deng 2009)

Convnets (Krizhevsky, 2012)

Compared to MLP, CNN is more competent.

- Require fewer params
- convolutions can be parallelized across gpu

## 7.1. fully connected to convolutions

### Invariance

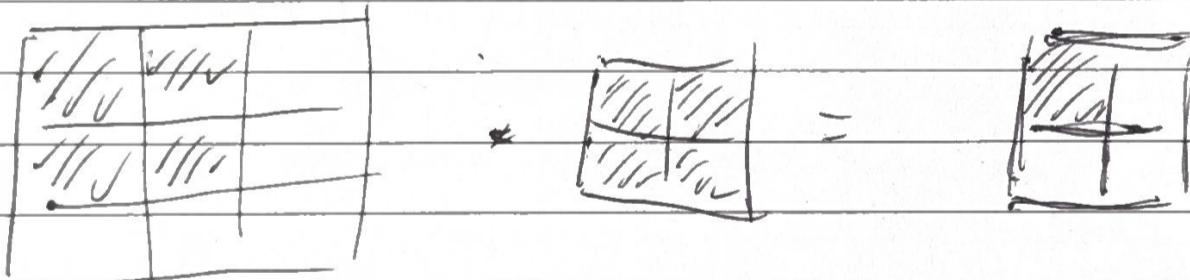
when recognizing objects in a image, we likely do not want to be too concerned w/ precise location in the image

## 7.2. Convolutions for images

Convolutions = cross correlation

- An input layer
- kernel

(~~is~~) one combined to produce an output tensor



Aggregates down shape

the two params of a convolution:

- the kernel
- A scalar bias

ernels are generally randomized

## Cross correlation vs convolution

A true convolution can be performed by flipping the kernel vertically & horizontally

Same trained kernel

Cross-correlation is the standard procedure

### 7.2.6 feature map & receptive field

Output of a convolution layer is called a feature map

it can be regarded as the learned representations (features) in the spatial dimensions ( $w, h$ )

the receptive field

( $\Leftrightarrow$  given an element/Pixel

the RF is the elements from the prev field that went into creating it

### 7.3 Padding & Stride

P & S offer more control over the size of the output

Assuming the kernel  $w$  or  $h \geq 1$ , we tend to "run out" of space to stride & the image shrinks

the shrunk size is lost info

Padding is the most popular tool to handle <sup>this</sup>

that sometimes we want to reduce dimensions

↳ Sided convolutions

Padding = retain size

Sided Convolutions = Reduce dimensions

#### 7.4 Multiple input & multiple output Channels

when adding in channels both the input & hidden layers become 3-D

$$3 \times n \times w$$

3 = channel dim

kernel needs to have same number of channels as input

(if more than 1 channel)

Multi-channels = multi-filters

(5)

each filter acts on each channel

↳ 3 outputs are then concat into 1

this happens for each filter so the output pairs from input to output remain the same

each output layer is a product of all the channels

### 7.4.2. Multi output channels

$1 \times 1$  convolutions have the effect of reducing the number of channels

## 7.5. Pooling

the deeper we go into the network, the larger the receptive field to which each hidden node is sensitive

Pooling layers

↳ Mitigating sensitivity of convo layers to location

↳ Spatially down sampling representations

### 7.5.1 Max & Average Pool

Similar to conv layer that is slid over

Contains no parameters

Instead deterministic

↳ Avg or max the kernel values

Akin to down-sampling

Can be adjusted w/ Padding & Stride  
to change output size

unlike conv layers

↳ Pooling acts on each channel

↳ Does not concat or aggregate

channels remain

### 7.6 Conv Neural Networks - (heNet)

fully functional CNN

high level =

- 2 conv layers
- Dense Block 3 MLPs

(7)

Image  $\rightarrow$  Conv  $\rightarrow$  Pool  $\rightarrow$  Conv  $\rightarrow$  Pool  $\rightarrow$  MCP  $\times 3$

Conv layer = Conv & sigmoid activation

Pool = Average

$\hookrightarrow$  ReLU & Max not discovered yet

Conv layers have 8 & 16 filters turning the 2-D input into having channels

Pooling is used to downsample by factor of four

2nd conv also has a downsample effect

output of conv is flattened to be pass MCP

$\hookrightarrow$  take minibatches & flatten

first conv uses padding to compensate to retain size

$\hookrightarrow$  2nd foregoes which is why decreases size