#### **CS6421: Deep Neural Networks**

#### **Gregory Provan**

#### Spring 2020 Lecture xx: Convolution Neural Networks

Based on notes from John Canny, Ismini Lourentzou

#### **Overview**

- Introduction
- Applications of CNNs
- CNN Operations
  - Convolution
  - Pooling
  - Classification

#### Introduction to CNNs

- A CNN is a feed-forward network that can extract topological properties from an image.
- Like almost every other neural network, they are trained with a version of the back-propagation algorithm.
- Convolutional Neural Networks are designed to recognize visual patterns directly from pixel images with minimal preprocessing.
- They can recognize patterns with extreme variability (such as handwritten characters).

#### Classification





#### Architecture: Fully-Connected vs CNN

- We know it is good to learn a small model.
- From this fully connected model, do we really need all the edges?
- Can some of these be shared?



## **CNN Topology**



#### **Overview**

#### Introduction

### Applications of CNNs

## CNN Operations

- Convolution
- Pooling
- Classification

## **Applications of CNNs**

- Game playing
- Speech
- Text classification

#### AlphaGo: Playing the Game Go



#### 19 x 19 matrix

- Black: 1
- white: -1

none: 0

#### Fully-connected feedforward network can be used

But CNN performs much better

### **CNN in speech recognition**





Source of image: http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.703.6858&rep=rep1&type=pdf

#### **Overview**

- Introduction
- Applications of CNNs
- CNN Operations
  - Convolution
  - Pooling
  - Classification

#### **Convolutional Neural Networks**

Subsampling

Convolutions



Convolutions

Image Credit: Yann LeCun, Kevin Murphy

Full connection

Full connection

Subsampling

Gaussian connections

#### Feature extraction layer or Convolution layer

Detect the same feature at different positions in the input image.



#### Feature extraction



#### **Feature extraction**

- Shared weights: all neurons in a feature share the same weights (but not the biases).
- In this way all neurons detect the same feature at different positions in the input image.
- Reduce the number of free parameters.



#### **Feature extraction**

If a neuron in the feature map fires, this corresponds to a match with the template.



#### **Convolutional Layer**



Slide Credit: Marc'Aurelio Ranzato

#### **Convolutional Layer**



**Preview:** ConvNet is a sequence of Convolutional Layers, interspersed with activation functions



Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

#### Extraction of Feature Maps at Multiple Levels





Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n



Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n

the subsampling layers reduce the spatial resolution of each feature map

Sy reducing the spatial resolution of the feature map, a certain degree of shift and distortion invariance is achieved.



## subsampling layers reduce the spatial resolution of each feature map



#### Subsampling layer





# weight sharing is also applied in subsampling layers.



weight sharing reduces the effect of noise and shift or distortion



#### **Process of CNN Inference**



30

#### Image- and Region-Specific Filters

- CNNs use specific filters
- Detect specific image properties
  - Multiple levels
    - Low-level: edge
    - Higher-level: face

#### Consider learning an image:

#### Some patterns are much smaller than the whole image

Can represent a small region with fewer parameters



### **CNN for Image Analysis**

Same pattern appears in different places: They can be compressed! What about training a lot of such "small" detectors and each detector must "move around".



#### A convolutional layer

- A CNN is a neural network with some convolutional layers (and some other layers).
- A convolutional layer has a number of filters that performs a convolution operation.



#### What's a convolution?

#### Basic idea:

- Pick a 3x3 matrix F of weights
- Slide this over an image and compute the "inner product" (similarity) of F and the corresponding field of the image, and replace the pixel in the center of the field with the output of the inner product operation

#### • Key point:

- Different convolutions extract different types of low-level "features" from an image
- All that we need to vary to generate these different features is the weights of F

#### Convolving an image with an ANN

## Note that the parameters in the matrix defining the convolution are **tied** across all places that it is used

#### input neurons

000000000000000000000000000000000000000	first hidden layer
000000	
	000000000000000000000000000000000000000
000000000000000000000000000000000000000	
000000000000000000000000000000000000000	000000000000000000000000000000000000000
000000000000000000000000000000000000000	
000000000000000000000000000000000000000	000000000000000000000000000000000000000
000000000000000000000000000000000000000	000000000000000000000000000000000000000
	000000000000000000000000000000000000000
000000000000000000000000000000000000000	000000000000000000000000000000000000000
000000000000000000000000000000000000000	

# How do we do many convolutions of an image with an ANN?

 $28 \times 28$  input neurons

first hidden layer:  $3\times 24\times 24$  neurons


### Example: 6 convolutions of a digit

#### http://scs.ryerson.ca/~aharley/vis/conv/



### CNNs typically alternate convolutions, nonlinearity, and then downsampling

Downsampling is usually averaging or (more common in recent CNNs) max-pooling

# Why do max-pooling?

- Saves space
- Reduces overfitting?

PROC. OF THE IEEE, NOVEMBER 1998

- Because I'm going to add more convolutions after it!
  - Allows the short-range convolutions to extend over larger subfields of the images
    - So we can spot larger objects
    - Eg, a long horizontal line, or a corner, or …



 $\overline{7}$ 

Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

### Another CNN visualization

https://cs.stanford.edu/people/karpathy/convnetjs/demo/mnist.html



### **Alternating Convolution and Sub-sampling**



5 layers up

The subfield in a large dataset that gives the strongest output for a neuron



### **Convolution filters**



6 x 6 image



Each filter detects a small pattern (3 x 3).





stride=1



6 x 6 image



If stride=2

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

-1 1 -1 -1 -1 1 -1 -1 1

#### Filter 1



6 x 6 image





6 x 6 image



### Filter 1





 -1
 1
 -1

 -1
 1
 -1

 -1
 1
 -1

#### Filter 2

#### stride=1



6 x 6 image

### Repeat this for each filter



Two 4 x 4 images Forming 2 x 4 x 4 matrix

### Color image: RGB 3 channels



### Convolution v.s. Fully Connected



Fullyconnected







### The whole CNN



# Max Pooling



Filter 1









### Subsampling pixels will not change the object bird



We can subsample the pixels to make image smaller

fewer parameters to characterize the image

- A CNN compresses a fully connected network in two ways
  - Reducing number of connections
  - Shared weights on the edges
  - Max pooling further reduces the complexity

## Max Pooling



6 x 6 image

2 x 2 image Each filter is a channel

# Why do max-pooling?

- Saves space
- Reduces overfitting?
- Because we add more convolutions after it!
  - Allows the short-range convolutions to extend over larger subfields of the images
    - So we can spot larger objects
    - Eg, a long horizontal line, or a corner, or ...



Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

PROC. OF THE IEEE, NOVEMBER 1998

# Why do max-pooling (2)?

PROC. OF THE IEEE, NOVEMBER 1998

- At some point the feature maps start to get very sparse and blobby
  - indicators of some semantic property, not a recognizable transformation of the image
- Then just use them as features in a "normal" ANN



 $\overline{7}$ 

Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.









# Flattening



### **CNN in Keras**

Only modified the *network structure* and *input format (vector -> 3-D tensor)* 



# Only modified the *network structure* and *input format (vector -> 3-D array)*



**CNN** in Keras



### Contents



Focus on 1 & 2-D signals but signal dimensionality is arbitrary (usually 1,2,3,4-D)

The 2-D discrete convolution of two signals *I* and *K* is defined as:

$$I * K)(i,j) = \sum_{m} \sum_{n} I(m,n)K(i-m,j-n)$$

$$=\sum_{m}\sum_{n}I(i - m, j - n)K(m, n)$$



Convolution

Image by Cmglee - Own work, CC BY-SA 3.0, https://commons.wikimedia.o rg/w/index.php?curid=20206 883

Where  $-\infty \le m, n \le \infty$ . Finite signals can be extended by adding zeros (more on this later)

## **Convolution vs Cross-correlation**

$$(I * K)(i,j) =$$
$$= \sum_{m} \sum_{n} I(m,n)K(i-m,j-n)$$

$$(I \star K)(i,j) = \sum_{m} \sum_{n} I(m,n)K(i+m,j+n)$$

Cross-correlation quantifies presence of I(i, j) in (shifted) K(i, j)

Most ML libraries implement convolutional layers as cross-correlation layers.

- Commutativity not important for most ConvNets
- Can avoid flipping one of the signals ⇒ easier to implement



Image by Cmglee - Own work, CC BY-SA 3.0, https://commons.wikimedia.org/w/index.php?curid=20206883
## **Continuous convolution**

https://en.wikipedia.org/wiki/Convolution

1-D 
$$(f * g)(t) \stackrel{\text{def}}{=} \int_{-\infty}^{\infty} f(\tau) g(t - \tau) d\tau$$
  
=  $\int_{-\infty}^{\infty} f(t - \tau) g(\tau) d\tau$ .



## **Continuous convolution**

https://en.wikipedia.org/wiki/Convolution

1-D 
$$(f * g)(t) \stackrel{\text{def}}{=} \int_{-\infty}^{\infty} f(\tau) g(t - \tau) d\tau$$
  
=  $\int_{-\infty}^{\infty} f(t - \tau) g(\tau) d\tau$ .















