Introduction to Machine Learning



Convolutional Neural Networks

Neural Networks

Universal Function Approximator:

 Neural networks with large-enough width (and bounded depth) can approximate any "nice" function



(Analogous result for bounded width and large-enough depth)

Are we done?

Universal Function Approximator:

Let's suppose we had an infinite compute: Can we just set the number of hidden units to 10^1000?

Are we done?

Universal Function Approximator:

Let's suppose we had an infinite compute: Can we just set the number of hidden units to 10^1000?

- The theorem says nothing about **sample efficiency**
- Sample efficiency (informal): Model A is more sample efficient than model B if A needs less training data to get the same test performance

Neural Network Architectures

We want to inject structural knowledge about the input domain into our neural network.



Neural Network Architectures

We want to inject structural knowledge about the input domain into our neural network.



If the neural network is "aware" of input domain structure, then it can learn faster and generalize better.

Neural Network Architectures

We want to inject structural knowledge about the input domain into our neural network.



If the neural network is "aware" of input domain structure, then it can learn faster and generalize better.











"0" if black, "1" if white



 $x_{22} \ x_{23} \ x_{24} \ x_{25}$

"Fully-Connected" Neural Networks

Every hidden unit is connected to the input!





Logistic Regression/ Neural Network Whether "H" is in image or not

"Fully-Connected" Neural Networks

Every hidden unit is connected to the input!



x_1	x_2	x_3	x_4	x_5
x_6	x_7	x_8	x_9	x_{10}
x_{11}	x_{12}	x_{13}	x_{14}	x_{15}
x_{16}	x_{17}	x_{18}	x_{19}	x_{20}
x_{21}	x_{22}	x_{23}	x_{24}	x_{25}

The model doesn't know apriori that x7 is close to x1.

Image Structure

What do we know about images?



Image Structure

What do we know about images?

Fur

Spatial Locality: things nearby in 2D space are more likely to be similar (or have the same source of information)



Image Structure

What do we know about images?

Translation Invariance: The same pattern of pixels has the same "interpretation" no matter where it occurs in the image.

Still a "samoyed" even though absolute position is different





1D sequence





1D sequence



1D sequence



1D sequence



1D sequence



1D sequence



1D sequence



1D sequence



1D sequence



Convolutional Filter (model parameters)

(Optional) Padding to preserve sequence length

-1

1D sequence



Convolutional Filter (model parameters)



Preactivation Layer

(Optional) Padding to preserve sequence length

-1 1

1D sequence



Convolutional Filter (model parameters)

with bias +1

Convolutional filter can have biases (as in fully connected layers)

1D sequence



Convolutional Filter (model parameters)

> Filter weights (and bias) will be learned via gradient descent

1D sequence



Convolutional Filter (model parameters) $w_1 \, w_2$ with bias b



Filters can have different widths

1D sequence



Filters can have different widths

1D sequence

Convolutional Filter (model parameters)

Preactivation Layer

ReLU nonlinearity

Convolutional Layer: Spatial Locality



What is this filter doing?





Convolutional Layer: Spatial Locality



What is this filter doing?



Looking for length-3 sequences things that it has high dot product with.

Convolutional Layer: Spatial Locality



What is this filter doing?



Looking for length-3 sequences things thatit has high dot product with.

Looks for **local** information! Spatial locality is encoded into the model

Convolutional Layer: Translational Invariance

0 0 1 1 1 0 1 1 1 0



-1 0 -1 0 -2 0 -1 0

Convolutional Layer: Translational Invariance





Same local region => same output!
Convolutional Layer: Translational Invariance





Doesn't matter where the feature occurs

Same local region => same output!

 $\mathbf{h} = \mathbf{W}\mathbf{x}$



 $\mathbf{h} = \mathbf{W}\mathbf{x}$

									0
h1		w1	w2	w3	0	0	0	0	X
h2		0	w1	w2	w3	0	0	0	xź
h3	=	0	0	w1	w2	w3	0	0	x
h4		0	0	0	w1	w2	w3	0	X4
h5		0	0	0	0	w1	w2	w3	X

0

-1	0
h1 w1 w2 w3 0 0 0 0 x	x1
h2 0 w1 w2 w3 0 0 0 x	x2
h3 0 0 w1 w2 w3 0 0 x	x3
h4 0 0 0 w1 w2 w3 0 x	x4
h5 0 0 0 0 w1 w2 w3 x	x5



Fully Connected Layer

0

x1

x2

x3

x4

x5

0

 $\mathbf{h} = \mathbf{W}\mathbf{x}$

	-							
h1		w1	w2	w3	w4	w5	w6	w7
h2		w8	w9	w10	w11	w12	w13	w14
h3	=	w15	w16	w17	w18	w19	w20	w21
h4		w22	w23	w24	w25	w26	w27	w28
h5		w29	w30	w31	w32	w33	w34	w35

These weights can be anything!

- 1D Convolution is just matrix multiplication with a structured matrix!
- Convolutional layers used shared weights to enforce spatial locality and translation invariance.
- They are less flexible of a model than fully connected layers (i.e., parameters).
- When would convolutional layers be a bad idea??





















- Tensor: generalization of a matrix
 - E.g. 1D: vector, 2D: matrix





- Tensor: generalization of a matrix
 - E.g. 1D: vector, 2D: matrix





- Tensor: generalization of a matrix
 - E.g. 1D: vector, 2D: matrix





- Tensor: generalization of a matrix
 - E.g. 1D: vector, 2D: matrix



- Input : height x width x depth
- Each fiter F produces height x width "hidden layer"
- Multiple filters ("filter bank") produce a layer of dimension height x width x number of filters

Specifying Each Layer

- Denote current layer with I
- Input: n^{I-1} x n^{I-1} x m^{I-1}
- Number of filters: m^l
- Size of filters: $k^{I} x k^{I} x m^{I-1}$
- Stride: s^I
- Output: $n^{I} \times n^{I} \times m^{I}$, where $n^{I} = n^{I-1} / s^{I}$



Study Question

- 1) How many weights are in a convolutional layer specified as below?
- 2) If we used a fully-connected layer with the same size inputs and outputs, how many weights would it have?

Denote current layer with I

Input: n^{I-1} x n^{I-1} x m^{I-1}

Number of filters: m^l

Size of filters: $k^{I}x k^{I}x m^{I-1}$

Stride: s^I

Output: $n^{I} \times n^{I} \times m^{I}$, where $n^{I} = n^{I-1} / s^{I}$



Output from the convolutional layer & ReLU:



Max pooling: returns max of its arguments

• E.g. size 3x3 ("size 3")







Output from the convolutional layer & ReLU:



Max pooling: returns max of its arguments

• E.g. size 3x3 ("size 3")

After max pooling:



Output from the convolutional layer & ReLU:



Max pooling: returns max of its arguments

• E.g. size 3x3 ("size 3")

After max pooling:



"Stride" 1 pooling







Output from the convolutional layer & ReLU:



Max pooling: returns max of its arguments

- E.g. size 3x3 ("size 3")
- E.g. stride 3

After max pooling:

"Stride" 3 pooling

Output from the convolutional layer & ReLU:



Max pooling: returns max of its arguments

- E.g. size 3x3 ("size 3")
- E.g. stride 3

After max pooling:

"Stride" 3 pooling

Output from the convolutional layer & ReLU:



Max pooling: returns max of its arguments

- E.g. size 3x3 ("size 3")
- E.g. stride 3

After max pooling: 0 1 1 0

"Stride" 3 pooling

No parameters in max pooling layer!

Flatten

Flattening transforms a two-dimensional matrix of features into a vector that can be fed into a fully connected neural network classifier


Putting it all together



Simple CNN: Forward Pass



$$Z_{i}^{1} = W^{1^{\mathsf{T}}} \cdot A_{[i-\lfloor k/2 \rfloor:i+\lfloor k/2 \rfloor}^{0}$$
$$A^{1} = \operatorname{ReLU}(Z^{1})$$
$$A^{2} = W^{2^{\mathsf{T}}}A^{1}$$
$$L(A^{2}, y) = (A^{2} - y)^{2}$$

Simple CNN: Backpropagation



$$Z_{i}^{1} = W^{1^{\mathsf{T}}} \cdot A_{[i-\lfloor k/2 \rfloor:i+\lfloor k/2 \rfloor}^{0}$$
$$A^{1} = \operatorname{ReLU}(Z^{1})$$
$$A^{2} = W^{2^{\mathsf{T}}}A^{1}$$
$$L(A^{2}, y) = (A^{2} - y)^{2}$$

∂loss	∂Z1	∂A^1	∂loss
$\partial W^1 =$	$\overline{\partial W^1}$	∂Z^1	∂A^1

Deep Learning Today

Used both for analysis and synthesis



Input image

Advances in Dataset Collection



IM GENET

An ontology of images based on WordNet

ImageNet currently has

- 13,000+ categories of visual concepts
- 10 million human-cleaned images (~700im/categ)
- 1/3+ is released online @ www.image-net.org



Deng, Dong, Socher, Li & Fei-Fei, CVPR 2009

Application to ImageNet



- ~14 million labeled images, 20k classes
- Images gathered from Internet
- Human labels via Amazon Turk

ImageNet Classification with Deep Convolutional Neural Networks [NIPS 2012]

Alex Krizhevsky University of Toronto kriz@cs.utoronto.ca Ilya Sutskever University of Toronto ilya@cs.utoronto.ca Geoffrey E. Hinton University of Toronto hinton@cs.utoronto.ca

Solving Different Tasks



Object Detection: What and Where





Object Detection Results



Object Detection Results



Object Detection Results



Segmentation, Depth, & More



monocular depth estimation (Liu et al. 2015)

boundary prediction (Xie & Tu 2015)