Convolutional Neural Network

Can the network be simplified by considering the properties of images?

Artificial Neural Networks

- Connectionist, PDP, etc. models
- A biologically-inspired approach for
 - intelligent computing machines
 - massive parallelism
 - distributed computing
 - learning, generalization, adaptivity
 - Tolerant of fault, uncertainty, imprecise info

Compared to Von Neumann

	Von Neumann	Biological
	computer	neural systems
Processor	complex, high	simple, low
	speed, few	speed, many
Memory	separate from	integrated into
	processor,	processor,
	non-content	content
	addressable	addressable
Computing	centralized,	distributed,
	sequential stored	parallel self-
	programs	learning
Reliability	vulnerable	fault tolerant

Biological Neural Networks

- soma (cell body)
- dendrites (receivers)
- axon (transmitters)
- synapses (connection points, axon-soma, axon-dendrite, axonaxon)
- Chemicals (neurotransmitters)
- 10¹¹ neurons
- each makes about $10^3 \sim 10^4$ connections
- with an operating speed of a few milliseconds
- one-hundred-step rule



Signal Generation

- Resting potential
 - Charge difference across neuron membrane approximately –70mV
- Graded potential
 - Stimulus across synapses of post-synaptic neuron
- Action potential
 - If accumulation of graded potential across neuron membrane over a short period of time is higher than ~15mV, action potential is generated and propagated across axon
 - Same form and amplitude regardless of stimulus, signal by frequency rather than amplitude

Signal Generation



Computational Neuron Model (McCulloch and Pitts)



What does a neuron do mathematically?

- W and b determine a "basis" function
- Unlike Fourier bases, these bases are learned to fit the data
- Pattern finder or classifier

- W and b determine a partitioning hyperplane
- They separate highdimensional feature spaces into regions
- Pattern classifier
- Again, the partitioning hyperplanes are learned

Why CNN for Image

Some patterns are much smaller than the whole image

A neuron does not have to see the whole image to discover the pattern.

Connecting to small region with less parameters



Why CNN for Image

• The same patterns appear in different regions.



Why CNN for Image

Subsampling the pixels will not change the object

bird



We can subsample the pixels to make image smaller

Less parameters for the network to process the image







Those are the network parameters to be learned.

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



6 x 6 image

Property 1 Each filter detects a small pattern (3 x 3).

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

Filter 1

stride=1

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



6 x 6 image



Filter 1

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



We set stride=1 below

6 x 6 image



Filter 1

stride=1



6 x 6 image





Filter 2

stride=1

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

6 x 6 image

Do the same process for every filter



CNN – Colorful image



Convolution v.s. Fully Connected



Fullyconnected

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









6 x 6 image

Less parameters!

Even less parameters!





CNN – Max Pooling







Filter 2





CNN – Max Pooling



6 x 6 image

Each filter is a channel

CNN Demo - Mnist

- Hand-written digits
- 60,000 training samples and 10,000 test samples
- 28x28 binary images

CNN Demo – CiFar10

- 10 classes (airplane, auto, bird, cat, dog, etc.)
- 50,000 training samples and 10,000 test samples
- 32x32 color images

CNN for Mnist

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

Loss Function

- Ground truth 1-hot vector of dimension nx1
 - N: number of categories
 - 1: true category (dog, cat, etc.)
- CNN output
 - Softmax post processing to make a probability distribution
- Loss:
 - Sum over mini-batch
 - 2-norm error (doesn't work)
 - Cross entropy

Softmax Function

$$\sigma(\mathbf{z})_j = rac{e^{z_j}}{\sum_{k=1}^K e^{z_k}} \quad ext{for } j = 1, ..., K ext{ and } \mathbf{z} = (z_1, \dots, z_K) \in \mathbb{R}^K$$

- Exponent CNN output
- Normalized by sum of exponents
- Output is a probability distribution
- Accentuate positives and suppress negatives

- Number of bits required to encode a lexicon
- Amount of info (surprise)
- Large Pi -> small –log(Pi) -> short codeword
- Small pi -> large –log(pi) -> large codeword
- Compared to ASCII (8-bit per char)
- Entropy is largest when all pi are the same (most uncertain)

Cross Entropy

$$H(p,q) = -\sum_{x\in\mathcal{X}} p(x)\,\log q(x)$$

- Use a codebook designed for one distribution for another
- Not symmetric, hence, not a distance function
- q(x) > p(x) => x has a short codeword => with small p(x)
- q(x) > p(x) => x has a long codeword => with large p(x)
- Smallest when p==q

<u>CNN in Keras</u>

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

Backpropagation Learning rule

Cost function

What does machine learn?

http://newsneakernews.wpengine.netdna-cdn.com/wpcontent/uploads/2016/11/rihanna-puma-creeper-velvet-release-date-02.jpg

Deconvolutional Network Map activations back to the input pixel space, show what input pattern originally caused a given activation.

Steps

- Compute activations at a specific layer.
- □ Keep one activation and set all others to zero.
- Unpooling and deconvolution
- □ Construct input image

Layer1

Layer2

Layer3

Layer4

Layer5

More Application: Playing Go

More Application: Playing Go

Why CNN for playing Go?

Some patterns are much smaller than the whole image

Alpha Go uses 5 x 5 for first layer

• The same patterns appear in different regions.

Why CNN for playing Go?

• Subsampling the pixels will not change the object

Max Pooling How to explain this???

Neural network architecture. The input to the policy network is a $19 \times 19 \times 48$ image stack consisting of 48 feature planes. The first hidden layer zero pads the input into a 23 \times 23 image, then convolves k filters of kernel size 5 \times 5 with stride 1 with the input image and applies a <u>rectifier nonlinearity</u>. Each of the subsequent hidden layers 2 to 12 zero pads the respective previous hidden layer into a 21×21 image, then convolves *k* filters of kernel size 3×3 with stride 1, again followed by a rectifier nonlinearity. The final layer convolves 1 filter of kernel size 1×1 with stride 1 with a different bies for each position and applies a softmax func-tion. The Alpha Go does not use Max Pooling Extended Data Table 3 additionally show the results of training with k = 128, 256 and 384 filters.

More Application: Speech

More Application: Text

