









Wow



false positives

no good filtr

so misclassified

what class

cool kernel

Goals

Build a classifier which is more powerful at representing complex functions *and* more suited to the learning problem.

What does this mean?

1. Assume that the *underlying data generating function* relies on a composition of factors.

2. Learn a feature representation that is specific to the dataset.

Neural Networks

- Basic building block for composition is a *perceptron* (Rosenblatt c.1960)
- Linear classifier vector of weights w and a 'bias' b





Mark 1 Perceptron c.1960

20x20 pixel camera feed

Universality

A single-layer network can learn any function:

- So long as it is differentiable
- To some approximation;
 More perceptrons = a better approximation

Visual proof (Michael Nielson):

http://neuralnetworksanddeeplearning.com/chap4.html

If a single-layer network can learn any function... ...given enough parameters...

...then why do we go deeper?

Intuitively, composition is efficient because it allows *reuse*.

Empirically, deep networks do a better job than shallow networks at learning such hierarchies of knowledge.

Composition



Layers that are in between the input and the output are called *hidden layers*, because we are going to *learn* their weights via an optimization process.

Interpretation of many layers

[0 0 1 0 0 0 0 1 0 0 1 1 0 0 1 0 ...] truck feature



Exponentially more efficient than a 1-of-N representation (a la k-means)



Interpretation

[1 1 0 0 0 1 0 1 0 0 0 0 1 1 0 1...] motorbike

[0 0 1 0 0 0 0 1 0 0 1 1 0 0 1 0 ...] truck





Interpretation



Ranzato

Lee et al. "Convolutional DBN's ..." ICML 2009

Activation functions: Rectified Linear Unit

• ReLU $f(x) = \max(0, x)$



Neural Networks: example

$$\begin{array}{c} x \\ \hline max(0, W^{1}x) \end{array} \xrightarrow{h^{1}} max(0, W^{2}h^{1}) \xrightarrow{h^{2}} W^{3}h^{2} \end{array} \xrightarrow{O}$$

- *x* input
- h^1 1-st layer hidden units
- h^2 2-nd layer hidden units
- *o* output

Example of a 2 hidden layer neural network (or 4 layer network, counting also input and output).



Does anyone pass along the weight without an activation function?

No – this is linear chaining.



Does anyone pass along the weight without an activation function?

No – this is linear chaining.



What is the relationship between SVMs and perceptrons?

SVMs attempt to learn the support vectors which maximize the margin between classes.



What is the relationship between SVMs and perceptrons?

SVMs attempt to learn the support vectors which maximize the margin between classes.

A perceptron does not.

Both of these perceptron classifiers are equivalent.

'Perceptron of optimal stability' is used in SVM:

Perceptron + optimal stability + kernel trick

= foundations of SVM



Outline

- Supervised Neural Networks
- Convolutional Neural Networks
- Examples





Images as input to neural networks





Images as input to neural networks



Images as input to neural networks

Example: 200x200 image 40K hidden units

~2B parameters!!!

- Spatial correlation is local
- Waste of resources + we have not enough training samples anyway..



Motivation

- Sparse interactions *receptive fields*
 - Assume that in an image, we care about 'local neighborhoods' only for a given neural network layer.
 - Composition of layers will expand local -> global.

Example: 200x200 image 40K hidden units Filter size: 10x10 4M parameters

Note: This parameterization is good when input image is registered (e.g., face recognition).



Motivation

- Sparse interactions *receptive fields*
 - Assume that in an image, we care about 'local neighborhoods' only for a given neural network layer.
 - Composition of layers will expand local -> global.
- Parameter sharing
 - 'Tied weights' use same weights for more than one perceptron in the neural network.
 - Leads to equivariant representation
 - If input changes (e.g., translates), then output changes similarly

Share the same parameters across different locations (assuming input is stationary):



Filtering reminder: Correlation (rotated convolution)



I[.,.]

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

h[.,.]

0	10	20	30	30	30	20	10	
0	20	40	60	60	60	40	20	
0	30	60	90	90	90	60	30	
0	30	50	80	80	90	60	30	
0	30	50	80	80	90	60	30	
0	20	30	50	50	60	40	20	
10	20	30	30	30	30	20	10	
10	10	10	0	0	0	0	0	

 $h[m,n] = \sum_{k,l} f[k,l] I[m+k,n+l]$

Credit: S. Seitz

Perceptron: output =
$$\begin{cases} 0 & \text{if } u \\ 1 & \text{if } u \end{cases}$$

 $egin{cases} 0 & ext{if} \ w\cdot x + b \leq 0 \ 1 & ext{if} \ w\cdot x + b > 0 \end{cases}$

$$w\cdot x\equiv \sum_j w_j x_j,$$

This is convolution!

Share the same parameters across different locations (assuming input is stationary):

Convolutions with learned kernels








































































Filter size: 10x10 10K parameters



Interpretation



Lee et al. "Convolutional DBN's ..." ICML 2009





n = layer number K = kernel size j = # channels (input) or # filters (depth)

 $h_1^{n-1} \longrightarrow h_1^n$ h_2^{n-1} h_3^{n-1} h_3^{n-1}



























Pooling Layer

Let us assume filter is an "eye" detector.

Q.: how can we make the detection robust to the exact location of the eye?



Pooling Layer

By *pooling* responses at different locations, we gain robustness to the exact spatial location of image features.



Pooling is similar to downsampling



...except sometimes we don't want to blur, as other functions might be better for classification.

Pooling Layer: Receptive Field Size



Pooling Layer: Examples

Max-pooling:

$$h_{j}^{n}(x, y) = max_{\bar{x} \in N(x), \bar{y} \in N(y)} h_{j}^{n-1}(\bar{x}, \bar{y})$$

Average-pooling:

$$h_{j}^{n}(x, y) = 1/K \sum_{\bar{x} \in N(x), \bar{y} \in N(y)} h_{j}^{n-1}(\bar{x}, \bar{y})$$

Max pooling

Single depth slice



Wikipedia

Pooling Layer: Examples

Max-pooling:

$$h_{j}^{n}(x, y) = max_{\bar{x} \in N(x), \bar{y} \in N(y)} h_{j}^{n-1}(\bar{x}, \bar{y})$$

Average-pooling:

$$h_{j}^{n}(x, y) = 1/K \sum_{\bar{x} \in N(x), \bar{y} \in N(y)} h_{j}^{n-1}(\bar{x}, \bar{y})$$

L2-pooling:

$$h_{j}^{n}(x, y) = \sqrt{\sum_{\bar{x} \in N(x), \bar{y} \in N(y)} h_{j}^{n-1}(\bar{x}, \bar{y})^{2}}$$

L2-pooling over features:

$$h_{j}^{n}(x, y) = \sqrt{\sum_{k \in N(j)} h_{k}^{n-1}(x, y)^{2}}$$



Pooling Layer: Receptive Field Size



If convolutional filters have size KxK and stride 1, and pooling layer has pools of size PxP, then each unit in the pooling layer depends upon a patch (at the input of the preceding conv. layer) of size:





Pooling Layer: Receptive Field Size



If convolutional filters have size KxK and stride 1, and pooling layer has pools of size PxP, then each unit in the pooling layer depends upon a patch (at the input of the preceding conv. layer) of size:













$$h^{i+1}(x, y) = \frac{h^{i}(x, y) - m^{i}(N(x, y))}{\sigma^{i}(N(x, y))}$$

Performed also across features and in the higher layers..

Effects:

- improves invariance
- improves optimization
- increases sparsity

Note: computational cost is negligible w.r.t. conv. layer.

70

Ranzate
ConvNets: Typical Stage

One stage (zoom)





ConvNets: Typical Architecture

One stage (zoom)









Conceptually similar to:

SIFT \rightarrow K-Means \rightarrow Pyramid Pooling \rightarrow SVM Lazebnik et al. "...Spatial Pyramid Matching..." CVPR 2006

SIFT \rightarrow Fisher Vect. \rightarrow Pooling \rightarrow SVM Sanchez et al. "Image classification with F.V.: Theory and practice" IJCV 2012



Yann LeCun's MNIST CNN architecture





32x32x3 image 32 32 3

5x5x3 filter

Convolution Layer







For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:



We stack these up to get a "new image" of size 28x28x6!



Ν



Output size: (N - F) / stride + 1

Our connectomics diagram

Auto-generated from network declaration by nolearn (for Lasagne / Theano)

Input 75x75x4



Reading architecture diagrams

Layers

- Kernel sizes
- Strides
- # channels
- # kernels
- Max pooling



[Krizhevsky et al. 2012]

AlexNet diagram (simplified)

Input size 227 x 227 x 3



Outline

- Supervised Neural Networks
- Convolutional Neural Networks
- Examples





- OCR / House number & Traffic sign classification





Ciresan et al. "MCDNN for image classification" CVPR 2012 Wan et al. "Regularization of neural networks using dropconnect" ICML 2013 Jaderberg et al. "Synthetic data and ANN for natural scene text recognition" arXiv 2014

- Scene Parsing



Farabet et al. "Learning hierarchical features for scene labeling" PAMI 201385Pinheiro et al. "Recurrent CNN for scene parsing" arxiv 2013Ranzato

- Segmentation 3D volumetric images



Ciresan et al. "DNN segment neuronal membranes..." NIPS 2012 Turaga et al. "Maximin learning of image segmentation" NIPS 2009



- Object detection



Sermanet et al. "OverFeat: Integrated recognition, localization, ..." arxiv 2013 Girshick et al. "Rich feature hierarchies for accurate object detection..." arxiv 2013 91 Szegedy et al. "DNN for object detection" NIPS 2013 Ranzato

- Face Verification & Identification



Taigman et al. "DeepFace..." CVPR 2014



Dataset: ImageNet 2012



- <u>S:</u> (n) <u>Eskimo dog</u>, husky (breed of heavy-coated Arctic sled dog)
 - o direct hypernym / inherited hypernym / sister term
 - S: (n) working dog (any of several breeds of usually large powerful dogs bred to work as draft animals and guard and guide dogs)
 - S: (n) dog, domestic dog, Canis familiaris (a member of the genus Canis (probably descended from the common wolf) that has been domesticated by man since prehistoric times; occurs in many breeds) "the dog barked all night"
 - S: (n) canine, canid (any of various fissiped mammals with nonretractile claws and typically long muzzles)
 - S: (n) carnivore (a terrestrial or aquatic flesh-eating mammal) "terrestrial carnivores have four or five clawed digits on each limb"
 - S: (n) placental, placental mammal, eutherian, eutherian mammal (mammals having a placenta; all mammals except monotremes and marsupials)
 - <u>S:</u> (n) mammal, mammalian (any warm-blooded vertebrate having the skin more or less covered with hair; young are born alive except for the small subclass of monotremes and nourished with milk)
 - S: (n) vertebrate, craniate (animals having a bony or cartilaginous skeleton with a segmented spinal column and a large brain enclosed in a skull or cranium)
 - S: (n) chordate (any animal of the phylum Chordata having a notochord or spinal column)
 - S: (n) animal, animate being, beast, brute, creature, fauna (a living organism characterized by voluntary movement)
 - S: (n) organism, being (a living thing that has (or can develop) the ability to act or function independently)
 - S: (n) living thing, animate thing (a living (or once living) entity)
 - S: (n) whole, unit (an assemblage of parts that is regarded as a single entity) "how big is that part compared to the whole?"; "the team is a unit"
 - <u>S</u> (n) <u>object</u>, <u>physical object</u> (a tangible and visible entity; an entity that can cast a shadow) "it was full of rackets, balls and other objects"
 - <u>S:</u> (n) physical entity (an entity that has physical existence)
 - <u>S</u>: (n) <u>entity</u> (that which is perceived or known or inferred to have its own distinct existence (living or nonliving))

Deng et al. "Imagenet: a large scale hierarchical image database" CVPR 2009



| grille | mushroom | grape | spider monkey |
|-------------|--------------------|------------------------|---------------|
| pickup | jelly fungus | elderberry | titi |
| beach wagon | gill fungus | ffordshire bullterrier | indri |
| fire engine | dead-man's-fingers | currant 🛛 | howler monkey |

Architecture for Classification



Results: ILSVRC 2012



Krizhevsky et al. "ImageNet Classification with deep CNNs" NIPS 2012

98

Ranzato

Phew!

- Friday:
- Network training

More ConvNet explanations

 <u>https://ujjwalkarn.me/2016/08/11/intuitive-</u> explanation-convnets/