Administration

CSCI567 Machine Learning (Fall 2025)

Haipeng Luo

University of Southern California

Sep 26, 2025

HW2 available now. Due on Oct 8th.

1 / 47

Review of Last Lecture

Outline

Outline

- Review of Last Lecture
- Neural Nets
- 3 Convolutional neural networks (ConvNets/CNNs)

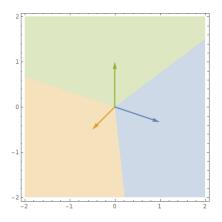
- Review of Last Lecture
- 2 Neural Nets
- 3 Convolutional neural networks (ConvNets/CNNs)

3 / 47

4 / 47

2 / 47

Linear models: from binary to multiclass



$$\mathbf{w}_1 = (1, -\frac{1}{3})$$

 $\mathbf{w}_2 = (-\frac{1}{2}, -\frac{1}{2})$
 $\mathbf{w}_3 = (0, 1)$

• Blue class:

 $\{\boldsymbol{x}: 1 = \operatorname{argmax}_k \boldsymbol{w}_k^{\mathrm{T}} \boldsymbol{x}\}$

Orange class:

 $\{\boldsymbol{x}: 2 = \operatorname{argmax}_{k} \boldsymbol{w}_{k}^{\mathrm{T}} \boldsymbol{x}\}$

Green class:

 $\{x: 3 = \operatorname{argmax}_k w_k^{\mathrm{T}} x\}$

$$\mathcal{F} = \left\{ f(oldsymbol{x}) = rgmax_{k \in [\mathsf{C}]} \ oldsymbol{w}_k^{\mathrm{T}} oldsymbol{x} \mid oldsymbol{w}_1, \dots, oldsymbol{w}_\mathsf{C} \in \mathbb{R}^\mathsf{D}
ight\}$$

Review of Last Lecture

Kernel functions

Definition: a function $k : \mathbb{R}^{D} \times \mathbb{R}^{D} \to \mathbb{R}$ is called a *kernel function* if there exists a function $\phi:\mathbb{R}^\mathsf{D} \to \mathbb{R}^\mathsf{M}$ so that for any $x,x' \in \mathbb{R}^\mathsf{D}$,

$$k(\boldsymbol{x}, \boldsymbol{x}') = \boldsymbol{\phi}(\boldsymbol{x})^{\mathrm{T}} \boldsymbol{\phi}(\boldsymbol{x}')$$

Examples we have seen

$$\begin{split} k(\boldsymbol{x}, \boldsymbol{x}') &= (\boldsymbol{x}^{\mathrm{T}} \boldsymbol{x}')^2 \\ k(\boldsymbol{x}, \boldsymbol{x}') &= \sum_{d=1}^{\mathsf{D}} \frac{\sin(2\pi(x_d - x_d'))}{x_d - x_d'} \\ k(\boldsymbol{x}, \boldsymbol{x}') &= (\boldsymbol{x}^{\mathrm{T}} \boldsymbol{x}' + c)^d & \text{(polynomial kernel)} \\ k(\boldsymbol{x}, \boldsymbol{x}') &= e^{-\frac{\|\boldsymbol{x} - \boldsymbol{x}'\|_2^2}{2\sigma^2}} & \text{(Gaussian/RBF kernel)} \end{split}$$

Softmax + MLE = minimizing cross-entropy loss

Maximize probability of see labels y_1, \ldots, y_N given x_1, \ldots, x_N

$$P(\boldsymbol{W}) = \prod_{n=1}^{\mathsf{N}} \mathbb{P}(y_n \mid \boldsymbol{x}_n; \boldsymbol{W}) = \prod_{n=1}^{\mathsf{N}} \operatorname{softmax}(\boldsymbol{W} \boldsymbol{x}_n)_{y_n} = \prod_{n=1}^{\mathsf{N}} \frac{e^{\boldsymbol{w}_{y_n}^{\mathsf{T}} \boldsymbol{x}_n}}{\sum_{k \in [\mathsf{C}]} e^{\boldsymbol{w}_k^{\mathsf{T}} \boldsymbol{x}_n}}$$

By taking **negative log**, this is equivalent to minimizing

$$F(\boldsymbol{W}) = \sum_{n=1}^{\mathsf{N}} \ln \left(\frac{\sum_{k \in [\mathsf{C}]} e^{\boldsymbol{w}_k^{\mathrm{T}} \boldsymbol{x}_n}}{e^{\boldsymbol{w}_{y_n}^{\mathrm{T}} \boldsymbol{x}_n}} \right) = \sum_{n=1}^{\mathsf{N}} \ln \left(1 + \sum_{k \neq y_n} e^{(\boldsymbol{w}_k - \boldsymbol{w}_{y_n})^{\mathrm{T}} \boldsymbol{x}_n} \right)$$

This is the multiclass logistic loss, a.k.a cross-entropy loss.

Review of Last Lecture

Kernelizing ML algorithms

Key idea: rewrite the algorithm so that its dependence on the transformed dataset Φ is only through the Gram matrix $K = \Phi \Phi^{\mathrm{T}}$.

Perceptron (primal form with ϕ)

issue: time/space linear in M

Initialize $w = 0 \in \mathbb{R}^{\mathsf{M}}$

Repeat:

- Pick a data point index n uniformly at random
- If $\operatorname{sgn}(\boldsymbol{w}^{\mathrm{T}}\boldsymbol{\phi}(\boldsymbol{x}_n)) \neq y_n$, update $\boldsymbol{w} \leftarrow \boldsymbol{w} + y_n \boldsymbol{\phi}(\boldsymbol{x}_n)$

Kernelized Perceptron

independent of M

Initialize $\alpha_m = 0$ for all $m \in [N]$ Repeat:

- Pick a data point index n uniformly at random
- If $\operatorname{sgn}(\sum_{m=1}^{N} \alpha_m k(\boldsymbol{x}_m, \boldsymbol{x}_n)) \neq y_n$, update $\alpha_n \leftarrow \alpha_n + y_n$

Outline

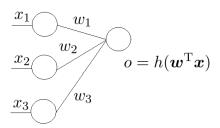
- Review of Last Lecture
- Neural Nets
 - Definition
 - Backpropagation
 - Preventing overfitting
- 3 Convolutional neural networks (ConvNets/CNNs)

9 / 47

Neural Nets

Definition

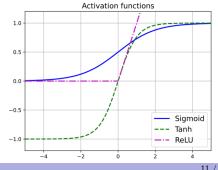
Linear model as a one-layer neural net



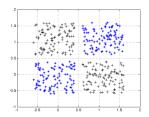
h(a) = a for linear model

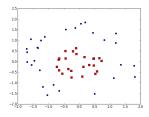
To create non-linearity, can use

- Rectified Linear Unit (ReLU):
 h(a) = max{0, a}
- sigmoid function: $h(a) = \frac{1}{1 + e^{-a}}$
- TanH: $h(a) = \frac{e^a e^{-a}}{e^a + e^{-a}}$
- many more



Linear models are not always adequate





We can use a nonlinear mapping as discussed:

$$oldsymbol{\phi}(oldsymbol{x}):oldsymbol{x}\in\mathbb{R}^{\mathsf{D}}
ightarrowoldsymbol{z}\in\mathbb{R}^{\mathsf{M}}$$

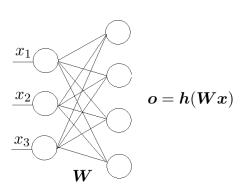
Instead of manually picking ϕ or kernel k, can we learn this nonlinear mapping from data?

The most popular nonlinear models nowadays: neural nets

10 / 47

Neural Nets Definition

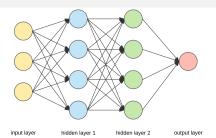
More output nodes



 $W \in \mathbb{R}^{4 \times 3}$, $h : \mathbb{R}^4 \to \mathbb{R}^4$ so $h(a) = (h_1(a_1), h_2(a_2), h_3(a_3), h_4(a_4))$

Can think of this as a nonlinear mapping: $\phi(x) = h(Wx)$

More layers



Becomes a network:

- each node is called a neuron
- h is called the activation function
 - can use h(a) = 1 for one neuron in each layer to *incorporate bias term*
 - output neuron can use h(a) = a
- #layers refers to #hidden_layers (plus 1 or 2 for input/output layers)
- deep neural nets can have many layers and *millions* of parameters
- this is a **feedforward**, **fully connected** neural net, there are many variants (convolutional nets, recurrent nets, transformers, etc.)

13 / 47

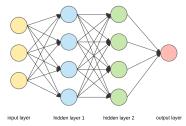
Neural Nets

Definition

Math formulation

An L-layer neural net can be written as

$$\boldsymbol{f}(\boldsymbol{x}) = \boldsymbol{h}_{\mathsf{L}} \left(\boldsymbol{W}_{\!L} \boldsymbol{h}_{\mathsf{L}-1} \left(\boldsymbol{W}_{\!L-1} \cdots \boldsymbol{h}_{1} \left(\boldsymbol{W}_{\!1} \boldsymbol{x} \right) \right) \right)$$



To ease notation, for a given input x, define recursively

$$o_0 = x, \qquad a_\ell = W_\ell o_{\ell-1}, \qquad o_\ell = h_\ell(a_\ell) \qquad \qquad (\ell = 1, \dots, L)$$

where

- $m{W}_\ell \in \mathbb{R}^{\mathsf{D}_\ell imes \mathsf{D}_{\ell-1}}$ is the weights between layer $\ell-1$ and ℓ
- $D_0 = D, D_1, \dots, D_L$ are numbers of neurons at each layer
- ullet $a_\ell \in \mathbb{R}^{\mathsf{D}_\ell}$ is input to layer ℓ
- $oldsymbol{o}_\ell \in \mathbb{R}^{\mathsf{D}_\ell}$ is output of layer ℓ
- $h_\ell: \mathbb{R}^{\mathsf{D}_\ell} \to \mathbb{R}^{\mathsf{D}_\ell}$ is activation functions at layer ℓ

How powerful are neural nets?

Universal approximation theorem (Cybenko, 89; Hornik, 91):

A feedforward neural net with a single hidden layer can approximate any continuous functions.

It might need a huge number of neurons though, and depth helps!

Designing network architecture is important and very complicated

• for feedforward network, need to decide number of hidden layers, number of neurons at each layer, activation functions, etc.

Demo: http://playground.tensorflow.org

14 / 47

Neural Nets Definition

Learning the model

No matter how complicated the model is, our goal is the same: minimize

$$F(\boldsymbol{W}_1,\ldots,\boldsymbol{W}_L) = \frac{1}{N} \sum_{n=1}^{N} F_n(\boldsymbol{W}_1,\ldots,\boldsymbol{W}_L)$$

where

$$F_n(\boldsymbol{W}_1,\dots,\boldsymbol{W}_{\mathsf{L}}) = egin{cases} \|\boldsymbol{f}(\boldsymbol{x}_n) - \boldsymbol{y}_n\|_2^2 & \text{for regression} \\ \ln\left(1 + \sum_{k \neq y_n} e^{f(\boldsymbol{x}_n)_k - f(\boldsymbol{x}_n)_{y_n}}
ight) & \text{for classification} \end{cases}$$

How to optimize such a complicated function?

Same thing: apply **SGD**! even if the model is *nonconvex*.

What is the gradient of this complicated function?

Chain rule is the only secret:

• for a composite function f(g(w))

$$\frac{\partial f}{\partial w} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial w}$$

• for a composite function $f(g_1(w), \ldots, g_d(w))$

$$\frac{\partial f}{\partial w} = \sum_{i=1}^{d} \frac{\partial f}{\partial g_i} \frac{\partial g_i}{\partial w}$$

the simplest example $f(g_1(w), g_2(w)) = g_1(w)g_2(w)$

17 / 47

Neural Nets

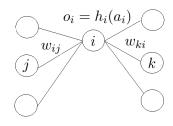
Backpropagation

Computing the derivative

Adding the subscript for layer:

$$\frac{\partial F_n}{\partial w_{\ell,ij}} = \frac{\partial F_n}{\partial a_{\ell,i}} o_{\ell-1,j}$$

$$\frac{\partial F_n}{\partial a_{\ell,i}} = \left(\sum_k \frac{\partial F_n}{\partial a_{\ell+1,k}} w_{\ell+1,ki}\right) h'_{\ell,i}(a_{\ell,i})$$



For the last layer, for square loss

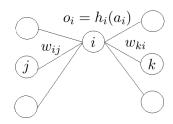
$$\frac{\partial F_n}{\partial a_{\mathsf{L},i}} = \frac{\partial (h_{\mathsf{L},i}(a_{\mathsf{L},i}) - y_{n,i})^2}{\partial a_{\mathsf{L},i}} = 2(h_{\mathsf{L},i}(a_{\mathsf{L},i}) - y_{n,i})h'_{\mathsf{L},i}(a_{\mathsf{L},i})$$

Exercise: try to do it for cross-entropy loss yourself.

Computing the derivative

Drop the subscript ℓ for layer for simplicity.

Find the derivative of F_n w.r.t. to w_{ij}



$$\frac{\partial F_n}{\partial w_{ij}} = \frac{\partial F_n}{\partial a_i} \frac{\partial a_i}{\partial w_{ij}} = \frac{\partial F_n}{\partial a_i} \frac{\partial (w_{ij}o_j)}{\partial w_{ij}} = \frac{\partial F_n}{\partial a_i} o_j$$

$$\frac{\partial F_n}{\partial a_i} = \frac{\partial F_n}{\partial o_i} \frac{\partial o_i}{\partial a_i} = \left(\sum_k \frac{\partial F_n}{\partial a_k} \frac{\partial a_k}{\partial o_i}\right) h_i'(a_i) = \left(\sum_k \frac{\partial F_n}{\partial a_k} w_{ki}\right) h_i'(a_i)$$

18 / 47

Neural Nets

Backpropagation

Computing the derivative

Using matrix notation greatly simplifies presentation and implementation:

$$\frac{\partial F_n}{\partial \boldsymbol{W}_{\ell}} = \frac{\partial F_n}{\partial \boldsymbol{a}_{\ell}} \boldsymbol{o}_{\ell-1}^{\mathrm{T}} \in \mathbb{R}^{\mathsf{D}_{\ell} \times \mathsf{D}_{\ell-1}}$$

$$\frac{\partial F_n}{\partial \boldsymbol{a}_{\ell}} = \begin{cases} \left(\boldsymbol{W}_{\ell+1}^{\mathrm{T}} \frac{\partial F_n}{\partial \boldsymbol{a}_{\ell+1}}\right) \circ \boldsymbol{h}'_{\ell}(\boldsymbol{a}_{\ell}) & \text{if } \ell < \mathsf{L} \\ 2(\boldsymbol{h}_{\mathsf{L}}(\boldsymbol{a}_{\mathsf{L}}) - \boldsymbol{y}_n) \circ \boldsymbol{h}'_{\mathsf{L}}(\boldsymbol{a}_{\mathsf{L}}) & \text{else} \end{cases}$$

where $v_1 \circ v_2 = (v_{11}v_{21}, \cdots, v_{1D}v_{2D})$ is the element-wise product (a.k.a. Hadamard product).

Verify yourself! (But no need to remember this formula.)

Putting everything into SGD

The backpropagation algorithm (Backprop)

Initialize $\boldsymbol{W}_1,\ldots,\boldsymbol{W}_{\mathsf{L}}$ randomly. Repeat:

- **1** randomly pick one data point $n \in [N]$
- **2 forward propagation**: for each layer $\ell = 1, ..., L$

ullet compute $oldsymbol{a}_\ell = oldsymbol{W}_\ell oldsymbol{o}_{\ell-1}$ and $oldsymbol{o}_\ell = oldsymbol{h}_\ell (oldsymbol{a}_\ell)$

 $(o_0 = x_n)$

- **3** backward propagation: for each $\ell = L, \ldots, 1$
 - compute

$$rac{\partial F_n}{\partial oldsymbol{a}_\ell} = egin{cases} \left(oldsymbol{W}_{\ell+1}^{
m T} rac{\partial F_n}{\partial oldsymbol{a}_{\ell+1}}
ight) \circ oldsymbol{h}_\ell'(oldsymbol{a}_\ell) & ext{ if } \ell < \mathsf{L} \ 2(oldsymbol{h}_\mathsf{L}(oldsymbol{a}_\mathsf{L}) - oldsymbol{y}_n) \circ oldsymbol{h}_\mathsf{L}'(oldsymbol{a}_\mathsf{L}) & ext{ else} \end{cases}$$

• update weights

$$\boldsymbol{W}_{\ell} \leftarrow \boldsymbol{W}_{\ell} - \eta \frac{\partial F_n}{\partial \boldsymbol{W}_{\ell}} = \boldsymbol{W}_{\ell} - \eta \frac{\partial F_n}{\partial \boldsymbol{a}_{\ell}} \boldsymbol{o}_{\ell-1}^{\mathrm{T}}$$

(Important: should W_{ℓ} be overwritten immediately in the last step?)

21 / 47

Neural Nets Backpropagation

SGD with momentum (a simple version)

Initialize w_0 and velocity v=0

For t = 1, 2, ...

- ullet form a stochastic gradient $oldsymbol{g}_t$
- update velocity $\boldsymbol{v} \leftarrow \alpha \boldsymbol{v} + \boldsymbol{g}_t$ for some discount factor $\alpha \in (0,1)$
- ullet update weight $oldsymbol{w}_t \leftarrow oldsymbol{w}_{t-1} \eta oldsymbol{v}$

Updates for first few rounds:

- $w_1 = w_0 \eta g_1$
- $w_2 = w_1 \alpha \eta g_1 \eta g_2$
- $w_3 = w_2 \alpha^2 \eta g_1 \alpha \eta g_2 \eta g_3$
- ..

Adam (most popular) \approx SGD + adaptive learning rate + momentum

Important tricks to optimize neural nets

Many important tricks on top on Backprop

- mini-batch: randomly sample a batch of examples to form a stochastic gradient (common batch size: 32, 64, 128, etc.)
- batch normalization: normalize the inputs of each neuron over the mini-batch (to zero-mean and one-variance; c.f. Lec 1)
- adaptive learning rate: scale the learning rate of each parameter based on some moving average of the magnitude of the gradients
- momentum: make use of previous gradients (taking inspiration from physics)
- · · ·

22 / 47

Neural Nets Pre

Preventing overfitting

Overfitting

Overfitting is very likely since neural nets are too powerful.

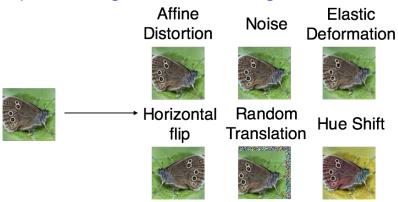
Methods to overcome overfitting:

- data augmentation
- regularization
- dropout
- early stopping
- • •

Data augmentation

Data: the more the better. How do we get more data?

Exploit prior knowledge to add more training data



25 / 47

Neural Nets

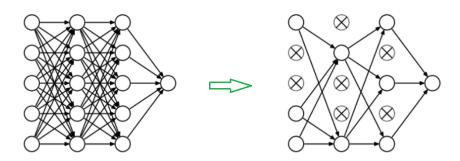
Preventing overfitting

Neural Nets

Preventing overfitting

Dropout

Independently delete each neuron with a fixed probability (say 0.5), during each iteration of Backprop (only for training, not for testing)



Very effective, makes training faster as well

Regularization

L2 regularization: minimize

$$F'(\mathbf{W}_1, \dots, \mathbf{W}_{\mathsf{L}}) = F(\mathbf{W}_1, \dots, \mathbf{W}_{\mathsf{L}}) + \lambda \sum_{\ell=1}^{\mathsf{L}} \|\mathbf{W}_{\ell}\|_2^2$$

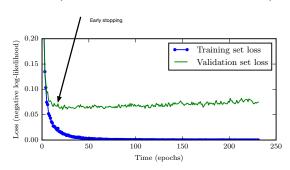
Simple change to the gradient:

$$\frac{\partial F'}{\partial w_{ij}} = \frac{\partial F}{\partial w_{ij}} + 2\lambda w_{ij}$$

Introduce weight decaying effect

Early stopping

Stop training when the performance on validation set stops improving



26 / 47

Conclusions for neural nets

Deep neural networks

- are hugely popular, achieving best performance on many problems
- do need a lot of data to work well
- take a lot of time to train (need GPUs for massive parallel computing)
- take some work to select architecture and hyperparameters
- are still not well understood in theory

29 / 4

Convolutional neural networks (ConvNets/CNNs)

Acknowledgements

Not much math, a lot of empirical intuitions

The materials borrow heavily from the following sources:

- Stanford Course CS231n: http://cs231n.stanford.edu/
- Dr. Ian Goodfellow's lectures on deep learning: http://deeplearningbook.org

Both website provides tons of useful resources: notes, demos, videos, etc.

Also, demo from https://poloclub.github.io/cnn-explainer/

Outline

- Review of Last Lecture
- 2 Neural Nets
- 3 Convolutional neural networks (ConvNets/CNNs)
 - Motivation
 - Architecture

30 / 47

Image Classification: A core task in Computer Vision



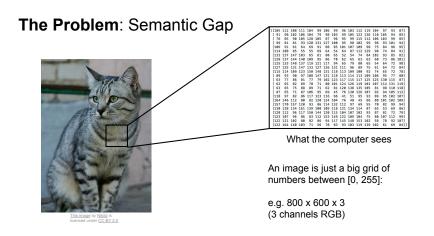
(assume given set of discrete labels) {dog, cat, truck, plane, ...}

→ cat

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 2 - 6

April 6, 2017

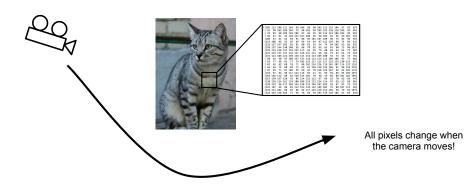


Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 2 - 7

April 6, 2017

Challenges: Viewpoint variation



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 2 - 8

April 6, 2017

Challenges: Illumination









Challenges: Deformation









Challenges: Occlusion







Challenges: Background Clutter





This image is CC0 1.0 public domain

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 2 - 11

April 6, 2017

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 2 - 12

April 6, 2017

Convolutional neural networks (ConvNets/CNNs) Motivation

Fundamental problems in vision

Challenges: Intraclass variation



Fei-Fei Li & Justin Johnson & Serena Yeung Lecture 2 - 13

April 6, 2017

The key challenge

How to train a model that can tolerate all those variations?

Main ideas

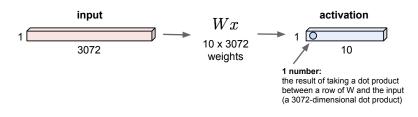
- need a lot of data that exhibits those variations
- need more specialized models to capture the invariance

Convolutional neural networks (ConvNets/CNNs)

Issues of standard NN for image inputs

Fully Connected Layer

32x32x3 image -> stretch to 3072 x 1



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 27

April 18, 2017

Spatial structure is lost!

Convolutional neural networks (ConvNets/CNNs)

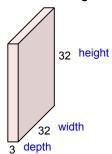
Architecture

Convolution layer

Arrange neurons as a **3D volume** naturally

Convolution Layer

32x32x3 image -> preserve spatial structure



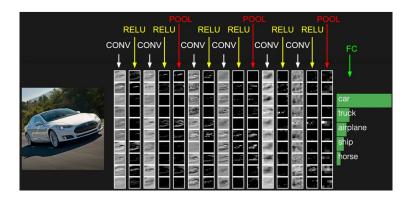
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 28 April 18, 2017

Solution: Convolutional Neural Net (ConvNet/CNN)

A special case of fully connected neural nets

- usually consist of convolution layers, ReLU layers, pooling layers, and regular fully connected layers
- key idea: learning from low-level to high-level features



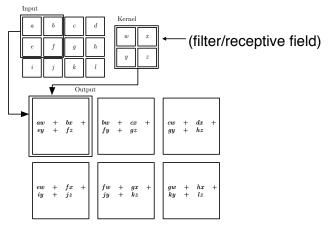
34 / 47

Convolutional neural networks (ConvNets/CNNs)

Architecture

Convolution

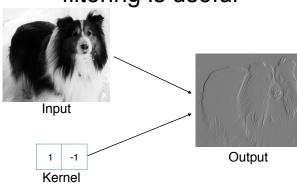
2D Convolution



Why convolution makes sense?

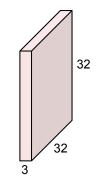
Main idea: if a filter is useful at one location, it should be useful at other locations.

A simple example why filtering is useful



Convolution Layer

32x32x3 image



5x5x3 filter



Convolve the filter with the image i.e. "slide over the image spatially, computing dot products"

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 29

April 18, 2017

Convolution Layer

32x32x3 image 32

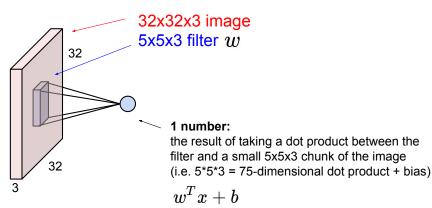
Filters always extend the full depth of the input volume

5x5x3 filter

Convolve the filter with the image i.e. "slide over the image spatially, computing dot products"

April 18, 2017

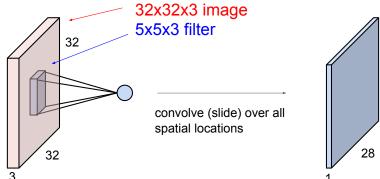
Convolution Layer



Fei-Fei Li & Justin Johnson & Serena Yeung

April 18, 2017 Lecture 5 - 31

Convolution Layer



Fei-Fei Li & Justin Johnson & Serena Yeung

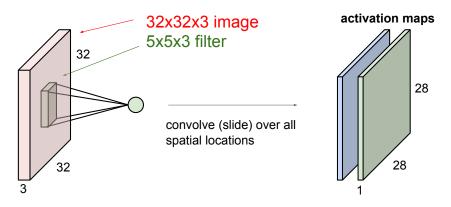
April 18, 2017 Lecture 5 - 32

activation map

28

Convolution Layer

consider a second, green filter

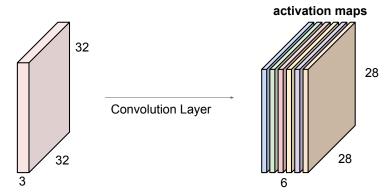


Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 33

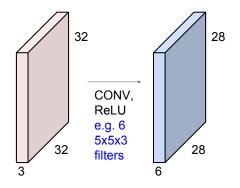
April 18, 2017

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!

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



Fei-Fei Li & Justin Johnson & Serena Yeung April 18, 2017 Lecture 5 - 34

Fei-Fei Li & Justin Johnson & Serena Yeung

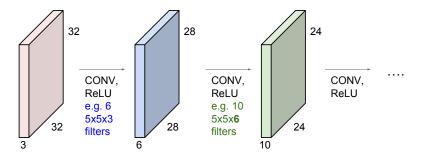
April 18, 2017 Lecture 5 - 35

Connection to fully connected NNs

A convolution layer is a special case of a fully connected layer:

• filter = weights with sparse connection

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



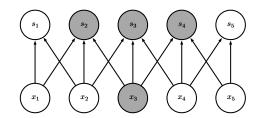
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 36 April 18, 2017

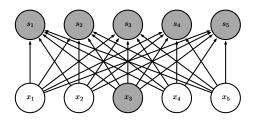
Demo: https://poloclub.github.io/cnn-explainer

Local Receptive Field Leads to Sparse Connectivity (affects less)

Sparse connections due to small convolution kernel

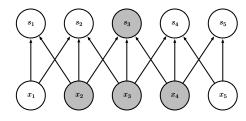


Dense connections



Sparse connectivity: being affected by less

Sparse connections due to small convolution kernel



Dense connections

(Goodfellow 2016)

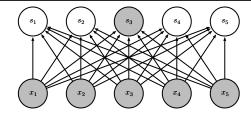


Figure 9.3

39 / 47

(Goodfellow 2016)

Connection to fully connected NNs

A convolution layer is a special case of a fully connected layer:

- filter = weights with sparse connection
- parameters sharing

40 / 47

Convolutional neural networks (ConvNets/CNNs) Architecture

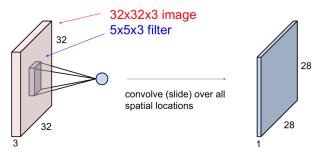
Connection to fully connected NNs

A convolution layer is a special case of a fully connected layer:

- filter = weights with sparse connection
- parameters sharing

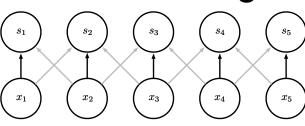
Much fewer parameters! Example (ignore bias terms):

- FC: $(32 \times 32 \times 3) \times (28 \times 28) \approx 2.4M$
- CNN: $5 \times 5 \times 3 = 75$



Parameter Sharing

Convolution shares the same parameters across all spatial locations



Traditional matrix multiplication does not share any parameters

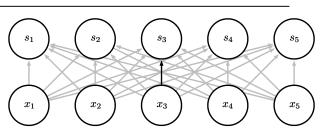


Figure 9.5

(Goodfellow 2016)

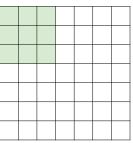
Convolutional neural networks (ConvNets/CNNs)

Spatial arrangement: stride and padding

7

A closer look at spatial dimensions:

7



7x7 input (spatially) assume 3x3 filter

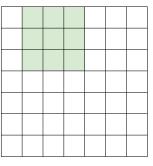
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 42

April 18, 2017

A closer look at spatial dimensions:

7



assume 3x3 filter

7x7 input (spatially)

Fei-Fei Li & Justin Johnson & Serena Yeung

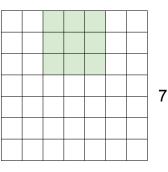
7

Lecture 5 - 43

April 18, 2017

A closer look at spatial dimensions:

7



7x7 input (spatially) assume 3x3 filter

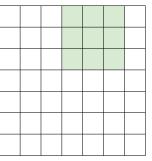
Fei-Fei Li & Justin Johnson & Serena Yeung

A closer look at spatial dimensions:

April 18, 2017 Lecture 5 - 44

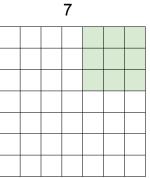
A closer look at spatial dimensions:

7



7x7 input (spatially) assume 3x3 filter

7



7x7 input (spatially) assume 3x3 filter

=> 5x5 output

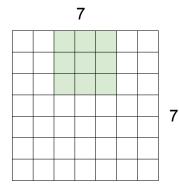
7

A closer look at spatial dimensions:

7

7x7 input (spatially) assume 3x3 filter applied with stride 2

A closer look at spatial dimensions:



7x7 input (spatially) assume 3x3 filter applied with stride 2

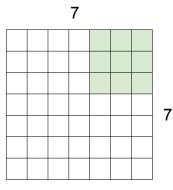
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 47 April 18, 2017

Fei-Fei Li & Justin Johnson & Serena Yeung

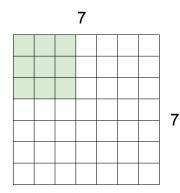
Lecture 5 - 48 April 18, 2017

A closer look at spatial dimensions:



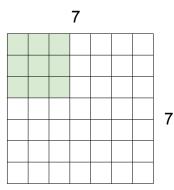
7x7 input (spatially) assume 3x3 filter applied with stride 2 => 3x3 output!

A closer look at spatial dimensions:



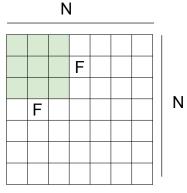
7x7 input (spatially) assume 3x3 filter applied with stride 3?

A closer look at spatial dimensions:



7x7 input (spatially) assume 3x3 filter applied with stride 3?

doesn't fit! cannot apply 3x3 filter on 7x7 input with stride 3.



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

e.g. N = 7, F = 3: stride 1 => (7 - 3)/1 + 1 = 5stride 2 => (7 - 3)/2 + 1 = 3stride 3 => (7 - 3)/3 + 1 = 2.33 :\

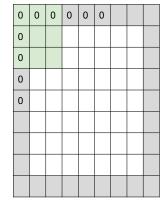
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 51 April 18, 2017

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 52 April 18, 2017

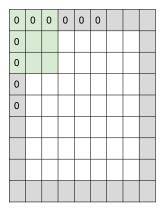
In practice: Common to zero pad the border



e.g. input 7x7
3x3 filter, applied with stride 1
pad with 1 pixel border => what is the output?

(recall:)
(N - F) / stride + 1

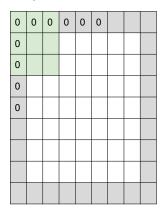
In practice: Common to zero pad the border



e.g. input 7x7
3x3 filter, applied with stride 1
pad with 1 pixel border => what is the output?

7x7 output!

In practice: Common to zero pad the border



e.g. input 7x7 3x3 filter, applied with stride 1 pad with 1 pixel border => what is the output?

7x7 output!

in general, common to see CONV layers with stride 1, filters of size FxF, and zero-padding with (F-1)/2. (will preserve size spatially)

e.g. $F = 3 \Rightarrow zero pad with 1$ $F = 5 \Rightarrow zero pad with 2$ $F = 7 \Rightarrow zero pad with 3$

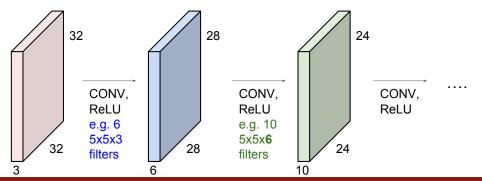
Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 55

April 18, 2017

Remember back to...

E.g. 32x32 input convolved repeatedly with 5x5 filters shrinks volumes spatially! (32 -> 28 -> 24 ...). Shrinking too fast is not good, doesn't work well.



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 56

April 18, 2017

Convolutional neural networks (ConvNets/CNNs) Architecture

Summary for convolution layer

Input: a volume of size $W_1 \times H_1 \times D_1$

Hyperparameters:

- K filters of size $F \times F$
- stride S
- amount of zero padding P (for one side)

Output: a volume of size $W_2 \times H_2 \times D_2$ where

•
$$W_2 = (W_1 + 2P - F)/S + 1$$

•
$$H_2 = (H_1 + 2P - F)/S + 1$$

•
$$D_2 = K$$

#parameters: $(F \times F \times D_1 + 1) \times K$ weights

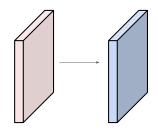
Common setting: F = 3, S = P = 1

Examples time:

Input volume: 32x32x3

10 5x5 filters with stride 1, pad 2

Output volume size: ?



Fei-Fei Li & Justin Johnson & Serena Yeung

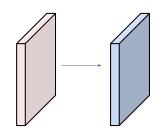
Lecture 5 - 57

April 18, 2017

Examples time:

Input volume: 32x32x3

10 5x5 filters with stride 1, pad 2



Output volume size:

(32+2*2-5)/1+1 = 32 spatially, so

32x32x10

Fei-Fei Li & Justin Johnson & Serena Yeung

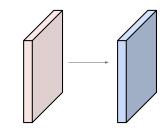
Lecture 5 - 58

April 18, 2017

Examples time:

Input volume: 32x32x3

10 5x5 filters with stride 1, pad 2



Number of parameters in this layer?

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 59

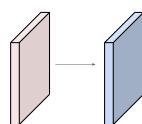
Architecture

April 18, 2017

Examples time:

Input volume: 32x32x3

10 5x5 filters with stride 1, pad 2



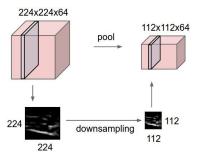
Number of parameters in this layer? each filter has 5*5*3 + 1 = 76 params (+1 for bias)

=> 76*10 = **760**

Convolutional neural networks (ConvNets/CNNs) Another element: pooling

Pooling layer

- makes the representations smaller and more manageable
- operates over each activation map independently:



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 5 - 72

April 18, 2017

Convolutional neural networks (ConvNets/CNNs)

Architecture

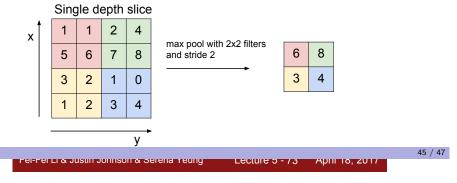
Pooling

Similar to a filter, except

- depth is always 1
- different operations: average, L2-norm, max
- no parameters to be learned

Max pooling with 2×2 filter and stride 2 is very common

MAX POOLING



Convolutional neural networks (ConvNets/CNNs)

Architecture

How to train a CNN?

How do we learn the filters/weights?

Essentially the same as FC NNs: apply SGD/backpropagation

Convolutional neural networks (ConvNets/CNNs)

Architectur

Putting everything together

Typical architecture for CNNs:

$$\mathsf{Input} \to [[\mathsf{Conv} \to \mathsf{ReLU}]^*\mathsf{N} \to \mathsf{Pool?}]^*\mathsf{M} \to [\mathsf{FC} \to \mathsf{ReLU}]^*\mathsf{Q} \to \mathsf{FC}$$

Common choices: $N \leq 5, Q \leq 2, M$ is large

Well-known CNNs: LeNet, AlexNet, ZF Net, GoogLeNet, VGGNet, etc.

All achieve excellent performance on image classification tasks.

Demo: https://poloclub.github.io/cnn-explainer

46 / 47