CS376 Computer Vision Lecture 20: Deep Learning Basics



Slide material from: http://cs231n.stanford.edu/slides/2018/cs231n_2018_lecture04.pdf http://cs231n.stanford.edu/slides/2018/cs231n_2018_lecture05.pdf

Today's Topic

- Neural Networks
 - Activation functions
 - Fully connected
 - Convolutional neural networks

• Optimization

Back-propagation

Neural Networks

From linear to multi-layer

(**Before**) Linear score function: f = Wx(Now) 2-layer Neural Network $f = W_2 \max(0, W_1 x)$ W2 h W1 Х S 10 100 3072 plane bird deer dog frog horse ship truck cat

From linear to multi-layer

(Before) Linear score function: f = Wx(Now) 2-layer Neural Network $f = W_2 \max(0, W_1 x)$ or 3-layer Neural Network

 $f=W_3\max(0,W_2\max(0,W_1x))$

Activation functions



Leaky ReLU $\max(0.1x, x)$



10

 $\begin{array}{l} \textbf{Maxout} \\ \max(w_1^T x + b_1, w_2^T x + b_2) \end{array}$



Neural Networks: Architectures



Summary

 We arrange neurons into fully-connected layers

• The abstraction of a layer has the nice property that it allows us to use efficient vectorized code (e.g. matrix multiplies)

• Neural networks are not really neural

Convolutional Neural Network

Convolutional Neural Networks

• A bit of history

Gradient-based learning applied to document recognition [LeCun, Bottou, Bengio, Haffner 1998]



LeNet-5

Convolutional Neural Networks

• A bit of history

ImageNet Classification with Deep Convolutional Neural Networks [Krizhevsky, Sutskever, Hinton, 2012]



Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

"AlexNet"

Recall fully connected layers

32x32x3 image -> stretch to 3072 x 1



Recall fully connected layers

32x32x3 image -> stretch to 3072 x 1



32x32x3 image -> preserve spatial structure







5x5x3 filter



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





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!

ConvNet

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



A typical ConvNet



A closer look at spatial dimensions



Activations



Output size

Ν



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

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

stride 2 => $(7 - 3)/2 + 1 = 3$
stride 3 => $(7 - 3)/3 + 1 = 2.33$:\

Padding is necessary

Padding



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 => zero pad with 1

F = 5 => zero pad with 2

F = 7 => zero pad with 3

Note that

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.



An example



Input volume: **32x32x3 10 5x5** filters with stride 1, pad 2

Number of parameters in this layer?

General form of a Conv layer

Common settings:

Summary. To summarize, the Conv Layer:

- Accepts a volume of size $W_1 imes H_1 imes D_1$
- Requires four hyperparameters:
 - Number of filters K,
 - $\circ\;$ their spatial extent F ,
 - the stride S,
 - the amount of zero padding P.
- Produces a volume of size $W_2 imes H_2 imes D_2$ where:

$$W_2 = (W_1 - F + 2P)/S + 1$$

- $H_2 = (H_1 F + 2P)/S + 1$ (i.e. width and height are computed equally by symmetry) • $D_2 = K$
- With parameter sharing, it introduces $F \cdot F \cdot D_1$ weights per filter, for a total of $(F \cdot F \cdot D_1) \cdot K$ weights and K biases.
- In the output volume, the d-th depth slice (of size $W_2 \times H_2$) is the result of performing a valid convolution of the d-th filter over the input volume with a stride of S, and then offset by d-th bias.

K = (powers of 2, e.g. 32, 64, 128, 512) F = 3, S = 1, P = 1 F = 5, S = 1, P = 2 F = 5, S = 2, P = ? (whatever fits) F = 1, S = 1, P = 0

An extreme case

(btw, 1x1 convolution layers make perfect sense)



The brain/neuron view of CONV layer





An activation map is a 28x28 sheet of neuron outputs:

- 1. Each is connected to a small region in the input
- 2. All of them share parameters

"5x5 filter" -> "5x5 receptive field for each neuron"

The brain/neuron view of CONV layer





E.g. with 5 filters, CONV layer consists of neurons arranged in a 3D grid (28x28x5)

There will be 5 different neurons all looking at the same region in the input volume

Reminder: Fully Connected Layer



Pooling layer

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



Max pooling

Single depth slice



У

Х

max pool with 2x2 filters and stride 2

6	8
3	4

General form

Common settings:

- Accepts a volume of size $W_1 imes H_1 imes D_1$
- Requires three hyperparameters:
 - their spatial extent F,
 - the stride S,
- Produces a volume of size $W_2 imes H_2 imes D_2$ where:
- Introduces zero parameters since it computes a fixed function of the input
- Note that it is not common to use zero-padding for Pooling layers

F = 2, S = 2F = 3, S = 2

Fully Connected Layer (FC Layer)

 Contains neurons that connect to the entire input volume, as in ordinary Neural Networks



Summary

- ConvNets stack CONV,POOL,FC layers
- Trend towards smaller filters and deeper architectures
- Trend towards getting rid of POOL/FC layers (just CONV)
- Typical architectures look like
- [(CONV-RELU)*N-POOL?]*M-(FC-RELU)*K,SOFTMAX where N is usually up to ~5, M is large, 0 <= K <= 2.
 - but recent advances such as ResNet/GoogLeNet challenge this paradigm

Optimization

Let us go back to the linear case

$$egin{aligned} s &= f(x;W) = Wx & ext{scores function} \ L_i &= \sum_{j
eq y_i} \max(0, s_j - s_{y_i} + 1) & ext{SVM loss} \ L &= rac{1}{N} \sum_{i=1}^N L_i + \sum_k W_k^2 & ext{data loss + regularization} \ ext{want} \ egin{aligned}
abla_W L & ext{} \end{aligned}$$

 $\vee WL$

Optimization





Vanilla Gradient Descent

while True:

weights_grad = evaluate_gradient(loss_fun, data, weights)
weights += - step_size * weights_grad # perform parameter update

Gradient descent

$$rac{df(x)}{dx} = \lim_{h o 0} rac{f(x+h) - f(x)}{h}$$

Numerical gradient: slow :(, approximate :(, easy to write :) Analytic gradient: fast :), exact :), error-prone :(

In practice: Derive analytic gradient, check your implementation with numerical gradient

Computtional graphs



Back-propagation

A toy example



Back-propagation



Back-propagation

- recursive application of the chain rule along a computational graph to compute the gradients of all inputs /parameters / intermediates
- implementations maintain a graph structure, where the nodes implement the forward() / backward() API
- forward: compute result of an operation and save any intermediates needed for gradient computation in memory
- backward: apply the chain rule to compute the gradient of the loss function with respect to the inputs