

Dr. Sedat Ozer





Announcements / Questions:

- HW1 questions
- Class schedule Fridays
- Midterm: Do we have a midterm?
- Paper presentations: **Due today**. Submit your presentation info at the link that was sent to you in my first email.
- Final project: Start working on your project early. Have frequent meetings with me especially if you are aiming



6x6x3

Multiple filters



Multiple filters



A Convolutional Layer of CNN (ConvNet)



Convolutional Layer Output: Activation_map = ReLU ((RGBImage * filter) + b) Output of a neuron in a FC NN: a = ReLU ($\mathbf{w}^{T}\mathbf{x} + b$) in FC NN

6

Input – Output Dimensions for *l*th Convolutional Layer

Input (One image): $n_{H}^{[l-1]} \ge n_{W}^{[l-1]} \ge n_{C}^{[l-1]}$

Padding size: $p^{[l]} x p^{[l]}$

Stride size: $s^{[l]} x s^{[l]}$

Filter size:
$$f^{[l]} \times f^{[l]} \times n_C^{[l-1]} \times n_C^{[l]} = One_Filter_{Dims} \times n_C^{[l]}$$

(Number of Weights = Filter size, Bias size: $n_C^{[l]}$)



The output dimensions for *m* input samples $A^{[l]}: m \ge n_H^{[l]} \ge n_W^{[l]} \ge n_W^{[l]}$

Output: $n_H^{[l]} \ge n_W^{[l]} \ge n_C^{[l]}$

Common Layer Types in a ConvNet

- Convolutional Layer (CONV)
- Fully Connected (FC, Dense)
- Locally Connected (not so commonly used!)
- Pooling (POOL)

- Activation Layer
 - (Sometimes, activation function is considered as a separate layer)

We have already seen both CONV and FC layers!

For now, we will use only those two layer types mentioned above.

These three layers have parameters to be learned!

No learned parameters for these layers!

Pooling layer: Max pooling



Pooling is typically applied on each channel separately.

9

2

4

Pooling layer: Max pooling







Pooling layer: Average pooling



Filter size= 2 x 2 Stride: s = 2



Pooling layer

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



Summary of pooling

Hyperparameters:

f : filter sizes : strideMax or Average pooling

Input:
$$n_H \ge n_W \ge n_C$$
 Output: $\left\lfloor \frac{n_H - f}{s} + 1 \right\rfloor \ge \left\lfloor \frac{n_W - f}{s} + 1 \right\rfloor \ge n_C$

Fully Connected Layer

32x32x3 image -> stretch to 3072 x 1



Fully Connected Layer

32x32x3 image -> stretch to 3072 x 1



(a 3072-dimensional dot product)

Convolution Layer



Convolution Layer

32x32x3 image



5x5x3 filter

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



Source: Fei-Feli Li & Justin Johnson & Serena Yeung

Convolution Layer



Convolution Layer



activation map



consider a second, green 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!

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



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



Parameter Sharing



One filter to rule them all!

One filter to rule them all, One w-vector to find them, One channel to bring them all and in the same-layer bind them.

- Gandalf the Grey



- Image Classification & Recognition: A main task in computer vision
 - PASCAL
 - IMAGENET

Pop Quiz

- How many parameters do we have in a conv layer with 100 filters that are 3x3x3 dimensional each?
- For one filter: we have 3x3x3 = 27 filter weights; and + 1 bias = **28** parameters.
- For all 100 filters: we have 28 x 100 = 2800 parameters to learn for that layer!

Pop Quiz

- Imagine that you are designing a CNN with 5 layers. The first 4 layers are convolutional and the last layer is FC layer with one neuron (using logistic reg).
 - The input dims are: 600x600x3.
 - In each of the 4 conv. layers:
 - We have 10 filters (for each, f = 3)
 - Stride is 2
 - Padding is "valid".

What is the output dims at the end of the 5th layer?

30 seconds to return your answers! ③

(Slightly Modified) LeNet-5



[LeCun et al., 1998. Gradient-based learning applied to document recognition]

(Slightly Modified) LeNet-5



	Activation shape	Activation Size	# parameters		Activation shape	Activation Size	# parameters
Input:	(32,32,1)	1,024	0	Input:	(32,32,3)	3,072	0
CONV1 (f=5, s=1)	(28,28,6)	4,704	150 + 6 bias = 156	CONV1 (f=5, s=1)	(28,28,6)	4,704	450 + 6 bias = 456
POOL1	(14,14,6)	1,176	0	POOL1	(14,14,6)	1,176	0
CONV2 (f=5, s=1)	(10,10,16)	1,600	$(5x5x6) \ge 16=$ 2400 2400 + 16 = 2416	CONV2 (f=5, s=1)	(10,10,16)	1,600	$(5x5x6) \ge 16=$ 2400 2400 + 16 = 2416
POOL2	(5,5,16)	400	0	POOL2	(5,5,16)	400	0
FC3	(120,1)	120	48,000 + 120 = 48,120	FC3	(120,1)	120	48,000 + 120 = 48,120
FC4	(84,1)	84	10,080 + 84 = 10,164	FC4	(84,1)	84	10,080 + 84 = 10,164
Softmax	(10,1)	10	84x10 + 10 = 850	Softmax	(10,1)	10	84x10 + 10 = 850

For: Colored (RGB) image

For: Grayscale image (single channel)

Variables: name (type shape) [size]

Tensorflow output for the (grayscale) model summarizing the model parameters: ^{conv2_biases:0} (float32_ref 16) [16, bytes: 64] Variable_2:0 (float32_ref 400x120) [48000, bytes: 192000] fc1_biases:0 (float32_ref 120) [120, bytes: 480] Variable_3:0 (float32_ref 120) [120, bytes: 480] fc2_biases:0 (float32_ref 84) [10080, bytes: 40320]

Variable:0 (float32_ref 5x5x1x6) [150, bytes: 600] conv1_biases:0 (float32_ref 5x5x1x6) [2400, bytes: 24] Variable_1:0 (float32_ref 5x5x6x16) [2400, bytes: 9600] conv2_biases:0 (float32_ref 16) [16, bytes: 64] Variable_2:0 (float32_ref 120) [120, bytes: 480] Variable_3:0 (float32_ref 120) [120, bytes: 480] Variable_3:0 (float32_ref 120) [10080, bytes: 40320] fc2_biases:0 (float32_ref 84) [84, bytes: 336] Variable_4:0 (float32_ref 84x10) [840, bytes: 3360] fc3_biases:0 (float32_ref 10) [10, bytes: 40] Total size of variables: 61706 Total bytes of variables: 246824 sedat@sedat-ThinkPad-P50:~/



[Krizhevsky et al., 2012. ImageNet classification with deep convolutional neural networks]

Lecture Notes for Computer Vision Sedat Ozer

- 8 layers total
- Trained on Imagenet dataset [Deng et al. CVPR'09]
- 18.2% top-5 error
- Our reimplementation: 18.1% top-5 error



Rob Fergus, NIPS 2013

- Remove top fully connected layer
 - Layer 7
- Drop 16 million parameters
- Only 1.1% drop in performance!



- Remove both fully connected layers
 - Layer 6 & 7
- Drop ~50 million parameters
- 5.7% drop in performance









Lecture Notes for Computer Vision Sedat Ozer

Residual Network - ResNet (Microsoft)



A "plain" (deep) network



A residual network (with 5 residual blocks)

[He et al., 2015. Deep residual networks for image recognition]

Residual Building Block

Skip Connection (Short Cut) $a^{[l]} \rightarrow a^{[l+1]} \rightarrow a^{[l+2]}$

$$a^{[l]} \xrightarrow{} z^{[l+1]} \xrightarrow{} \text{ReLU} \xrightarrow{} a^{[l+1]} \xrightarrow{} z^{[l+2]} \xrightarrow{} \xrightarrow{} \text{ReLU} \xrightarrow{} a^{[l+2]}$$

 $z^{[l+1]} = W^{[l+1]} a^{[l]} + b^{[l+1]} a^{[l+1]} = g(z^{[l+1]}) z^{[l+2]} = W^{[l+2]} a^{[l+1]} + b^{[l+2]} a^{[l+2]} = g(z^{[l+2]})$

With skip connection:
$$a^{[l+2]} = g(z^{[l+2]} + a^{[l]})$$

[He et al., 2015. Deep residual networks for image recognition]

Residual Network



ResNet



[Source: He et al., 2015. Deep residual networks for image recognition]

Inception Network



(GoogleNet)

Main motivation for inception network



[Szegedy et al. 2014. Going deeper with convolutions]

What does a 1 × 1 convolution do?



[Lin et al., 2013. Network in network]

45

1×1 convolutions can shrink the layer dims



Helps reducing the number of parameters!

[Lin et al., 2013. Network in network]

The problem of computational cost



- One filter volume: 5x5x192 for one voxel in the output volume (in the activation map).
- Since we have 28x28x32 dimensional volume at the output, we need: (5x5x192) x (28x28x32) multiplications! (~120Million)

47

Using 1×1 convolution for reducing the computation



 $28 \times 28 \times 192$

Multiplications needed: (28x28x16) x (1x1x192) = ~ 2.4 Million (28x28x32) x (5x5x16) = ~ 10 Million

Total: 12.4 Million << 120Million

Inception module



Inception Network (Google)

GooLeNet or GoogLeNet





Local Response Normalization (LRN)

[Szegedy et al., 2014, Going Deeper with Convolutions]

• This lecture has materials from CS231n (Stanford), Andrew Ng and Ulas Bagci.