

CSCI 599: Deep Learning and its Applications

Lecture 5

Fall 2017
Joseph J. Lim

Help from Shao-Hua Sun, Youngwoon Lee, and Te-Lin Wu

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Speakers visiting USC

- Phillip Isola (OpenAI) on October 19th
- Jimmy Ba (Univ. of Toronto) on November 9th
- Carl Vondrick (Google) on November 28th

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Disclaimer

- This course is taught for the 1st time @ USC. This course is 599, and thus an **experimental** course.
- The syllabus, course policy, and grading details **may change** over the semester (**check website!**)
- If you prefer a well-structured course, this is **NOT** a course for you, and I encourage you to take the course next year. We really mean this.
- But, it will be **fun** and **challenging**!

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Important Dates

- **Midterm: week 7**
- Assignment 1: **week 8** (out: week 5)
- Assignment 2: week 11
- Project
 - Team formation: week 4
 - Project proposal: week 8
 - **Project meeting with TA #1: between week 4 - week 8**
 - Project meeting with Instructor #1: week 8 (M-W)
 - Project mid-report: week 12
 - Project meeting with TA #2: between week 8 - week 12
 - Project meeting with Instructor #2: week 11 (M-W)
 - Project report + Final presentation: week 15 (5-9:30pm) **4.5 hours**
 - Project meeting with TA #3: between week 12 - week 15

Subject to change!

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Assignment 1

- Assignment 1 is released on the website (<https://csci599-dl.github.io/>)
- DUE October 11th, 2017 (week 8)
- Collaboration is OK! Please list all collaborators' names.

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Project

- Report: a blog post (rather than a paper format)
- Lots of “what if we fail” questions

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Project

Please take this course to do a fun project, not to get an A

- Report: a blog post (rather than a paper format)
- Lots of “what if we fail” questions

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Project meeting with TA #1

- We posted on Piazza your TA assignment
 - Temporary until the proposal due date
- Please sign up for your team on your TA's doodle link.
- Please be available — if not, we may not be able to schedule one with you.

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Project meeting with Instructor #1

- We will follow up on this (Monday to Wednesday of week 8).
- Please be available — if not, we may not be able to schedule one with you.
- I am meeting 58 teams over 3 days



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Note

- **No office hour next Wednesday (9/27).**
- Potentially, also a week after (10/4) — This will be announced on Piazza.

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Today's agenda

- Part 1
 - Convolutional Neural Networks
- Part 2
 - Training Neural Networks

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Today's agenda

- Part 1
 - Convolutional Neural Networks
- Part 2
 - Training Neural Networks

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Different NN models

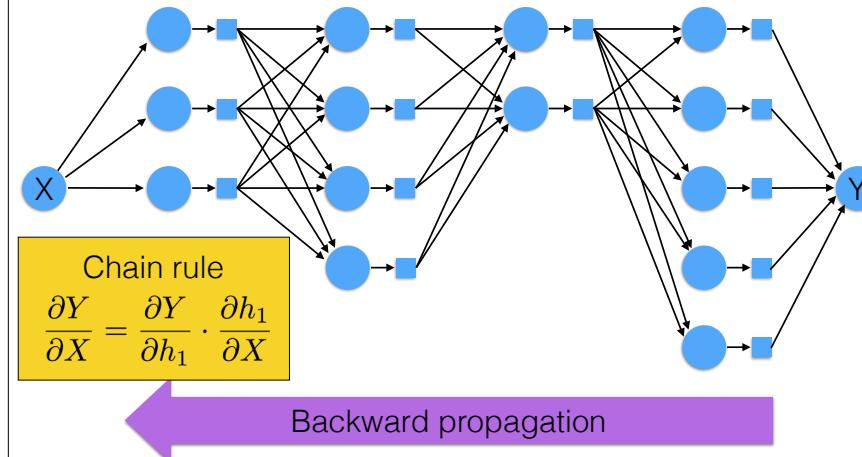
- Feedforward networks (CNNs, ...)
- Recurrent networks (RNNs, ...)
- Memory networks
- ...

Joseph J. Lim

CSCI 599 @ USC

Lecture 4

Recap: Backward Propagation



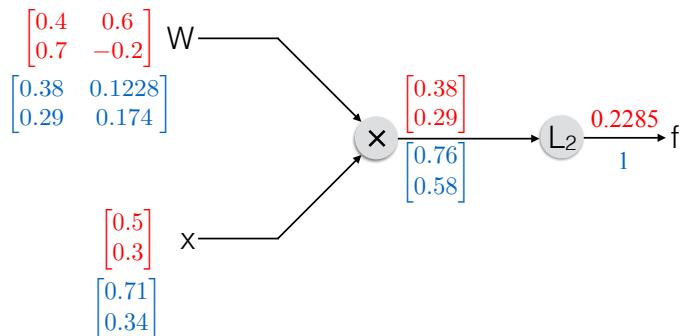
Joseph J. Lim

CSCI 599 @ USC

Lecture 4

Recap: Computational Graph

A vectorized example: $f(x, W) = \|W \cdot x\|^2 = \sum_{i=1}^n (W \cdot x)_i^2$



Joseph J. Lim

CSCI 599 @ USC

Lecture 4

All useless

Geoff Hinton (the godfather of deep learning) said we need to start over.

On 9/15 (two days after our back-prop lecture), he said he is “**deeply suspicious**” about back-propagation.



Geoffrey Hinton: Artificial intelligence pioneer says we need to start over

In 1986, Geoffrey Hinton co-authored a paper that, three decades later, is central to the explosion of artificial intelligence. But Hinton says his breakthrough method should be dispensed with, and a new path to AI found.

Speaking with Axios on the sidelines of an AI conference in Toronto on Wednesday, Hinton, a professor emeritus at the University of Toronto and a Google researcher, said he is now “deeply suspicious” of back-propagation, the workhorse method that underlies most of the advances we are seeing in the AI field today, including the capacity to sort through photos and talk to Siri. “My view is throw it all away and start again,” he said.

The bottom line: Other scientists at the conference said back-propagation still has a core role in AI’s future. But Hinton said that, to push forward, entirely new methods will probably have to be invented. “Max Planck said, ‘Science progresses one funeral at a time.’ The future depends on some graduate student who is deeply suspicious of everything I have said.”

How it works: In back propagation, labels or “weights” are used to represent a photo or voice within a brain-like neural layer. The weights are then adjusted and readjusted, layer by layer, until the network can perform an intelligent function with the fewest possible errors.

But Hinton suggested that, to get where neural networks are able to become intelligent on their own, what is known as “unsupervised learning,” “I suspect that means getting rid of back-propagation.”

“I don’t think it’s how the brain works,” he said. “We clearly don’t need all the labeled data.”

<https://wwwaxios.com/ai-pioneer-advocates-starting-over-2485537027.html>

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

All useless

Geoff Hinton (the godfather)

Artificial intelligence pioneer says we need to start over



More seriously...

watch his capsule talk @ <https://www.youtube.com/watch?v=TFIMqt0yT2I>

But Hinton suggested that, to get to where neural networks are able to become intelligent on their own, what is known as "unsupervised learning." "I suspect that means getting rid of back-propagation."

"I don't think it's how the brain works," he said. "We clearly don't need all the labeled data."

<https://wwwaxios.com/ai-pioneer-advocates-starting-over-2485537027.html>

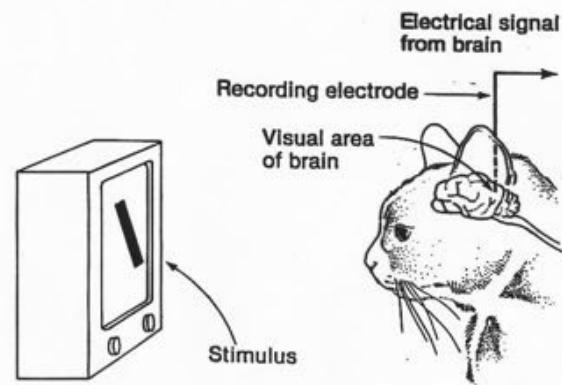
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

History

Hubel & Wiesel



Understanding how the system constructs complex representation from simple stimulus features.

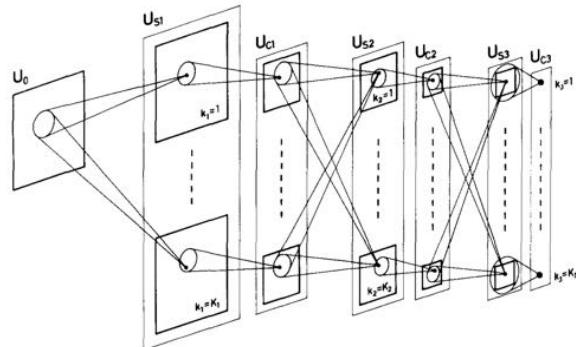
Joseph J. Lim

CSCI 599 @ USC

Lecture 1

History

Neocognitron



Fukushima, Kunihiko. Neocognitron: A self-organizing neural network model for a mechanism of visual pattern recognition. Competition and cooperation in neural nets 1982.

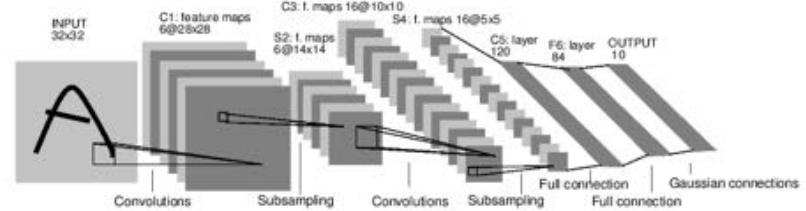
Joseph J. Lim

CSCI 599 @ USC

Lecture 1

History

LeNet



Y. LeCun, et. al. Handwritten digit recognition with a back-propagation network. NIPS 1989.

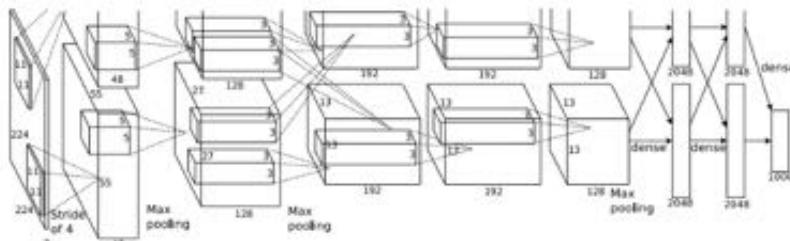
Joseph J. Lim

CSCI 599 @ USC

Lecture 1

History

AlexNet



A Krizhevsky, et. al. ImageNet Classification with Deep Convolutional Neural Networks. NIPS 2012.

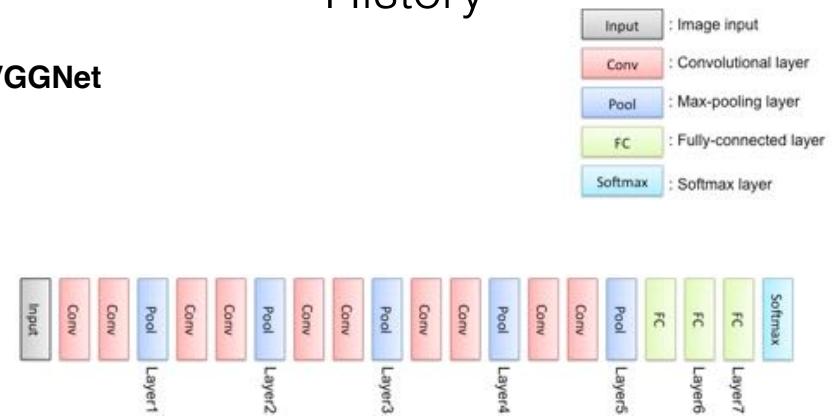
Joseph J. Lim

CSCI 599 @ USC

Lecture 1

History

VGGNet



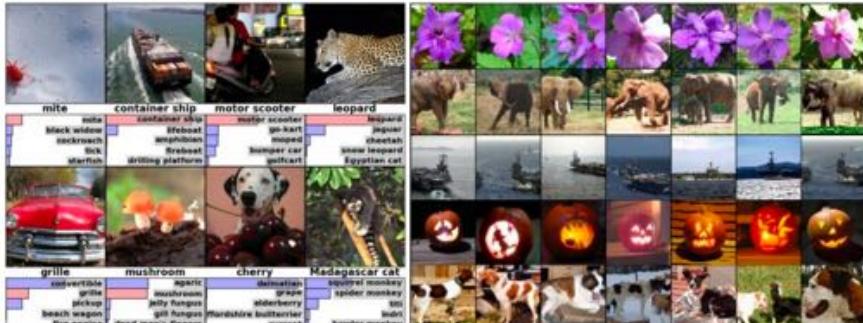
Joseph J. Lim

CSCI 599 @ USC

Lecture 1

Convolutional NNs

Classification

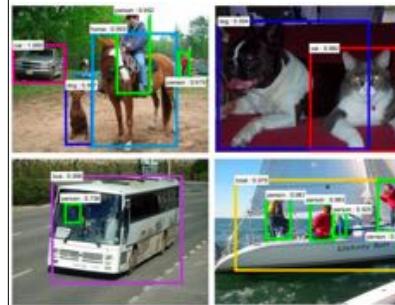


Figures copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

Retrieval

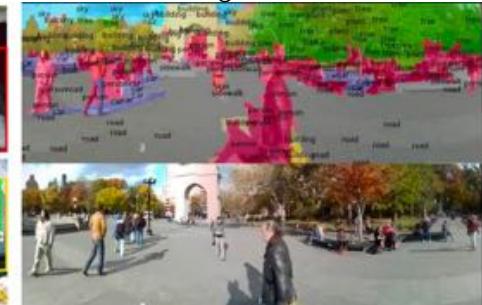
Convolutional NNs

Detection



Figures copyright Shaoqing Ren, Kaiming He, Ross Girshick, Jian Sun, 2015. Reproduced with permission.
[Faster R-CNN: Ren, He, Girshick, Sun 2015]

Segmentation



Figures copyright Clement Farabet, 2012. Reproduced with permission.
[Farabet et al., 2012]

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Slide credit: CS 231N

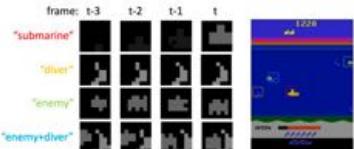
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Convolutional NNs

Game AI



[Guo et al. 2014]



Figures copyright Xiaoxiao Guo, Satinder Singh, Honglak Lee, Richard Lewis, and Xiaoshi Wang, 2014. Reproduced with permission.



Original image is CC0 public domain
Starry Night and Tree Roots by Van Gogh are in the public domain
Bokesh image is in the public domain
Stylized images copyright Justin Johnson, 2017;
reproduced with permission.

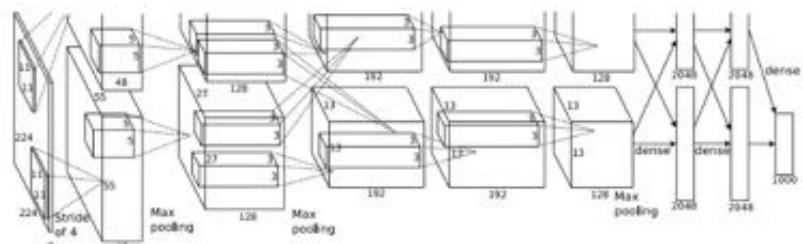
Arts

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

AlexNet in details



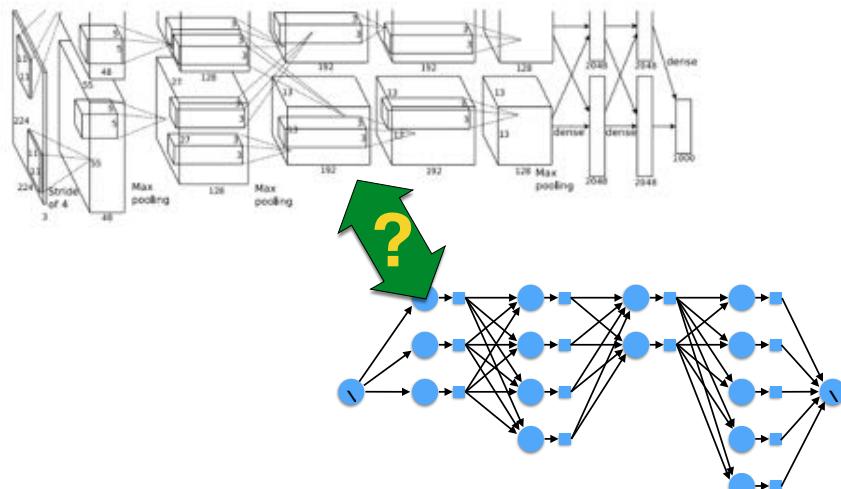
A Krizhevsky, et. al. ImageNet Classification with Deep Convolutional Neural Networks. NIPS 2012.

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

AlexNet in details

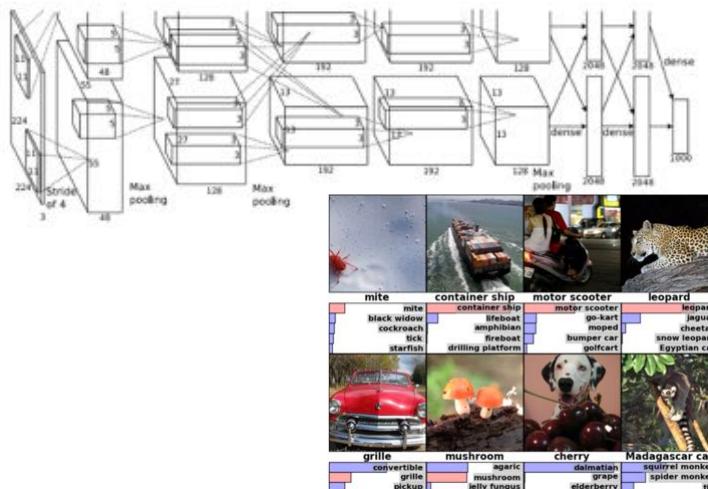


Joseph J. Lim

CSCI 599 @ USC

Lecture 4

AlexNet in details

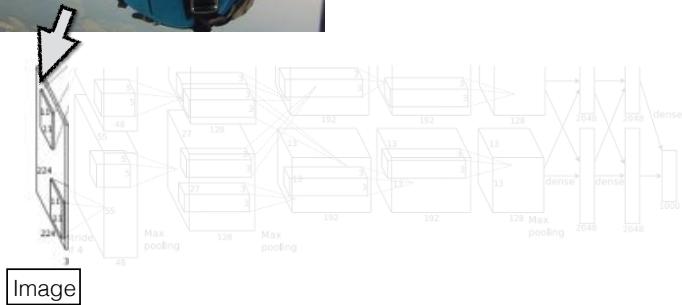


Joseph J. Lim

CSCI 599 @ USC

Lecture 4

AlexNet in details

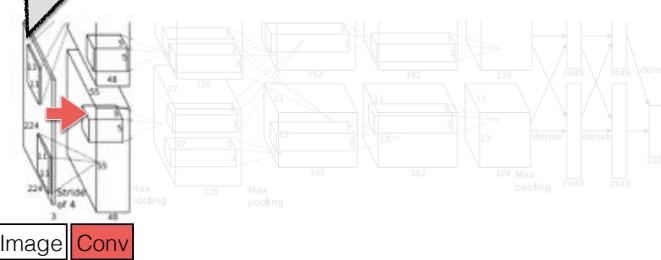


Joseph J. Lim

CSCI 599 @ USC

PC: Daniel Yang
Lecture 5

AlexNet in details

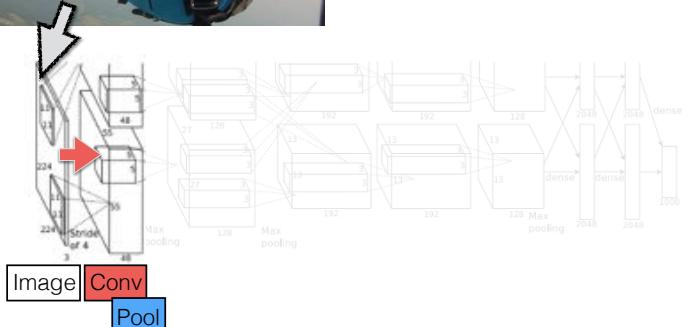


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

AlexNet in details

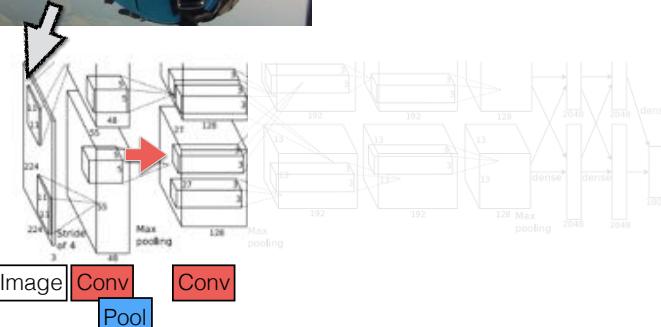


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

AlexNet in details

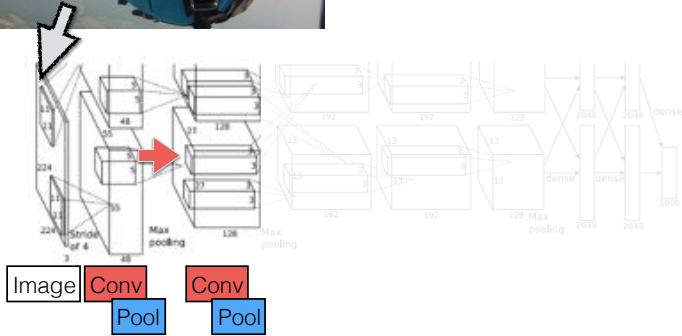


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

AlexNet in details

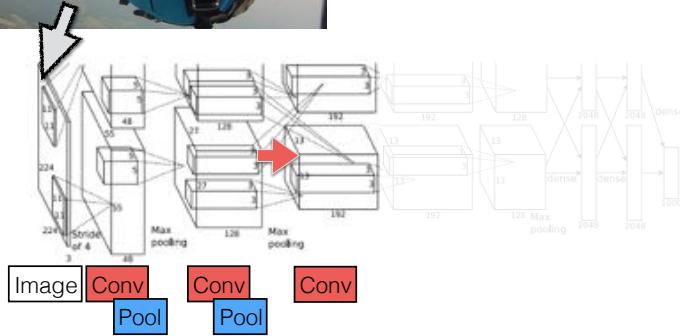


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

AlexNet in details

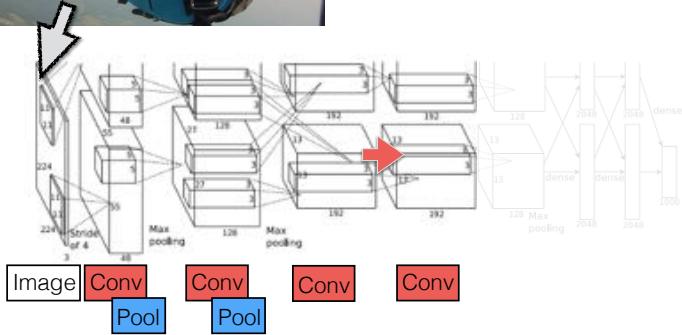


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

AlexNet in details

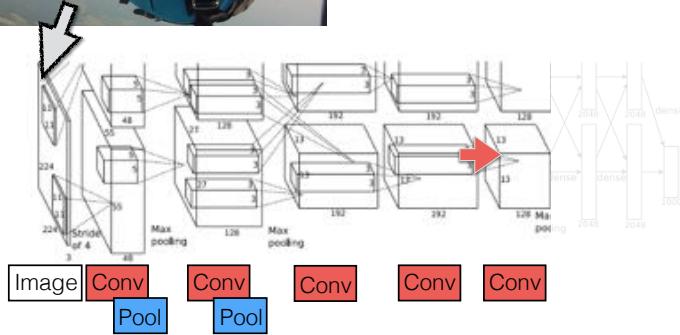


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

AlexNet in details

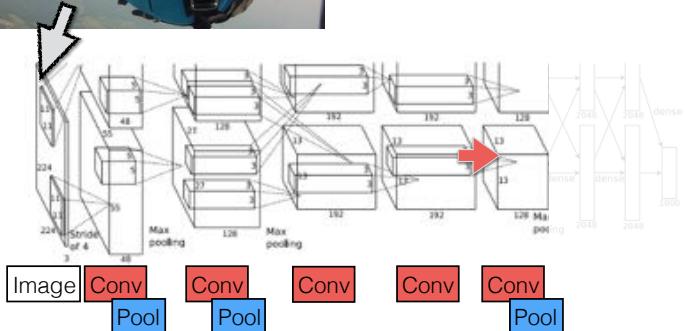


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

AlexNet in details

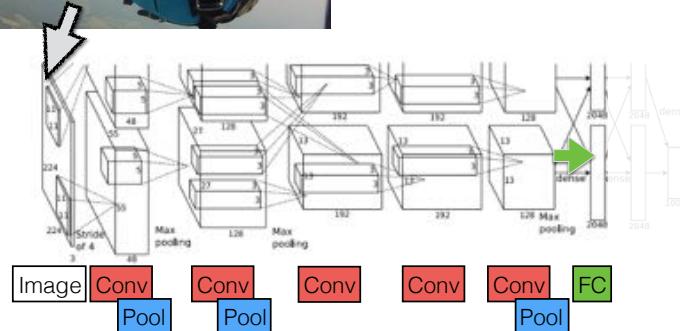


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

AlexNet in details

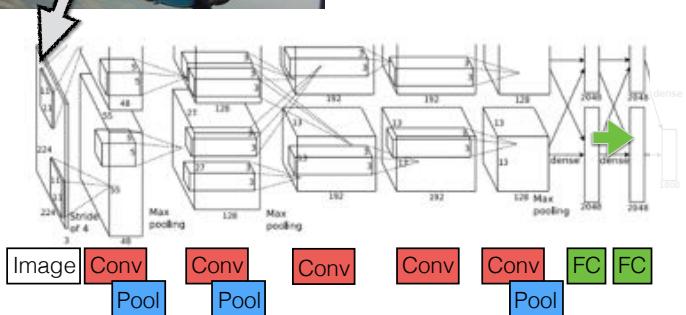


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

AlexNet in details

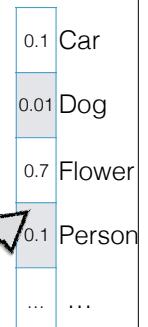
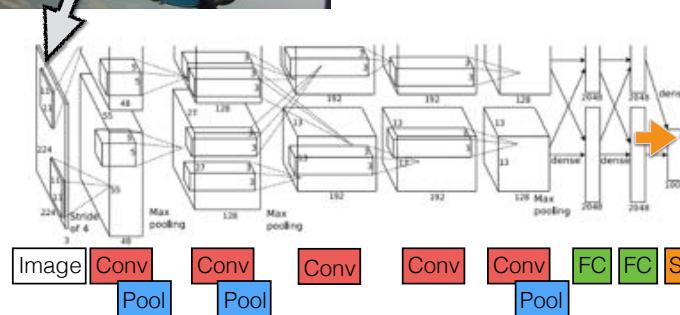


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

AlexNet in details

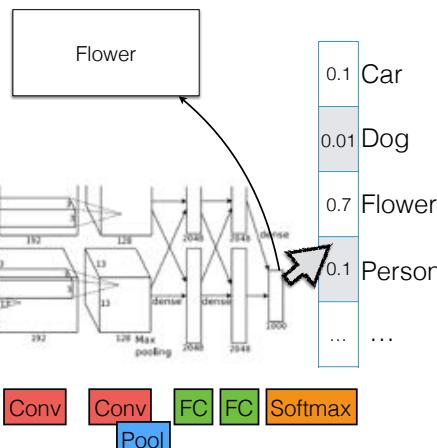


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

AlexNet in details



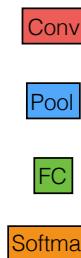
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

What are these layers?

- Convolutional layer
- Pool layer
- Fully connected layer
- Softmax layer



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

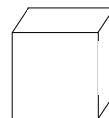
What are these layers?

- Convolutional layer
- Pool layer
- Fully connected layer
- Softmax layer



Layer = Function

Fully connected layer (FC)



M x N x P

Joseph J. Lim

CSCI 599 @ USC

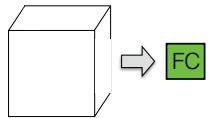
Lecture 5

Joseph J. Lim

CSCI 599 @ USC

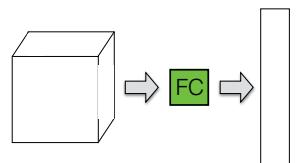
Lecture 5

Fully connected layer (FC)



$M \times N \times P$

Fully connected layer (FC)



$M \times N \times P$ $Q \times 1$

Joseph J. Lim

CSCI 599 @ USC

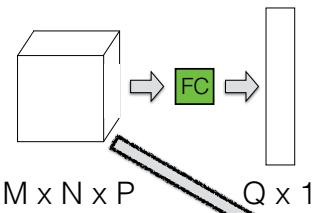
Lecture 5

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Fully connected layer (FC)



$M \times N \times P$ $Q \times 1$

FC operation takes an input, stretches it.

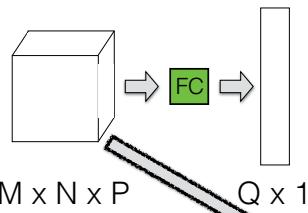
$(M^*N^*P) \times 1$

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Fully connected layer (FC)



$M \times N \times P$ $Q \times 1$

FC operation takes an input, stretches it.



$Q \times (M^*N^*P)$ $(M^*N^*P) \times 1$

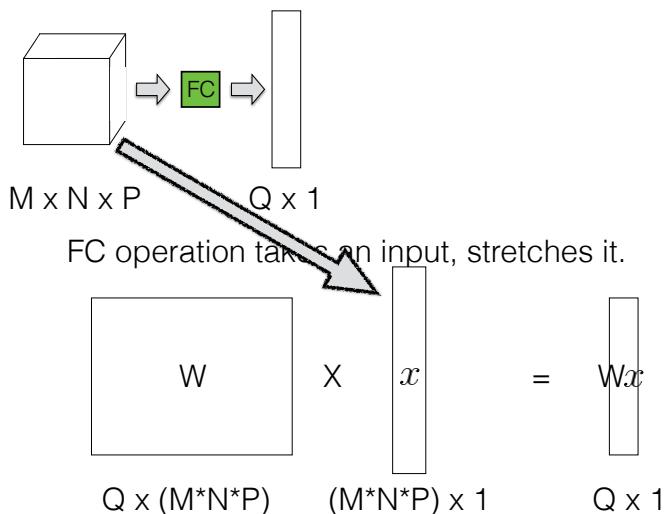
Joseph J. Lim

CSCI 599 @ USC

CSCI 599 @ USC

Lecture 5

Fully connected layer (FC)

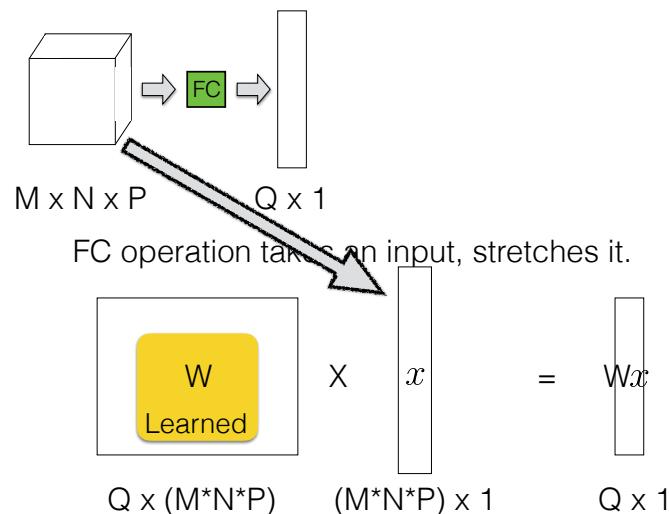


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Fully connected layer (FC)

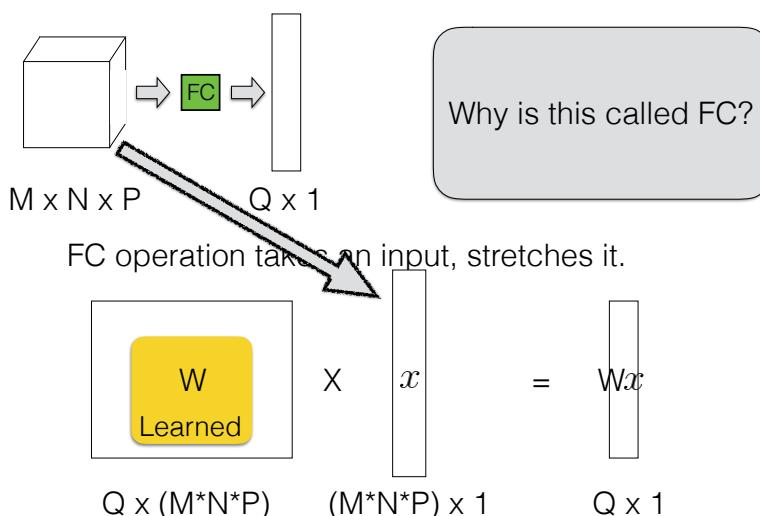


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Fully connected layer (FC)

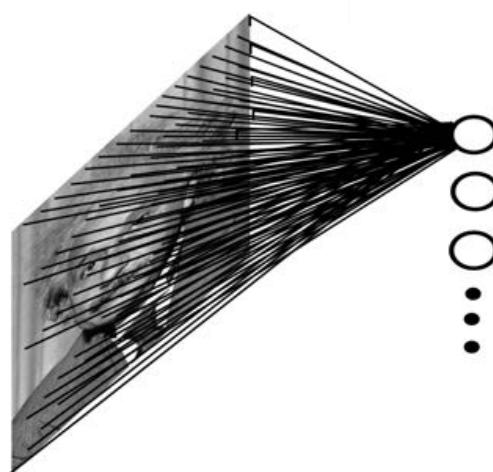


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Fully connected layer (FC)



Slide credit: Marc'Aurelio Ranzato

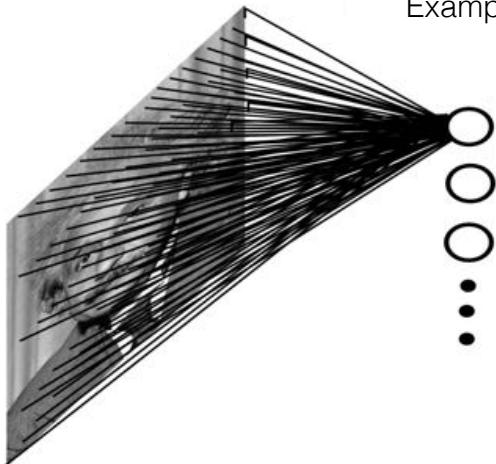
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Fully connected layer (FC)

Example: 200x200 image



Joseph J. Lim

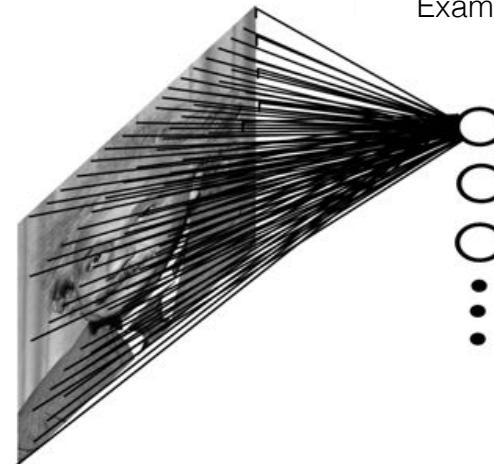
CSCI 599 @ USC

Lecture 5

Slide credit: Marc'Aurelio Ranzato

Fully connected layer (FC)

Example: 200x200 image
40k hidden units



Joseph J. Lim

CSCI 599 @ USC

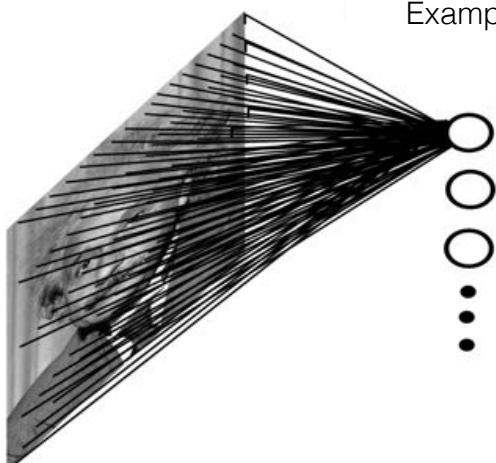
Lecture 5

Slide credit: Marc'Aurelio Ranzato

Fully connected layer (FC)

Example: 200x200 image
40k hidden units
~2B parameters!!!

Why?



Joseph J. Lim

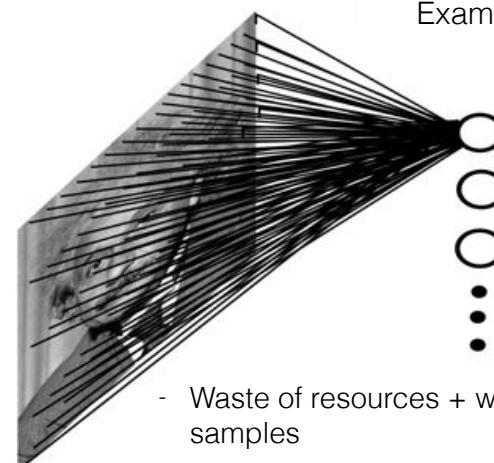
CSCI 599 @ USC

Lecture 5

Slide credit: Marc'Aurelio Ranzato

Fully connected layer (FC)

Example: 200x200 image
40k hidden units
~2B parameters!!!



- Waste of resources + we have not enough training samples

Joseph J. Lim

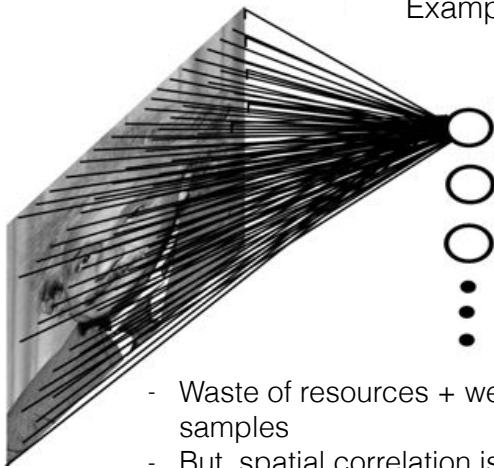
CSCI 599 @ USC

Lecture 5

Slide credit: Marc'Aurelio Ranzato

Fully connected layer (FC)

Example: 200x200 image
40k hidden units
~2B parameters!!!



- Waste of resources + we have not enough training samples
- But, spatial correlation is local!

Slide credit: Marc'Aurelio Ranzato

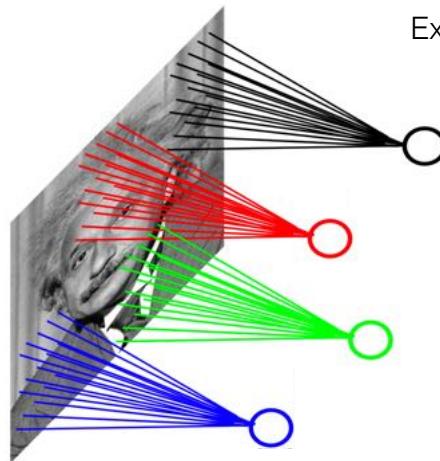
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Locally connected layer

Example: 200x200 image
40k hidden units
Filter size: 10x10



Slide credit: Marc'Aurelio Ranzato

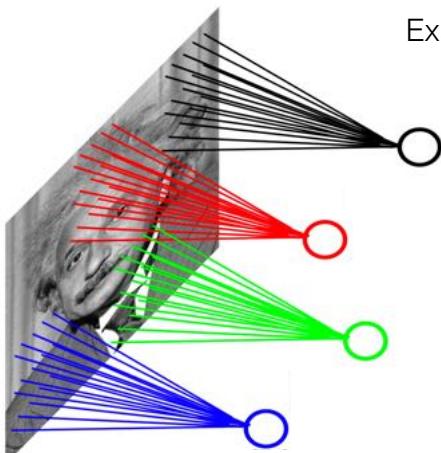
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Locally connected layer

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



Slide credit: Marc'Aurelio Ranzato

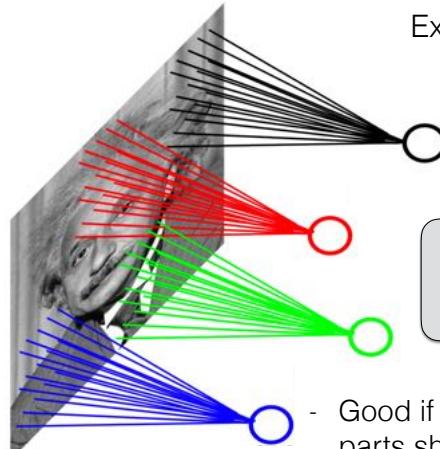
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Locally connected layer

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



What can we do?

- Good if images are registered (i.e. parts show up at the same locations)

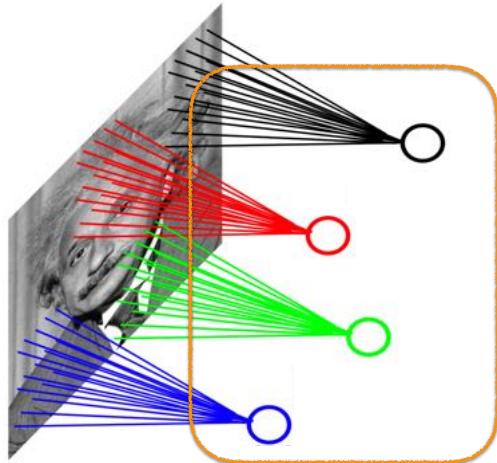
Slide credit: Marc'Aurelio Ranzato

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Locally connected layer



Share all local weights

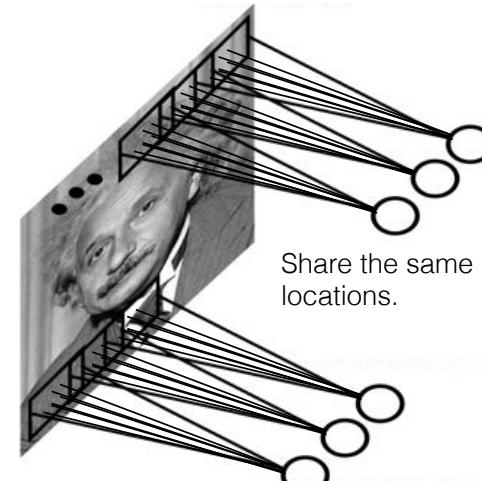
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Slide credit: Marc'Aurelio Ranzato

Convolutional layer



Share the same parameters across different locations.

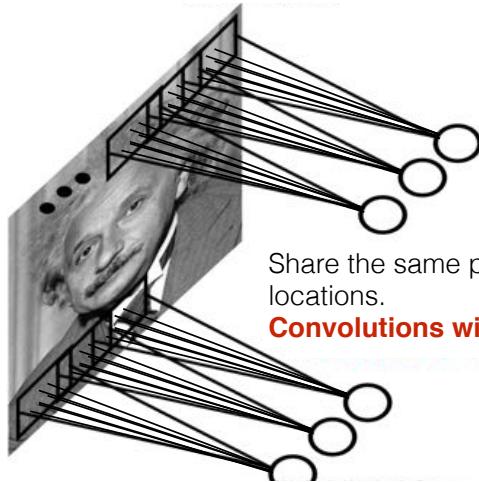
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Slide credit: Marc'Aurelio Ranzato

Convolutional layer



Share the same parameters across different locations.

Convolutions with learned kernels (weights)

Joseph J. Lim

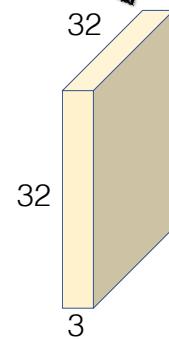
CSCI 599 @ USC

Lecture 5

Slide credit: Marc'Aurelio Ranzato

Convolutional layer

32x32x3 input



5x5x3 filter



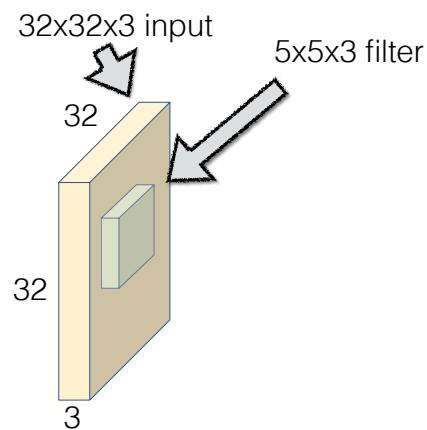
Slide credit: CS 231N

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Convolutional layer

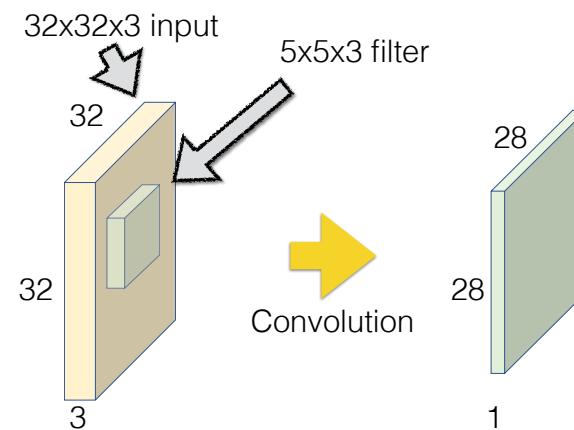


Joseph J. Lim

CSCI 599 @ USC

Slide credit: CS 231N
Lecture 5

Convolutional layer



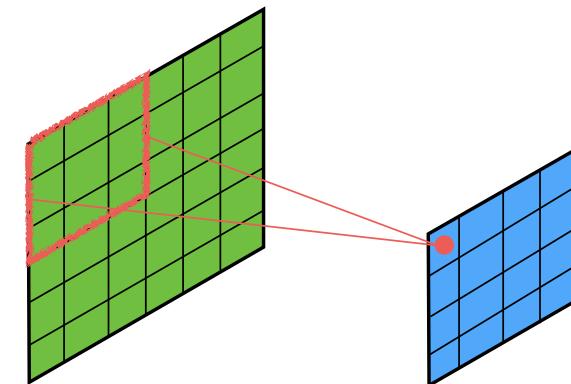
Joseph J. Lim

CSCI 599 @ USC

Slide credit: CS 231N
Lecture 5

Convolution?

Convolutional layer



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

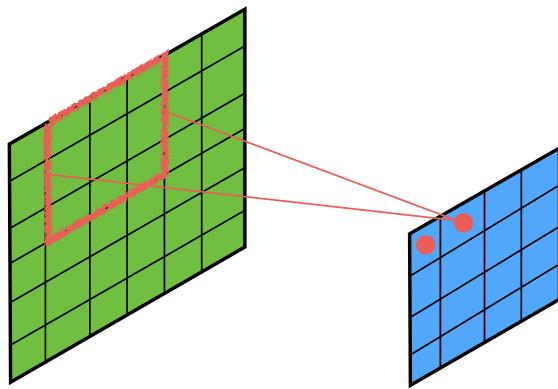
Inspired by Marc'Aurelio Ranzato's slides

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Convolutional layer



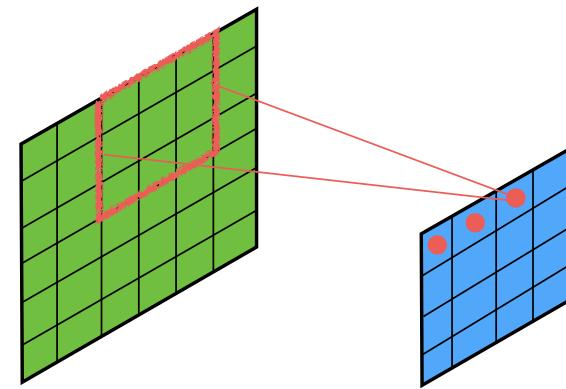
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Inspired by Marc'Aurelio Ranzato's slides

Convolutional layer



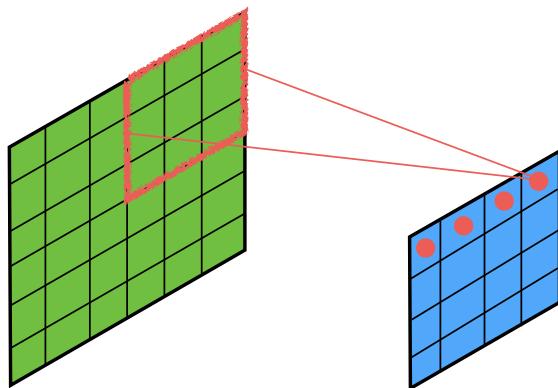
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Inspired by Marc'Aurelio Ranzato's slides

Convolutional layer



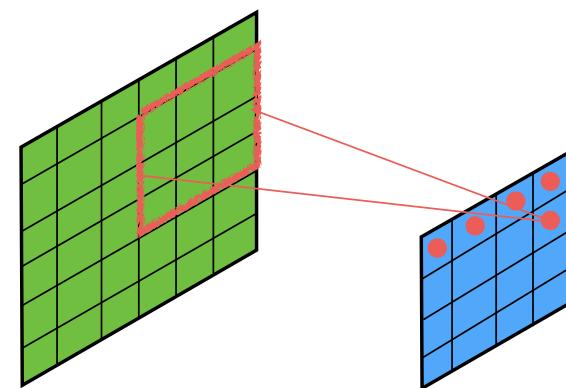
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Inspired by Marc'Aurelio Ranzato's slides

Convolutional layer



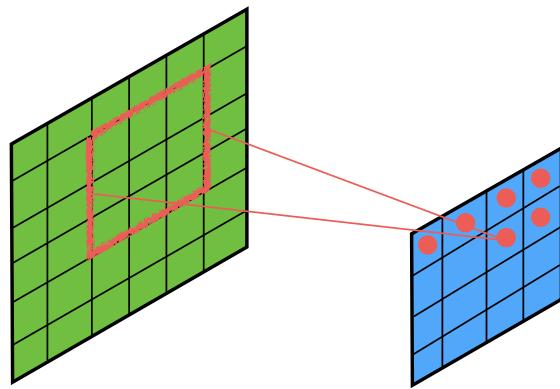
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Inspired by Marc'Aurelio Ranzato's slides

Convolutional layer



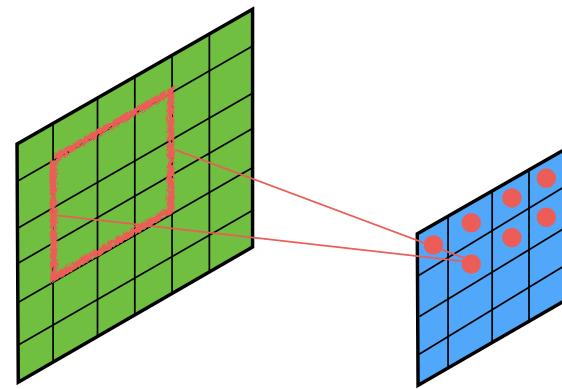
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Inspired by Marc'Aurelio Ranzato's slides

Convolutional layer



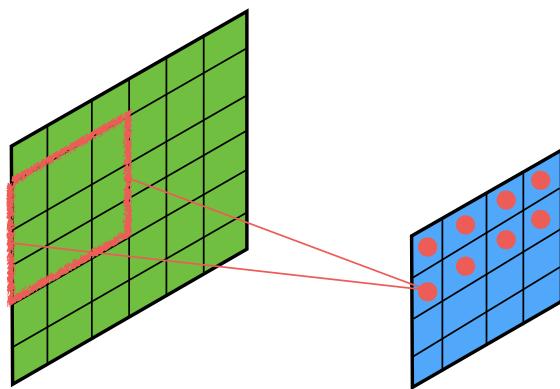
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Inspired by Marc'Aurelio Ranzato's slides

Convolutional layer



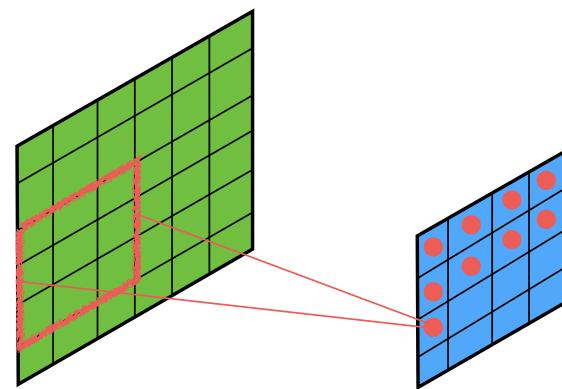
Inspired by Marc'Aurelio Ranzato's slides

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Convolutional layer



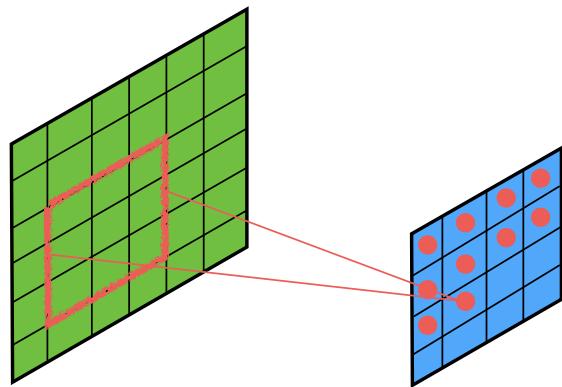
Inspired by Marc'Aurelio Ranzato's slides

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Convolutional layer



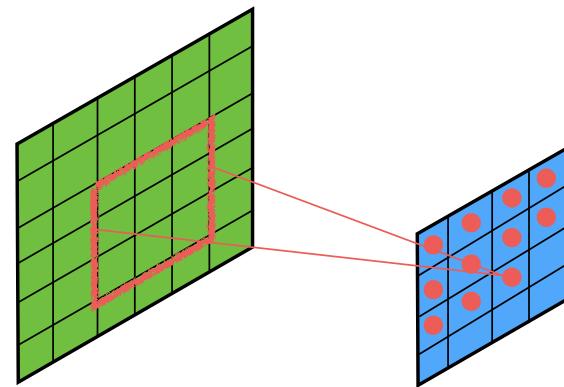
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Inspired by Marc'Aurelio Ranzato's slides

Convolutional layer



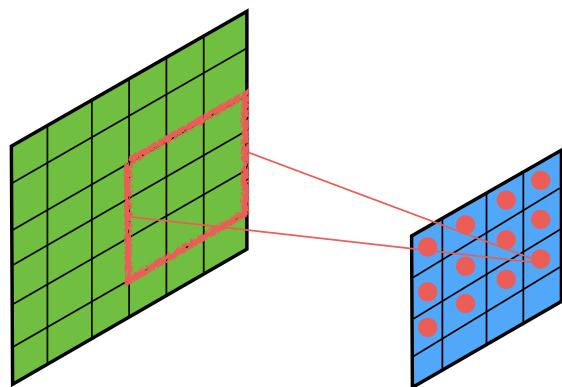
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Inspired by Marc'Aurelio Ranzato's slides

Convolutional layer



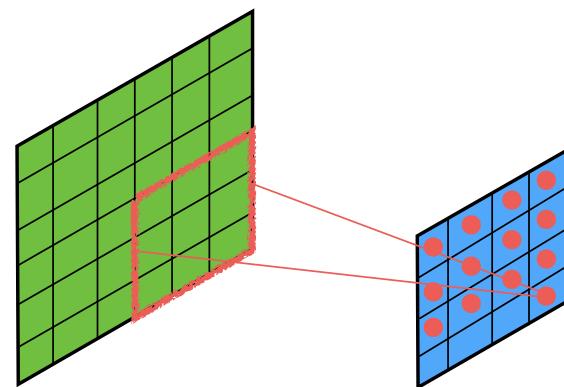
Inspired by Marc'Aurelio Ranzato's slides

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Convolutional layer



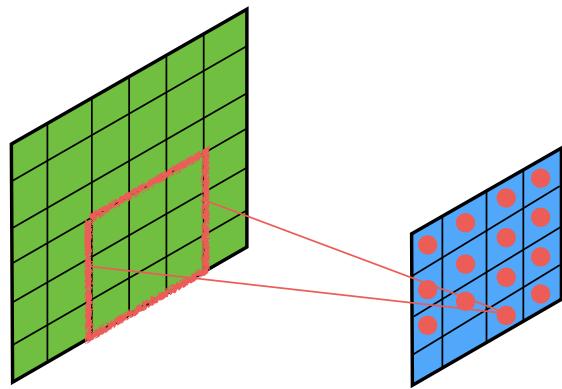
Inspired by Marc'Aurelio Ranzato's slides

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Convolutional layer



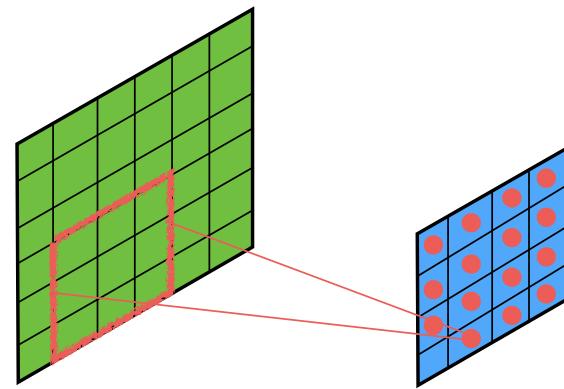
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Inspired by Marc'Aurelio Ranzato's slides

Convolutional layer



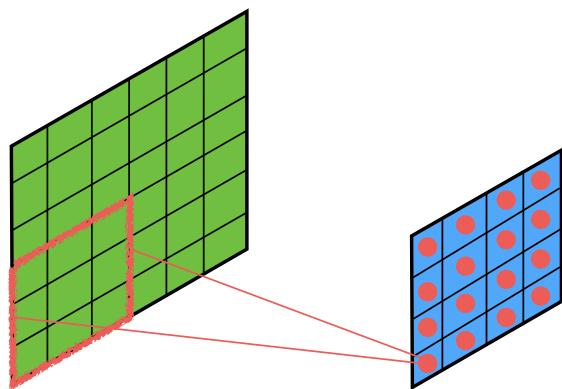
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Inspired by Marc'Aurelio Ranzato's slides

Convolutional layer



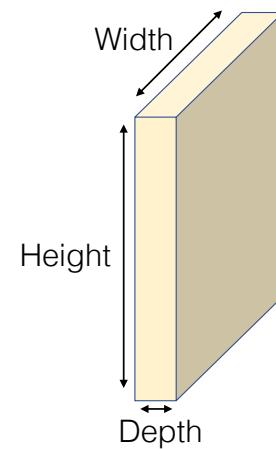
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Inspired by Marc'Aurelio Ranzato's slides

Terminologies

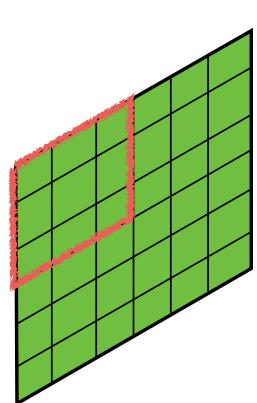


Joseph J. Lim

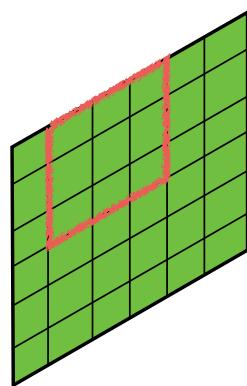
CSCI 599 @ USC

Lecture 5

Terminologies



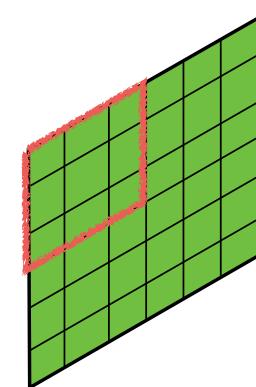
Step 1



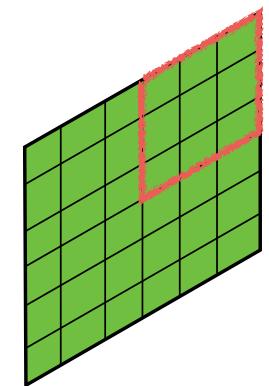
Stride 1

Step 2

Terminologies



Step 1



Stride 3

Step 2

Terminologies

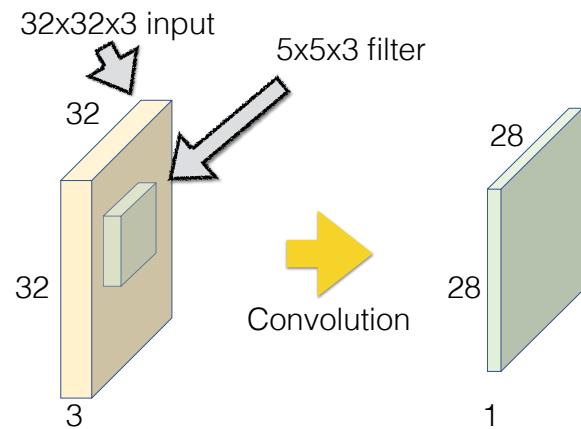
a	b	c	d	e
f	g	h	i	j
k	l	m	n	o
p	q	r	s	t
u	v	w	x	y

Terminologies

0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0
0	0	a	b	c	d	e	0	0
0	0	f	g	h	i	j	0	0
0	0	k	l	m	n	o	0	0
0	0	p	q	r	s	t	0	0
0	0	u	v	w	x	y	0	0
0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0

Zero-padding: 2

Convolutional layer

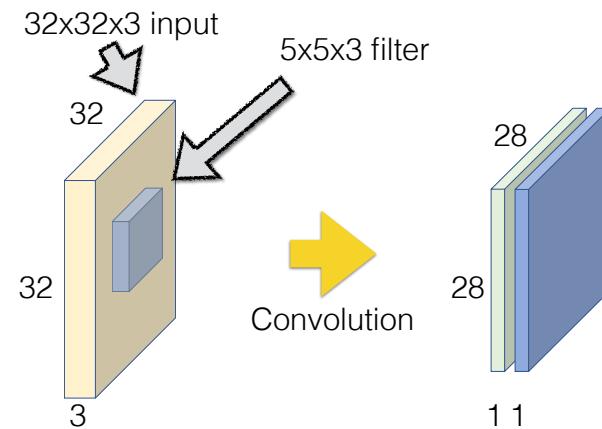


Joseph J. Lim

CSCI 599 @ USC

Slide credit: CS 231N
Lecture 5

Convolutional layer



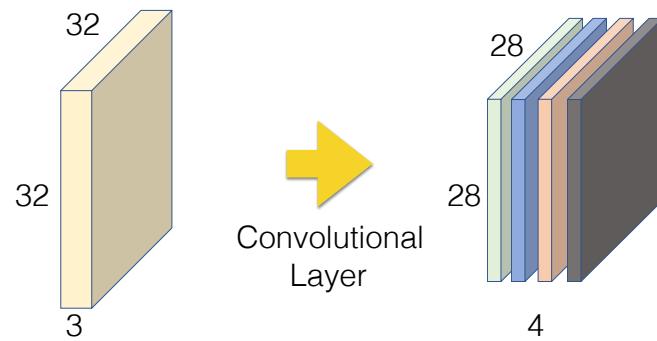
Joseph J. Lim

CSCI 599 @ USC

Slide credit: CS 231N
Lecture 5

Convolutional layer

With 4 5x5x3 filters, we will get 4 activation outputs.
Hence, the final output is 28x28x4.



Slide credit: CS 231N
Lecture 5

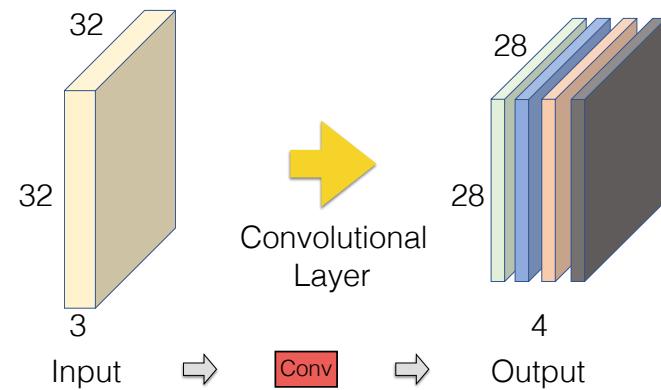
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Convolutional layer

With 4 5x5x3 filters, we will get 4 activation outputs.
Hence, the final output is 28x28x4.



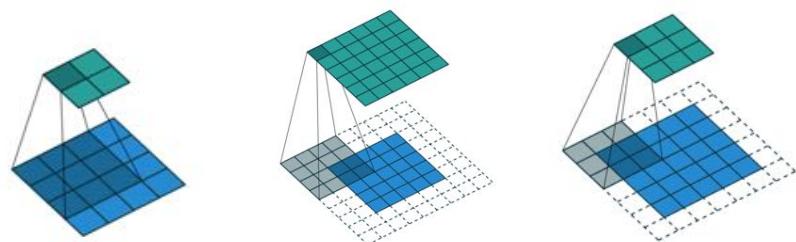
Slide credit: CS 231N
Lecture 5

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Convolutional layer



padding 0, stride 1

padding 2, stride 1

padding 1, stride 2

Note: it can have **depth** (3rd dimension).

Image credit: vdumoulin (github): https://github.com/vdumoulin/conv_arithmetic

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Deconvolutional layer

Deconvolutional layer

Warning: this is not about **de**-convolution!

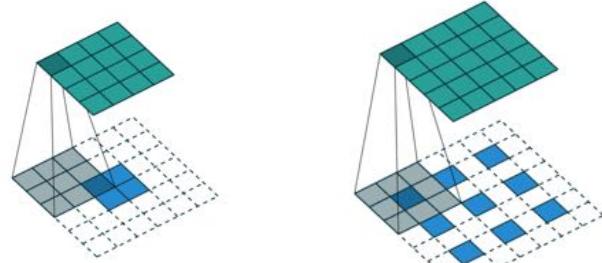
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Deconvolutional layer

Warning: this is not about **de**-convolution!



padding 0, stride 1

padding 1, stride 1

Why do we need this?

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Image credit: vdumoulin (github): https://github.com/vdumoulin/conv_arithmetic

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Fully connected vs Convolutional

- Fully connected layer
- Locally connected layer
- Convolutional layer
- Which layer is the best? Why?

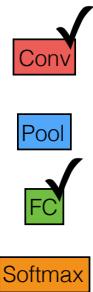
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

What are these layers?

- Convolutional layer
- Pool layer
- Fully connected layer
- Softmax layer

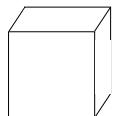


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Max pooling layer



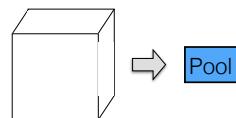
$M \times N \times P$

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Max pooling layer



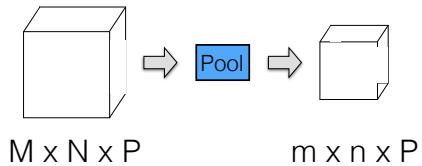
$M \times N \times P$

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Max pooling layer

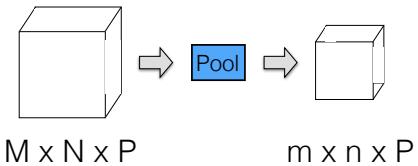


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Max pooling layer



- Shrinks the feature smaller (compact + manageable)

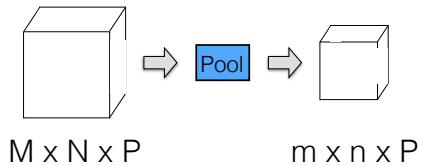
Slide credit: CS 231N

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Max pooling layer



- Shrinks the feature smaller (compact + manageable)
- Each activation map (3rd dim) is handled separately

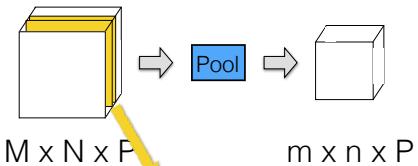
Slide credit: CS 231N

Joseph J. Lim

CSCI 599 @ USC

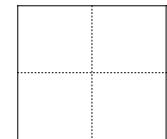
Lecture 5

Max pooling layer



1	3	2	7	2
5	1	3	0	5
4	8	5	0	6
3	2	6	0	8
0	7	4	6	9

max pool with
3x3 filters & stride 2

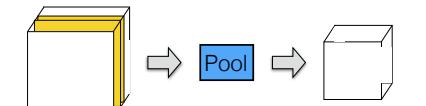


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Max pooling layer



$M \times N \times P$ $m \times n \times P$

1	3	2	7	2
5	1	3	0	5
4	8	5	0	6
3	2	6	0	8
0	7	4	6	9

max pool with
3x3 filters & stride 2

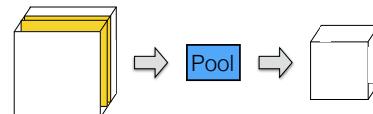
8	

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Max pooling layer



$M \times N \times P$ $m \times n \times P$

1	3	2	7	2
5	1	3	0	5
4	8	5	0	6
3	2	6	0	8
0	7	4	6	9

max pool with
3x3 filters & stride 2

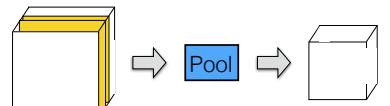
8	
	7

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Max pooling layer



$M \times N \times P$ $m \times n \times P$

1	3	2	7	2
5	1	3	0	5
4	8	5	0	6
3	2	6	0	8
0	7	4	6	9

max pool with
3x3 filters & stride 2

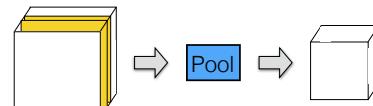
8	7
	8

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Max pooling layer



$M \times N \times P$ $m \times n \times P$

1	3	2	7	2
5	1	3	0	5
4	8	5	0	6
3	2	6	0	8
0	7	4	6	9

max pool with
3x3 filters & stride 2

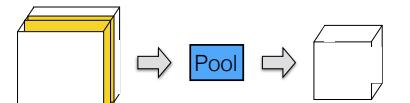
8	7
8	9

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Max pooling layer



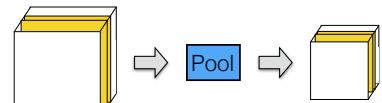
$M \times N \times P$ $m \times n \times P$

1	3	2	7	2
5	1	3	0	5
4	8	5	0	6
3	2	6	0	8
0	7	4	6	9

max pool with
3x3 filters & stride 2

8		7
8		9

Max pooling layer



$M \times N \times P$ $m \times n \times P$

1	3	2	7	2
5	1	3	0	5
4	8	5	0	6
3	2	6	0	8
0	7	4	6	9

max pool with
3x3 filters & stride 2

8		7
8		9

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Pooling layers

- **Max pooling** (most common)
- Average pooling
- L2-norm pooling
- ...

Why pooling layers?

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Why pooling layers?

- Reduces # of parameters in the following layers

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Why pooling layers?

- Reduces # of parameters in the following layers
 - Hence, less overfitting!

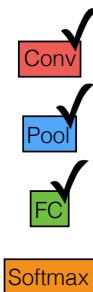
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

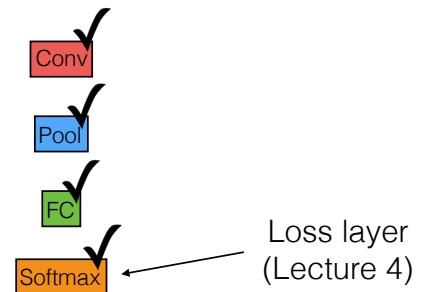
What are these layers?

- Convolutional layer
- Pool layer
- Fully connected layer
- Softmax layer



What are these layers?

- Convolutional layer
- Pool layer
- Fully connected layer
- Softmax layer



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

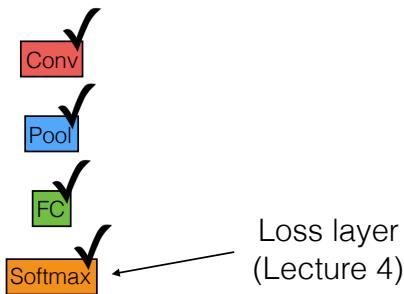
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

All these are linear!!!

- Convolutional layer
- Pool layer
- Fully connected layer
- Softmax layer

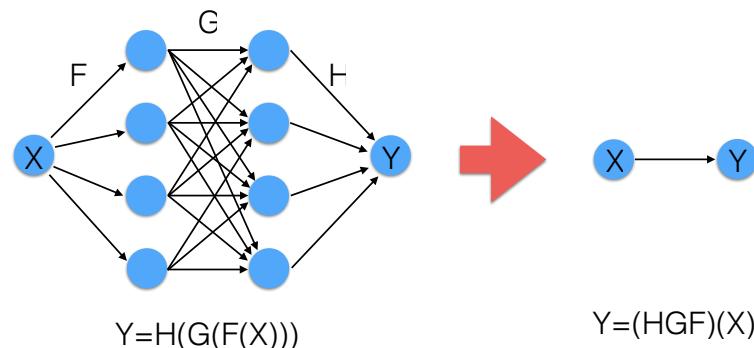


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Activation functions



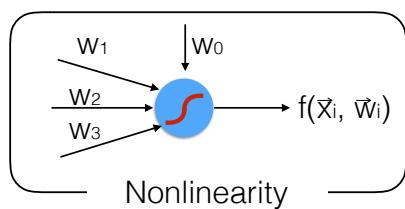
Without nonlinear activation,
multiple layers are reduced to one matrix multiplication.

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Activation functions



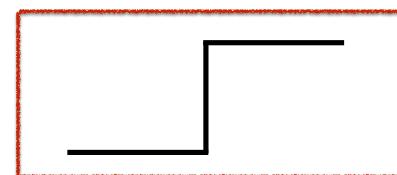
Nonlinearity makes a neural network deeper

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Step function



- Firing after a threshold
- Is not differentiable
- Cannot apply back-propagation

Anything else we can use?

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Activation functions

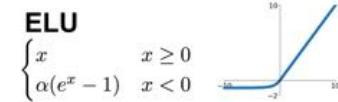
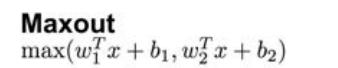
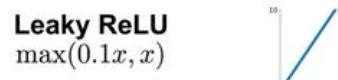
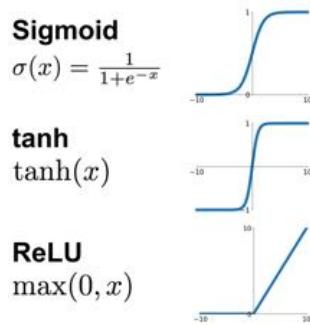


Figure from Stanford cs231n lecture slides

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Activation functions

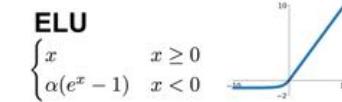
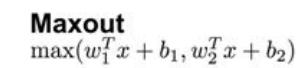
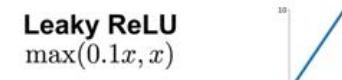
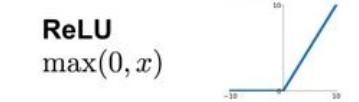
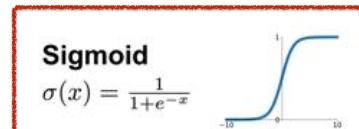


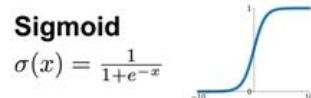
Figure from Stanford cs231n lecture slides

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Sigmoid



- Similar to step functions, but differentiable
- (-) easily saturated (no gradients)

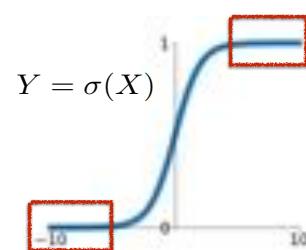
Figure from Stanford cs231n lecture slides

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Gradient vanishing problem



Chain rule
$$\frac{\partial Y}{\partial X} = \frac{\partial Y}{\partial \sigma} \frac{\partial \sigma}{\partial X}$$

is almost 0 for most of X

Thus, $\frac{\partial Y}{\partial X}$ also goes to 0.

Figure from Stanford cs231n lecture slides

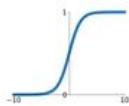
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Activation functions

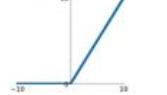
Sigmoid
 $\sigma(x) = \frac{1}{1+e^{-x}}$



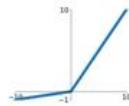
tanh
 $\tanh(x)$



ReLU
 $\max(0, x)$



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



Maxout

$\max(w_1^T x + b_1, w_2^T x + b_2)$

ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$

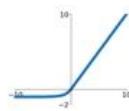


Figure from Stanford cs231n lecture slides

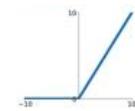
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

ReLU

ReLU
 $\max(0, x)$



- Prevent gradient vanishing when $x > 0$
- Computationally efficient
- Biologically plausible
- (-) Still loose gradient when $x < 0$

Figure from Stanford cs231n lecture slides

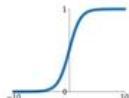
Joseph J. Lim

CSCI 599 @ USC

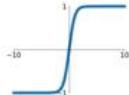
Lecture 5

Activation functions

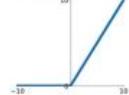
Sigmoid
 $\sigma(x) = \frac{1}{1+e^{-x}}$



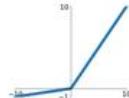
tanh
 $\tanh(x)$



ReLU
 $\max(0, x)$



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



Maxout

$\max(w_1^T x + b_1, w_2^T x + b_2)$

ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$

Figure from Stanford cs231n lecture slides

Joseph J. Lim

CSCI 599 @ USC

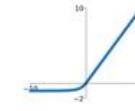
Lecture 5

Leaky ReLU/ELU

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



ELU
 $\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$



- No gradient vanishing problem
- Can be used instead of ReLU (e.g. Leaky ReLU for deconvolution)

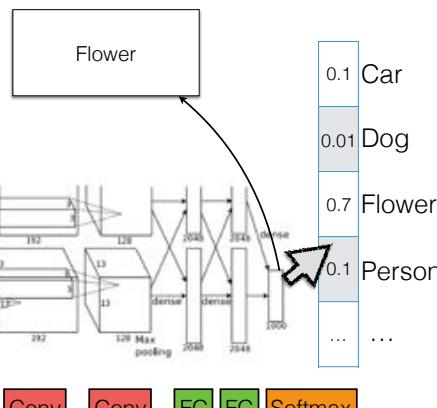
Figure from Stanford cs231n lecture slides

Joseph J. Lim

CSCI 599 @ USC

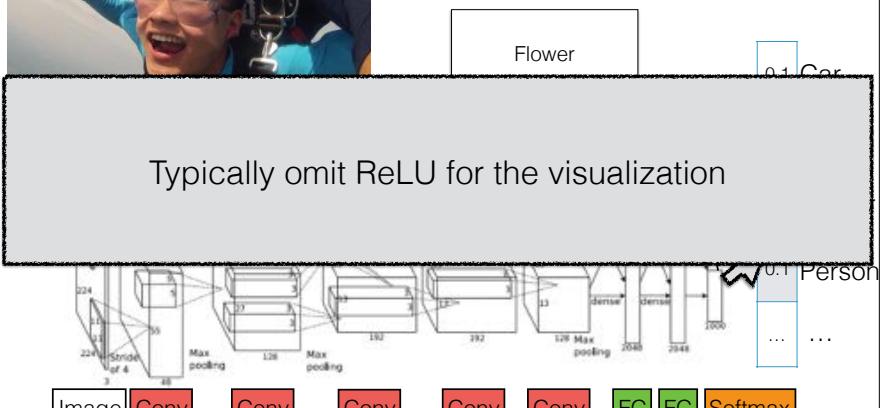
Lecture 5

AlexNet in details



Joseph J. Lin, ReLU ReLU SCI 599 @ USC ReLU Lecture 5

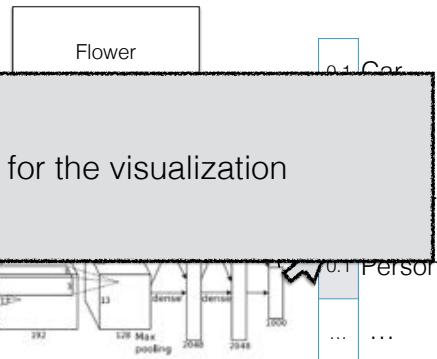
AlexNet in details



Typically omit ReLU for the visualization

Joseph J. Lin, ReLU ReLU SCI 599 @ USC ReLU Lecture 5

AlexNet in details



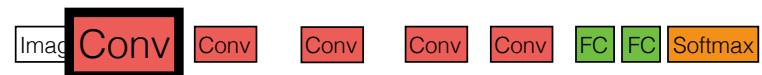
Typically omit ReLU for the visualization

Joseph J. Lim

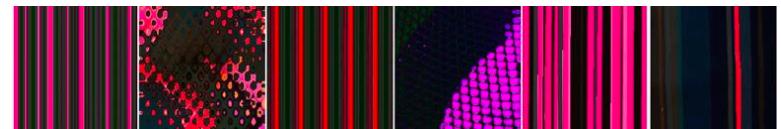
CSCI 599 @ USC

Lecture 5

What did AlexNet learn?



red

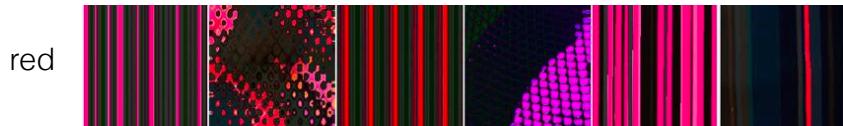
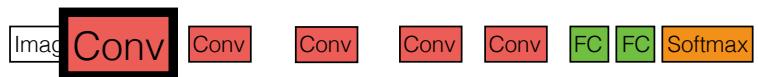


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

What did AlexNet learn?



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Bau, Zhou, and et. al. Network Dissection: Quantifying Interpretability of Deep Visual Representations. CVPR 2017.

What did AlexNet learn?



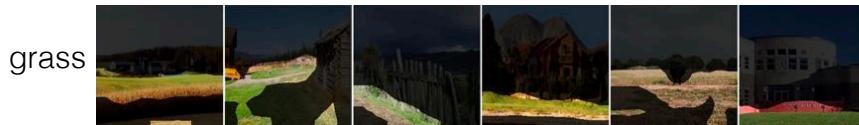
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Bau, Zhou, and et. al. Network Dissection: Quantifying Interpretability of Deep Visual Representations. CVPR 2017.

What did AlexNet learn?



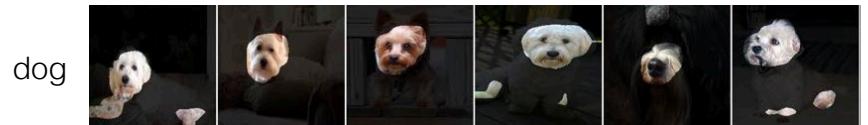
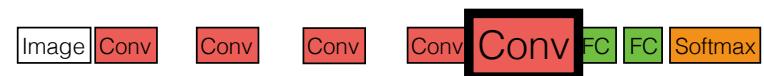
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Bau, Zhou, and et. al. Network Dissection: Quantifying Interpretability of Deep Visual Representations. CVPR 2017.

What did AlexNet learn?



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Bau, Zhou, and et. al. Network Dissection: Quantifying Interpretability of Deep Visual Representations. CVPR 2017.

What did AlexNet learn?



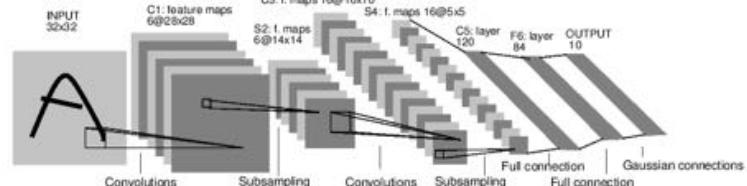
Bau, Zhou, and et. al. Network Dissection: Quantifying Interpretability of Deep Visual Representations. CVPR 2017.

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

CNNs: LeNet



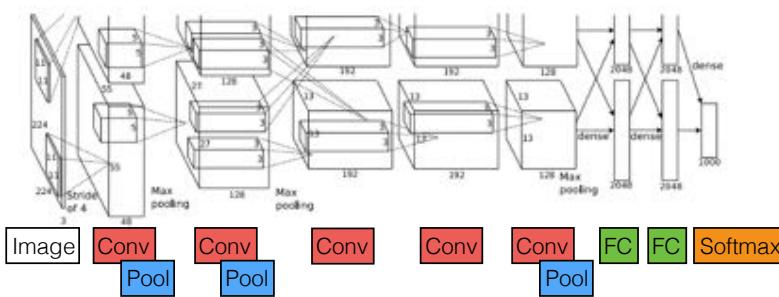
LeCun et. al. Handwritten digit recognition: Applications of neural net chips and automatic learning. IEEE Communication 1989.

Joseph J. Lim

CSCI 599 @ USC

Lecture 1

CNNs: AlexNet



A Krizhevsky, et. al. ImageNet Classification with Deep Convolutional Neural Networks. NIPS 2012.

Joseph J. Lim

CSCI 599 @ USC

Lecture 1

What made AlexNet AlexNet?

- GPU?
- Data?
- ReLU?
- Convolutional?
- Hyperparameter?

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

CNNs: VGGNet



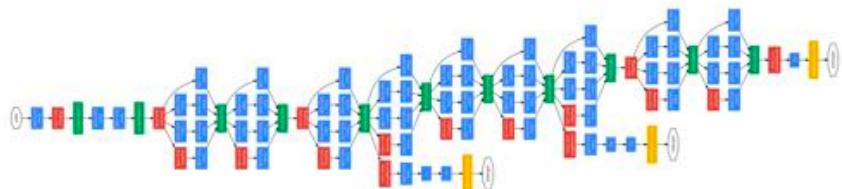
Simonyan and Zisserman. Very Deep Convolutional Networks for Large-scale Image Recognition. ICLR 2015.

Joseph J. Lim

CSCI 599 @ USC

Lecture 1

CNNs: GoogLeNet



Szegedy, et. al. Going Deeper with Convolution. CVPR 2015.

Joseph J. Lim

CSCI 599 @ USC

Lecture 1

CNNs: ResNet



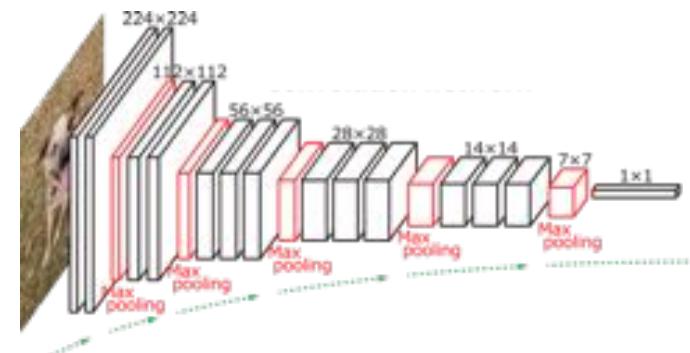
He, et. al. Deep **Residual** Learning for Image Recognition. CVPR 2016.

Joseph J. Lim

CSCI 599 @ USC

Lecture 1

Encoder



Encoding (into a feature / representation)

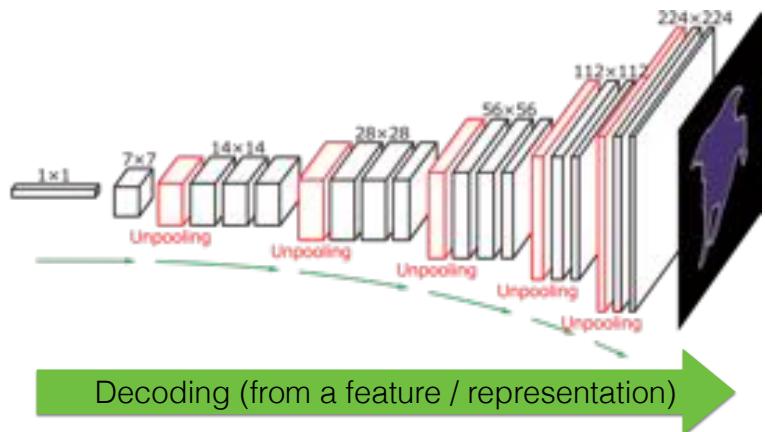
Image credit: Noh, et. al. Learning Deconvolution Network for Semantic Segmentation. ICCV 2015.

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Decoder



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Congratulations!

All basic vocabbs are covered.

Let's revisit CV tasks!

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

CNN applications

- Image classification
- Object detection
- Object segmentation
- Image resolution
- Human pose detection
- Action classification
- Image captioning

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Image Classification



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Image Classification



Car



Image Conv Conv Pool Conv Pool FC FC Softmax

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Object Detection



Object Detection



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Object Detection



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Object Detection



1. Object proposal

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Object Detection



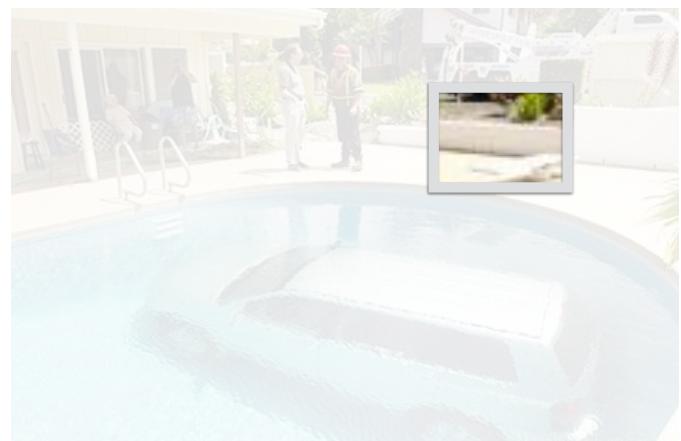
2. Evaluate each proposal

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Object Detection



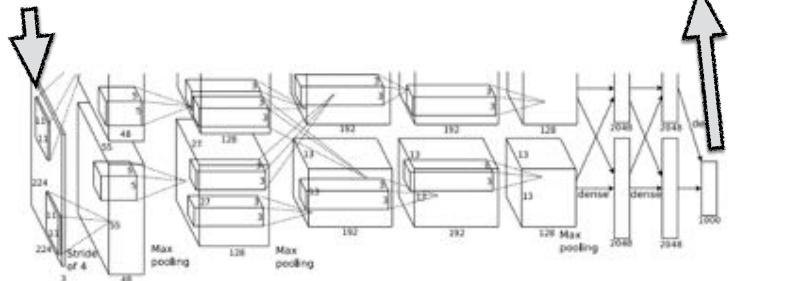
2. Evaluate each proposal

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Object Detection



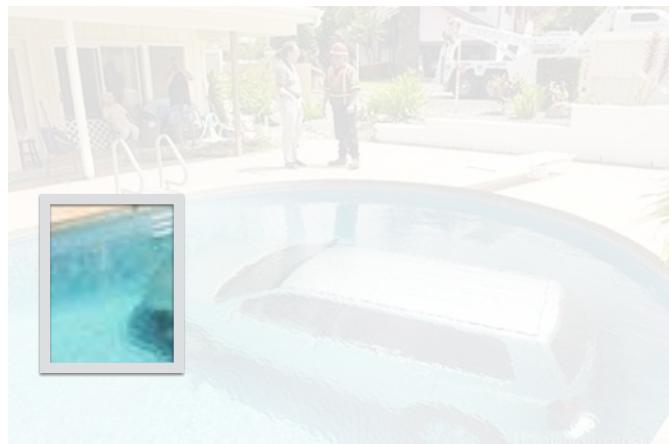
2. Evaluate each proposal

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Object Detection



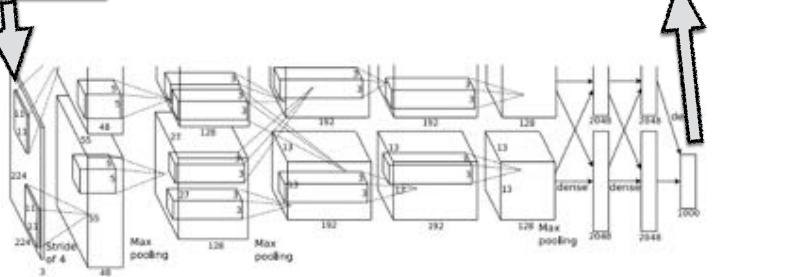
2. Evaluate each proposal

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Object Detection



2. Evaluate each proposal

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Object Detection



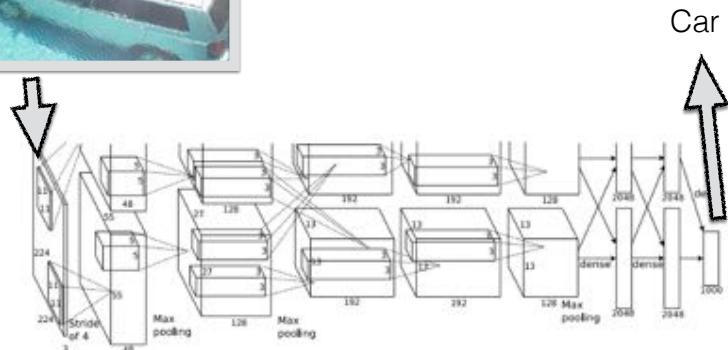
2. Evaluate each proposal

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Object Detection



2. Evaluate each proposal

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Object Detection



For more:

R-CNN (<https://github.com/rbgirshick/rcnn>) and Fast R-CNN (<https://github.com/rbgirshick/fast-rcnn>)

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Semantic Segmentation

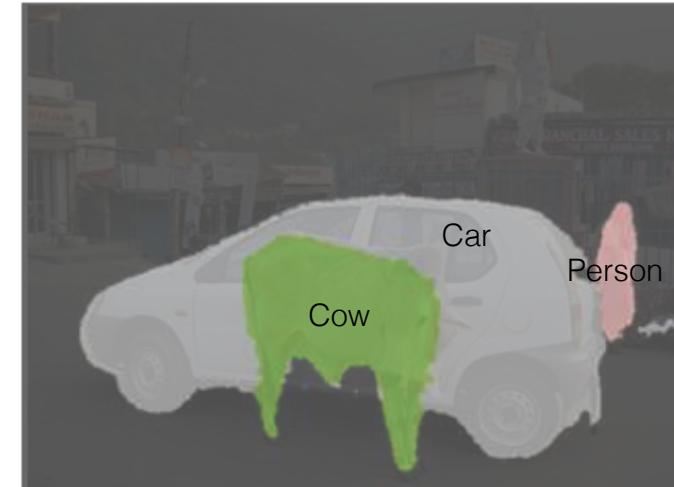


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Semantic Segmentation



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Semantic Segmentation

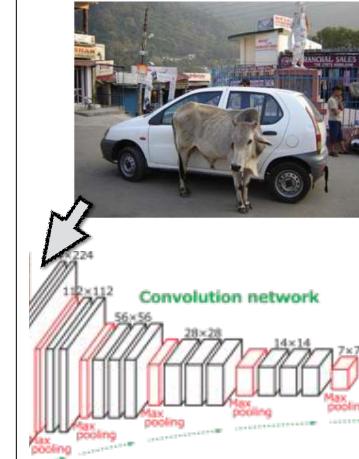


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Semantic Segmentation



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Semantic Segmentation

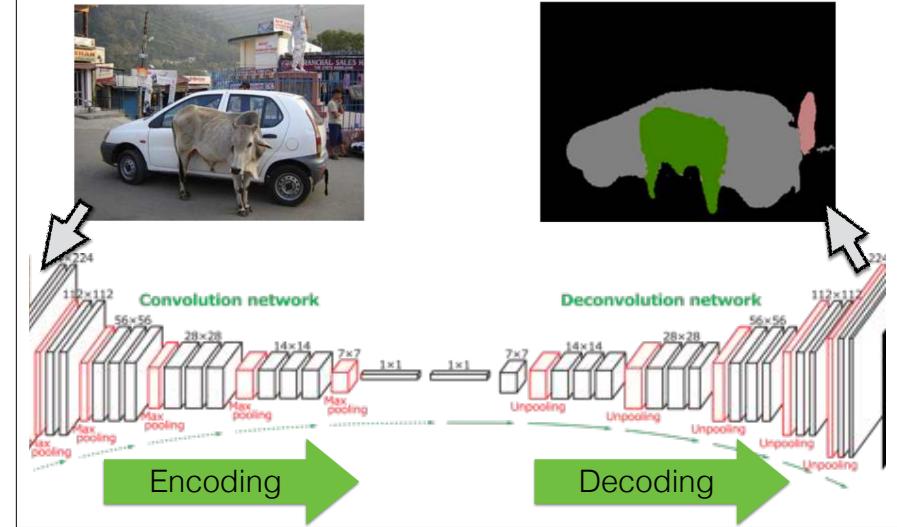


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Semantic Segmentation

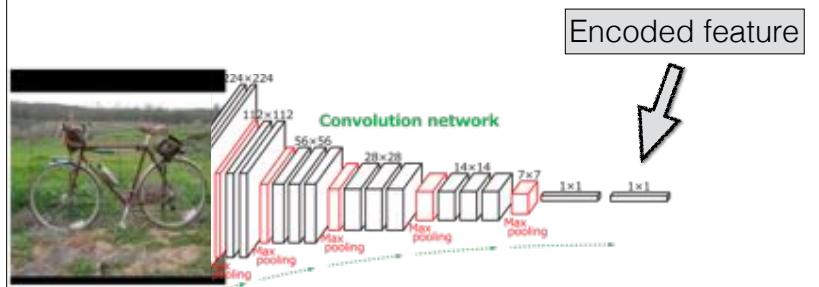


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Semantic Segmentation

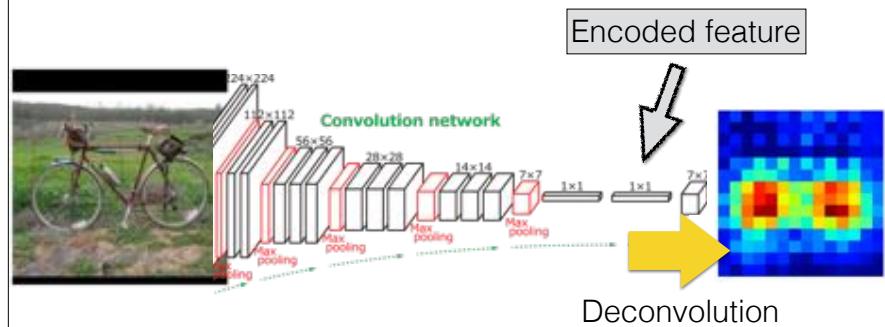


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Semantic Segmentation

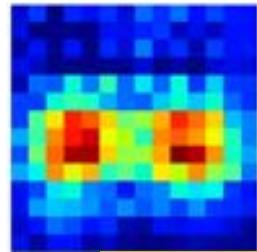


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Semantic Segmentation



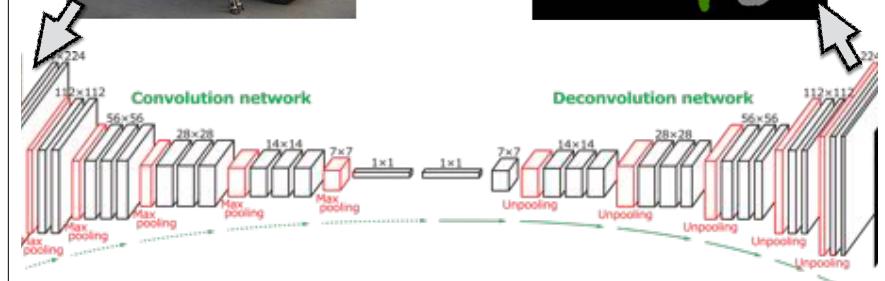
Deconvolution

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Semantic Segmentation



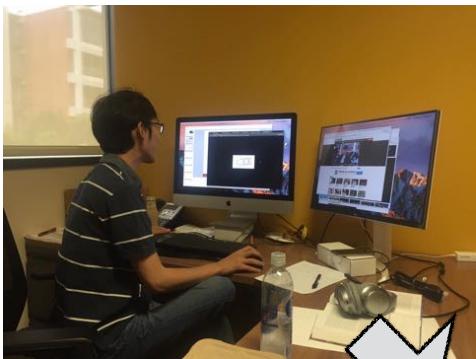
For more: Noh, et.al. Learning Deconvolution Network for Semantic Segmentation. ICCV 2015.

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Image captioning



Hao Li (USC prof.) said
Josepo ah! Let me tell you!
Hard working Deep Nearning
Propressor at USC university!

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Image captioning



Encoding

<Image Feature>

Image captioning



Encoding

<Image Feature>

$f(x; W)$

Decoding

Josepo ah! Let me tell you!
Hard working Deep Nearning
Propressor at USC university!

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Demo

- 3D Face Reconstruction from a Single Image

<http://www.cs.nott.ac.uk/~psxasj/3dme/>

Joseph J. Lim

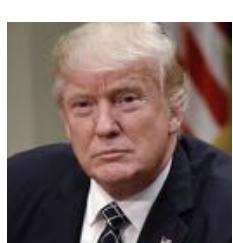
CSCI 599 @ USC

Lecture 5

Demo

- 3D Face Reconstruction from a Single Image

<http://www.cs.nott.ac.uk/~psxasj/3dme/>



Joseph J. Lim

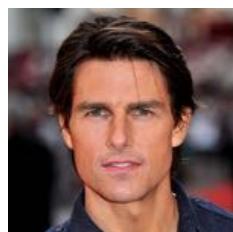
CSCI 599 @ USC

Lecture 5

Demo

- 3D Face Reconstruction from a Single Image

<http://www.cs.nott.ac.uk/~psxasj/3dme/>



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Demo

- 3D Face Reconstruction from a Single Image

<http://www.cs.nott.ac.uk/~psxasj/3dme/>



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Demo

- 3D Face Reconstruction from a Single Image

<http://www.cs.nott.ac.uk/~psxasj/3dme/>



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Summary

- Fully connected layer
- Convolutional layer
- Deconvolutional layer
- Pooling layers
- Activation layers (e.g. ReLU)

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Break Time

See you in 20 mins!

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Today's agenda

- Part 1
 - Convolutional Neural Networks
- Part 2
 - Training Neural Networks

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

All inputs are just numbers



Image



a	b	c	d	e
f	g	h	i	j
k	l	m	n	o
p	q	r	s	t
u	v	w	x	y



Sound



a	b	c	d	e
---	---	---	---	---

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Data preprocessing

- Zero-center data
- Normalize data
- PCA Whitening

Joseph J. Lim

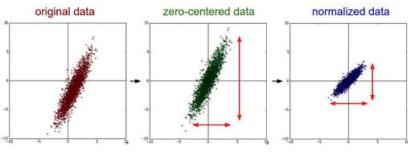
CSCI 599 @ USC

Lecture 5

Data preprocessing

Why do we need preprocessing?

- Zero-center data
- Normalize data
- PCA Whitening



Activation functions (ReLU) works around **zero**.

If input data is biased toward positive or negative, randomly initialized layers produce biased output.

Figure from Stanford cs231n lecture slides

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Data preprocessing

Why do we need preprocessing?

- Zero-center data
- Normalize data
- PCA Whitening

Joseph J. Lim

CSCI 599 @ USC

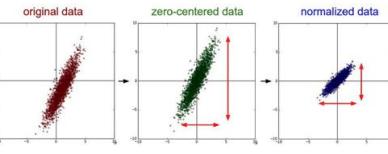
Lecture 5

Data preprocessing

Data preprocessing

Why do we need preprocessing?

- Zero-center data
- Normalize data
- PCA Whitening



A feature with small scale has a negligible effect on backprop.

Figure from Stanford cs231n lecture slides

Joseph J. Lim

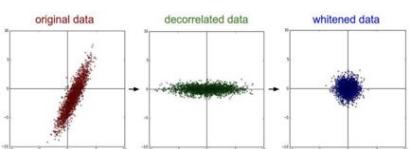
CSCI 599 @ USC

Lecture 5

Data preprocessing

Why do we need preprocessing?

- Zero-center data
- Normalize data
- PCA Whitening



Remove correlation between input features

Figure from Stanford cs231n lecture slides

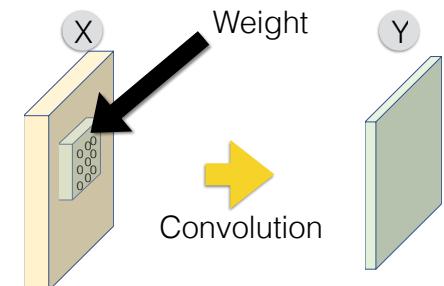
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Weight initialization

If all weights are 0?



Slide credit: CS 231N

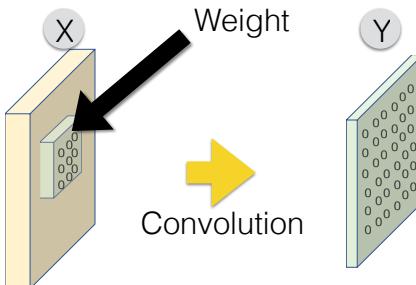
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Weight initialization

If all weights are 0?



Gradient vanishing problem

Output \mathbf{Y} will be always 0 and then gradients will go to 0.

- Small random numbers
- Xavier
- Xavier / 2 (He et. al.)

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Slide credit: CS 231N

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Small Random Numbers

- Sample from a gaussian distribution with 0 mean and given standard deviation (such as 0.01)

$$W \sim \mathcal{N}(\mu, \sigma^2)$$

- Doesn't work for deeper networks

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Small Random Numbers

If the weights are too small,

the activation signals shrink quickly

If the weights are too large,

the activation signals explode

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Xavier Initialization

- Set the variance according to # of input neurons

$$W \sim \mathcal{N}(\mu, \sigma^2) \quad \text{Var}(W) = \frac{1}{n_{in}}$$

- Consider # output neurons for back-prop

$$\text{Var}(W) = \frac{2}{n_{in} + n_{out}}$$

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

He et. al. Initialization

Random Numbers of dimension (D_{in}, D_{out})

$$\sqrt{D_{in}/2}$$

- Default choice of training a deep network

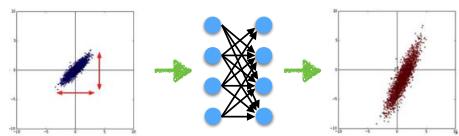
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Normalization

- After each layer, the distribution of activation signals changes (Internal Covariate Shift)



- As a network becomes deeper, distribution shifts more.
- Recall why we need data preprocessing

Figure from Stanford cs231n lecture slides

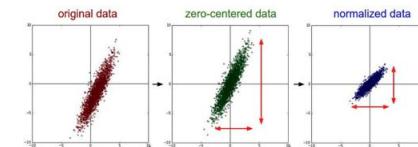
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Normalization

- Normalize the distribution of activations to have zero mean and unit variance



Same as data preprocessing

$$\hat{x}^{(k)} = \frac{x^{(k)} - \mathbb{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

Figure from Stanford cs231n lecture slides

Joseph J. Lim

CSCI 599 @ USC

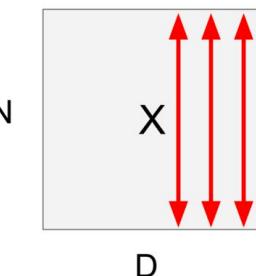
Lecture 5

Batch normalization

- Why do we need normalization?
 - Internal Covariate Shift
 - To train deeper networks on larger data
 - Make each dimension Unit Gaussian

Batch normalization

- Compute the empirical mean and variance independently for each dimension and normalize.



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Batch normalization

- Vanilla Normalization:

$$\hat{x}^{(k)} = \frac{x^{(k)} - \text{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Batch normalization

- Squashing the range:

$$y^{(k)} = \gamma^{(k)} \hat{x}^{(k)} + \beta^{(k)}$$

- Network can learn:

$$\gamma^{(k)} = \sqrt{\text{Var}[x^{(k)}]}$$

$$\beta^{(k)} = \text{E}[x^{(k)}]$$

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Batch normalization

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_{1\dots m}\}$;
Parameters to be learned: γ, β
Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\begin{aligned}\mu_{\mathcal{B}} &\leftarrow \frac{1}{m} \sum_{i=1}^m x_i && // \text{mini-batch mean} \\ \sigma_{\mathcal{B}}^2 &\leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 && // \text{mini-batch variance} \\ \hat{x}_i &\leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} && // \text{normalize} \\ y_i &\leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) && // \text{scale and shift}\end{aligned}$$

[Ioffe and Szegedy, 2015]

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Batch normalization

- Why is it good?
 - Improves gradient flow
 - Allows higher learning rates
 - Reduces strong dependence on initialization
 - Regularization

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Batch normalization

- At test time
 - The mean and variance are not computed based on the batch.
 - Empirical mean of activations during training is used.

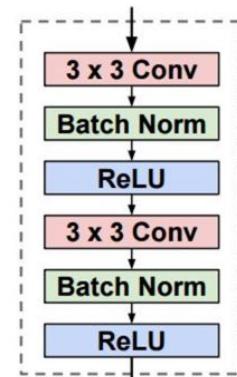
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Batch normalization

- Usually added after convolutional layers or fully connected layers
- before non-linearities
- Question: Which is better?
 - Before or after non-linearities



Johnson, Justin, Alexandre Alahi, and Li Fei-Fei, "Perceptual losses for real-time style transfer and super-resolution," *European Conference on Computer Vision*. Springer International Publishing, 2016.

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Training details

- Learning rate
- Batch size
- Hyperparameter
- Loss curve

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Learning Rate

How **fast** or **slow** we update the model parameters

$$\theta_{t+1} = \theta_t - \eta_t \nabla f(\theta_t)$$

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Learning Rate

How **fast** or **slow** we update the model parameters

$$\theta_{t+1} = \theta_t - \eta_t \boxed{\nabla f(\theta_t)}$$

Gradients

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Learning Rate

Too large

- Explode

Too small

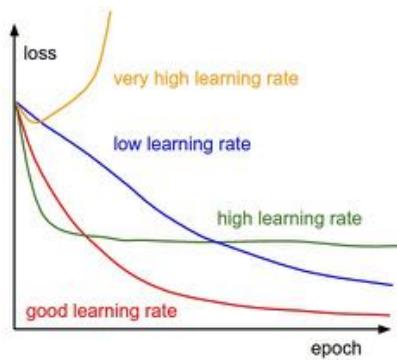
- Get stuck in local minimum

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Learning Rate



ref: <http://cs231n.github.io/neural-networks-3/#anneal>

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Learning Rate

Too large

- Explode

Too small

- Get stuck in local minimum

Anneal the learning rate

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Learning Rate

- Constant
- Step decay
- Exponential decay
- Time-based decay
- Etc

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Learning Rate

Constant

$$lr = lr_0$$



Dataset: CIFAR-10

ref: <https://goo.gl/xnrb3N>

Joseph J. Lim

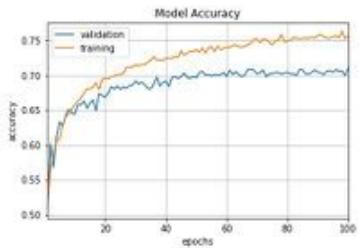
CSCI 599 @ USC

Lecture 5

Learning Rate

Time-based decay

$$lr = \frac{lr_0}{1 + k * iter}$$



Dataset: CIFAR-10

ref: <https://goo.gl/xnrb3N>

Joseph J. Lim

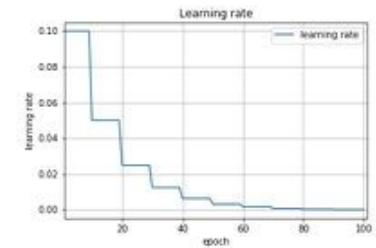
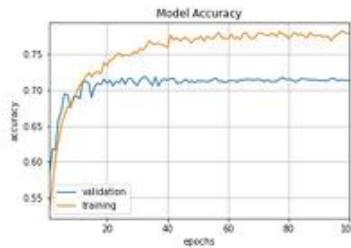
CSCI 599 @ USC

Lecture 5

Learning Rate

Step decay

$$lr = lr_0 * \text{floor}\left(\frac{\text{epoch}}{\text{epoch_drops}}\right)$$



Dataset: CIFAR-10

ref: <https://goo.gl/xnrb3N>

Joseph J. Lim

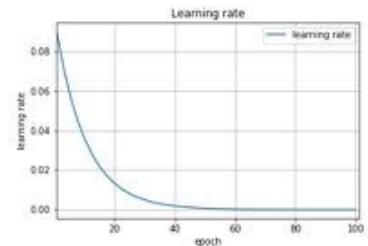
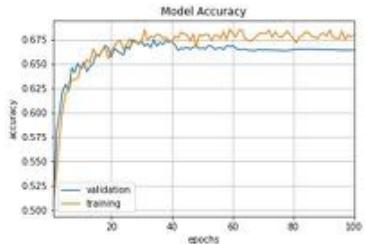
CSCI 599 @ USC

Lecture 5

Learning Rate

Exponential decay

$$lr = \frac{lr_0}{\exp^{k*iter}}$$



Dataset: CIFAR-10

Joseph J. Lim

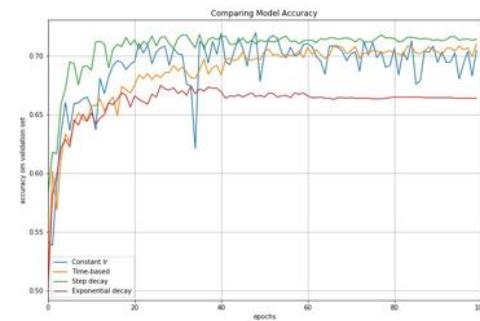
ref: <https://goo.gl/xnrb3N>

CSCI 599 @ USC

Lecture 5

Learning Rate

Comparisons



Dataset: CIFAR-10

Joseph J. Lim

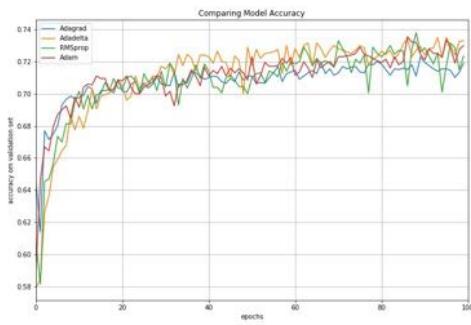
ref: <https://goo.gl/xnrb3N>

CSCI 599 @ USC

Lecture 5

Learning Rate

Adaptive learning rate methods



Dataset: CIFAR-10

Joseph J. Lim

ref: <https://goo.gl/xnrb3N>

CSCI 599 @ USC

Lecture 5

Batch Size

Number of samples in one forward/backward pass

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Batch Size

Too large

- Memory inefficient
- Computation inefficient
- Converge to sharp minimizers of the training function

Too small

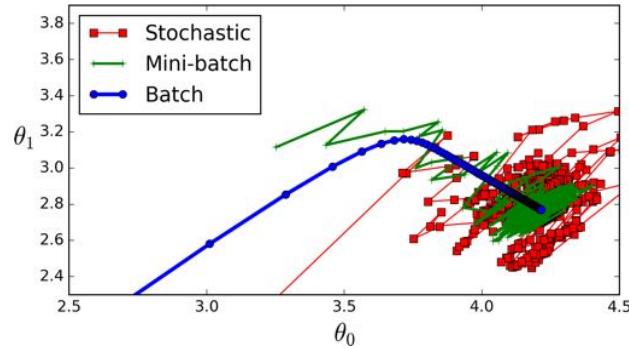
- Less accurate gradients at each iteration
- Converge to flat minimizers of the training function

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Batch Size



Stochastic: batch size is 1

ref: <https://stats.stackexchange.com/questions/153531/what-is-batch-size-in-neural-network>

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Optimizers and Algorithms

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Gradient Descent

- Gradient descent
- Stochastic gradient descent
- Mini-batch gradient descent
- Momentum
- Nesterov accelerated gradient
- Etc

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Gradient Descent

Computes the gradient of the cost function w.r.t. to the parameters θ for **the entire training dataset at each iteration**

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} f(\theta_t)$$

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Stochastic Gradient Descent

Computes the gradient of the cost function w.r.t. to the parameters θ for **each data point at each iteration**

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} f(\theta_t; x^{(i)})$$

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Mini-batch Gradient Descent

Computes the gradient of the cost function w.r.t. to the parameters θ for performs an update for **every mini-batch of n training examples at each iteration**

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} f(\theta_t; x^{(i:i+n)})$$

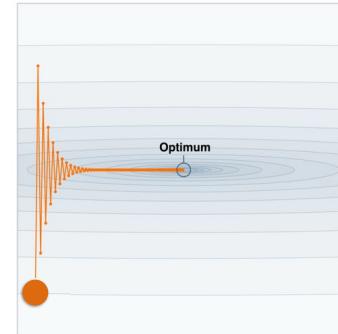
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Momentum

GD family updates parameters inefficiency



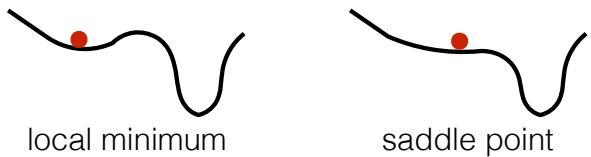
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Momentum

GD family gets stuck at local minimums and saddle points easily



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Momentum

Consider the history:
Adding a fraction γ of **the update vector of the past time step** to the current update vector

$$\nu_t = \gamma \nu_{t-1} + \eta \nabla_{\theta_t} f(\theta_t)$$

$$\theta_{t+1} = \theta_t - \nu_t$$

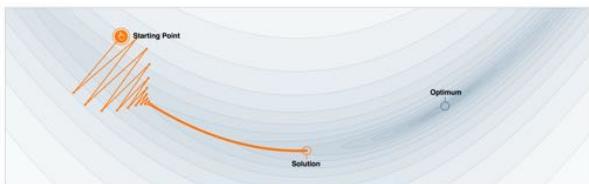
Joseph J. Lim

CSCI 599 @ USC

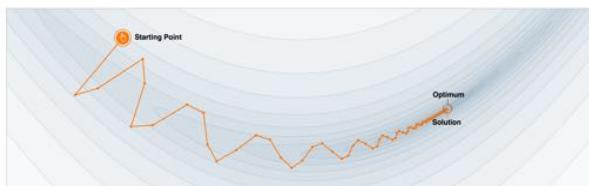
Lecture 5

Momentum

Without momentum



With momentum



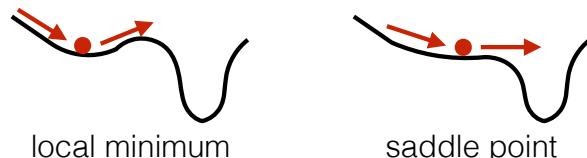
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Momentum

GD family gets stuck at local minimums and saddle points easily



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Nesterov accelerated gradient

Calculate the gradient not w.r.t. to our current parameters but w.r.t. **the approximate future position of our parameters**:

$$\nu_t = \gamma\nu_{t-1} + \eta\nabla_{\theta_t} f(\theta_t - \gamma\nu_{t-1})$$

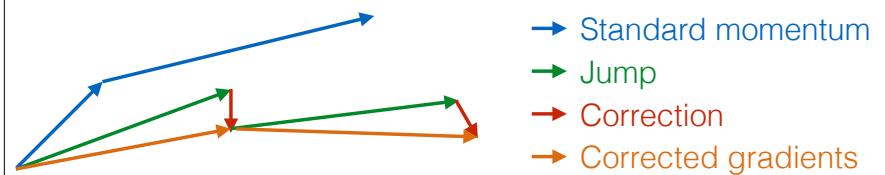
$$\theta_{t+1} = \theta_t - \nu_t$$

Nesterov accelerated gradient

Calculate the gradient not w.r.t. to our current parameters but w.r.t. **the approximate future position of our parameters**:

$$\nu_t = \gamma\nu_{t-1} + \eta\nabla_{\theta_t} f(\theta_t - \gamma\nu_{t-1})$$

$$\theta_{t+1} = \theta_t - \nu_t$$



Regularization

- Extra terms
- Dropout
- Data augmentation

Regularization

$$L = L_{original} + \lambda R(\theta)$$

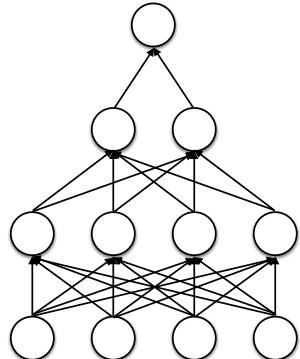
L2 norm regularizer $R(\theta) = \sum_i \sum_j \theta_{i,j}^2$

L1 norm regularizer $R(\theta) = \sum_i \sum_j |\theta_{i,j}|$

L1+L2 norm regularizer $R(\theta) = \sum_i \sum_j |\theta_{i,j}| + \alpha \theta_{i,j}^2$

Dropout

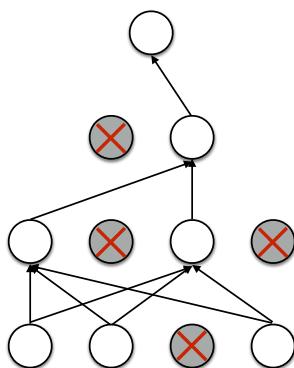
At each iteration, randomly shutdown some neurons with a certain probability (usually 0.5)



Joseph J. Lim

CSCI 599 @ USC

Lecture 5



Dropout

Intuition

- Forces the network to have a redundant representation
- Training a bunches of small models sharing parameters
 - Each binary mask is one model

Joseph J. Lim

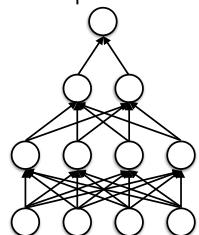
CSCI 599 @ USC

Lecture 5

Dropout

Testing phase

- All neurons are active always
- Scale the activations so that for each neuron
 - Testing output = expected output at training time



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Data Augmentation

Deep networks need to be trained on a huge number of training images



What if... we don't have enough amount of data?

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Data Augmentation

Some tricks for image data augmentation

- Flip
- Scale
- Rotate
- Crop
- Color jitter
- etc

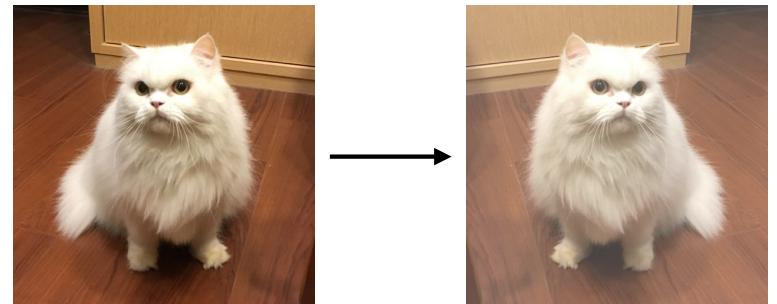
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Data Augmentation

Flip



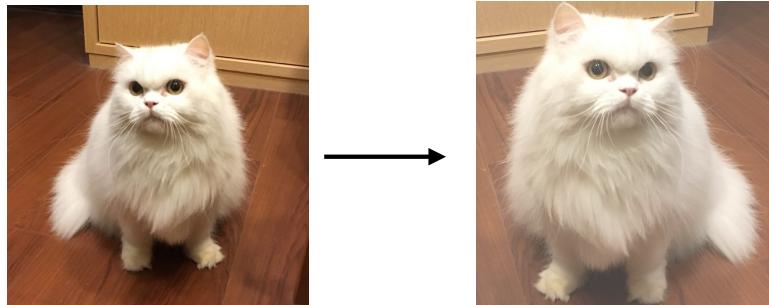
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Data Augmentation

Scale



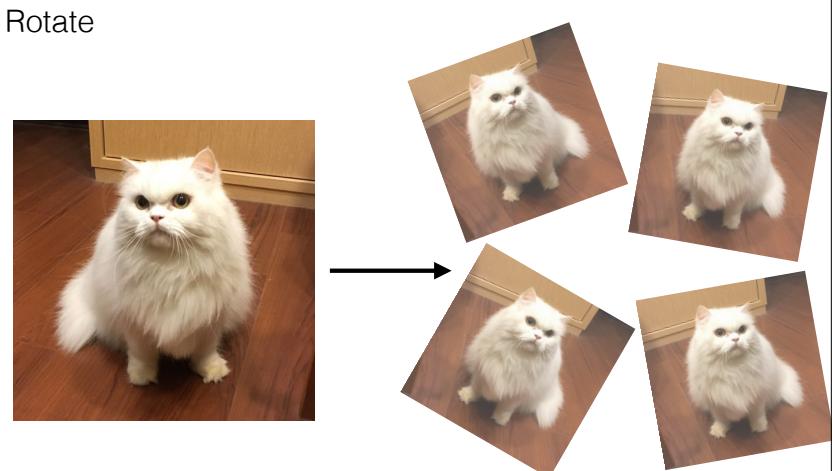
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Data Augmentation

Rotate



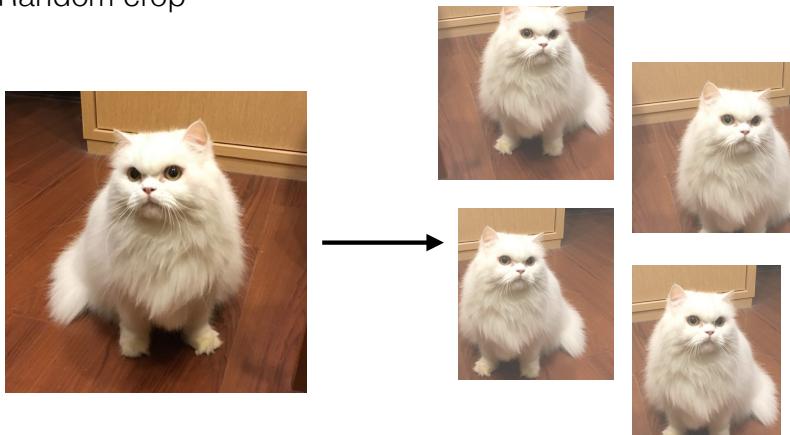
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Data Augmentation

Random crop

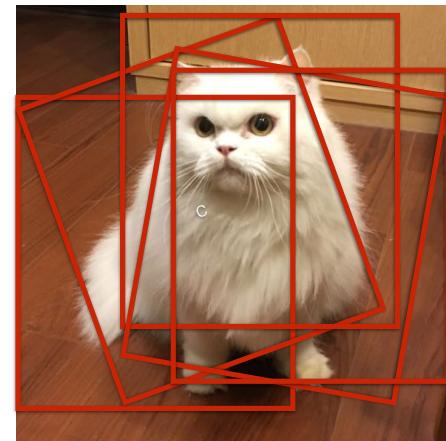


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Data Augmentation



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Summary

Deep learning is all about getting everything right...

- Data
- Model
 - Architecture
- Optimizer
- Hyperparameters
- Training tricks
 - Initialization, normalization, regularization, etc

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Next week

- Deep learning software
- Cloud services (Google cloud, Amazon AWS)

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Todo

- Midterm on week 7
- Assignment #1 & Project proposal due week 8
- Project Meeting with TA #1 by week 8

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Questions?

Joseph J. Lim

CSCI 599 @ USC

Lecture 5