Network Architectures

Computer Vision (UW EE/CSE 576)



Richard Szeliski Facebook & UW Lecture 9 – Apr 28, 2020



Class calendar

Date	Торіс	Slides	Reading	Homework
April 9	Filters and convolutions	Google Slides	<u>Szeliski</u> , Chapter 3	HW1 due, <u>HW2</u> assigned
April 14	Interpolation and Optimization	pdf, pptx	<u>Szeliski</u> , Chapter 4	
April 16	Machine Learning	pdf, pptx	Szeliski, Chapter 5.1-5.2	
April 21	Deep Neural Networks	pdf, pptx	<u>Szeliski</u> , Chapter 5.3	
April 23	Convolutional Neural Networks		<u>Szeliski</u> , Chapter 5.4	HW2 due, HW3 assigned
April 28	Network Architectures		<u>Szeliski</u> , Chapter 5.4	
April 30	Object Detection		<u>Szeliski</u> , Chapter 6.3	
May 5	Detection and Instance Segmentation		<u>Szeliski</u> , Chapter 6.4	
May 7	Edges, features, matching, RANSAC		<u>Szeliski</u> , Chapter 7.1-7.2, 8.1-8.2	HW3 due, HW4 assigned

References







https://d2l.ai/

Chapter 5 **Deep Learning**

5.3	Deep r	neural networks
	5.3.1	Weights and layers
	5.3.2	Activation functions
	5.3.3	Regularization and normalization
	5.3.4	Loss functions
	5.3.5	Backpropagation
	5.3.6	Training and optimization
5.4	Convo	lutional neural networks
	5.4.1	Pooling and unpooling 296
	5.4.2	Application: Digit classification
	5.4.3	Network architectures
	5.4.4	Visualizing weights and activations
	5.4.5	Adversarial examples
	5.4.6	Pre-training and fine-tuning networks
5.5	More of	complex networks
	5.5.1	Three-dimensional CNNs
	5.5.2	Spatio-temporal models
		Recurrent networks and LSTMs
	5.5.4	Generative models

Readings

Richard Szeliski

Network architectures

- Batch normalization
- CNN architectures
- Visualization
- Adversarial examples
- Deconvolution and U-Nets
- Deep learning software





As before, I'm borrowing slides from

EECS 498-007 / 598-005 Deep Learning for Computer Vision Fall 2019

Course Description

UNIVERSITY OF MICHIGAN

Computer Vision has become ubiquitous in our society, with applications in search, image understanding, apps, mapping, medicine, drones, and self-driving cars. Core to many of these applications are visual recognition tasks such as image classification and object detection. Recent developments in neural network approaches have greatly advanced the performance of these state-of-the-art visual recognition systems. This course is a deep dive into details of neural-network based deep learning methods for computer vision. During this course, students will learn to implement, train and debug their own neural networks and gain a detailed understanding of cutting-edge research in computer vision. We will cover learning algorithms, neural network architectures, and practical engineering tricks for training and fine-tuning networks for visual recognition tasks.

Instructor Graduate Student Instructors





EECS 498-007 / 598-005 Deep Learning for Computer Vision Fall 2019

Lecture 7: Convolutional Networks

Justin Johnson



Components of a Convolutional Network

Fully-Connected Layers



Activation Function



Convolution Layers



Pooling Layers



Normalization



Justin Johnson



Idea: "Normalize" the outputs of a layer so they have zero mean and unit variance. Why?

Why? Helps reduce "internal covariate shift", improves optimization

We can normalize a batch of activations like this:



This is a **differentiable function**, so we can use it as an operator in our networks and backprop through it!

loffe and Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift", ICML 2015

Justin Johnson



Activation and weight scaling



Figure 5.54 Simple two hidden unit network with a ReLU activation function and no bias parameters for regressing the function $y = |x_1 + 1.1x_2|$: (a) can you guess a set of weights would fit this function? (b) a reasonable set of starting weights; (c) a poorly scaled set of weights.

Ex 5.1: Activation and weight scaling Consider the two hidden unit network shown in Figure 5.54 which uses ReLU activation functions and has no additive bias parameters. Your task is to find a set of weights that will fit the function

$$y = |x_1 + 1.1x_2|. \tag{5.91}$$

1. Can you guess a set of weights that will fit this function?

- 2. Starting with the weights shown in column b, compute the activations for the hidden and final units as well as the regression loss for the nine input values $(x_1, x_2) \in \{-1, 0, 1\} \times \{-1, 0, 1\}$.
- 3. Now compute the gradients of the squared loss with respect to all six weights using the backpropagation chain rule equations (5.79-5.82) and sum them up across the training samples to get a final gradient.
- 4. What step size should you take in the gradient direction, and what would your update squared loss become?
- 5. Repeat this exercise for the initial weights in column (c) of Figure 5.54.
- 6. Given this new set of weights, how much worse is your error decrease, and how many iterations would you expect it to take to achieve a reasonalbe solution?
- 7. Would batch normalization help in this case?

Richard Szeliski

UW CSE 576 - Networks Architectures

Activation and weight scaling

<u>x1</u>	<u>x2</u>	<u>w11</u>	<u>w12</u>	<u>w21</u>	<u>w22</u>	<u>z1</u>	<u>z2</u>	<u>w31</u>	<u>w32</u>	Y	target	L2 loss
-1	-1	-1	-1	1	1	2	0	1	1	2	2.1	0.01
-1	0					1	0			1	1	0
-1	1					0	0			0	0.1	0.01
0	-1					1	0			1	1.1	0.01
0	0					0	0			0	0	0
0	1					0	1			1	1.1	0.01
1	-1					0	0			0	0.1	0.01
1	0					0	1			1	1	0
1	1					0	2			2	2.1	0.01
												0.06
<u>x1</u>	<u>x2</u>	<u>d w11</u>	<u>d w12</u>	<u>d w21</u>	d w22	<u>d z1</u>	<u>d z2</u>	d w31	<u>d w32</u>	<u>g y</u>		
-1	-1	0.1	0.1	0	0	-0.1	-0.1	-0.2	0	-0.1		
-1	0	0	0	0	0	0	0	0	0	0		
-1	1	0	0	0	0	0	0	0	0	-0.1		
0	-1	0	0.1	0	0	-0.1	-0.1	-0.1	0	-0.1		
0	0	0	0	0	0	0	0	0	0	0		
0	1	0	0	0	-0.1	-0.1	-0.1	0	-0.1	-0.1		
1	-1	0	0	0	0	0	0	0	0	-0.1		
1	0	0	0	0	0	0	0	0	0	0		
1	1	0	0	-0.1	-0.1	-0.1	-0.1	0	-0.2	-0.1		
grad-:		-0.1	-0.2	0.1	0.2			0.3	0.3			
• • • • • • • • • • • • • • • • • • •	Balanc	ed Imb	alanced	(+)				•				



Figure 5.54 Simple two hidden unit network with a ReLU activation function and no bias parameters for regressing the function $y = |x_1 + 1.1x_2|$: (a) can you guess a set of weights would fit this function? (b) a reasonable set of starting weights; (c) a poorly scaled set of weights.

Activation and weight scaling



<u>x1</u>	x	2 w1	1	<u>w12</u>	<u>w21</u>	<u>w22</u>	<u>z1</u>	<u>z2</u>	<u>w3</u>	1	<u>w32</u>	Y		targ	et L2	2 loss		$_{x_1}$	/	\sum_{x_2}	$_{x_1}$	```	x_2	r ₁	$_{x_2}$
	-1	-1	-1	-1	. 1	. 1		2	0	1		1	2	2	2.1	0.0	1		(a)		(1	b)		(c)	
	-1	0						v1	x2	w11		w12	w21		w22	71		7)	w31	w32	v		target	12 000	und no bias
	-1	1						-1	-	1 -	100	-100)	100	10	0	200	0	0.0	0.0	1 1	2	2.1	0.01	of weights
	0	-1					-	-1		-	100	200		100			100	0	0.0.		-	- 1	1	0	caled set of
	0	0						-1		1							0	0				0	0.1	0.01	-
	0	1					· · · ·	0	-	1							100	0				1	1.1	0.01	
	1	-1						0		-							0	0				0	0	0	
	1	0						0		1			_				0	100				1	1.1	0.01	
	1	1						1	-	1							0	0			_	0	0.1	0.01	-
								1		0							0	100				1	1	0	-
<u>x1</u>	x	<u>2 d v</u>	<mark>v11</mark>	<u>d w12</u>	<u>d w21</u>	<u>d w22</u>	<u>d z1</u>	1		1							0	200				2	2.1	0.01	
	-1	-1	0.1	0.1	. 0	0																		0.06	
	-1	0	0	0	0 0	0		x1	x2	d w11	1	d w12	d w2	21	d w22	d z1	L	d z2	d w31	d w32	gу				
	-1	1	0	0	0 0	0		-1	-:	1 0.	.001	0.001	L	0		0 -	0.001	-0.001	-20)	0	-0.1			-
	0	-1	0	0.1	0	0		-1	(0	0	()	0		0	0	0	()	0	0			
	0	0	0	0	0 0	0		-1		1	0	()	0		0	0	0	()	0	-0.1			
	0	1	0	0	0	-0.1		0	-:	1	0	0.001	L	0		0 -	0.001	-0.001	-10)	0	-0.1			
	1	-1	0	0	0	0		0		0	0	()	0		0	0	0	()	0	0			
	1	0	0	0	0	0		0		1	0	()	0	-0.00)1 -	0.001	-0.001	() -:	LO	-0.1			
	1	1	0	0	-0.1	-0.1		1	-:	1	0	()	0		0	0	0	()	0	-0.1			
grad-:			-0.1	-0.2	2 0.1	0.2		1		0	0	()	0		0	0	0	()	0	0			
		Delement	lue le	-	0			1		1	0	0) -(0.001	-0.00)1 -	0.001	-0.001	() -:	20	-0.1			
4 F.		Balanced	Imp	alanced	+			grad-:		-0.	.001	-0.002	2 0	0.001	0.00)2			30) :	30				
								<	Balan	ced	lmba	alanced	(Ð											[



loffe and Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift", ICML 2015

Justin Johnson





Ioffe and Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift", ICML 2015

Justin Johnson



Input: $x: N \times D$

Learnable scale and shift parameters:

 $\gamma,\beta:D$

Learning $\gamma = \sigma$, $\beta = \mu$ will recover the identity function!

$$\begin{split} \mu_{j} &= \frac{1}{N} \sum_{i=1}^{N} x_{i,j} & \text{Per-channel} \\ \text{mean, shape is D} \\ \sigma_{j}^{2} &= \frac{1}{N} \sum_{i=1}^{N} (x_{i,j} - \mu_{j})^{2} & \text{Per-channel} \\ \text{std, shape is D} \\ \hat{x}_{i,j} &= \frac{x_{i,j} - \mu_{j}}{\sqrt{\sigma_{j}^{2} + \varepsilon}} & \text{Normalized x,} \\ \gamma_{j} &= \gamma_{j} \hat{x}_{i,j} + \beta_{j} & \text{Output,} \\ \text{Shape is N x D} \\ \end{split}$$

Justin Johnson

Batch Normalization: Test-Time minibatch; can't do this at test-time

Input: $x: N \times D$

Learnable scale and shift parameters:

 $\gamma,\beta:D$

Learning $\gamma = \sigma$, $\beta = \mu$ will recover the identity function!

$$\begin{split} \mu_j &= \frac{1}{N} \sum_{i=1}^N x_{i,j} & \text{Per-channel} \\ \text{mean, shape is D} \\ \sigma_j^2 &= \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2 & \text{Per-channel} \\ \text{std, shape is D} \\ \hat{x}_{i,j} &= \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}} & \text{Normalized x,} \\ \text{Shape is N x D} \end{split}$$

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + eta_j$$
 Output,
Shape is N x D

Batch Normalization: Test-Time

Input: $x: N \times D$

Learnable scale and shift parameters:

 $\gamma,\beta:D$

Learning $\gamma = \sigma$, $\beta = \mu$ will recover the identity function!

(Running) average of Per-channel $\mu_j = \text{values seen during}$ mean, shape is D training (Running) average of $\sigma_j^2 =$ values seen during Per-channel std, shape is D training $\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$ Normalized x, Shape is N x D $y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$ Output, Shape is N x D

Batch Normalization: Test-Time

Input: $x: N \times D$

Learnable scale and shift parameters:

 $\gamma,\beta:D$

During testing batchnorm becomes a linear operator. Can be fused with the previous fully-connected or conv layer

(Running) average of Per-channel $\mu_j = \text{values seen during}$ mean, shape is D training (Running) average of $\sigma_j^2 =$ values seen during Per-channel std, shape is D training $\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$ Normalized x, Shape is N x D Output, $y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$ Shape is N x D

Batch Normalization for ConvNets

Batch Normalization for **fully-connected** networks

Batch Normalization for **convolutional** networks (Spatial Batchnorm, BatchNorm2D)







Usually inserted after Fully Connected or Convolutional layers, and before nonlinearity.



loffe and Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift", ICML 2015

Justin Johnson





- Makes deep networks **much** easier to train!
- Allows higher learning rates, faster convergence
- Networks become more robust to initialization
- Acts as regularization during training
- Zero overhead at test-time: can be fused with conv!



loffe and Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift", ICML 2015

Justin Johnson



-

_

-



- Makes deep networks **much** easier to train!
- Allows higher learning rates, faster convergence
- Networks become more robust to initialization
- Acts as regularization during training
- Zero overhead at test-time: can be fused with conv!
- Not well-understood theoretically (yet)
- Behaves differently during training and testing: this is a very common source of bugs!

loffe and Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift", ICML 2015



Comparison of Normalization Layers



Wu and He, "Group Normalization", ECCV 2018

Justin Johnson



Group Normalization



Wu and He, "Group Normalization", ECCV 2018

Justin Johnson



Components of a Convolutional Network

Convolution Layers



Pooling Layers



Fully-Connected Layers



Activation Function



Normalization

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

Justin Johnson



Components of a Convolutional Network

Problem: What is the right way to combine all these components?



	list	in	hn	SO	n
J	ust			30	<u> </u>



Network architectures

- Batch normalization
- CNN architectures
- Visualization
- Adversarial examples
- Deconvolution and U-Nets
- Deep learning software







EECS 498-007 / 598-005 Deep Learning for Computer Vision Fall 2019

Lecture 8: CNN Architectures

Justin Johnson



ImageNet Classification Challenge



Justin Johnson



ImageNet Classification Challenge



Justin Johnson



AlexNet



227 x 227 inputs
5 Convolutional layers
Max pooling
3 fully-connected layers
ReLU nonlinearities

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



AlexNet



227 x 227 inputs
5 Convolutional layers
Max pooling
3 fully-connected layers
ReLU nonlinearities

Used "Local response normalization"; Not used anymore

Trained on two GTX 580 GPUs – only 3GB of memory each! Model split over two GPUs

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





Interesting trends here!



	Inpu	t size		Laye	er		Outp	out size			
Layer	С	H / W	filters	kernel	stride	pad	С	H / W	memory (KB)	params (k)	flop (M)
conv1	3	227	64	11	4	2	<mark>2</mark> 64	- 56	784	23	73
pool1	64	56		3	2	. C	<mark>)</mark> 64	- 27	182	0	0
conv2	64	27	192	5	1	. 2	2 192	. 27	547	307	224
pool2	192	27		3	2	. C) 192	13	127	0	0
conv3	192	13	384	3	1	. 1	L 384	. 13	254	664	112
conv4	384	13	256	3	1	. 1	L 256	5 13	169	885	145
conv5	256	13	256	3	1	. 1	L 256	5 13	169	590	100
pool5	256	13		3	2	. C) 256	6	36	0	0
flatten	256	6					9216		36	0	0
fc6	9216		4096				4096		16	37,749	38
fc7	4096		4096				4096		16	16,777	17
fc8	4096		1000				1000		4	4,096	4

Justin Johnson



AlexNet

Most of the **memory usage** is in the early convolution layers

Memory (KB)



Nearly all **parameters** are in the fully-connected layers

Params (K)

\dense 2048 2048 192 128 128 dense 192 128 Max 2048 2048 pooling Max Max 128 pooling poolina

Most **floating-point ops** occur in the convolution layers

MFLOP





Justin Johnson



ImageNet Classification Challenge



Justin Johnson



ImageNet Classification Challenge



Justin Johnson


ZFNet: A Bigger AlexNet

ImageNet top 5 error: 16.4% -> 11.7%



AlexNet but: CONV1: change from (11x11 stride 4) to (7x7 stride 2) CONV3,4,5: instead of 384, 384, 256 filters use 512, 1024, 512 More trial and error =(

Zeiler and Fergus, "Visualizing and Understanding Convolutional Networks", ECCV 2014



The really cool part

- Visualizing features with backprojection (``deconvnet'')
- How could you do this?
- What about max pooling?



Visualizing features through backprojection



Zeiler and Fergus, "Visualizing and Understanding Convolutional Networks", ECCV 2014

gray: max stimulus, color: patch

Richard Szeliski

Visualizing features through backprojection



Zeiler and Fergus, "Visualizing and Understanding Convolutional Networks", ECCV 2014

gray: max stimulus, color: patch

Richard Szeliski

UW CSE 576 - Networks Architectures

Visualizing features through backprojection



Attribution: saliency via occlusions



ImageNet Classification Challenge



Justin Johnson



ImageNet Classification Challenge



Justin Johnson



VGG: Deeper Networks, Regular Design

VGG Design rules:

All conv are 3x3 stride 1 pad 1 All max pool are 2x2 stride 2 After pool, double #channels

Network has 5 convolutional **stages**: Stage 1: conv-conv-pool Stage 2: conv-conv-pool Stage 3: conv-conv-pool Stage 4: conv-conv-conv-[conv]-pool Stage 5: conv-conv-conv-[conv]-pool

(VGG-19 has 4 conv in stages 4 and 5)

Softmax FC 1000 FC 4096 FC 4096 Pool

x3 conv, 25

Pool

Pool

1x11 conv. 9 Input

AlexNet



Simonyan and Zissermann, "Very Deep Convolutional Networks for Large-Scale Image Recognition", ICLR 2015

Lecture 8 - 50

VGG19

VGG: Deeper Networks, Regular Design

VGG Design rules:

All conv are 3x3 stride 1 pad 1

All max pool are 2x2 stride 2 After pool, double #channels

<u>Option 1:</u> Conv(5x5, C -> C)	<u>Option 2:</u> Conv(3x3, C -> C) Conv(3x3, C -> C)
Params: 25C ²	Params: 18C ²
FLOPs: 25C ² HW	FLOPs: 18C ² HW

	, negulai		511
	Two 3x3 conv hat receptive field as conv, but has few and takes less co	s same a single 5x ver parame mputation	5 ters
		Softmax	
		FC 1000	
		FC 4096	
		FC 4096	
		Pool	
->	()	3x3 conv, 256	
-/	C)	3x3 conv, 384	
->	()	Pool	
-		3x3 conv, 384	
		Pool	
		5x5 conv, 256	
2		11x11 conv, 96	

Input

AlexNet



Simonyan and Zissermann, "Very Deep Convolutional Networks for Large-Scale Image Recognition", ICLR 2015

AlexNet vs VGG-16: Much bigger network!





AlexNet vs VGG-16 (MFLOPs)



Justin Johnson

Lecture 8 - 52

Fall 2019

ImageNet Classification Challenge



Justin Johnson



ImageNet Classification Challenge



Justin Johnson



GoogLeNet: Focus on Efficiency

Many innovations for efficiency: reduce parameter count, memory usage, and computation



Fall 2019

Szegedy et al, "Going deeper with convolutions", CVPR 2015

Justin Johnson

GoogLeNet: Aggressive Stem

Stem network at the start aggressively downsamples input (Recall in VGG-16: Most of the compute was at the start)



3x3+2(5

LocalRespNorm

LocalRespNorm

MaxPool 3x3+2(S

input

Szegedy et al, "Going deeper with convolutions", CVPR 2015

Justin Johnson





Szegedy et al, "Going deeper with convolutions", CVPR 2015

Justin Johnson

Lecture 8 - 59

Fall 2019



Szegedy et al, "Going deeper with convolutions", CVPR 2015

Justin Johnson

Lecture 8 - 60

Fall 2019

GoogLeNet: Global Average Pooling

No large FC layers at the end! Instead uses **global average pooling** to collapse spatial dimensions, and one linear layer to produce class scores (Recall VGG-16: Most parameters were in the FC layers!)

	Input	size		Layer				t size			
Layer	С	H/W	filters	kernel	stride	pad	С	H/W	memory (KB)	params (k)	<mark>flop (M)</mark>
avg-pool	1024	7		7	1	0	1024	1	4	0	0
fc	1024		1000				1000		0	1025	1



Justin Johnson



GoogLeNet: Global Average Pooling

No large FC layers at the end! Instead uses "global average pooling" to collapse spatial dimensions, and one linear layer to produce class scores (Recall VGG-16: Most parameters were in the FC layers!)

	Input size			Laye	er		Outpu	t size			
Layer	С	H/W	filters	kernel	stride	pad	С	H/W	memory (KB)	params (k)	flop (M)
avg-pool	1024	7		7	1	0	1024	1	4	0	0
fc	1024		1000				1000		0	1025	1

Compare with VGG-16:

Layer	С	H/W	filters	kernel	stride	pad	С	H/W	memory (KB)	params (K)	flop (M)
flatten	512	7					25088		98		
fc6	25088			4096			4096		16	102760	103
fc7	4096			4096			4096		16	16777	17
fc8	4096			1000			1000		4	4096	4



Justin Johnson

GoogLeNet: Auxiliary Classifiers

Training using loss at the end of the network didn't work well: Network is too deep, gradients don't propagate cleanly

As a hack, attach "auxiliary classifiers" at several intermediate points in the network that also try to classify the image and receive loss

GoogLeNet was before batch normalization! With BatchNorm no longer need to use this trick



Fall 2019

ImageNet Classification Challenge



Justin Johnson



ImageNet Classification Challenge



Justin Johnson



Once we have Batch Normalization, we can train networks with 10+ layers. What happens as we go deeper?

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

Justin Johnson



Once we have Batch Normalization, we can train networks with 10+ layers. What happens as we go deeper?

Deeper model does worse than shallow model!

Initial guess: Deep model is **overfitting** since it is much bigger than the other model



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

Justin Johnson



Once we have Batch Normalization, we can train networks with 10+ layers. What happens as we go deeper?



In fact the deep model seems to be **underfitting** since it also performs worse than the shallow model on the training set! It is actually **underfitting**

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

Justin Johnson



A deeper model can <u>emulate</u> a shallower model: copy layers from shallower model, set extra layers to identity

Thus deeper models should do at least as good as shallow models

Hypothesis: This is an <u>optimization</u> problem. Deeper models are harder to optimize, and in particular don't learn identity functions to emulate shallow models

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

A deeper model can <u>emulate</u> a shallower model: copy layers from shallower model, set extra layers to identity

Thus deeper models should do at least as good as shallow models

Hypothesis: This is an <u>optimization</u> problem. Deeper models are harder to optimize, and in particular don't learn identity functions to emulate shallow models

Solution: Change the network so learning identity functions with extra layers is easy!

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



Solution: Change the network so learning identity functions with extra layers is easy!



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

Justin Johnson



Solution: Change the network so learning identity functions with extra layers is easy!



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

Justin Johnson



A residual network is a stack of many residual blocks

Regular design, like VGG: each residual block has two 3x3 conv

Network is divided into **stages**: the first block of each stage halves the resolution (with stride-2 conv) and doubles the number of channels



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



Uses the same aggressive **stem** as GoogleNet to downsample the input 4x before applying residual blocks:

	Ir	put					Οι	utput			
	S	size	Layer					size			
									params	flop	
Layer	С	H/W	filters	kernel	stride	pad	С	H/W	memory (KB)	(k)	(M)
conv	3	224	64	7	2	3	64	112	3136	9	<mark>118</mark>
max-pool	64	112		3	2	1	64	56	784	0	2

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

Justin Johnson



Fall 2019

Like GoogLeNet, no big fully-connected-layers: instead use **global average pooling** and a single linear layer at the end



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

Justin Johnson

Fall 2019

ResNet-18:

Stem: 1 conv layer Stage 1 (C=64): 2 res. block = 4 conv Stage 2 (C=128): 2 res. block = 4 conv Stage 3 (C=256): 2 res. block = 4 conv Stage 4 (C=512): 2 res. block = 4 conv Linear

ImageNet top-5 error: 10.92 GFLOP: 1.8

He et al, "Deep Residual Learning for Image Recognition", CVPR 2016 Error rates are 224x224 single-crop testing, reported by <u>torchvision</u>



Justin Johnson

ResNet-18:

```
Stem: 1 conv layer
Stage 1 (C=64): 2 res. block = 4 conv
Stage 2 (C=128): 2 res. block = 4 conv
Stage 3 (C=256): 2 res. block = 4 conv
Stage 4 (C=512): 2 res. block = 4 conv
Linear
```

ResNet-34:

Stem: 1 conv layer Stage 1: 3 res. block = 6 conv Stage 2: 4 res. block = 8 conv Stage 3: 6 res. block = 12 conv Stage 4: 3 res. block = 6 conv Linear

ImageNet top-5 error: 10.92 GFLOP: 1.8

ImageNet top-5 error: 8.58 GFLOP: 3.6



He et al, "Deep Residual Learning for Image Recognition", CVPR 2016 Error rates are 224x224 single-crop testing, reported by <u>torchvision</u>

Justin Johnson



ResNet-18:

```
Stem: 1 conv layer
Stage 1 (C=64): 2 res. block = 4 conv
Stage 2 (C=128): 2 res. block = 4 conv
Stage 3 (C=256): 2 res. block = 4 conv
Stage 4 (C=512): 2 res. block = 4 conv
Linear
```

ImageNet top-5 error: 10.92 GFLOP: 1.8

He et al, "Deep Residual Learning for Image Recognition", CVPR 2016 Error rates are 224x224 single-crop testing, reported by <u>torchvision</u> ResNet-34:

Stem: 1 conv layer Stage 1: 3 res. block = 6 conv Stage 2: 4 res. block = 8 conv Stage 3: 6 res. block = 12 conv Stage 4: 3 res. block = 6 conv Linear

ImageNet top-5 error: 8.58 GFLOP: 3.6

VGG-16:

ImageNet top-5 error: 9.62 GFLOP: 13.6



Justin Johnson



Residual Networks: Basic Block



"Basic" Residual block

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

Justin Johnson


Residual Networks: Basic Block



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

Justin Johnson



Residual Networks: Bottleneck Block





"Bottleneck" Residual block

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

Justin Johnson



Residual Networks: Bottleneck Block



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

Justin Johnson

Lecture 8 - 82

			Stag	ge 1	Stag	ge 2	Stag	ge 3	Stag	ge 4			
	Block	Stem									FC		ImageNet
	type	layers	Blocks	Layers	Blocks	Layers	Blocks	Layers	Blocks	Layers	layers	GFLOP	top-5 error
ResNet-18	Basic	1	2	4	2	4	2	4	2	4	1	1.8	10.92
ResNet-34	Basic	1	3	6	4	8	6	12	3	6	1	3.6	8.58

Softmax FC 1000 Pool 3x3 conv, 512 3x3 conv, 512, /2 3x3 conv, 128 3x3 conv. 128 3x3 conv, 64 3x3 conv. 64 3x3 conv. 64 3x3 conv, 64 3x3 conv, 64 3x3 conv, 64 Input

He et al, "Deep Residual Learning for Image Recognition", CVPR 2016 Error rates are 224x224 single-crop testing, reported by <u>torchvision</u>

Justin Johnson



ResNet-50 is the same as ResNet-34, but replaces Basic blocks with Bottleneck Blocks. This is a great baseline architecture for many tasks even today!

			Stag	ge 1	Stag	ge 2	Stag	ge 3	Stag	ge 4			
	Block	Stem									FC		ImageNet
	type	layers	Blocks	Layers	Blocks	Layers	Blocks	Layers	Blocks	Layers	layers	GFLOP	top-5 error
ResNet-18	Basic	1	2	4	2	4	2	4	2	4	. 1	. 1.8	10.92
ResNet-34	Basic	1	3	6	4	8	6	12	3	6	1	3.6	8.58
ResNet-50	Bottle	1	3	9	4	12	6	18	3	9	1	3.8	7.13



He et al, "Deep Residual Learning for Image Recognition", CVPR 2016 Error rates are 224x224 single-crop testing, reported by <u>torchvision</u>

Justin Johnson



Deeper ResNet-101 and ResNet-152 models are more accurate, but also more computationally heavy

			Stag	ge 1	Stag	ge 2	Sta	ge 3	Stag	ge 4			
	Block	Stem									FC		ImageNet
	type	layers	Blocks	Layers	Blocks	Layers	Blocks	Layers	Blocks	Layers	layers	GFLOP	top-5 error
ResNet-18	Basic	1	2	4	2	4	2	4	2	4	1	. 1.8	10.92
ResNet-34	Basic	1	3	6	4	8	6	12	3	6	1	. 3.6	8.58
ResNet-50	Bottle	1	3	9	4	12	6	18	3	9	1	. 3.8	7.13
ResNet-101	Bottle	1	3	9	4	12	23	69	3	9	1	. 7.6	6.44
ResNet-152	Bottle	1	3	9	8	24	36	108	3	9	1	. 11.3	5.94



He et al, "Deep Residual Learning for Image Recognition", CVPR 2016 Error rates are 224x224 single-crop testing, reported by <u>torchvision</u>

Justin Johnson

- Able to train very deep networks
- Deeper networks do better than shallow networks (as expected)
- Swept 1st place in all ILSVRC and COCO 2015 competitions
- Still widely used today!

MSRA @ ILSVRC & COCO 2015 Competitions

1st places in all five main tracks

- ImageNet Classification: "Ultra-deep" (quote Yann) 152-layer nets
- ImageNet Detection: 16% better than 2nd
- ImageNet Localization: 27% better than 2nd
- COCO Detection: 11% better than 2nd
- COCO Segmentation: 12% better than 2nd

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

Justin Johnson



Improving Residual Networks: Block Design

Original ResNet block

"Pre-Activation" ResNet Block



He et al, "Identity mappings in deep residual networks", ECCV 2016

Justin Johnson



Improving Residual Networks: Block Design

Original ResNet block

"Pre-Activation" ResNet Block



Slight improvement in accuracy (ImageNet top-1 error)

ResNet-152: 21.3 vs **21.1** ResNet-200: 21.8 vs **20.7**

Not actually used that much in practice

Lecture 8 - 88



Fall 2019

He et al, "Identity mappings in deep residual networks", ECCV 2016

Justin Johnson





Fall 2019

BN: batch normalized version

Canziani et al, "An analysis of deep neural network models for practical applications", 2017

Justin Johnson



Inception-v4: Resnet + Inception!



BN: batch normalized version

Canziani et al, "An analysis of deep neural network models for practical applications", 2017

Justin Johnson



VGG: Highest memory, most operations





BN: batch normalized version

Canziani et al, "An analysis of deep neural network models for practical applications", 2017

Justin Johnson





GoogLeNet: Very efficient!



BN: batch normalized version

Canziani et al, "An analysis of deep neural network models for practical applications", 2017

Justin Johnson





AlexNet: Low compute, lots of parameters



BN: batch normalized version

Canziani et al, "An analysis of deep neural network models for practical applications", 2017

Justin Johnson







BN: batch normalized version

Canziani et al, "An analysis of deep neural network models for practical applications", 2017

Justin Johnson



ImageNet Classification Challenge



Justin Johnson



ImageNet Classification Challenge



Justin Johnson



ImageNet 2016 winner: Model Ensembles

Multi-scale ensemble of Inception, Inception-Resnet, Resnet, Wide Resnet models

	Inception- v3	Inception- v4	Inception- Resnet-v2	Resnet- 200	Wrn-68-3	Fusion (Val.)	Fusion (Test)
Err. (%)	4.20	4.01	3.52	4.26	4.65	2.92 (-0.6)	2.99

Shao et al, 2016



Improving ResNets



Justin Johnson

Lecture 8 - 98



Xie et al, "Aggregated residual transformations for deep neural networks", CVPR 2017

Justin Johnson

Lecture 8 - 99



Xie et al, "Aggregated residual transformations for deep neural networks", CVPR 2017

Justin Johnson

Lecture 8 - 100



<u>Convolution with groups=1</u>: Normal convolution

Input: $C_{in} \times H \times W$ Weight: $C_{out} \times C_{in} \times K \times K$ Output: $C_{out} \times H' \times W'$ FLOPs: $C_{out}C_{in}K^2HW$

All convolutional kernels touch all C_{in} channels of the input





<u>Convolution with groups=1</u>: Normal convolution

Input: C_{in} x H x W Weight: C_{out} x C_{in} x K x K Output: C_{out} x H' x W' FLOPs: C_{out}C_{in}K²HW

All convolutional kernels touch all C_{in} channels of the input

Convolution with groups=2: Two parallel convolution layers that work on half the channels Input: C_{in} x H x W Split Group 1: Group 2: $(C_{in} / 2) \times H \times W$ $(C_{in} / 2) \times H \times W$ $Conv(K \times K, C_{in}/2 \rightarrow C_{out}/2)$ $Conv(K \times K, C_{in}/2 \rightarrow C_{out}/2)$ Out 1: Out 2: $(C_{out} / 2) \times H' \times W'$ $(C_{out} / 2) \times H' \times W'$ → Concat Output: C_{out} x H' x W'

Justin Johnson



<u>Convolution with groups=1</u>: Normal convolution

Input: C_{in} x H x W Weight: C_{out} x C_{in} x K x K Output: C_{out} x H' x W' FLOPs: C_{out}C_{in}K²HW

All convolutional kernels touch all C_{in} channels of the input

<u>Convolution with groups=G</u>: G parallel conv layers; each "sees" C_{in}/G input channels and produces C_{out}/G output channels

Input: $C_{in} \times H \times W$ Split to $G \times [(C_{in} / G) \times H \times W]$ Weight: $G \times (C_{out} / G) \times (C_{in} \times G) \times K \times K$ G parallel convolutions Output: $G \times [(C_{out} / G) \times H' \times W']$ Concat to $C_{out} \times H' \times W'$ FLOPs: $C_{out}C_{in}K^2HW/G$

Justin Johnson



<u>Convolution with groups=1</u>: Normal convolution

Input: C_{in} x H x W Weight: C_{out} x C_{in} x K x K Output: C_{out} x H' x W' FLOPs: C_{out}C_{in}K²HW

All convolutional kernels touch all C_{in} channels of the input

Depthwise Convolution

Special case: $G=C_{in}$, $C_{out} = nC_{in}$ Each input channel is convolved with n different K x K filters to produce n output channels <u>Convolution with groups=G</u>: G parallel conv layers; each "sees" C_{in}/G input channels and produces C_{out}/G output channels

Input: $C_{in} \times H \times W$ Split to $G \times [(C_{in} / G) \times H \times W]$ Weight: $G \times (C_{out} / G) \times (C_{in} \times G) \times K \times K$ G parallel convolutions Output: $G \times [(C_{out} / G) \times H' \times W']$ Concat to $C_{out} \times H' \times W'$ FLOPs: $C_{out}C_{in}K^2HW/G$

Justin Johnson



Grouped Convolution in PyTorch

PyTorch convolution gives an option for groups!

Conv2d

CLASS torch.nn.Conv2d(in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1, groups=1, bias=True, [SOURCE] padding_mode='zeros')





Justin Johnson



ResNeXt: Maintain computation by adding groups!

Model	Groups	Group width	Top-1 Error
ResNet-50	1	64	23.9
ResNeXt-50	2	40	23
ResNeXt-50	4	24	22.6
ResNeXt-50	8	14	22.3
ResNeXt-50	32	4	22.2

Model	Groups	Group width	Top-1 Error
ResNet-101	1	64	22.0
ResNeXt-101	2	40	21.7
ResNeXt-101	4	24	21.4
ResNeXt-101	8	14	21.3
ResNeXt-101	32	4	21.2

Adding groups improves performance with same computational complexity!

Xie et al, "Aggregated residual transformations for deep neural networks", CVPR 2017



ImageNet Classification Challenge



Justin Johnson



Squeeze-and-Excitation Networks

Adds a "Squeeze-and-excite" branch to each residual block that performs global pooling, full-connected layers, and multiplies back onto feature map

Adds **global context** to each residual block!

Won ILSVRC 2017 with ResNeXt-152-SE



SE-ResNet Module

Justin Johnson



ImageNet Classification Challenge



Justin Johnson

Densely Connected Neural Networks

Dense blocks where each layer is connected to every other layer in feedforward fashion

Alleviates vanishing gradient, strengthens feature propagation, encourages feature reuse



Dense Block

Justin Johnson





MobileNets: Tiny Networks (For Mobile Devices)

Standard Convolution Block

Total cost: 9C²HW

Depthwise Separable Convolution

Total cost: $(9C + C^2)HW$



Howard et al, "MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications", 2017

Justin Johnson



MobileNets: Tiny Networks (For Mobile Devices)

Depthwise Separable Convolution

Fall 2019

Total cost: $(9C + C^2)HW$



Howard et al, "MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications", 2017

Justin Johnson

Neural Architecture Search

Designing neural network architectures is hard – let's automate it!

- One network (controller) outputs network architectures
- Sample child networks from controller and train them
- After training a batch of child networks, make a gradient step on controller network (Using **policy gradient**)
- Over time, controller learns to output good architectures!

Lecture 8 - 115



Fall 2019

Zoph and Le, "Neural Architecture Search with Reinforcement Learning", ICLR 2017

Justin Johnson
Neural Architecture Search

Designing neural network architectures is hard – let's automate it!

- One network (controller) outputs network architectures
- Sample child networks from controller and train them
- After training a batch of child networks, make a gradient step on controller network (Using **policy gradient**)
- Over time, controller learns to output good architectures!
- VERY EXPENSIVE!! Each gradient step on controller requires training a batch of child models!
- Original paper trained on 800 GPUs for 28 days!
- Followup work has focused on efficient search



Zoph and Le, "Neural Architecture Search with Reinforcement Learning", ICLR 2017



Neural Architecture Search

Neural architecture search can be used to find efficient CNN architectures!



Zoph et al, "Learning Transferable Architectures for Scalable Image Recognition", CVPR 2018

Justin Johnson

Lecture 8 - 117



CNN Architectures Summary

Early work (AlexNet -> ZFNet -> VGG) shows that **bigger networks work better**

GoogLeNet one of the first to focus on **efficiency** (aggressive stem, 1x1 bottleneck convolutions, global avg pool instead of FC layers)

ResNet showed us how to train extremely deep networks – limited only by GPU memory! Started to show diminishing returns as networks got bigger

After ResNet: **Efficient networks** became central: how can we improve the accuracy without increasing the complexity?

Lots of **tiny networks** aimed at mobile devices: MobileNet, ShuffleNet, etc

Neural Architecture Search promises to automate architecture design

Which Architecture should I use?

Don't be a hero. For most problems you should use an off-the-shelf architecture; don't try to design your own!

If you just care about accuracy, **ResNet-50** or **ResNet-101** are great choices

If you want an efficient network (real-time, run on mobile, etc) try **MobileNets** and **ShuffleNets**



Network architectures

- Batch normalization
- CNN architectures
- Visualization
- Adversarial examples
- Deconvolution and U-Nets
- Deep learning software







EECS 498-007 / 598-005 Deep Learning for Computer Vision Fall 2019

Lecture 14: Visualizing CNNs

(see slides on Web site)



Justin Johnson



Intermediate Features via (guided) backprop



Maximally activating patches (Each row is a different neuron)



Guided Backprop

Zeiler and Fergus, "Visualizing and Understanding Convolutional Networks", ECCV 2014 Springenberg et al, "Striving for Simplicity: The All Convolutional Net", ICLR Workshop 2015 Figure copyright Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, Martin Riedmiller, 2015; reproduced with permission

Justin Johnson



(Guided) backprop:

Find the part of an image that a neuron responds to

Gradient ascent:

Generate a synthetic image that maximally activates a neuron





$$\arg\max_{I} \frac{S_c(I) - \lambda \|I\|_2^2}{2}$$

1. Initialize image to zeros

score for class c (before Softmax)



Repeat:

- 2. Forward image to compute current scores
- 3. Backprop to get gradient of neuron value with respect to image pixels
- 4. Make a small update to the image

Justin Johnson



$$\arg\max_{I} S_c(I) - \frac{\lambda \|I\|_2^2}{2}$$

Simple regularizer: Penalize L2 norm of generated image

Simonyan, Vedaldi, and Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and

Saliency Maps", ICLR Workshop 2014.

Figures copyright Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, 2014; reproduced with permission.

Justin Johnson



$$\arg\max_{I} S_c(I) - \frac{\lambda \|I\|_2^2}{2}$$

Simple regularizer: Penalize L2 norm of generated image









Simonyan, Vedaldi, and Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014. Figures copyright Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, 2014; reproduced with permission.

bell pepper

lemon

husky

Justin Johnson



$$\arg\max_{I} S_c(I) - \frac{\lambda \|I\|_2^2}{2}$$

Simple regularizer: Penalize L2 norm of generated image



Simonyan, Vedaldi, and Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.

Figures copyright Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, 2014; reproduced with permission.

goose

ostrich

limousine

Justin Johnson



$$\arg\max_{I} S_c(I) - \lambda \|I\|_2^2$$

Better regularizer: Penalize L2 norm of image; also during optimization periodically

- 1. Gaussian blur image
- 2. Clip pixels with small values to 0
- 3. Clip pixels with small gradients to 0

Yosinski et al, "Understanding Neural Networks Through Deep Visualization", ICML DL Workshop 2014.



$$\arg\max_{I} S_c(I) - \lambda \|I\|_2^2$$

Better regularizer: Penalize L2 norm of image; also during optimization periodically

- 1. Gaussian blur image
- 2. Clip pixels with small values to 0
- 3. Clip pixels with small gradients to 0



Flamingo







Pelican



Indian Cobra

Yosinski et al, "Understanding Neural Networks Through Deep Visualization", ICML DL Workshop 2014.



$$\arg\max_{I} S_c(I) - \lambda \|I\|_2^2$$

Better regularizer: Penalize L2 norm of image; also during optimization periodically

- 1. Gaussian blur image
- 2. Clip pixels with small values to 0
- 3. Clip pixels with small gradients to 0



Hartebeest



Station Wagon



Billiard Table



Black Swan

Yosinski et al, "Understanding Neural Networks Through Deep Visualization", ICML DL Workshop 2014.

Justin Johnson



Use the same approach to visualize intermediate features



Yosinski et al, "Understanding Neural Networks Through Deep Visualization", ICML DL Workshop 2014.

Justin Johnson

Adding "multi-faceted" visualization gives even nicer results: (Plus more careful regularization, center-bias)

Reconstructions of multiple feature types (facets) recognized by the same "grocery store" neuron



Corresponding example training set images recognized by the same neuron as in the "grocery store" class



Nguyen et al, "Multifaceted Feature Visualization: Uncovering the Different Types of Features Learned By Each Neuron in Deep Neural Networks", ICML Visualization for Deep Learning Workshop 2016.

Justin Johnson





Nguyen et al, "Multifaceted Feature Visualization: Uncovering the Different Types of Features Learned By Each Neuron in Deep Neural Networks", ICML Visualization for Deep Learning Workshop 2016.

Justin Johnson







leaf beetle



badger



toaster





cloak

lawn mower







candle French loaf barn table lamp lemon D pool table aircraft carrier entertainment ctr chest running shoe water jug broom cellphone jean

Nguyen et al, "Synthesizing the preferred inputs for neurons in neural networks via deep generator networks," NIPS 2016

Justin Johnson



Feature visualization by optimization

How neural networks build up their understanding of images



Edges (layer conv2d0)

Textures (layer mixed3a)

Parts (layers mixed4b & mixed4c)

Objects (layers mixed4d & mixed4e)

Feature visualization allows us to see how GoogLeNet [1], trained on the ImageNet [2] dataset, builds up its understanding of images over many layers. Visualizations of all channels are available in the appendix.

AUTHORS	AFFILIATIONS	PUBLISHED	DOI
Chris Olah	Google Brain Team	Nov. 7, 2017	10.23915/distill.00007
Alexander Mordvintsev	Google Research		
Ludwig Schubert	Google Brain Team		

https://distill.pub/2017/feature-visualization/

Adversarial Examples

 Start from an arbitrary image
Pick an arbitrary category
Modify the image (via gradient ascent) to maximize the class score
Stop when the network is fooled



Adversarial Examples

African elephant



koala



schooner







Difference

10x Difference



10x Difference



Boat image is CCO public domain Elephant image is CCO public domain

Justin Johnson



Network architectures

- Batch normalization
- CNN architectures
- Visualization
- Adversarial examples
- Deconvolution and U-Nets
- Deep learning software







EECS 498-007 / 598-005 Deep Learning for Computer Vision Fall 2019

Lecture 16: Semantic Segmentation

Justin Johnson



Computer Vision Tasks: Semantic Segmentation



Semantic Segmentation

Label each pixel in the image with a category label

Don't differentiate instances, only care about pixels



This image is CC0 public domain

Justin Johnson

Lecture 16 - 141

Grass

Sky

Semantic Segmentation Idea: Sliding Window



Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013 Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

Justin Johnson



Semantic Segmentation Idea: Sliding Window



Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013 Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

Justin Johnson

between overlapping patches



Design a network as a bunch of convolutional layers to make predictions for pixels all at once!



Design a network as a bunch of convolutional layers to make predictions for pixels all at once!



Input:**Problem #1**: Effective receptive3 x H x Wfield size is linear in number of
conv layers: With L 3x3 conv
layers, receptive field is 1+2L

Long et al, "Fully convolutional networks for semantic segmentation", CVPR 2015

Lecture 16 - 145

Fall 2019

Design a network as a bunch of convolutional layers to make predictions for pixels all at once!



Input:Problem #1: Effective receptive3 x H x Wfield size is linear in number of
conv layers: With L 3x3 conv
layers, receptive field is 1+2L

Problem #2: Convolution on high res images is expensive! Recall ResNet stem aggressively downsamples

Long et al, "Fully convolutional networks for semantic segmentation", CVPR 2015

Justin Johnson



Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!



Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015 Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

Justin Johnson



Downsampling: Pooling, strided convolution



Input: 3 x H x W Design network as a bunch of convolutional layers, with downsampling and upsampling inside the network!

Upsampling: ???





Predictions: HxW

Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015 Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

Justin Johnson



In-Network Upsampling: "Unpooling"

Bed of Nails



C x 2 x 2 C x 4 x 4

Justin Johnson



In-Network Upsampling: "Unpooling"

Bed of Nails

Nearest Neighbor



In-Network Upsampling: Bilinear Interpolation



Input: C x 2 x 2 Output: C x 4 x 4

 $f_{x,y} = \sum_{i,j} f_{i,j} \max(0, 1 - |x - i|) \max(0, 1 - |y - j|) \quad i \in \{\lfloor x \rfloor - 1, \dots, \lceil x \rceil + 1\}$ Use two closest neighbors in x and y to construct linear approximations $j \in \{\lfloor y \rfloor - 1, \dots, \lceil y \rceil + 1\}$

Justin Johnson


In-Network Upsampling: Bicubic Interpolation



Input: C x 2 x 2

Output: C x 4 x 4

Use **three** closest neighbors in x and y to construct **cubic** approximations (This is how we normally resize images!)

Justin Johnson



In-Network Upsampling: "Max Unpooling"

Max Pooling: Remember which position had the max

Max Unpooling: Place into remembered positions





Pair each downsampling layer with an upsampling layer

Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

Justin Johnson



Recall: Normal 3 x 3 convolution, stride 1, pad 1

Input: 4 x 4

Output: 4 x 4

Justin Johnson



Recall: Normal 3 x 3 convolution, stride 1, pad 1



Input: 4 x 4

Output: 4 x 4

Justin Johnson

Recall: Normal 3 x 3 convolution, stride 1, pad 1



Input: 4 x 4

Output: 4 x 4

Justin Johnson

Recall: Normal 3 x 3 convolution, <u>stride 2</u>, pad 1



Input: 4 x 4

Output: 2 x 2

Fall 2019

Justin Johnson

Recall: Normal 3 x 3 convolution, <u>stride 2</u>, pad 1



Input: 4 x 4

Output: 2 x 2

Justin Johnson



Recall: Normal 3 x 3 convolution, <u>stride 2</u>, pad 1





Output: 2 x 2

Justin Johnson



Recall: Normal 3 x 3 convolution, <u>stride 2</u>, pad 1



Convolution with stride > 1 is "Learnable Downsampling" Can we use stride < 1 for "Learnable Upsampling"?

Dot product between input and filter



Input: 4 x 4

Output: 2 x 2

Justin Johnson



3 x 3 convolution transpose, stride 2



Justin Johnson

3 x 3 convolution transpose, stride 2



Weight filter by input value and copy to output



Output: 4 x 4

Input: 2 x 2

Justin Johnson



3 x 3 **convolution transpose**, stride 2

Filter moves 2 pixels in <u>output</u> for every 1 pixel in <u>input</u>



Weight filter by input value and copy to output



Output: 4 x 4

Input: 2 x 2

3 x 3 convolution transpose, stride 2

Filter moves 2 pixels in <u>output</u> for every 1 pixel in <u>input</u>



Weight filter by input value and copy to output



Output: 4 x 4

Input: 2 x 2

Justin Johnson



Justin Johnson

Lecture 16 - 165

Weight filter by input value and copy to output



Output: 4 x 4

3 x 3 convolution transpose, stride 2

This gives 5x5 output – need to trim one pixel from top and left to give 4x4 output



Input: 2 x 2

Learnable Upsampling: Transposed Convolution

output overlaps

Fall 2019

Transposed Convolution: 1D example

Filter Output Input ах ay Х a az**+**bx Y b by Ζ bz

Output has copies of filter weighted by input

Stride 2: Move 2 pixels output for each pixel in input

Sum at overlaps



Transposed Convolution: 1D example



This has many names:

- Deconvolution (bad)!
- Upconvolution
- Fractionally strided convolution
- Backward strided convolution
- <u>Transposed Convolution</u> (best name)

Justin Johnson



We can express convolution in terms of a matrix multiplication

$$\vec{x} * \vec{a} = X\vec{a}$$

$$\begin{bmatrix} x & y & x & 0 & 0 & 0 \\ 0 & x & y & x & 0 & 0 \\ 0 & 0 & x & y & x & 0 \\ 0 & 0 & 0 & x & y & x \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ay + bz \\ ax + by + cz \\ bx + cy + dz \\ cx + dy \end{bmatrix}$$

Example: 1D conv, kernel size=3, stride=1, padding=1

Justin Johnson



We can express convolution in terms of a matrix multiplication

Transposed convolution multiplies by the transpose of the same matrix:

$$\vec{x} * \vec{a} = X\vec{a} \qquad \qquad \vec{x} *^{T} \vec{a} = X^{T}\vec{a}$$

$$\begin{bmatrix} x & y & x & 0 & 0 & 0 \\ 0 & x & y & x & 0 & 0 \\ 0 & 0 & x & y & x & 0 \\ 0 & 0 & 0 & x & y & x \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ay + bz \\ ax + by + cz \\ bx + cy + dz \\ cx + dy \end{bmatrix} \begin{bmatrix} x & 0 & 0 & 0 \\ y & x & 0 & 0 \\ z & y & x & 0 \\ 0 & z & y & x \\ 0 & 0 & z & y \\ 0 & 0 & 0 & z \end{bmatrix} \begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix} = \begin{bmatrix} ax \\ ay + bx \\ az + by + cx \\ bz + cy + dx \\ cz + dy \\ dz \end{bmatrix}$$

Example: 1D conv, kernel size=3, stride=1, padding=1

When stride=1, transposed conv is just a regular conv (with different padding rules)

Justin Johnson

We can express convolution in terms of a matrix multiplication

 $\tau T \rightarrow$

$$\vec{x} *^T \vec{a} = X^T \vec{a}$$

$$\begin{aligned} x * a &= X a \\ \begin{bmatrix} x & y & x & 0 & 0 & 0 \\ 0 & 0 & x & y & x & 0 \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ay + bz \\ bx + cy + dz \end{bmatrix} \end{aligned}$$

Example: 1D conv, kernel size=3, <u>stride=2</u>, padding=1



We can express convolution in terms of a matrix multiplication

$$\vec{x} * \vec{a} = X\vec{a}$$

$$\begin{bmatrix} x & y & x & 0 & 0 & 0 \\ 0 & 0 & x & y & x & 0 \end{bmatrix} \begin{bmatrix} 0 \\ a \\ b \\ c \\ d \\ 0 \end{bmatrix} = \begin{bmatrix} ay + bz \\ bx + cy + dz \end{bmatrix}$$

Transposed convolution multiplies by the transpose of the same matrix:

$$\vec{x} *^{T} \vec{a} = X^{T} \vec{a}$$

$$\begin{bmatrix} x & 0 \\ y & 0 \\ z & x \\ 0 & y \\ 0 & z \\ 0 & 0 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} ax \\ ay \\ az + bx \\ by \\ bz \\ 0 \end{bmatrix}$$

Example: 1D conv, kernel size=3, <u>stride=2</u>, padding=1

When stride>1, transposed convolution cannot be expressed as normal conv

Justin Johnson



Semantic Segmentation: Fully Convolutional Network

Downsampling: Pooling, strided convolution



Input: 3 x H x W Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!



Upsampling: linterpolation, transposed conv



Predictions: H x W

Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015 Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

Loss function: Per-Pixel cross-entropy

Justin Johnson



<u>U-Nets</u>

- Skip connections between encoding and decoding stages
- Widely used in image denoising, inpainting, flow, stereo, ...



Figure 5.41 (a) The deconvolution network of Noh, Hong, and Han (2015) © 2015 IEEE and (b-c) the U-Net of Ronneberger, Fischer, and Brox (2015), redrawn using the PlotNeuralNet LaTeX package by Matt Dietke. In addition to the fine-to-coarse-to-fine bottleneck used in (a), the U-Net also has skip connections between encoding and decoding layers at the same resolution.



EECS 498-007 / 598-005 Deep Learning for Computer Vision Fall 2019



Justin Johnson

Lecture 9 - 174



Summary: Network architectures

- Batch normalization
- CNN architectures
- Visualization
- Adversarial examples
- Deconvolution and U-Nets
- Deep learning software



