#### 6.S094: Deep Learning for Self-Driving Cars 2018

https://selfdrivingcars.mit.edu

Lex Fridman



## Lecture 4: Computer Vision



January

# 最专业报告分享群:

#### •每日分享5+科技行业报告

- 同行业匹配,覆盖人工智能、大数据、机器人、 智慧医疗、智能家居、物联网等行业。
- 高质量用户,同频的人说同样的话

扫描右侧二维码, 或直接搜索关注公众号: 智东西(zhidxcom) 回复"报告群"加入



### **Computer Vision** is Deep Learning

Technology



https://selfdrivingcars.mit.edu lex.mit.edu

internal state

January 2018

### **Images are Numbers**



- **Regression:** The output variable takes continuous values
- Classification: The output variable takes class labels
  - Underneath it may still produce continuous values such as probability of belonging to a particular class.

Computer Vision with Deep Learning: Our intuition about what's "hard" is flawed (in complicated ways)

Visual perception:540,000,000 years of dataBipedal movement:230,000,000 years of dataAbstract thought:100,000 years of data



Prediction: Dog

+ Distortion

Prediction: Ostrich

"Encoded in the large, highly evolve sensory and motor portions of the human brain is a **billion years of experience** about the nature of the world and how to survive in it.... Abstract thought, though, is a new trick, perhaps less than **100 thousand years** old. We have not yet mastered it. It is not all that intrinsically difficult; it just seems so when we do it." - Hans Moravec, Mind Children (1988)



References: [6, 7, 11, 68]

## **Neuron:** Biological Inspiration for Computation



• Neuron: computational building block for the brain



 (Artificial) Neuron: computational building block for the "neural network"

#### Differences (among others):

- **Parameters:** Human brains have ~10,000,000 times synapses than artificial neural networks.
- **Topology:** Human brains have no "layers". Topology is complicated.
- **Async:** The human brain works asynchronously, ANNs work synchronously.
- Learning algorithm: ANNs use gradient descent for learning. Human brains use ... (we don't know)
- **Processing speed**: Single biological neurons are slow, while standard neurons in ANNs are fast.
- **Power consumption:** Biological neural networks use very little power compared to artificial networks
- Stages: Biological networks usually don't stop / start learning. ANNs have different fitting (train) and prediction (evaluate) phases.

#### Similarity (among others):

• Distributed computation on a large scale.



#### Retinal Ganglion Cell Activity



## Human Vision

Its structure is instructive and inspiring!

Thalamocortical System Simulation: 8 million cortical neurons + 2 billion synapses:





References: [118]

## Visual Cortex

(Its Structure is Instructive and Inspiring)





Lex Fridman lex.mit.edu

### **Deep Learning is Hard:** Illumination Variability



Massachusetts Institute of Technology

For the full updated list of references visit: https://selfdrivingcars.mit.edu/references



MIT 6.S094: Deep Learning for Self-Driving Cars https://selfdrivingcars.mit.edu

Lex Fridman lex.mit.edu

January 2018

### **Deep Learning is Hard: Pose Variability**



### Figure 1. The deformable and truncated cat. Cats exhibit (al-

Parkhi et al. "The truth about cats and dogs." 2011.



For the full updated list of references visit: https://selfdrivingcars.mit.edu/references

[69]

MIT 6.S094: Deep Learning for Self-Driving Cars https://selfdrivingcars.mit.edu

### Deep Learning is Hard: Intra-Class Variability















Parkhi et al. "Cats and dogs." 2012.



For the full updated list of references visit: https://selfdrivingcars.mit.edu/references



MIT 6.S094: Deep Learning for Self-Driving Cars https://selfdrivingcars.mit.edu Lex Fridman January lex.mit.edu 2018

## Occlusion





MIT 6.S094: Deep Learning for Self-Driving Cars Lex Fridman https://selfdrivingcars.mit.edu lex.mit.edu

January <mark>2018</mark>

## Occlusion





January <mark>2018</mark>

## Occlusion





January

2018

## Philosophical Ambiguity: "Image Classification" is not (yet) "Understanding"





References: [121]

Lex Fridman lex.mit.edu January 2018

## **Image Classification Pipeline**



References: [81, 89]

## **Famous Computer Vision Datasets**



#### **MNIST:** handwritten digits



#### CIFAR-10(0): tiny images



#### ImageNet: WordNet hierarchy



#### Places: natural scenes

January

2018

## Let's Build an Image Classifier for CIFAR-10

airplane automobile bird cat deer dog frog horse ship truck

test image					training image				pixe	pixel-wise absolute value differences				nces
56	32	10	18		10	20	24	17		46	12	14	1	C
90	23	128	133		8	10	89	100		82	13	39	33	
24	26	178	200		12	16	178	170	=	12	10	0	30	>
2	0	255	220		4	32	233	112		2	32	22	108	



January

2018

## Let's Build an Image Classifier for CIFAR-10

1	test image					training image				pix	pixel-wise absolute value differences				
	56	32	10	18	12	10	20	24	17		46	12	14	1	
	90	23	128	133		8	10	89	100		82	13	39	33	
	24	26	178	200	-	12	<mark>1</mark> 6	178	170	=	12	10	0	30	-
	2	0	255	220		4	32	233	<mark>11</mark> 2		2	32	22	108	



### Accuracy

Random: 10% Our image-diff (with L1): 38.6% Our image-diff (with L2): 35.4%

- 456

### K-Nearest Neighbors: Generalizing the Image-Diff Classifier



Tuning (hyper)parameters:





Massachusetts Institute of Technology

References: [89]

Lex Fridman January lex.mit.edu 2018

### K-Nearest Neighbors: Generalizing the Image-Diff Classifier



. . .



### Accuracy

Random: **10%** Training and testing on the same data: **35.4%** 7-Nearest Neighbors: **~30%** Human: **~95%** 

Convolutional Neural Networks: ~97.75%

## *Reminder:* Weighing the Evidence



lassachusetts

Institute of

[echnology



## Reminder: "Learning" is Optimization of a Function



Ground truth for "6": $y(x) = (0,0,0,0,0,0,0,0,0,0,0)^T$ 

"Loss" function:

$$C(w,b)\equiv rac{1}{2n}\sum_x \|y(x)-a\|^2$$



## Classify and Image of a Number

**Input:** (28x28)



January

2018

## **Convolutional Neural Networks**

Regular neural network (fully connected):



Convolutional neural network:



Each layer takes a 3d volume, produces 3d volume with some smooth function that may or may not have parameters.



## **Convolutional Neural Networks: Layers**

- **INPUT** [32x32x3] will hold the raw pixel values of the image, in this case an image of width 32, height 32, and with three color channels R,G,B.
- **CONV** layer will compute the output of neurons that are connected to local regions in the input, each computing a dot product between their weights and a small region they are connected to in the input volume. This may result in volume such as [32x32x12] if we decided to use 12 filters.
- **RELU** layer will apply an elementwise activation function, such as the *max(0,x)* thresholding at zero. This leaves the size of the volume unchanged ([32x32x12]).
- **POOL** layer will perform a downsampling operation along the spatial dimensions (width, height), resulting in volume such as [16x16x12].
- FC (i.e. fully-connected) layer will compute the class scores, resulting in volume of size [1x1x10], where each of the 10 numbers correspond to a class score, such as among the 10 categories of CIFAR-10. As with ordinary Neural Networks and as the name implies, each neuron in this layer will be connected to all the numbers in the previous volume.



Layers **highlighted in blue** have learnable parameters.

January 2018

## Dealing with Images: Local Connectivity



Same neuron. Just more focused (narrow "receptive field").

### The parameters on a each filter are spatially "shared" (if a feature is useful in one place, it's useful elsewhere)



## ConvNets: Spatial Arrangement of Output Volume



- Depth: number of filters
- Stride: filter step size (when we "slide" it)
- **Padding:** zero-pad the input



Technology

Filte	er W	'1 (3x3x3)
w1 [	:,	:,0]
-1	0	0
1	-1	-1
0	0	-1
w1[	:,	:,1]
-1	0	1
1	-1	1
-1	0	1
w1[	:,	:,2]
-1	-1	-1
-1	1	-1
0	1	-1

Bias b1 (1x1x1)

b1[:,:,0]

0

Out	put V	/olu	me (3	x3x2)
0[:	, : ,	0]		
-3	-1	4		
-2	-7	-4		
1	-1	1		
0[:	,:,	1]		
-7	3	1		
-7	-11	-1		
-4	-2	-4		

MIT 6.S094: Deep Learning for Self-Driving Cars

https://selfdrivingcars.mit.edu

Lex Fridman lex.mit.edu



Filte	Filter W1 (3x3x3)						
w1[	:,:	:,0]					
-1	0	0					
1	-1	-1					
0	0	-1					
w1 [	:,:	:,1]					
-1	0	1					
1	-1	1					
-1	0	1					
w1[	:,:	,2]					
-1	-1	-1					
-1	1	-1					
0	1	-1					

Output	Volume	(3x3x2)
--------	--------	---------

o(:,:,0]

1	-1	1
0[:	,:,	1]
-7	3	1
-7	-11	-1
-4	-2	-4

-1 4

-7 -4

-3

-2

Bias b1 (1x1x1) b1[:,:,0] 0



References: [95]







Output Volume (3x3x2) • [:,:,0] -3 -1 4 -2 -7 -4 1 -1 1 • [:,:,1] -7 3 1 -7 -11 -1 -4 -2 -4

MIT 6.S094: Deep Learning for Self-Driving Cars

https://selfdrivingcars.mit.edu

References: [95]

Massachusetts

Institute of

Technology



Out	put V	/olum	e (3x3x2)
0[:	,:,	0]	
-3	-1	4	
-2	-7	-4	
1	-1	1	
0[:	,:,	1]	
-7	3	1	
-7	-11	-1	
-4	-2	-4	

MIT 6.S094: Deep Learning for Self-Driving Cars

https://selfdrivingcars.mit.edu

Massachusetts References: [95]

Institute of

Technology



Output Volume (3x3x2) • [:,:,0] -3 -1 4 -2 -7 -4 1 -1 1 • [:,:,1] -7 3 1 -7 -11 -1

MIT 6.S094: Deep Learning for Self-Driving Cars

https://selfdrivingcars.mit.edu

-4

References: [95]

Massachusetts

Institute of

Technology



Out	put V	/olu	me (3x3x2)
0[:	,:,	0]	
-3	-1	4	
-2	-7	-4	
1	-1	1	
o[:	,:,	1]	
-7	3	1	
-7	-11	-1	
-4	-2	-4	

MIT 6.S094: Deep Learning for Self-Driving Cars

https://selfdrivingcars.mit.edu

Massachusetts Institute of Technology

References: [95]



Output Volume (3x3x2) • [:,:,0] -3 -1 4 -2 -7 -4 1 -1 1 • [:,:,1] -7 3 1 -7 -11 -1 -4 -2 -4

MIT 6.S094: Deep Learning for Self-Driving Cars

https://selfdrivingcars.mit.edu


Output Volume (3x3x2) 0[:,:,0] -3 -1 4 -2 -7 -4 1 -1 1 0[:,:,1] -7 3 1 -7 -11 -1 -4 -2 -4

MIT 6.S094: Deep Learning for Self-Driving Cars

https://selfdrivingcars.mit.edu

Institute of

Technology

### Convolution

100	1000	0
les		7
		ŧ.
		3.

Operation	Filter	Convolved Image
Identity	$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$	
Edge detection	$\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix}$	
	$\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$	
	$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$	



References: [124]

#### Convolution



Massachusetts Institute of Technology References: [124]

#### **Convolution:** Representation Learning



References: [124]

n January 2018

### **ConvNets:** Pooling

Single depth slice



У

max pool with 2x2 filters and stride 2

6	8
3	4



Institute of

Technology

January

#### Same Architecture, Many Applications



This part might look different for:

- Different image classification domains
- Image captioning with recurrent neural networks
- Image object localization with bounding box
- Image segmentation with fully convolutional networks
- Image segmentation with deconvolution layers



#### Object Recognition Case Study: ImageNet







#### What is ImageNet?

- ImageNet: dataset of 14+ million images (21,841 categories)
- Let's take the high level category of **fruit** as an example:
  - Total 188,000 images of fruit
  - There are 1206 Granny Smith apples:







#### What is ImageNet?



# **ILSVRC** Challenge Evaluation for Classification

- Top 5 error rate:
  - You get 5 guesses to get the correct label



- ~20% reduction in accuracy for Top 1 vs Top 5
- Human annotation is a binary task: "apple" or "not apple"

echnology



- Human error: 5.1%
  - Surpassed in 2015

- AlexNet (2012): First CNN (15.4%)
  - 8 layers
  - 61 million parameters
- ZFNet (2013): 15.4% to 11.2%
  - 8 layers
  - More filters. Denser stride.
- VGGNet (2014): 11.2% to 7.3%
  - Beautifully uniform: 3x3 conv, stride 1, pad 1, 2x2 max pool
  - 16 layers
  - 138 million parameters
- GoogLeNet (2014): 11.2% to 6.7%
  - Inception modules
  - 22 layers
  - 5 million parameters (throw away fully connected layers)
- ResNet (2015): 6.7% to 3.57%
  - More layers = better performance
  - 152 layers
- CUImage (2016): 3.57% to 2.99%
  - Ensemble of 6 models
- SENet (2017): 2.99% to 2.251%
  - Squeeze and excitation block: network is allowed to adaptively adjust the weighting of each feature map in the convolutional block.

lassachusetts

Institute of

Technology



#### ImageNet Classification Error (Top 5)

- AlexNet (2012): First CNN (15.4%)
  - 8 layers
  - 61 million parameters
- ZFNet (2013): 15.4% to 11.2%
  - 8 layers
  - More filters. Denser stride.

#### • VGGNet (2014): 11.2% to 7.3%

- Beautifully uniform: 3x3 conv, stride 1, pad 1, 2x2 max pool
- 16 layers
- 138 million parameters
- GoogLeNet (2014): 11.2% to 6.7%
  - Inception modules
  - 22 layers
  - 5 million parameters (throw away fully connected layers)
- ResNet (2015): 6.7% to 3.57%
  - More layers = better performance
  - 152 layers
- CUImage (2016): 3.57% to 2.99%
  - Ensemble of 6 models

in January J 2018





- AlexNet (2012): First CNN (15.4%)
  - 8 layers
  - 61 million parameters
- ZFNet (2013): 15.4% to 11.2%
  - 8 layers
  - More filters. Denser stride.
- VGGNet (2014): 11.2% to 7.3%
  - Beautifully uniform: 3x3 conv, stride 1, pad 1, 2x2 max pool
  - 16 layers
  - 138 million parameters
- GoogLeNet (2014): 11.2% to 6.7%
  - Inception modules
  - 22 layers
  - 5 million parameters (throw away fully connected layers)
- ResNet (2015): 6.7% to 3.57%
  - More layers = better performance
  - 152 layers
- CUImage (2016): 3.57% to 2.99%
  - Ensemble of 6 models

Krizhevsky et al. "Imagenet classification with deep convolutional neural networks." Advances in neural information processing systems. 2012.

lassachusetts

Institute of

Technology

January 2018



Simonyan et al. "Very deep convolutional networks for large-scale image recognition." 2014.

lassachusetts

Institute of

Technology

- AlexNet (2012): First CNN (15.4%)
  - 8 layers
  - 61 million parameters
- ZFNet (2013): 15.4% to 11.2%
  - 8 layers
  - More filters. Denser stride.
- VGGNet (2014): 11.2% to 7.3%
  - Beautifully uniform: 3x3 conv, stride 1, pad 1, 2x2 max pool
  - 16 layers
  - 138 million parameters
- GoogLeNet (2014): 11.2% to 6.7%
  - Inception modules
  - 22 layers
  - 5 million parameters (throw away fully connected layers)
- ResNet (2015): 6.7% to 3.57%
  - More layers = better performance
  - 152 layers
- CUImage (2016): 3.57% to 2.99%
  - Ensemble of 6 models

References: [128]



Szegedy et al. "Going deeper with convolutions." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2015. AlexNet (2012): First CNN (15.4%)

- 8 layers •
- 61 million parameters •
- ZFNet (2013): 15.4% to 11.2%
  - 8 layers ٠
  - More filters. Denser stride.

#### VGGNet (2014): 11.2% to 7.3% ٠

- Beautifully uniform: ٠ 3x3 conv, stride 1, pad 1, 2x2 max pool
- 16 layers •
- 138 million parameters
- GoogLeNet (2014): 11.2% to 6.7% •
  - Inception modules •
  - 22 layers ٠
  - 5 million parameters (throw away fully connected layers)
- ResNet (2015): 6.7% to 3.57% ٠
  - More layers = better performance
  - 152 layers
- CUImage (2016): 3.57% to 2.99% ٠
  - Ensemble of 6 models

assachusetts Institute of Technology

# Inception Module



- Process: do different size convolutions, and concatenate
- Convolution sizes: ullet
  - Smaller convolutions: local features
  - Larger convolutions: high-abstracted features
- **Result**: Fewer parameters and better performance



He, Kaiming, et al. "Deep residual learning for image recognition." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2016.

• AlexNet (2012): First CNN (15.4%)

- 8 layers
- 61 million parameters
- ZFNet (2013): 15.4% to 11.2%
  - 8 layers
  - More filters. Denser stride.
- VGGNet (2014): 11.2% to 7.3%
  - Beautifully uniform: 3x3 conv, stride 1, pad 1, 2x2 max pool
  - 16 layers
  - 138 million parameters
- GoogLeNet (2014): 11.2% to 6.7%
  - Inception modules
  - 22 layers
  - 5 million parameters (throw away fully connected layers)
- ResNet (2015): 6.7% to 3.57%
  - More layers = better performance
  - 152 layers
- CUImage (2016): 3.57% to 2.99%
  - Ensemble of 6 models

Massachusetts Institute of Technology References: [130]



#### Initial Observation:

- Network depth often increases representation power, but is harder to train.
- Residual Block:
  - Repeat a simple network block (think: RNN)
  - Pass input along without transformation: help ensure that each layer learns something new





# SENet: Squeeze-and-Excitation Networks



- **Content-aware channel weighting:** Add parameters to each channel of a convolutional block so that the network can adaptively adjust the weighting of each feature map
- This approach is simple and can be added to any model
  - **Takeaway for thought:** Parameterize everything (that's cost-effective) including higher-order hyper-parameters.



### Capsule Networks (Hinton)



- A CNN see both images as the same. The problem:
  - Internal data representation of a convolutional neural network does not take into account important spatial hierarchies between simple and complex objects.
- See upcoming online-only lecture on capsule networks.



#### Same Architecture, Many Applications



This part might look different for:

- Different image classification domains
- Image captioning with recurrent neural networks
- Image object localization with bounding box
- Image segmentation with fully convolutional networks
- Image segmentation with deconvolution layers



#### **Object Detection**



#### **R-CNN:** Regions with CNN features





### **Fully Convolutional Networks**

- Goal: Classify every pixel in an image.
- Difficulty: Hard
- Why?
  - When precise boundaries of objects matter (medical, driving)
  - Useful for fusing with other sensors (LIDAR)





For the full updated list of references visit: https://selfdrivingcars.mit.edu/references Lex Fridman January lex.mit.edu 2018

## FCN (Nov 2014)

Paper: "Fully Convolutional Networks for Semantic Segmentation"

- **Repurpose Imagenet pretrained nets** •
- Upsample using deconvolution
- Skip connections to improve coarseness of upsampling



sachusetts istitute of chnology

For the full updated list of references visit: https://selfdrivingcars.mit.edu/references Lex Fridman January lex.mit.edu 2018

#### SegNet (Nov 2015)

Paper: "SegNet: A Deep Convolutional Encoder-Decoder Architecture for Image Segmentation"

Maxpooling indices transferred to decoder to improve the • segmentation resolution.





## Dilated Convolutions (Nov 2015)

Paper: "Multi-Scale Context Aggregation by Dilated Convolutions"

- Since pooling decreases resolution:
  - Added "dilated convolution layer"
- Still interpolate up from 1/8 of original image size





January 2018

# DeepLap v1, v2 (Jun 2016)

Paper: "DeepLab: Semantic Image Segmentation with Deep Convolutional Nets, Atrous Convolution, and Fully Connected CRFs"

- Added fully-connected Conditional Random Fields (CRFs) as a post-processing step
  - Smooth segmentation based on the underlying image intensities



# Key Aspects of Segmentation

- Fully convolutional networks (FCNs) replace fully-connected layers with convolutional layers
  - Deeper, updated models (now ResNet) consistent with ImageNet Challenge object classification tasks.
- **Conditional Random Fields (CRFs)** to capture both local and long-range dependencies within an image to refine the prediction map.
- Dilated convolution (aka Atrous convolution) maintain computational cost, increase resolution of intermediate feature maps



# ResNet-DUC (Nov 2017)

Paper: "Understanding Convolution for Semantic Segmentation"

- Dense upsampling convolution (DUC) instead of bilinear upsampling
  - Learnable: Learn the upscaling filters
- Hybrid dilated convolution (HDC)
  - Use a different dilation rate





# FlowNet (May 2015)

Paper: "FlowNet: Learning Optical Flow with Convolutional Networks "

- Learn flow from image-pair, end to end.
  - FlowNetS stacks two images as input
  - FlowNetC convolute separately, combine with correlation layer



Fig. 1





For the full updated list of references visit: <u>https://selfdrivingcars.mit.edu/references</u> [177]

January 2018

# FlowNet 2.0 (Dec 2016)

Paper: "FlowNet 2.0: Evolution of Optical Flow Estimation with Deep Networks"

- Stack FlowNetS and FlowNetC
- Improvement over FlowNet
  - Smooth flow fields •
  - Preserves fine-motion detail ٠
  - Runs at 8-140fps ٠

- **Observations:** 
  - Stacking networks as an approach
  - Order of training dataset matters





For the full updated list of references visit: https://selfdrivingcars.mit.edu/references



Lex Fridman January lex.mit.edu 2018

# Original Video

#### cars.mit.edu/segfuse

Massachusetts Institute of Technology

For the full updated list of references visit: https://selfdrivingcars.mit.edu/references MIT 6.S094: Deep Learning for Self-Driving Cars https://selfdrivingcars.mit.edu

Lex Fridman lex.mit.edu January



#### cars.mit.edu/segfuse

Massachusetts Institute of Technology

For the full updated list of references visit: https://selfdrivingcars.mit.edu/references MIT 6.S094: Deep Learning for Self-Driving Cars https://selfdrivingcars.mit.edu

Lex Fridman lex.mit.edu January

# Segmentation

#### cars.mit.edu/segfuse

Massachusetts Institute of Technology

For the full updated list of references visit: https://selfdrivingcars.mit.edu/references MIT 6.S094: Deep Learning for Self-Driving Cars https://selfdrivingcars.mit.edu Lex Fridman lex.mit.edu January



#### cars.mit.edu/segfuse

Massachusetts Institute of Technology

For the full updated list of references visit: https://selfdrivingcars.mit.edu/references MIT 6.S094: Deep Learning for Self-Driving Cars https://selfdrivingcars.mit.edu Lex Fridman lex.mit.edu January









#### cars.mit.edu/segfuse



For the full updated list of references visit: https://selfdrivingcars.mit.edu/references MIT 6.S094: Deep Learning for Self-Driving Cars https://selfdrivingcars.mit.edu Lex Fridman lex.mit.edu January
## Thank You

## Tomorrow: Waymo



## Next lecture: Deep Learning for Human Sensing



## **Upcoming online-only lectures:**

- Capsule networks
- Generative adversarial networks

