



comp150dl: Deep Learning for Computer Vision

Instructor: Genevieve Patterson







Ridiculously Brief History of Computer Vision





Artificial Intelligence Group Vision Memo. No. 100.

Seymour Papert

The summer vision project is an attempt to use our summer workers effectively in the construction of a significant part of a visual system. The particular task was chosen partly because it can be segmented into sub-problems which will allow individuals to work independently and yet participate in the construction of a system complex enough to be a real landmark in the development of "pattern recognition".

MASSACHUSETTS INSTITUTE OF TECHNOLOGY PROJECT MAC

July 7, 1966

THE SUMMER VISION PROJECT





Parts-and-shape models

- Model:
 - Object as a set of parts
 - Relative locations between parts
 - Appearance of part





Figure from [Fischler & Elschlager 73]



Pictorial structure model

Fischler and Elschlager(73), Felzenszwalb and Huttenlocher(00)





Eigenfaces (Turk & Pentland, 1991)







Local features for object instance recognition - SIFT





D. Lowe (1999, 2004) cor









Carefully Considered Features

 Histogram of Oriented Gradients







• Self-Similarity

HOG [1]

Inverse (Us)

Original













Canonical Challenges





Classification Challenge



is there a cat?



11

Detection Challenge







Segmentation Challenge







Training Pipeline





Library of Classifiers



Sliding window approaches



- Turk and Pentland, 1991 ullet
- Belhumeur, Hespanha, & Kriegman, 1997
- Schneiderman & Kanade 2004
- Viola and Jones, 2000 \bullet







Spatial pyramid representation

- Extension of a bag of features \bullet
- \bullet





Locally orderless representation at several levels of resolution



Lazebnik, Schmid & Ponce (CVPR 2006)



Caltech101 dataset



Multi-class classification results (30 training images per class)

	Weak feat	ures (16)	Strong features (200)		
Level	Single-level	Pyramid	Single-level	Pyramid	
0	15.5 ± 0.9		41.2 ± 1.2		
1	31.4 ± 1.2	32.8 ± 1.3	55.9 ± 0.9	57.0 ± 0.8	
2	47.2 ± 1.1	49.3 ± 1.4	63.6 ± 0.9	64.6 ±0.8	
3	52.2 ± 0.8	54.0 ± 1.1	60.3 ± 0.9	$64.6\pm\!0.7$	



17

- TV/monitor)
- Three challenges: image?)
 - every X)
 - Segmentation challenge



The PASCAL Visual Object Classes Challenge 2009 (VOC2009) • Twenty object categories (aeroplane to

- Classification challenge (is there an X in this

- Detection challenge (draw a box around



Slides from Noah Snavely



Discriminatively trained part-based



















P. Felzenszwalb, R. Girshick, D. McAllester, D. Ramanan, **"Object Detection with Discriminatively Trained** Part-Based Models," PAMI 2009









Why Deep Networks?





Global (End-to-End) Learning: Energy-Based Models.



Making every single module in the system trainable.

Y LeCun

Every module is trained simultaneously so as to optimize a global loss function.

Includes the feature extractor, the recognizer, and the contextual post-processor (graphical model)

Problem: back-propagating gradients through the graphical model.

Components of End-to-End Learning





Texture representations vs CNNs



Subhransu Maji (UMass Amherst)



Logistic Function



$$ah \phi(v_i) = anhigl(eta_1 + eta_0 \sum_j v_{i,j} x_jigr)$$

ReLU

$$f(x) = egin{cases} 0 & ext{for} & x < 0 \ x & ext{for} & x \ge 0 \end{cases}$$



Feedforward Neural Network





24

Neurons



Very bad (slow to train)

Perceptron

Very good (quick to train)

Rectified Linear Unit (ReLU)

Figure from Karpathy 2015



Convolutional Network







Filter Bank +non-linearity

Pooling

Filter Bank +non-linearity

Pooling

Filter Bank +non-linearity





Filter W1 (3x3x3					
wil		,0	J		
-1	1	-1			
1	-1	1			
-1	0	1			
w1[:,:	,1]		
1	1	0			
-1	1	1			
-1	0	1			
w1[:,:,2]					
0	0	1			
-1	1	-1			
-1	0	-1			

Output Volume $(3x3x2)$ o[:,:,0]					
1	-5	1			
0	0	5			
1	0	2			
0[:,:,1]					
4	-3	1			
5	7	0			
6	4	0			

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

How Convolution Works

Figure from Karpathy 2015



Image Filters





The Dot Product

- Also called 'scalar product'
- Sum of the product of each element of two sequences
- $(1,2,3) \cdot (3,4,5) = 1*3+2*4+3*5 = 24$



- b = (3,1)
- a b = 5



• Dot product is the length of **a** projected on **b**







Original











Original





Filtered (no change)







Original











Original





Shifted left By 1 pixel





Box Filter

What does it do?

- Replaces each pixel with an average of its neighborhood
- Achieve smoothing effect (remove sharp features)





Slide credit: David Lowe (UBC)





Smoothing with box filter





Image filtering



0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0



Credit: S. Seitz




0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0







0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0







0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0







0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

		-								
0	0									
0	0		0	10	20	30	30			
0	0									
0	0									
0	0									
0	0									
0	0									
0	0									
0	0									
0	0		C							
COL	np15	50dl		S					Cr	redit.





0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	00	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

0	0									
0	0		0	10	20	30	30			
0	0									
0	0									
0										
0	0									
0	0					?				
0	0									
0	0									
0	0									
0	0									
CO	mp15	50dl		S					Cr	redit•





0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

0	0									
0	0		0	10	20	30	30			
0	0									
0	0							?		
0	0									
0	0									
0	0					50				
0	0									
0	0									
0	0		C							
CO	mp15	50dl		S					Cr	redit•





0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

comp150dl





	0	10	20	30	30	30	20	10	
	0	20	40	60	60	60	40	20	
	0	30	60	90	90	90	60	30	
	0	30	50	80	80	90	60	30	
	0	30	50	80	80	90	60	30	
	0	20	30	50	50	60	40	20	
	10	20	30	30	30	30	20	10	
	10	10	10	0	0	0	0	0	
8	6								



What else is possible with Filters?

- Really important for photo editing!
 - Enhance images
 - Denoise, resize, increase contrast, etc.
 - Extract information from images
 - Texture, edges, distinctive points, etc.
 - - Template matching

Detect patterns —-> Convolutional Networks!!





Practice with linear filters



0	0	0
0	2	0
0	0	0

(Note that filter sums to 1)

Original







Source: D. Lowe



Practice with linear filters



0	0	0
0	2	0
0	0	0

Original







Sharpening filter

- Accentuates differences with local average

Source: D. Lowe



Sharpening



before



after



Source: D. Lowe





Noise

Gaussian filter







Median filters

- selecting the median intensity in the window.
- over a mean filter?

• A Median Filter operates over a window by

• What advantage does a median filter have



Slide by Steve Seitz



Comparison: salt and pepper noise



Gaussian

Median





Slide by Steve Seitz



Other filters







Sobel



Vertical Edge (absolute value)





Other filters





Sobel



Horizontal Edge (absolute value)





Key properties of linear filters

Linearity: $filter(f_1 + f_2) = filter(f_1) + filter(f_2)$

Shift invariance: same behavior regardless of pixel location filter(shift(f)) = shift(filter(f))



Source: S. Lazebnik



More properties

- Commutative: *a* * *b* = *b* * *a*
 - Conceptually no difference between filter and signal
- Associative: a * (b * c) = (a * b) * c

 - This is equivalent to applying one filter: $a * (b_1 * b_2 * b_3)$
- Distributes over addition:
- Scalars factor out: ka * b =
- Identity: unit impulse *e* = [0, 0, 1, 0, 0], $a^* e = a$

• But particular filtering implementations might break this equality

Often apply several filters one after another: (((a * b₁) * b₂) * b₃)

$$a^{*}(b + c) = (a^{*}b) + (a^{*}c)$$

$$= a * kb = k (a * b)$$



Source: S. Lazebnik



Practical matters

- What about near the edge?
 - the filter window falls off the edge of the image
 - need to extrapolate
 - methods:
 - clip filter (black)
 - wrap around
 - copy edge
 - reflect across edge





Source: S. Marschner



Take-home messages about filters





Be aware of details for filter size, extrapolation, cropping



• Linear filtering is sum of dot product at each position - Can smooth, sharpen, translate (among many other uses)







Kernels: Layer 1 (11x11)

Layer 1: 3x96 kernels, RGB->96 feature maps, 11x11 Kernels, stride 4





Current Computer Vision









Fig. 8. More results using our multiscale convolutional network and a flat CRF on the Stanford Background Dataset.

Learning Hierarchical Features for Scene Labeling

Clement Farabet, Camille Couprie, Laurent Najman, Yann LeCun [PAMI '13]



This time though, the reviewers were particularly clueless, or negatively biased, or both. I was very sure that this paper was going to get good reviews because: 1) it has two simple and generally applicable ideas for segmentation ("purity tree" and "optimal cover"); 2) it uses no hand-crafted features (it's all learned all the way through. Incredibly, this was seen as a negative point by the reviewers!); 3) it beats all published results on 3 standard datasets for scene parsing; 4) it's an order of magnitude faster than the competing methods.

If that is not enough to get good reviews, just don't know what is.

So, I'm giving up on submitting to computer vision conferences altogether.







AlexNet Architecture - 7 hidden weight layers



3 Fully connected layers

The output of the last fully-connected layer is fed to a 1000-way softmax which produces a distribution over the 1000 class labels The ReLU non-linearity is applied to the output of every convolutional and fully-connected layer.

Detection

Fast R-CNN

- convolve once
- project + detect



Faster R-CNN

- end-to-end proposals and detection
- 200 ms / image inference
- fully convolutional Region Proposal Net + Fast R-CNN

<u>arXiv</u> and <u>code</u> for Fast R-CNN

Ross Girshick, Shaoqing Ren, Kaiming He, Jian Sun



Pixelwise Prediction

Fully convolutional networks for pixel prediction applied to semantic segmentation

- end-to-end learning
- efficient inference and learning 150 ms per-image prediction
- multi-modal, multi-task

Further applications

- depth
- boundaries
- flow + more

CVPR15 <u>arXiv</u> and <u>pre-release</u>





Jon Long* & Evan Shelhamer*, **Trevor Darrell**





Sequences

Recurrent Net and Long Short Term Memory LSTM are sequential models

- video
- language
- dynamics

learned by backpropagation through time.

LRCN: Long-term Recurrent Convolutional Network

- activity recognition
- image captioning
- video captioning

CVPR15 <u>arXiv</u> and <u>project site</u>





A group of young men playing a game of soccer.

Jeff Donahue et al.







Pre-trained Models

Lots of Data



image by Andrej Karpathy



comp150dl **Tufts**

© kaggle.com







IM GENET Large Scale Visual Recognition Challenge (ILSVRC) 2015

Large-scale recognition



Participation in ILSVRC over the years



Year





IM GENET Large Scale Visual Recognition Challenge



Year 2012 SuperVision



[Krizhevsky NIPS 2012]



<u>Year 2015</u>







Our results on COCO – too many objects, let's check carefully!



Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. "Deep Residual Learning for Image Recognition". arXiv 2015. Shaoqing Ren, Kaiming He, Ross Girshick, & Jian Sun. "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks". NIPS 2015.

Microsoft Research

*the original image is from the COCO dataset





Visual Turing test for computer vision systems

Donald Geman^a, Stuart Geman^{b,1}, Neil Hallonquist^a, and Laurent Younes^a

3618-3623 PNAS March 24, 2015 vol. 112 no. 12



	1. Q: Is there a person in the blue region?	A: yes
	2. Q: Is there a unique person in the blue region?	A: yes
	(Label this person 1)	
	3. Q: Is person 1 carrying something?	A: yes
1	4. Q: Is person 1 female?	A: yes
l.	5. Q: Is person 1 walking on a sidewalk?	A: yes
	6. Q: Is person 1 interacting with any other object?	A: no
	:	
1	9. Q: Is there a unique vehicle in the yellow region?	A: yes
	(Label this vehicle 1)	
l	10. Q: Is vehicle 1 light-colored?	A: yes
۱	11. Q: Is vehicle 1 moving?	A: no
ł,	12. Q: Is vehicle 1 parked and a car?	A: yes
2	1	
	14. Q: Does vehicle 1 have exactly one visible tire?	A: no
1	15. Q: Is vehicle 1 interacting with any other object?	A: no
1	17. Q: Is there a unique person in the red region?	A: no
1	18. Q: Is there a unique person that is female in the red region?	A: no
	19. Q: Is there a person that is standing still in the red region?	A: yes
	20. Q: Is there a unique person standing still in the red region?	A: yes
	(Label this person 2)	
	1	
	23. Q: Is person 2 interacting with any other object?	A: yes
	24. Q: Is person 1 taller than person 2?	A: amb
	25. Q: Is person 1 closer (to the camera) than person 2?	A: no
	26. Q: Is there a person in the red region?	A: yes
	27. Q: Is there a unique person in the red region?	A: yes
	(Label this person 3)	
	!	
	36. Q: Is there an interaction between person 2 and person 3?	A: yes
	37. Q: Are person 2 and person 3 talking?	A: yes







Deep Learning IRL





Self-Driving Cars





DeepTesla

72
Product Search



















Auto-tagging







Clarifai

74





Medical Research



PathAl

75

Image Generation



Generated results



Scribbler





What is this class about?





Course Description

- Learned Representations \bullet
- Object Proposals
- CNN detection and segmentation Co-attention models
- Weakly Supervised and Unsupervised CNNs
- Ensemble methods Recurrent Neural Nets and Long-Short Term Memory Networks
- Generative Networks

• Siamese / Ranking / Triplet Networks

Residual Nets

Reinforcement Learning





Preparation

- Programming Experience
 - Python, Matlab
- Math
- Machine Learning
- Computer Vision

• Linear Algebra, Basic Calculus, Probability





Step 1: Datasets





flick from YAHOO!	
Home You - Organize & C	reate
Search	Phot
Everyone's Uploads] indig

Sort: Relevant Recent Interesting

From Steve...





From OwimanSA





From hart_curt



From Buzzie82



From Christian.



From tomelizab ...









From Dave 2x

















From Birds&.



From Bird Man...







From dwaynejava







From Jim Adams...



From dwaynejava



From Jim Adams...



From Bird Man...



From owisblood





View: Small Medium Detail



From KirkH1









From dmarshman



From DansPhotoArt





From iceberg_c...



From tanagergirl



From MoGov



From Dave 2x









From tomelizab.

From Dan and,





From MomOnTheR ...













6000 images from flickr.com









Building datasets





Is there an Indigo bunting in the image?

Annotators

amazonmechanical turk Artificial Artificial Intelligence

100s of training images









Slide credit: Welinder et al





✓ Instance segmentation✓ Non-iconic Images



- 330,000 images
- >2 million instances (700k people)
- Every instance is segmented
- 7.7 instances per image (3.5 categories)



0k people) nted (3.5 categories

Beyond detection

✓ Sentences

two giraffe standing next to each other in front of a wooden fence. two giraffes standing in the dirt near a gate. two giraffes stand by a food box awaiting the goods. two giraffes are standing next to a wooden fence. two giraffes standing alone by a picket fence.



Collecting Image Annotations Using Amazon's Mechanical Turk, C. Rashtchian, P. Young, M. Hodosh, J. Hockenmaier, NAACL HLT Workshop on Creating Speech and Language Data with Amazon's Mechanical Turk, 2010

Beyond detection

✓ Keypoints (provided by Facebook)





MS COCO Challenges at ICCV 2015





DetectionSegmentation









Evaluation Metrics

Average Precision	(AP) :				
AP					
ΔpIOU=.50		90	AP	at	IoU
ΔDIOU=.75		00	AP	at	IoU
AL		00	AP	at	IoU

Challenges Score: AP

- AP is averaged over multiple loU values between 0.5 and 0.95 (and categories, size).
- More comprehensive metric than the traditional AP at a fixed IoU value (0.5 for Pascal).

i=.50:.05:.95 (determines challenge winner)
i=.50 (PASCAL VOC metric)
i=.75 (strict metric)







Evaluation Metrics

AP	Across	Scales:				
	AP ^{small}		Q	ΛD	for	<u> </u>
	⊿ pmedium		Ō	AP	TOT	SIII
	nlarge		00	AP	for	me
	AP-aryo		olo	AP	for	la

Other Scores: Size AP

large (A > 96 x 96) instances of objects.

32x32 < A < 96x96



< 32x32



all objects: area < 32² dium objects: 32^2 < area < 96^2 rge objects: area > 96²

• AP is averaged over small ($A < 32 \times 32$), medium ($32 \times 32 < A < 96 \times 96$) and

>96x96







Evaluation Metrics

Average Recall (AR):				
AR ^{max=1}	olo	AR	give	en
AR ^{max=10}	00	AR	given	
AR ^{max=100}	00	AR	₹ given	
AR Across Scales:				
AR ^{small}	00	AR	for	sm
AR ^{medium}	00	AR	for	me
AR ^{large}	olo	AR	for	la

Other Scores: AR

- allowed in the image of 1, 10, 100.
- large $(A > 96 \times 96)$ instances of objects.

```
1 detection per image
10 detections per image
100 detections per image
```

```
all objects: area < 32<sup>2</sup>
dium objects: 32^2 < area < 96^2
rge objects: area > 96<sup>2</sup>
```

```
• Measures the maximum recall over a fixed number of detections
```

```
• AR is averaged over small (A < 32 \times 32), medium (32 \times 32 < A < 96 \times 96) and
```



Detection Leaderboard (II)



Object Localization can improve





BBox detections and IoU



Also hard for humans

IoU = 0.5



loU = 0.7

IoU = 0.95





BBox detections and IoU



Also hard for humans

IoU = 0.5



IoU = 0.75



IoU = 0.95





Performance Breakdown (I)



COCO AP varies across supercategories and size





Bounding Box Detection Errors (I)

What type of errors are algorithms making?





Super-category FP removed



Background FP removed

Category FP removed



FN errors are removed

MSRA













0.2

0.4

0.6

recall

0.8

0

0

Bounding Box Detection Errors (II)

Super-category FP removed





Background FP removed

FN errors are removed

indoor-book-medium







Some success cases ...



Results from FAIRCNN team.







Results from FAIRCNN team.

... and some failures

