Lecture 18: Introduction to Recognition

COMP 590/776: Computer Vision Instructor: Soumyadip (Roni) Sengupta TA: Mykhailo (Misha) Shvets



Course Website: Scan Me!

Image classification demo



https://cloud.google.com/vision/docs/drag-and-drop

See also:

https://aws.amazon.com/rekognition/

https://www.clarifai.com/

https://azure.microsoft.com/en-us/services/cognitive-services/computer-vision/

Next few slides adapted from Li, Fergus, & Torralba's excellent <u>short course</u> on category and object recognition



• Verification: is that a lamp?



• Verification: is that a lamp?

• Detection: where are the people?



- Verification: is that a lamp?
- Detection: where are the people?
- Identification: is that Potala Palace?



- Verification: is that a lamp?
- Detection: where are the people?
- Identification: is that Potala Palace?
- Object categorization



- Verification: is that a lamp?
- Detection: where are the people?
- Identification: is that Potala Palace?
- Object categorization
- Scene and context categorization



- Verification: is that a lamp?
- Detection: where are the people?
- Identification: is that Potala Palace?
- Object categorization
- Scene and context categorization
- Activity / Event Recognition



Recognition: What type of output?

Image classification



Semantic segmentation

Object detection



Instance segmentation





• And beyond: depth/3D structure prediction, image description, etc.

Object recognition: Is it really so hard?

Find the chair in this image



Output of normalized correlation



This is a chair



Object recognition: Is it really so hard?



Find the chair in this image





Pretty much garbage: Simple template matching is not going to do the trick

Object recognition: Is it really so hard?



Find the chair in this image



A "popular method is that of template matching, by point to point correlation of a model pattern with the image pattern. These techniques are inadequate for three-dimensional scene analysis for many reasons, such as occlusion, changes in viewing angle, and articulation of parts." Nivatia & Binford, 1977.

Why not use SIFT matching for everything?

• Works well for object *instances* (or distinctive images such as logos)



• Not great for generic object categories



And it can get a lot harder



Brady, M. J., & Kersten, D. (2003). Bootstrapped learning of novel objects. J Vis, 3(6), 413-422



Svetlana Lazebnik



Variation Makes Recognition Hard

 The same class of object can appear very differently in different images





The Semantic Gap



PTOTOTIT TTOTOTOTO(J0J0J0JJJJ0JJJ0JJ(PICOIIICOI I CI I(POT TTTOT TOOTOOTO] PODIJIJ TOODIJOO]]07007007007770007] 70000 J000JJJ0JJJ0(0770000007777][٦.) JOJJOOJ OJOOJJ J(<u>, U J' J</u> JTTT TOTOTO TOTTT]]00070707007070770] 7007 7777000070 77(JOTTTTOOOOTTTT TTO]

What we see

What the computer sees

Image Classifiers in a Nutshell

- Input: an image
- Output: the class label for that image
- Label is generally one or more of the discrete labels used in training
 - e.g. {cat, dog, cow, toaster, apple, tomato, truck, ... }

def classifier(image):
//Do some stuff
return class_label;



The Problem is Under-constrained

- Distinct realities can produce the same image...
- We generally can't compute the "right" answer, but we can compute the most likely one...
- We need some kind of prior to condition on. We can learn this prior from data:

$$f(x) = \underset{\ell_x}{\operatorname{argmax}} P(\ell_x | data)$$





What Matters in Recognition?

• Data

- More is always better (as long as it is good data)
- Annotation is the hard part
- Representation
 - Low level: SIFT, HoG, GIST, edges
 - Mid level: Bag of words, sliding window, deformable model
 - High level: Contextual dependence
 - Deep learned features
- Learning Techniques
 - E.g. choice of classifier or inference method

What Matters in Recognition?

• Data

- More is always better (as long as it is good data)
- Annotation is the hard part

Representation

- Low level: SIFT, HoG, GIST, edges
- Mid level: Bag of words, sliding window, deformable model
- High level: Contextual dependence
- Deep learned features
- Learning Techniques
 - E.g. choice of classifier or inference method

24 Hrs in Photos

Flickr Photos From 1 Day in 2011



https://www.kesselskramer.com/project/24-hrs-in-photos/

Data Sets

- PASCAL VOC
 - *Not* Crowdsourced, bounding boxes, 20 categories
- ImageNet
 - Huge, Crowdsourced, Hierarchical, *Iconic* objects
- SUN Scene Database, Places
 - Not Crowdsourced, 397 (or 720) scene categories
- LabelMe (Overlaps with SUN)
 - Sort of Crowdsourced, Segmentations, Open ended
- SUN Attribute database (Overlaps with SUN)
 - Crowdsourced, 102 attributes for every scene
- OpenSurfaces
 - Crowdsourced, materials
- Microsoft COCO
 - Crowdsourced, large-scale objects

The PASCAL Visual Object Classes Challenge 2009 (VOC2009)

- 20 object categories (aeroplane to TV/monitor)
- Three challenges:
 - Classification challenge (is there an X in this image?)
 - Detection challenge (draw a box around every X)
 - Segmentation challenge (which class is each pixel?)



Large Scale Visual Recognition Challenge (ILSVRC) IM GENET

20 object classes22,591 images1000 object classes1,431,167 images

Image: Constraint of the second s

http://image-net.org/challenges/LSVRC/{2010,2011,2012}

Variety of object classes in ILSVRC

flamingo







bottle



car

ILSVRC



cock





quail



partridge



ruffed grouse



beer bottle wine bottle water bottle pop bottle ... pill bottle



race car wagon







cab

. .

bottles

cars

birds

Variety of object classes in ILSVRC



What Matters in Recognition?

• Data

- More is always better (as long as it is good data)
- Annotation is the hard part
- Representation
 - Low level: SIFT, HoG, GIST, edges
 - Mid level: Bag of words, sliding window, deformable model
 - High level: Contextual dependence
 - Deep learned features
- Learning Techniques
 - E.g. choice of classifier or inference method

Recall: Origins of computer vision







(b) Differentiated picture.



(a) Original picture.

(c) Line drawing

(d) Rotated view.



PICTURE

Pattern Classification and Scene Analysis Richard O. Duda and Peter E. Hart



Hough, 1959



Rosenfeld, 1969

Duda & Hart, 1972

History of recognition: Geometric alignment





Perkins (1978)



Grimson & Lozano-Perez (1984)



Lowe (1985)











Huttenlocher & Ullman (1987)

Ayache & Faugeras (1986)

(b)

History of recognition: Hierarchies of parts



History of recognition: Deformable templates

0.000	12345+7890123455799012345573901234557890
1.	00433334 0042 FEERFF PFFFFFFF44999348498
2	
3	B#####################################
4	**************************************
5	***************************************
5	833###################################
7	SERVICES CONSERBORNAGE BAREFFEELDEDTAT
6	EELETETTOBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
9	Eters # + + + + + + + + + + + + + + + + + +
10	EFFEFHHIBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBBB
11	
12	**************************************
13	
14	\$\$\$ { { { { { { { { { { { { { { { { { {
15	EF#FEE+#888##88x) = = +1 *88881244434Hd
16	••••••••••••••••••••••••••••••••••••••
17	> 응유유유류의단상행수소 ™21+ = 1 A 명월방 7 A A 자귀하는 색
18	FRH = 254 8 - 1 X Z 1 + - 1 A 8 8 X X X A 4989 A
19	€€968MM28#489#A)==1)))2 x821 XXXAAM8MA
20	<pre>###### AA#############################</pre>
21	BHEMAAXX FAMER HE +ZMMER XXZ 1ZZ XXAAMAA
22	EPPMAAX A# AXX ZZXA=ZZ1) 112 XX1) Z XXAAMAA
23	66644XX464X1))22-1+ -+)241+)2XXXX444X
24	05************************************
25	0444 AXAAAAXX1) X+)1- +)XA+))111112222
26	€2*MAXXAAMX2)) =x 1x 1+ =+1AX++++)) 111112
27	MMAAXAXAXXXXZ2XAXXXZ1)+ZMZ++++)))))1111
28	AAAXXXXXXAAXA* 4XXXZZ })X91++++)))11111
29	AxxxxxxXZAExxtenquex(1A=1+++))))11111
30	XXX222221X#MXAAXXX2)2A±5?X)))11112222
31	XXX2222222 MBEA>2X2112A000 AX1272222XXX
32	XXXZZZX#+#####AZ1)1ZA#MA#A###AXXXXXXXAA
33	XXX7ZABB-PESSEAX XAGEYXAB78895474964A4
34	XXXXAGGA X86000000000000XXAM1830P3M393493
35	XXABBEBE) PMAMMAXXXXX AB) ABBEEEFX BEHMM

123456789012345678901234567890123456789

Original picture.

123456769012345678901234567490123456739 884 A3638 AYPEAA68A16X12XYYEsYXE +#24244 #38818###237 X*E87 48#*#2###18X##7146818888 E6+ERE#8*8*8xX 78+8x8=18X#+8887 16×844#44 8H8284+6H1488628XYHX88883**+2X3488887*XX L BZ BYEYBBEED XERBARBEE ZBABBEE A= X BY AB X A X Z 1 @ 88/ M848) 54 (78898 YBUG8 + KR848 + (141 / 884 / 88+8%+>H88888888888888888888848848141814181448 13 8=+8A8)F6f+8F6+1- +7) A*#L6*E8+*#+8#* €=5×2A8+88#→882= # 2 212-8888++ 82+2+ 888X4824984091))=-)1 +=6)881X=XX412+ 17 HZAXIAXBEB-16)A)-) - ZBUBXYMA MANKE Y8YX€ZL==88=-AX A1 1F844141#A#K4A 19 682A##X88#48#) X + 2A##)#)28X38X61 AZM+ ABBBB 42808842++62 4644178014764 4X 21 Z@XAZAAAHA##11#A =1Z@Bx@1= AXE-x2@H1@ ALSHMENMENEMEMYABX1+21 1M411AXZOXXX144M 22 23 MBZ MZ #7 @+A211 Z= Z11-M +#1 81 #XA4+4 BEEZZE== MY1 2+1) 2- MZ=+)1444)++ 78 24 25 8A-)816966AX *X++Z1) =+#)* +A XXX 828*18AA=84-1221 A42X X 4= 2411 XX1X1+ "8888)#+)XX "61 XZ4H-+ Z98+ X/ X. -9 +4 1x1A))21=M@MAM@X11)X A Z)1)) A = xx-A ++++ M) Z / X)-#88 ×8421X944+# 14 X AX+ ZAA=) A X== A 281)++) A1 A8 YY+ 4 A))+1=14-X 30)ZAMXZZAM#68-MBA# A+M888XE=X 1M))+) ASA ARBIELSEZ) 1) MEMA-EX99AYX ACCENS 1 EZX-8. #A88 #A8 24*-E14E*Starest##Z 33 ZIA MBR#1@18984%A"X "WAAAEEXX"ERXBEBAS 34 35 +8188× 218292×1×1*+€+×#8088×-444888 123456789012345676901234567890123456789

Noisy picture (sensed scene) as used in experiment.

	1234567890123456784012345678901234557990
1	1111XXXX7717X77711+++1)))XA*AZZXX111
2	111 XAA XA X17 X217 MAX 771117X MAAXXXXX M111
3	111 HOAM 4 AAAXXMXA 494 XX7774 HHHAXA XA111
4	111 MMMMAAAMXXXXXXXXXXXXXXXXXXXXXXXXXXXX
5	111MHHAMAX7XAXX7AXAXXAXXAAXAAAMMMX111
6	1116MMXAXX717XX711X6AXX77XX7XMMMM4111
ž	111 ***********************************
8	11144444447 8444444444444444444444444444
9	1111XXX7XXX7XA99 YOB COO YWWWX77XX7X111
10	
11	
12	1114M4X74/F966 F8888888888899999999999999111
13	111966FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF
14	
15	111 AMAAMMAMMAMAAAAAAAAAAAAAAAAAAAAAAAA
16	111AXA44XAX77771111+11)1764#A57777111
17	111 XA XXXAXX71+1111+ -7A 404 4AXXAA111
18	111XXXA1 4X7111+111+ = 17 AHR4XAAAXA111
10	111 AM AM ##############################
20	111 XMERESSAGE VAAVPAMEE OBAMAXXMA111
21	111X AAMMA ARE THAA AMAR BER MAAMAXXMMM9111
22	111 × 7 ΔΔ× ×******************************
23	1117717474X7XX7/1111177X1744X4HM4 4111
24	111 x 7 7 x 7 117 1 + + + 1 7 1 1 1 1 1 1 x x 7 x 4 4 A X X 4 1 1 1
25	111 AAMMWM AX111+==++++11111777X1111111
26	11140000004Y1X11771711+++1171177X17111
27	11100 MAXX7 177 XAAXX22 1++ 11 11117 12X111
23	1112x2221))2Ax2v A77)+=+)+++)))111111
29	11121+1)+1122X0H44X2272211+++1+11=111
30	111111++17++))7)11717xxxxx1)+)+ ++111
31	1112))11 XA*AX1))122XXAPM271)++= + 111
32	11171117 PRFX11711XXXMM8AX27+)17771111
33	111+)) 1 XAAAAAZ++1 Z1 X X Z Z X Z X AZ 1 Z Z Z 1111
34	111))) ZZX ZXX ZZZX)) 11) 12122 ZZZ XZ XAL111
35	1111127+11+111111+++112XXXAAAAA11111
	1234567890123456789012345678901234567890

L(EV)A for hair. (Density at a point is proportional

to probability that hair is present at that loca-



HAIR WAS LOCATED AT (11, 21) L/EDGE WAS LOCATED AT (25, 11) R/EDGE WAS LOCATED AT (25, 24) L/EYE WAS LOCATED AT (21, 15) R/EYE WAS LOCATED AT (21, 21) NOSE WAS LOCATED AT (22, 18) MOUTH WAS LOCATED AT (29, 17)

M. Fischler and R. Elschlager, <u>The representation and matching of pictorial structures</u>, IEEE Trans. on Computers, 1973

tion.)

History of recognition: Appearance-based models



M. Turk and A. Pentland, <u>Face recognition using</u> <u>eigenfaces</u>, CVPR 1991





H. Murase and S. Nayar, <u>Visual learning and recognition of 3-d</u> objects from appearance, IJCV 1995

History of recognition: Features and classifiers


History of recognition: Deformable templates

Pictorial structures revisited



Felzenszwalb & Huttenlocher (2000)

Discriminatively trained deformable part-based models



Felzenszwalb et al. (2008)

History of recognition: Constellation models



Weber, Welling & Perona (2000), Fergus, Perona & Zisserman (2003)

History of recognition: Bags of keypoints



Csurka et al. (2004), Willamowski et al. (2005), Grauman & Darrell (2005), Sivic et al. (2003, 2005)

Inspired by Bag of Words features from NLP

Spatial pyramids

• Orderless pooling of local features over a coarse grid



Lazebnik, Schmid & Ponce (CVPR 2006)

Spatial pyramids

• Caltech101 classification results:



	Weak features (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	$64.6\pm\!0.7$	

History of recognition: Neural networks

Perceptrons







Minsky & Papert (1969)

LeNet-5



Rumelhart, Hinton & Williams (1986)







AlexNet



Krizhevsky et al. (2012)

LeCun et al. (1998)

What Matters in Recognition?

• Data

- More is always better (as long as it is good data)
- Annotation is the hard part
- Representation
 - Low level: SIFT, HoG, GIST, edges
 - Mid level: Bag of words, sliding window, deformable model
 - High level: Contextual dependence
 - Deep learned features
- Learning Techniques
 - E.g. choice of classifier or inference method

Training & Testing a Classifier



Training & Testing a Classifier



Classifiers

- Nearest Neighbor
- kNN ("k-Nearest Neighbors")
- Linear Classifier
- Neural Network
- Deep Neural Network
- •

Classifiers

- Nearest Neighbor
- kNN ("k-Nearest Neighbors")
- Linear Classifier
- Neural Network

•

Deep Neural Network

Nearest Neighbor (NN) Classifier

• Train

- Remember all training images and their labels
- Predict
 - Find the closest (most similar) training image
 - Predict its label as the true label



CIFAR-10 and NN results

Example dataset: CIFAR-10 10 labels 50,000 training images 10,000 test images.

airplane	Sand .	X	-	X	¥	-	2	1		-
automobile		No.	C		-			٩,	100	*
bird	19	5	1	R		4	1	Y	-	4
cat	1			Se.		1	Z	A.	No.	1
deer	1	40	Ľ.	R		Y	Ŷ	K	T	<u>\$</u>
dog	376	1:	10	3	1	-	9	1	1	14
frog		1	-	S-	- 🐐	٠)	and the second s			32
horse	m.	-	P	7	3	TAB	18	t.	6	N.
ship	-	6	a inte	-	144	-	2	10	-	-
truck		No.	1					4		da

For every test image (first column), examples of nearest neighbors in rows



Slides from Andrej Karpathy and Fei-Fei Li http://vision.stanford.edu/teaching/cs231n/

k-nearest neighbor

- Find the k closest points from training data
- Take majority vote from K closest points



What does this look like?



What does this look like?



K-Nearest Neighbors: Distance Metric

L1 (Manhattan) distance $d_1(I_1, I_2) = \sum |I_1^p - I_2^p|$



K = 1

L2 (Euclidean) distance

$$d_2(I_1,I_2) = \sqrt{\sum_p \left(I_1^p - I_2^p
ight)^2}$$



K = 1

Demo: <u>http://vision.stanford.edu/teaching/cs231n-demos/knn/</u>

Hyperparameters

- What is the **best distance** to use?
- What is the **best value of k** to use?
- These are **hyperparameters**: choices about the algorithm that we set rather than learn
- How do we set them?
 - One option: try them all and see what works best

Setting Hyperparameters

Idea #1: Choose hyperparameters that work best on the data

BAD: K = 1 always works perfectly on training data

Your Dataset

Idea #2: Split data into train and test, chooseBAD: No idea how algorithmhyperparameters that work best on test datawill perform on new data

train test

Idea #3: Split data into train, val, and test; choose	Betterl
hyperparameters on val and evaluate on test	Detter

|--|

Setting Hyperparameters

Your Dataset

Idea #4: Cross-Validation: Split data into folds, try each fold as validation and average the results

fold 1	fold 2	fold 3	fold 4	fold 5	test
fold 1	fold 2	fold 3	fold 4	fold 5	test
fold 1	fold 2	fold 3	fold 4	fold 5	test

Useful for small datasets, but not used too frequently in deep learning

kNN -- Complexity and Storage

- N training images, M test images
- Training: O(1)
- Testing: O(MN)
- We often need the opposite:
 - Slow training is ok
 - Fast testing is necessary



k-Nearest Neighbors: Summary

- In image classification we start with a training set of images and labels, and must predict labels on the test set
- The K-Nearest Neighbors classifier predicts labels based on nearest training examples
- Distance metric and K are **hyperparameters**
- Choose hyperparameters using the validation set; only run on the test set once at the very end!

Problems with KNN: Distance Metrics

- terrible performance at test time
- distance metrics on level of whole images can be very unintuitive



(all 3 images have same L2 distance to the one on the left)

Problems with KNN: The Curse of Dimensionality

- As the number of dimensions increases, the same amount of data becomes more sparse.
- Amount of data we need ends up being exponential in the number of dimensions



Classifiers

- Nearest Neighbor
- kNN ("k-Nearest Neighbors")
- Linear Classifier
- Neural Network

•

Deep Neural Network

Linear Classifiers



This image is <u>CC0 1.0 public</u> domain

Linear Classification vs. Nearest Neighbors

- Nearest Neighbors
 - Store every image
 - Find nearest neighbors at test time, and assign same class



Linear Classification vs. Nearest Neighbors

- Nearest Neighbors
 - Store every image
 - Find nearest neighbors at test time, and assign same class
- Linear Classifier
 - Store hyperplanes that best separate different classes
 - We can compute continuous class score by calculating (signed) distance from hyperplane



We can interpret this as a linear "score function" for each class.

Score functions



class scores

Slide adapted from Andrej Karpathy and Fei-Fei Li http://vision.stanford.edu/teaching/cs231n/

Parametric Approach

image parameters f(x,W)

10 numbers, indicating class scores

[32x32x3] = 3072array of numbers 0...1 (3072 numbers total)

Slide adapted from Andrej Karpathy and Fei-Fei Li http://vision.stanford.edu/teaching/cs231n/

Parametric Approach: Linear Classifier



Parametric Approach: Linear Classifier



Linear Classifier



Interpretation: Algebraic

Example with an image with 4 pixels, and 3 classes (cat/dog/ship)



Interpretation: Geometric

 Parameters define a hyperplane for each class:

$$f(x_i, W, b) = Wx_i + b$$

• We can think of each class score as defining a distribution that is proportional to distance from the corresponding hyperplane



Interpretation: Template matching

• We can think of the rows in $W\,$ as templates for each class



Rows of W in $f(x_i, W, b) = Wx_i + b$
Linear Classifier: Three Viewpoints

f(x,W) = Wx

Algebraic Viewpoint



<u>Visual Viewpoint</u> One template per class



Hyperplanes cutting up space

Geometric Viewpoint



Hard Cases for a Linear Classifier

Class 1:

First and third quadrants

Class 2

Second and fourth quadrants

Class 1: 1 <= L2 norm <= 2

Class 2: Everything else



Class 1: Three modes

Class 2: Everything else



So far: Defined a (linear) <u>score function</u> f(x,W) = Wx + b

Example class scores for 3 images for some W:

How can we tell whether this W is good or bad?

Cat image by <u>Nikita</u> is licensed under <u>CC-BY 2.0</u> Car image is <u>CCO 1.0</u> public domain <u>Frog image</u> is in the public domain





airplane	-3.45	-0.51	3.42
automobile	-8.87	6.04	4.64
bird	0.09	5.31	2.65
cat	2.9	-4.22	5.1
deer	4.48	-4.19	2.64
dog	8.02	3.58	5.55
frog	3.78	4.49	-4.34
horse	1.06	-4.37	-1.5
ship	-0.36	-2.09	-4.79
truck	-0.72	-2.93	6.14

Recap

- Learning methods
 - k-Nearest Neighbors
 - Linear classification
- Classifier outputs a score function giving a score to each class
- How do we define how good a classifier is based on the training data? (Spoiler: define a *loss function*)

Linear classification



airplane	-3.45	-0.51	3.42
automobile	-8.87	6.04	4.64
bird	0.09	5.31	2.65
cat	2.9	-4.22	5.1
deer	4.48	-4.19	2.64
dog	8.02	3.58	5.55
frog	3.78	4.49	-4.34
horse	1.06	-4.37	-1.5
ship	-0.36	-2.09	-4.79
truck	-0.72	-2.93	6.14

Cat image by Nikita is licensed under CC-BY 2.0; Car image is CC0 1.0 public domain; Frog image is in the public domain

Output scores

TODO:

- Define a loss function that quantifies our unhappiness with the scores across the training data.
- 2. Come up with a way of efficiently finding the parameters that minimize the loss function.
 (optimization)

Loss functions

3.2

5.1

-1.7

cat

car

frog

Suppose: 3 training examples, 3 classes. With some W the scores f(x, W) = Wx are:



1.3

4.9

2.0

2.2

2.5

-3.1

A **loss function** tells how good our current classifier is

Given a dataset of examples $\{(x_i, y_i)\}_{i=1}^N$

Where $oldsymbol{x_i}$ is image and $oldsymbol{y_i}$ is (integer) label

Loss over the dataset is a sum of loss over examples:

$$L = \frac{1}{N} \sum_{i} L_i(f(x_i, W), y_i)$$

Simpler example: binary classification

- Two classes (e.g., "cat" and "not cat")
 - AKA "positive" and "negative" classes









not cat

Linear classifiers

- Find linear function (*hyperplane*) to separate positive and negative examples
 - $\mathbf{x}_i \text{ positive}: \quad \mathbf{x}_i \cdot \mathbf{w} + b \ge 0$ $\mathbf{x}_i \text{ negative}: \quad \mathbf{x}_i \cdot \mathbf{w} + b < 0$

Which hyperplane is best? We need a **loss function** to decide



What is a good loss function?

- One possibility: Number of misclassified examples
 - Problems: discrete, can't break ties
 - We want the loss to lead to good generalization
 - We want the loss to work for more than 2 classes



Softmax classifier

 Interpret Scores as unnormalized log probabilities of classes

$$f(x_i, W) = Wx_i$$
 (score function)



softmax function

Squashes values into *probabilities* ranging from 0 to 1

$$P(y_i \mid x_i; W)$$

Example with three classes:

 $[1,-2,0] \rightarrow [e^1,e^{-2},e^0] = [2.71,0.14,1] \rightarrow [0.7,0.04,0.26]$

Softmax classifier

Example with an image with 4 pixels, and 3 classes (cat/dog/ship)



Cross-entropy loss

 $f(x_i, W) = Wx_i$ (score function)

Cross-entropy loss

 $f(x_i, W) = W x_i$ (score function)



Cross-entropy loss

 $f(x_i, W) = W x_i$ (score function)



Losses

- Cross-entropy loss is just one possible loss function
 - One nice property is that it reinterprets scores as probabilities, which have a natural meaning
- SVM (max-margin) loss functions also used to be popular
 - But currently, cross-entropy is the most common classification loss

Summary

- Have score function and loss function
 - Currently, score function is based on linear classifier
 - Next, will generalize to convolutional neural networks
- Find W and b to minimize loss



What's Still Hard?

- Fine-grain classification
 - How do we distinguish between more subtle class differences?

Animal->Bird->Oriole...



Baltimore Oriole



Hooded Oriole



Scott Oriole

What's Still Hard?

- Few shot learning
 - How do we generalize from only a small number of examples?





Slide Credits

- <u>CS5670, Introduction to Computer Vision</u>, Cornell Tech, by Noah Snavely.
- CS 543 Computer Vision, by Stevlana Lazebnik, UIUC.
- EECS 442 <u>Computer Vision</u>, by Justin Johnson & David Fouhey, U Michigan.