Outline

- Motivation
- How to define problem?
- Case example: finding people
- Statistical classification

Zoom-out to the big picture



(We'll talk more about mid-level grouping toward end of semester)

Formalizing recognition



What are outputs of interest? How should space of outputs be represented?

yes-no answers: does this contain a car?



Image classification is this a beach?



Detection Does this contain a car? Where?



Attributes

Object poses + materials



Activities What are these people doing?



Instance recognition

Recognition meet geometry





Formalizing recognition

Human vision experiment: what can people describe when looking at images?

Tricky question to answer....

Rapid scene catgorization





People can distinguish high-level concepts (animal/transport) in under 150ms (Thorpe)



Appears to suggest feed-forward computations suffice (or at least dominate)

What do we perceive in a glance of a real-world scene?





PT = 107 ms

This is outdoors. A black, furry dog is running/walking towards the right of the picture. His tail is in the air and his mouth is open. Either he had a ball in his mouth or he was chasing after a ball. (Subject EC)

PT = 500 ms

I saw a black dog carrying a gray frisbee in the center of the photograph. The dog was walking near the ocean, with waves lapping up on the shore. It seemed to be a gray day out. (Subject JB)



Inside a house, like a living room, with chairs and sofas and tables, no ppl. (Subject HS) A room full of musical instruments. A piano in the foreground, a harp behind that, a guitar hanging on the wall (to the right). It looked like there was also a window behind the harp, and perhaps a bookcase on the left. (Subject RW)



PT = 107 ms

This is outdoors. A black, furry dog is running/walking towards the right of the picture. His tail is in the air and his mouth is open. Either he had a ball in his mouth or he was chasing after a ball. (Subject EC)

PT = 500 ms

I saw a black dog carrying a gray frisbee in the center of the photograph. The dog was walking near the ocean, with waves lapping up on the shore. It seemed to be a gray day out. (Subject JB)



Inside a house, like a living room, with chairs and sofas and tables, no ppl. (Subject HS) A room full of musical instruments. A piano in the foreground, a harp behind that, a guitar hanging on the wall (to the right). It looked like there was also a window behind the harp, and perhaps a bookcase on the left. (Subject RW)

Should language be the right output?

"The more you look, the more you see"

Hochstein & Ahissar 02



Vision at a glance (feedforward)

Rapid scene categorization

Vision with scrutiny (+feedback)

Fine-grained recognition Spatial localization for manipulation

Hilbert problems of vision



Vigorous discussion

Hierarchical annotation

Image Retrieval using Scene Graphs

Justin Johnson¹, Ranjay Krishna¹, Michael Stark², Li-Jia Li^{3,4}, David A. Shamma³, Michael S. Bernstein¹, Li Fei-Fei¹

¹Stanford University, ²Max Planck Institute for Informatics, ³Yahoo Labs, ⁴Snapchat



Other hiearchies for organizing recognition outputs

Partonomies (people are made out of arms and legs)

Taxonomies (dogs and cats are types of mammals)

Spatiotemporal Relations (people stand on floors)







Associative memory

Big-data philoshopy: ask not "what is this", but "what is this like"?



Malisiewicz et al, "Exemplar SVMs"

Visual Question Answering (VQA)



What color are her eyes? What is the mustache made of?



How many slices of pizza are there? Is this a vegetarian pizza?



Is this person expecting company? What is just under the tree?



Does it appear to be rainy? Does this person have 20/20 vision?

Generate questions via Amazon MTurk:

"We have built a smart robot. It understands a lot about images. It can recognize and name all the objects, it knows where the objects are, it can recognize the scene (e.g., kitchen, beach), people's expressions and poses, and properties of objects (e.g., color of objects, their texture). Your task is to stump this smart robot!".

My (current) favorite answer: semantic image segmentation



A "detail": instance segmentation

Augment semantic labels with an instance ID



Let's give it a try



Notes on image annotation

Adela Barriuso, Antonio Torralba Computer Science and Artificial Intelligence Laboratory (CSAIL), Massachusetts Institute of Technology

"I can see the ceiling, a wall and a ladder, but I do not know how to annotate what is on the right side of the picture. Maybe I just need to admit that I can not solve this picture in an easy and fast way. But if I was forced"

Semantic blindspots

Can we define a canonical list of objects, attributes, actions, materials....?



ImageNet (cf. WordNet, VerbNet, FrameNet,..)

Crowdsourced dataset construction



Task: select images that have WRONG object contour for toothbrush. Examples: Right Object Contour



Wrong Object Contour (not toothbrush, only contains parts of visible object contour, or multiple objects)



Tips: use n and b keys to move between rows of image.



"If you know what can be done with a ... object, what it can be used for, you can call it whatever you please" J. J Gibson [14]



Outline

- Motivation
- How to define problem?
- Case example: finding people
- Statistical classification

Goal: detecting pedestrians



Thought experiment: let's build a person detector (HW4). Why is this difficult?



variation in illumination



variation in appearance



variation in pose, viewpoint



occlusion & clutter

Classic "nuisance factors" for general object recognition

Main idea: use "invariant features"

edges!





with factor of the starts of a

e, filters and deformation c





(Simplified) HOG construction









- mage is partitioned
- n each block we co
- HAVAFIANT to chang



ineh prisetations

mations, etc.

Ve compute features at different resolutions (pyranida)

What should be the angle range of each bin?

[H x W] ->[H x W x 9] "orientation channel array"







Count up orientation bins over 8x8 pixel neighborhoods. (im2col) Get some spatial invariance (sort of)...





- mage is part in the 8x8 pixel blocks is part in the 8x8 pixel blocks is part in the second se
- Invariant to changes in lighting, small deformations, etc.
- Ve compute features at different resolutions (pyramid))

Get some lighting invariance (sort of)...

e, filters and deformation c





Recall SIFT

http://en.wikipedia.org/wiki/Scale-invariant_feature_transform



One can interpret HOG precisely as dense grid of SIFT descriptors (of size 2x2x9), computed at grid points of 8x8 pixel shifts

Template scoring



How can we fix output from naive thresholding? ³⁹

Nonmaximal suppression



- 1. Find highest scoring location
- 2. Zero-out responses that overlap
- 3. Repeat until highest remaining score is below a threshold



Pedestrian detection



Dalal and Triggs "Histograms of Gradients"

Face detection



template



Training

of images with labeled bounding boxe



Train

filters and defo



W







ed bounding boxes

nd deformation costs

arch over scales







Outline

- Motivation
- How to define problem?
- Case example: finding people
- Statistical classification

ata consists of images with labeled bounding boxes

arn the model structure, filters and deformation costs



Statistical classification (20 min review!)

Why?

- This is the world in which we live statistical models from data overpower classic "hand-designed" models
- Basic linear classification forms basis for nonlinear models (deep learning)







Recall N ~ 10x5x9 for typical HOG templates

Given training points (xi,yi), learn function f(x) that predicts a label {-1,1}

Statistical classification





Version 0: nearest neighbor classification

Train time:0 Test time: expensive Trivially handles multiple classes Is surprisingly powerful!

Statistical classification

Ask not what is this, but "what is this like"?







Exemplar



(in my view) this is the heart of all data-driven learning



Given training points (xi,yi), learn function f(x) that predicts a label {-1,1}

Version 0: nearest neighbor classification Version 1: linear classification

What's the best line?





$$\min_{w} \sum_{i} [y_i \neq thresh(w \cdot x_i)]$$

thresh(x) =
$$\begin{cases} 1, & \text{if } x > 0\\ -1, & \text{otherwise} \end{cases}$$

Find w that minimizes mistakes on training data [Hard to optimize]

Aside: doesn't this restrict line to go through origin?

What's the best line?



Easy to optimize - least squares!

The first term is a *regularizer* (prevents overfitting and makes optimization easier)

Alternate visualization: heightfeild





Unified notation: regularized loss minimization





$$\begin{array}{l} \textbf{Birds-eye view of ML}\\ \underset{w}{\min} \lambda R(w) + \sum_{i} loss(y_{i}f_{w}(x_{i}))\\ R_{L2}(w) = \frac{1}{2}||w||^{2} \qquad (L2 \ \text{regularization})\\ R_{L1}(w) = \sum_{j} |w_{j}| \qquad (L1/\text{sparse regularization})\\ loss_{squared}(m) = (1-m)^{2} \qquad (Linear \ \text{regression})\\ loss_{log}(m) = \log(1+e^{-m}) \qquad (Logistic \ \text{regression})\\ loss_{hinge}(m) = \max(0, 1+m) \qquad (Support \ \text{vector machine}) \end{array}$$

$$f_w(x_i) = w^T x_i$$
 Linear classifier $f_w(x_i) = CNN_w(x_i)$ Nonlinear classifier

Learning with losses



Bottom line: other than 01 and squared loss, most loss functions look similar (to me)

Learning with
gradient descent
$$\min_{w} L(w)$$
 where $L(w) = \frac{1}{2}||w||^{2} + \sum_{i}(f_{w}(x_{i}) - y_{i})^{2}$
 $w := w - step * \left(w + \sum_{i}(f_{w}(x_{i}) - y_{i})\frac{\partial f_{w}(x_{i})}{\partial w}\right)$



Early 2000's: obsession with convex L(w)



Trivially "out-of-core": most contemporary models are trained in this manner



Positive examples that score better than 1 and negative examples that score less than 1 do not affect the objective function

Implication: we can through them away without changing the solution

Large-scale learning

pos

neg





Our test set distribution is highly imbalanced; so should be the training set (hundreds of positives, hundreds of millions of negatives)

Support vector machines (SVMs) are attractive because they generate sparse learning problems (turn continuous search over parameters into combinatorial search over data)

Large-scale learning for SVMs Lots of large-scale solvers for quadratic programs (SVMS)

Two flavors

Batch: Require access to all training data (guarantees on convergence)

Online: Require access to on-the-fly training data (usually stochastic in practice)

In-between: Support-vectors fit in memory, but data doesn't (Relatively unexplored!)

Large-scale learning for SVMs

In practice, can get near-optimal models with a single pass of "dual coordinate descent" through large datasets

A. Bordes, L. Bottou, P. Gallinari, and J. Weston. Solving multiclass support vector machines with LaRank. In *ICML*, pages 89–96. ACM, 2007. 7

A. Bordes, S. Ertekin, J. Weston, and L. Bottou. Fast kernel classifiers with online and active learning. *The Journal of Machine Learning Research*, 6:1579–1619, 2005. 8

L. Bottou and O. Bousquet. The tradeoffs of large scale learning. Advances in neural information processing systems, 20:161–168, 2008. 1

http://www.csie.ntu.edu.tw/~cjlin/liblinear/

What the presentive weights interm? $w \cdot x > 0$

 $(w_{pos} - w_{neg}) \cdot x > 0$

 $W_{pos} \cdot X > W_{neg} \cdot X$

Pedestrian template



Pedestrian background template

Right approach is to compete pedestrian, pillar, doorway... models Background class is hard to model - easier to penalize particular vertical edges Historically, model-based approaches tend not to model negative set

But what if we actually try to model positives and negatives?



1. P(y): Flip (biased) coin to select class 'k'

2. P(x|y): Sample from Gaussian (mu_k, Sigma_k)

A look ahead





Simple linear discriminant analysis gets us 90% of the way there...

Class-conditional Gaussians

If we assume all classes have same Sigma.... (For notational ease, let's assume priors for class y in {0,1} are equal)



$$p(y = 1|x) = \frac{p(x|y = 1)p(y = 1)}{p(x|y = 1)p(y = 1) + p(x|y = 0)p(y = 0)}$$

Plug in the following and simplify:

$$p(y=1) = .5$$
$$p(x|y) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} e^{-\frac{1}{2}(x-\mu_y)^T \Sigma^{-1} (x-\mu_y)}$$

Class-conditional Gaussians

If we assume all classes have same Sigma....

$$p(y = 1|x) = sigmoid(w \cdot x + b)$$





p(y=1|x) > .5 when $f_w(x) = w \cdot x + b > 0$, $w = \Sigma^{-1}(\mu_1 - \mu_0)$

Class-conditional Gaussians



What happens if Sigmas for the two classes are different?

Alternative: discriminative fitting (logistic regression)



2

 (X_i^m)

 $p(y = 1|x) = sigmoid(w \cdot x + b)$

Given data $\{(x_i,y_i)\}$, directly fit parameters 'w,b' to maximize

$$\sum_{i} \log p(y_i | x_i)$$

This is sometimes called cross-entropy minimization, and is equivalent to miniming with log-loss!

Extension to multiple classes

Same derivation of class-conditional Gaussians for K classes (assuming all have same covariance)



Known as the *softmax* function

When parameters are fit to maximize $\log p(y|x)$, known as softmax cross-entropy minimmization

Let's build up intuition with generative model



SVM





Gaussian model

 $w = \Sigma^{-1}(\mu_1 - \mu_0)$

Centered model

 $w = \mu_1 - \mu_0$

Learn templates with generic (de)correlation model

Hariharan, Malik, Ramanan ECCV 12
What's the covariance matrix capturing?



The numbers



Simple gaussian model gets us 90% of the way there...

Parting thoughts on class-conditional Gaussians

1. Also known as Linear Discriminant Analysis (LDA) or Fischer Discriminant Analysis (FDA) when derived using other criteria (maximizing ratio of between to within-class variances)

2. One can also obtain the LDA / FDA solution by discriminative learning with a squared error loss

(Hastie et al, Elements of Statistical Learning)

Implies the distinction between generative and discriminative models can be blurred...



Parting thoughts on statistical classification

- Loss functions: hinge, log-loss, squared loss
- SVMs generate sparse optimization problems
- Generative models are promising, but current state-of-the-art relies on discrimative loss minimization

Things that will appear later:

Cross-entropy loss: minimizing cross-entropy of binary prediction is equivalent to log-loss Soft-max loss: minimizing cross-entropy of K-way prediction is equivalent to soft-max loss

Outline

- Motivation
- How to define problem?
- Case example: finding people
- Statistical classification