



Visual words

- Map high-dimensional descriptors to tokens/words by quantizing the feature space

- Quantize via clustering, let cluster centers be the prototype "words"
- Determine which word to assign to each new image region by finding the closest cluster center.

2
Kristen Grauman

Visual words

- Example: each group of patches belongs to the same visual word

3
Figure from Sivic & Zisserman, ICCV 2003
Kristen Grauman

Inverted file index

Database images

Word #	Image #
1	3
2	
...	
7	1, 2
8	3
9	
10	
...	
91	2

- Database images are loaded into the index mapping words to image numbers

Kristen Grauman

Inverted file index

New query image

Word #	Image #
1	3
2	
7	1, 2
8	3
9	
10	
...	
91	2

- New query image is mapped to indices of database images that share a word.

Kristen Grauman

Bags of visual words

- Summarize entire image based on its distribution (histogram) of word occurrences.
- Analogous to bag of words representation commonly used for documents.

6

Comparing bags of words

- Rank frames by normalized scalar product between their (possibly weighted) occurrence counts---nearest neighbor search for similar images.

[1 8 1 4]

\vec{d}_j

[5 1 1 0]

\vec{q}

$$sim(d_j, q) = \frac{\langle d_j, q \rangle}{\|d_j\| \|q\|}$$

$$= \frac{\sum_{i=1}^V d_j(i) * q(i)}{\sqrt{\sum_{i=1}^V d_j(i)^2} * \sqrt{\sum_{i=1}^V q(i)^2}}$$

for vocabulary of V words

Kristen Grauman

Application: Large-Scale Retrieval

Query Results from 5k Flickr images (demo available for 100k set)

8 [Philbin CVPR'07]

Spatial Verification: RANSAC

What else can we borrow from text retrieval?

INDEX

Apple 187, 190, 200, 210, 220, 230, 240, 250, 260, 270, 280, 290, 300, 310, 320, 330, 340, 350, 360, 370, 380, 390, 400, 410, 420, 430, 440, 450, 460, 470, 480, 490, 500, 510, 520, 530, 540, 550, 560, 570, 580, 590, 600, 610, 620, 630, 640, 650, 660, 670, 680, 690, 700, 710, 720, 730, 740, 750, 760, 770, 780, 790, 800, 810, 820, 830, 840, 850, 860, 870, 880, 890, 900, 910, 920, 930, 940, 950, 960, 970, 980, 990, 1000

China is forecasting a trade surplus of \$90bn (£51bn) to \$100bn this year, a threefold increase on 2004's \$32bn. The Commerce Ministry said the surplus would be predicted to rise in 2010 by a further 30% to \$120bn. China's deliberate policy to curb the surplus is one factor that has kept the yuan's value low. The yuan has stayed within a narrow band since July 2005 and permitted it to move more freely. However, Beijing has made it clear that it will take its time and tread carefully allowing the yuan to rise further in value.

China, trade, surplus, commerce, exports, imports, US, yuan, bank, domestic, foreign, increase, trade, value

tf-idf weighting

- Term frequency – inverse document frequency
- Describe frame by frequency of each word within it, downweight words that appear often in the database
- (Standard weighting for text retrieval)

Number of occurrences of word i in document d

Number of words in document d

$$t_i = \frac{n_{id}}{n_d} \log \frac{N}{n_i}$$

Total number of documents in database

Number of documents word i occurs in, in whole database

11
Kristen Grauman

Query expansion

Query: **golf green**

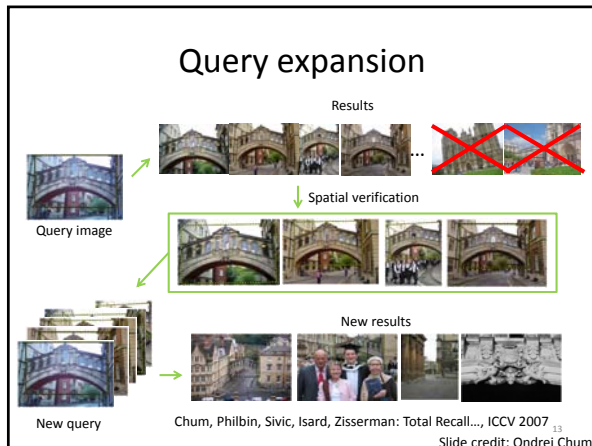
Results:

- How can the grass on the **greens** at a **golf** course be so perfect?
- For example, a skilled **golfer** expects to reach the **green** on a par-four hole in ...
- Manufactures and sells synthetic **golf** putting **greens** and mats.

Irrelevant result can cause a 'topic drift':

- Volkswagen **Golf**, 1999, **Green**, 2000cc, petrol, manual, hatchback, 94000miles, 2.0 GTI, 2 Registered Keepers, HPI Checked, Air-Conditioning, Front and Rear Parking Sensors, ABS, Alarm, Alloy

12
Slide credit: Andrei Chum



Recognition via local-feature based alignment



Pros:

- Effective when we are able to find reliable features within clutter
- Great results for matching specific instances

Cons:

- Spatial verification as post-processing – not seamless, expensive for large-scale problems
- Not suited for generic category recognition

14
Kristen Grauman



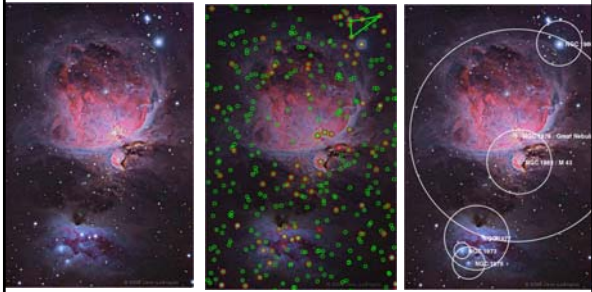
Making the Sky Searchable: Fast Geometric Hashing for Automated Astrometry

Sam Roweis, Dustin Lang & Keir Mierle
University of Toronto

David Hogg & Michael Blanton
New York University

15

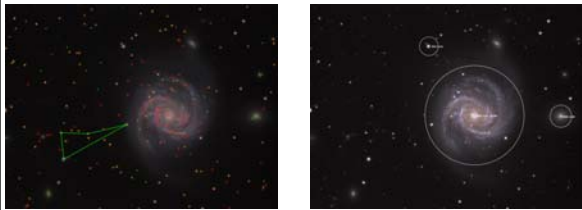
Example



A shot of the Great Nebula, by Jerry Lodriguss (c.2006), from [astropix.com](http://astrometry.net/gallery.html)
<http://astrometry.net/gallery.html>

16

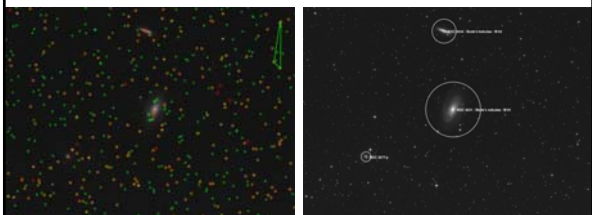
Example



An amateur shot of M100, by Filippo Ciferri (c.2007) from [flickr.com](http://astrometry.net/gallery.html)
<http://astrometry.net/gallery.html>

17

Example



A beautiful image of Bode's nebula (c.2007) by Peter Bressler, from [starglhfrieseid.de](http://astrometry.net/gallery.html)
<http://astrometry.net/gallery.html>

18

Today

- Generic object recognition

19

What does recognition involve?



Source: Fei-Fei Li, Rob Fergus, Antonio Torralba

20

Verification: is that a lamp?



Source: Fei-Fei Li, Rob Fergus, Antonio Torralba

21

Detection: are there people?



Source: Fei-Fei Li, Rob Fergus, Antonio Torralba. 22

Identification: is that Potala Palace?



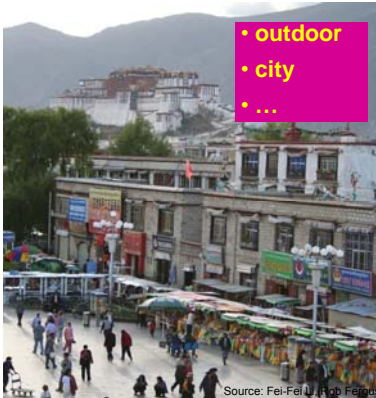
Source: Fei-Fei Li, Rob Fergus, Antonio Torralba. 23

Object categorization



Source: Fei-Fei Li, Rob Fergus, Antonio Torralba. 24

Scene and context categorization



Source: Fei-Fei Li, Rob Fergus, Antonio Torralba.

Instance-level recognition problem



John's car

26


Generic categorization problem



27

Object Categorization

- Task Description
 - “Given a small number of training images of a category, recognize a-priori unknown instances of that category and assign the correct category label.”
- Which categories are feasible visually?



“Fido” German shepherd dog animal living being

K. Grauman, B. Leibe 28

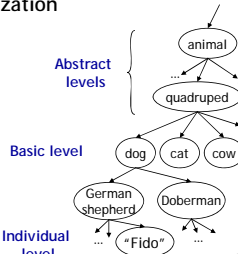
Visual Object Categories

- Basic Level Categories in human categorization [Rosch 76, Lakoff 87]
 - The highest level at which category members have similar perceived shape
 - The highest level at which a single mental image reflects the entire category
 - The level at which human subjects are usually fastest at identifying category members
 - The first level named and understood by children
 - The highest level at which a person uses similar motor actions for interaction with category members

K. Grauman, B. Leibe 29

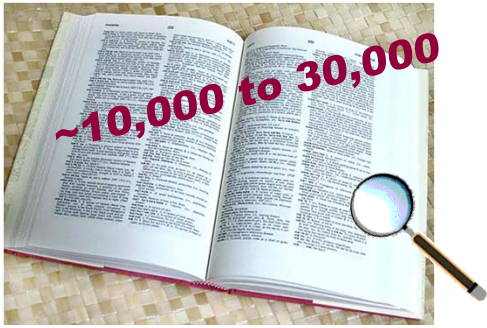
Visual Object Categories

- Basic-level categories in humans seem to be defined predominantly visually.
- There is evidence that humans (usually) start with basic-level categorization *before* doing identification.



K. Grauman, B. Leibe 30

How many object categories are there?



Source: Fei-Fel Li, Rob Fergus, Antonio Torralba. Biederman 1987



Other Types of Categories

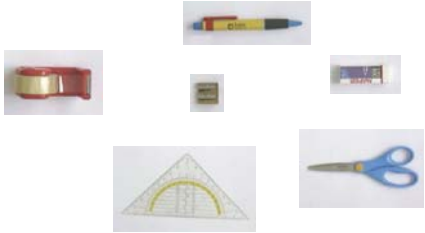
- Functional Categories
 - e.g. chairs = "something you can sit on"



K. Grauman, B. Leibe 33

Other Types of Categories

- Ad-hoc categories
 - e.g. "something you can find in an office environment"



Visual Object Recognition Tutorial

K. Grauman, B. Leibe


34

Why recognition?

- Recognition a fundamental part of perception
 - e.g., robots, autonomous agents
- Organize and give access to visual content
 - Connect to information
 - Detect trends and themes

35

Posing visual queries



Yeh et al., MIT

Digital Field Guides Eliminate the Guesswork

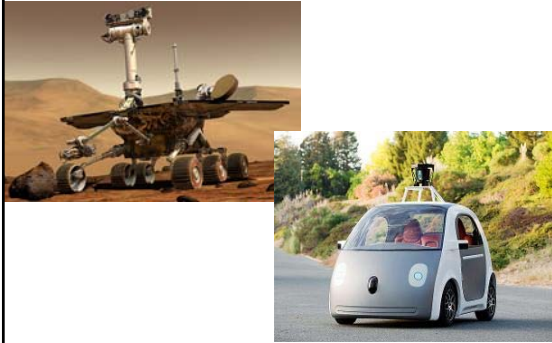
Belhumeur et al.

suptell part of AOL

Kooba, Bay & Quack et al.

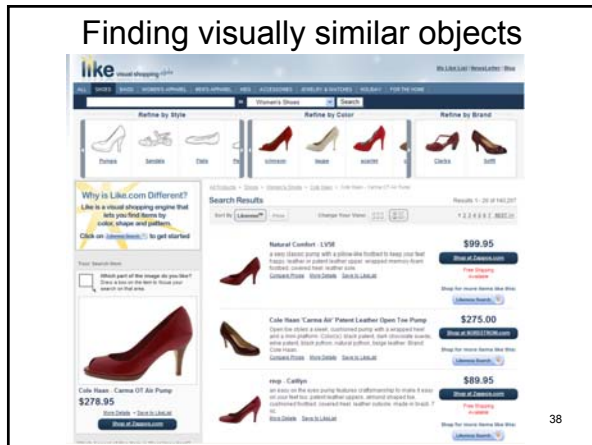
36

Autonomous agents able to detect objects



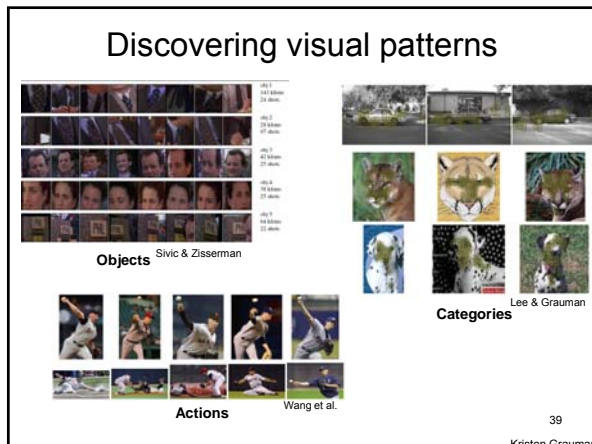
37

Finding visually similar objects



38

Discovering visual patterns



39

Kristen Grauman

Auto-annotation



Gammeter et al.



T. Berg et al.

President George W. Bush under a...
 British cleric from...
 American...
 Gammeter et al.
 T. Berg et al.

What are the challenges?



What we see



0	3	2	5	4	7	6	9	8
3	0	1	2	3	4	5	6	7
2	1	0	3	2	5	4	7	6
5	2	3	0	1	2	3	4	5
4	3	2	1	0	3	2	5	4
7	4	5	2	3	0	1	2	3
6	5	4	3	2	1	0	3	2
9	6	7	4	5	2	3	0	1
8	7	6	5	4	3	2	1	0

What a computer sees

Challenges: robustness

Illumination Object pose Clutter

Occlusions Intra-class appearance Viewpoint

43
Kristen Grauman

Challenges: robustness

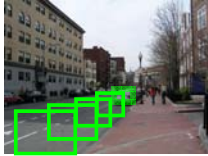
Realistic scenes are crowded, cluttered, have overlapping objects.

44

Challenges: importance of context

45
slide credit: Fei-Fei, Fergus & Torralba

Challenges: importance of context



46

Challenges: complexity

flickr
10 billion images

facebook
250 billion images

imgur
the simple image shorer
1 billion images served daily

YouTube
300 hours uploaded per minute

photobucket
10 billion images

CISCO

Almost 90% of web traffic is visual!

47

Challenges: complexity

- Thousands to millions of pixels in an image
- 30+ degrees of freedom in the pose of articulated objects (humans)
- About half of the cerebral cortex in primates is devoted to processing visual information [Felleman and van Essen 1991]

48

Kristen Grauman

Challenges: learning with minimal supervision

← Less | More →

Unlabeled, multiple objects

Classes labeled, some clutter

Cropped to object, parts and classes labeled

49
Kristen Grauman

What works today

- Reading license plates, zip codes, checks

3 6 8 / 7 9 6 6 9 1
6 7 5 7 8 6 3 4 8 5
2 1 7 9 7 / 2 3 4 5
4 8 1 9 0 1 8 8 9 4
7 6 1 8 6 4 1 5 6 0
7 5 9 2 6 5 8 1 9 7
2 2 2 2 3 4 4 8 0
0 2 3 8 0 7 3 8 5 7
0 1 4 6 4 6 0 2 4 3
7 / 2 8 9 6 9 8 6 1

50
Source: Lana Lazebnik

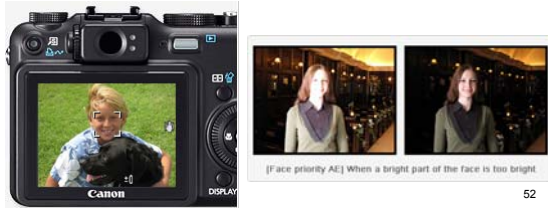
What works today

- Reading license plates, zip codes, checks
- Fingerprint recognition

51
Source: Lana Lazebnik

What works today

- Reading license plates, zip codes, checks
- Fingerprint recognition
- Face detection



52
Source: Lana Lazebnik

What works today

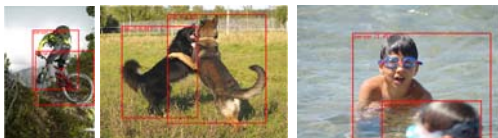
- Reading license plates, zip codes, checks
- Fingerprint recognition
- Face detection
- Recognition of flat textured objects (CD covers, book covers, etc.)



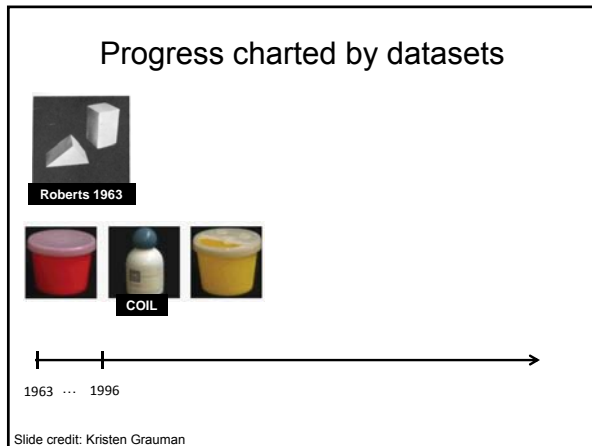
53
Source: Lana Lazebnik

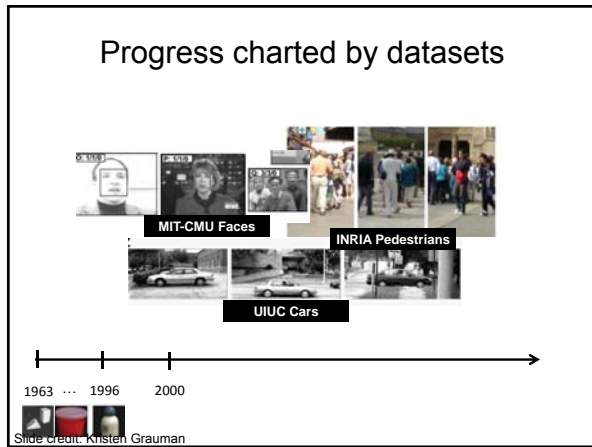
What works today

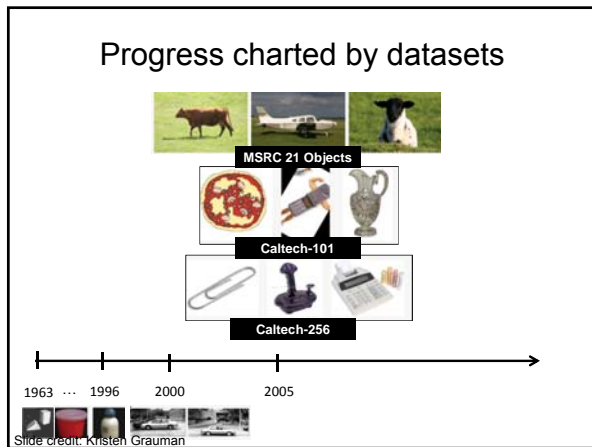
- Reading license plates, zip codes, checks
- Fingerprint recognition
- Face detection
- Recognition of flat textured objects (CD covers, book covers, etc.)
- Recognition of generic categories(*)!

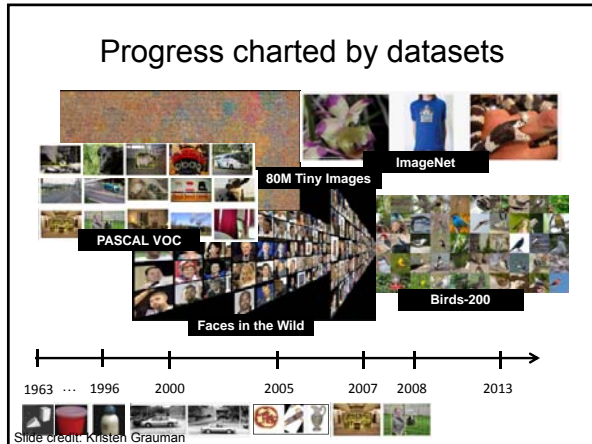


54









Evolution of methods

- Hand-crafted models
- 3D geometry
- Hypothesize and align
- Hand-crafted features
- Learned models
- Data-driven
- “End-to-end” learning of features and models*,**

_____→

* Labeled data availability
** Architecture design decisions, parameters.

Kristen Grauman

Supervised classification

- Given a collection of *labeled* examples, come up with a function that will predict the labels of new examples.

“four”					
“nine”					

Training examples Novel input

- How good is the function we come up with to do the classification?
- Depends on
 - Mistakes made
 - Cost associated with the mistakes

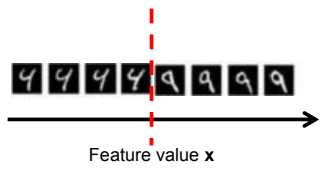
60
Kristen Grauman

Supervised classification

- Given a collection of *labeled* examples, come up with a function that will predict the labels of new examples.
- Consider the two-class (binary) decision problem
 - $L(4 \rightarrow 9)$: Loss of classifying a 4 as a 9
 - $L(9 \rightarrow 4)$: Loss of classifying a 9 as a 4
- Risk** of a classifier s is expected loss:
 $R(s) = \Pr(4 \rightarrow 9 | \text{using } s)L(4 \rightarrow 9) + \Pr(9 \rightarrow 4 | \text{using } s)L(9 \rightarrow 4)$
- We want to choose a classifier so as to minimize this total risk

61
Kristen Grauman

Supervised classification



Feature value x

Optimal classifier will minimize total risk.

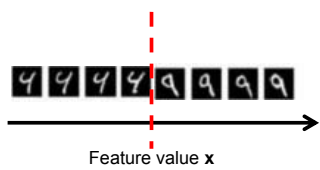
At decision boundary, either choice of label yields same expected loss.

If we choose class "four" for point x at boundary, expected loss is:
 $= P(\text{class is } 9 | \mathbf{x}) L(9 \rightarrow 4) + P(\text{class is } 4 | \mathbf{x}) L(4 \rightarrow 4)$

If we choose class "nine" for point x at boundary, expected loss is:
 $= P(\text{class is } 4 | \mathbf{x}) L(4 \rightarrow 9)$

62
Kristen Grauman

Supervised classification



Feature value x

Optimal classifier will minimize total risk.

At decision boundary, either choice of label yields same expected loss.

So, best decision boundary is at point x where
 $P(\text{class is } 9 | \mathbf{x}) L(9 \rightarrow 4) = P(\text{class is } 4 | \mathbf{x}) L(4 \rightarrow 9)$

To classify a new point, choose class with lowest expected loss; i.e., choose "four" if
 $P(4 | \mathbf{x}) L(4 \rightarrow 9) > P(9 | \mathbf{x}) L(9 \rightarrow 4)$

63
Kristen Grauman

Supervised classification

Optimal classifier will minimize total risk.

At decision boundary, either choice of label yields same expected loss.

So, best decision boundary is at point x where

$$P(\text{class is } 9 | x) L(9 \rightarrow 4) = P(\text{class is } 4 | x) L(4 \rightarrow 9)$$

To classify a new point, choose class with lowest expected loss: i.e., choose "four" if

$$P(4 | x) L(4 \rightarrow 9) < P(9 | x) L(9 \rightarrow 4)$$

How to evaluate these probabilities? Kristen Grauman

Probability

Basic probability

- X is a random variable
- $P(X)$ is the probability that X achieves a certain value

- $0 \leq P(X) \leq 1$
- $\int_{-\infty}^{\infty} P(X) dX = 1$ or $\sum P(X) = 1$
continuous X discrete X
- Conditional probability: $P(X | Y)$
- probability of X given that we already know Y

65 Source: Steve Seitz


Example: learning skin colors

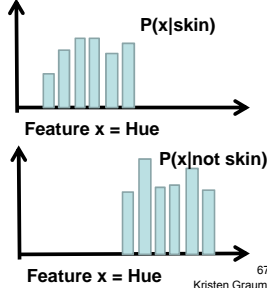
- We can represent a class-conditional density using a histogram (a "non-parametric" distribution)

66 Kristen Grauman

Example: learning skin colors

- We can represent a class-conditional density using a histogram (a "non-parametric" distribution)





Now we get a new image, and want to label each pixel as skin or non-skin.
 What's the probability we care about to do skin detection?

67
Kristen Grauman

Bayes rule

$$P(\text{skin} | x) = \frac{\overbrace{P(x | \text{skin})}^{\text{likelihood}} \overbrace{P(\text{skin})}^{\text{prior}}}{\underbrace{P(x)}^{\text{posterior}}}$$

Where does the prior come from?
 Why use a prior?

68

Example: classifying skin pixels

Now for every pixel in a new image, we can estimate probability that it is generated by skin.




Brighter pixels →
higher probability
of being skin

Classify pixels based on these probabilities

- if $p(\text{skin} | \mathbf{x}) > \theta$, classify as skin
- if $p(\text{skin} | \mathbf{x}) < \theta$, classify as not skin

69

Example: classifying skin pixels



Figure 6: A video image and its flesh probability image



Figure 7: Orientation of the flesh probability distribution marked on the source video image

Gary Bradski, 1998

70

Example: classifying skin pixels

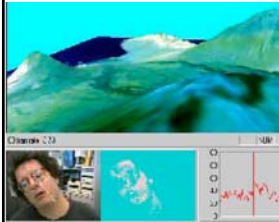


Figure 13: CAMSHIFT-based face tracker used to play over a 3D graphic's model of Hawaii



Figure 12: CAMSHIFT-based face tracker used to play Quake 2 hands free by inserting control variables into the mouse queue

Using skin color-based face detection and pose estimation as a video-based interface

Gary Bradski, 1998

71

Supervised classification

- Want to minimize the expected misclassification
- Two general strategies
 - Use the training data to build representative probability model; separately model class-conditional densities and priors (*generative*)
 - Directly construct a good decision boundary, model the posterior (*discriminative*)

72

Coming up

- Face detection
- Categorization with local features and part-based models
- Deep convolutional neural networks

73

Questions?

See you Tuesday!

74
