## **Object Recognition**

CS 554 – Computer Vision Pinar Duygulu Bilkent University

(Source:Kristen Grauman)

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



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

- 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:

 $R(s) = \Pr(4 \to 9 \mid \text{using } s)L(4 \to 9) + \Pr(9 \to 4 \mid \text{using } s)L(9 \to 4)$ 

• We want to choose a classifier so as to minimize this total risk



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

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

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

$$P(4 \mid \mathbf{x})L(4 \rightarrow 9) > P(9 \mid \mathbf{x})L(9 \rightarrow 4)$$



**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

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

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

How to evaluate these probabilities?

**Kristen Grauman** 

## Basic probability

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



- Conditional probability: P(X | Y)
  - probability of X given that we already know Y

Source: Steve Seitz

#### Example: learning skin colors

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



## 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?



#### Bayes rule



#### $P(skin | x) \alpha P(x | skin) P(skin)$

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

#### 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

Kristen Grauman

#### Example: classifying skin pixels





Figure 13: CAMSHIFT-based face tracker used to 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

Kristen Grauman

- 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*)

# Discriminative classifiers for image recognition

nearest neighbors (+ scene match app)

- support vector machines (+ gender, person app)

#### Nearest Neighbor classification

 Assign label of nearest training data point to each test data point

Black = negative Red = positive



Novel test example

Closest to a positive example from the training set, so classify it as positive.

from Duda *et al.* Voronoi partitioning of feature space for 2-category 2D data

#### K-Nearest Neighbors classification

- For a new point, find the k closest points from training data
- Labels of the k points "vote" to classify k = 5

Black = negative Red = positive



A nearest neighbor recognition example

#### Where in the World?



[Hays and Efros. **im2gps**: Estimating Geographic Information from a Single Image. CVPR 2008.] Slides: James Hays

#### Where in the World?



#### Where in the World?



#### 6+ million geotagged photos by 109,788 photographers



Annotated by Flickr users

#### Which scene properties are relevant?

#### **Spatial Envelope Theory of Scene Representation** Oliva & Torralba (2001)



## A scene is a single surface that can be represented by global (statistical) descriptors

Slide Credit: Aude Olivia

#### Global texture: capturing the "Gist" of the scene

## Capture global image properties while keeping some spatial information



Oliva & Torralba IJCV 2001, Torralba et al. CVPR 2003

#### Which scene properties are relevant?

- Gist scene descriptor
- **Color Histograms** L\*A\*B\* 4x14x14 histograms
- Texton Histograms 512 entry, filter bank based
- Line Features Histograms of straight line stats

#### Scene Matches







Croatia





Latvia











europe





France



Macau







Malta

Austria

[Hays and Efros. **im2gps**: Estimating Geographic Information from a Single Image. CVPR 2008.]







#### Scene Matches





[Hays and Efros. **im2gps**: Estimating Geographic Information from a Single Image. CVPR 2008.]







[Hays and Efros. im2gps: Estimating Geographic Information from a Single Image. CVPR-2008.]

#### Scene Matches



[Hays and Efros. **im2gps**: Estimating Geographic Information from a Single Image. CVPR 2008.]





[Hays and Efros. **im2gps**: Estimating Geographic Information from a Single Image. CVPR 2008.]

#### The Importance of Data



[Hays and Efros. **im2gps**: Estimating Geographic Information from a Single Image. CVPR 2008.] Slides: James Hays

#### **Feature Performance**



**Slides: James Hays** 

#### Nearest neighbors: pros and cons

#### • Pros:

- Simple to implement
- Flexible to feature / distance choices
- Naturally handles multi-class cases
- Can do well in practice with enough representative data

#### • Cons:

- Large search problem to find nearest neighbors
- Storage of data
- Must know we have a meaningful distance function

#### Linear classifiers




### Lines in R<sup>2</sup>



Let 
$$\mathbf{W} = \begin{bmatrix} a \\ c \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} x \\ y \end{bmatrix}$$

ax + cy + b = 0  $\mathbf{1}$   $\mathbf{w} \cdot \mathbf{x} + b = 0$ 



Let 
$$\mathbf{W} = \begin{bmatrix} a \\ c \end{bmatrix} \quad \mathbf{X} = \begin{bmatrix} x \\ y \end{bmatrix}$$

ax + cy + b = 0  $\downarrow$   $\mathbf{w} \cdot \mathbf{x} + b = 0$ 



Let  $\mathbf{W} = \begin{vmatrix} a \\ c \end{vmatrix} \quad \mathbf{X} = \begin{vmatrix} x \\ y \end{vmatrix}$ 

ax + cy + b = 0 $\mathbf{w} \cdot \mathbf{x} + b = 0$ 

distance from point to line



#### Linear classifiers

Find linear function to separate positive and negative examples



 $\mathbf{x}_i$  positive:  $\mathbf{x}_i \cdot \mathbf{w} + b \ge 0$  $\mathbf{x}_i$  negative:  $\mathbf{x}_i \cdot \mathbf{w} + b < 0$ 

Which line is best?

## Support Vector Machines (SVMs)



- Discriminative classifier based on optimal separating line (for 2d case)
- Maximize the *margin* between the positive and negative training examples

#### Support vector machines Want line that maximizes the margin.



 $\mathbf{x}_i$  positive  $(y_i = 1)$ : $\mathbf{x}_i \cdot \mathbf{w} + b \ge 1$  $\mathbf{x}_i$  negative  $(y_i = -1)$ : $\mathbf{x}_i \cdot \mathbf{w} + b \le -1$ For support, vectors, $\mathbf{x}_i \cdot \mathbf{w} + b = \pm 1$ 

C. Burges, <u>A Tutorial on Support Vector Machines for Pattern Recognition</u>, Data Mining and Knowledge Discovery, 1998

#### Support vector machines Want line that maximizes the margin.



 $\mathbf{x}_i$  positive $(y_i = 1)$ :  $\mathbf{x}_i \cdot \mathbf{w} + b \ge 1$  $\mathbf{x}_i$  negative  $(y_i = -1)$ :  $\mathbf{x}_i \cdot \mathbf{w} + b \leq -1$  $\mathbf{x}_i \cdot \mathbf{w} + b = \pm 1$ For support, vectors,  $|\mathbf{x}_i \cdot \mathbf{w} + b|$ Distance between point and line:  $\|\mathbf{w}\|$ For support vectors:  $\frac{\mathbf{w}^{T}\mathbf{x}+b}{\|\mathbf{w}\|} = \frac{\pm 1}{\|\mathbf{w}\|} \qquad M = \left|\frac{1}{\|\mathbf{w}\|} - \frac{-1}{\|\mathbf{w}\|}\right| = \frac{2}{\|\mathbf{w}\|}$ 

#### Support vector machines Want line that maximizes the margin.



 $\mathbf{x}_i$  positive  $(y_i = 1)$ : $\mathbf{x}_i \cdot \mathbf{w} + b \ge 1$  $\mathbf{x}_i$  negative  $(y_i = -1)$ : $\mathbf{x}_i \cdot \mathbf{w} + b \le -1$ For support, vectors, $\mathbf{x}_i \cdot \mathbf{w} + b = \pm 1$ Distance between point<br/>and line: $\frac{|\mathbf{x}_i \cdot \mathbf{w} + b|}{||\mathbf{w}||}$ Therefore, the margin is  $2 / ||\mathbf{w}||$ 

## Finding the maximum margin line

- 1. Maximize margin  $2/||\mathbf{w}||$
- 2. Correctly classify all training data points:

 $\mathbf{x}_{i} \text{ positive}(y_{i} = 1): \qquad \mathbf{x}_{i} \cdot \mathbf{w} + b \ge 1$  $\mathbf{x}_{i} \text{ negative}(y_{i} = -1): \qquad \mathbf{x}_{i} \cdot \mathbf{w} + b \le -1$ 

• Quadratic optimization problem: • Minimize  $\frac{1}{2}\mathbf{w}^T\mathbf{w}$ Subject to  $y_i(\mathbf{w}\cdot\mathbf{x}_i+b) \ge 1$ 

# Finding the maximum margin line



Finding the maximum margin line

• Solution:  $\mathbf{w} = \sum_{i} \alpha_{i} y_{i} \mathbf{x}_{i}$ 

$$b = y_i - \mathbf{w} \cdot \mathbf{x}_i \quad \text{(for any support} \\ \text{vector)} \quad \mathbf{w} \cdot \mathbf{x} + b = \sum_i \alpha_i y_i \mathbf{x}_i \cdot \mathbf{x} + b$$

• Classification function:

$$f(x) = \operatorname{sign} (\mathbf{w} \cdot \mathbf{x} + \mathbf{b})$$
$$= \operatorname{sign} \left( \sum_{i} \alpha_{i} \mathbf{x}_{i} \cdot \mathbf{x} + \mathbf{b} \right)$$

If f(x) < 0, classify as negative, if f(x) > 0, classify as positive

C. Burges, <u>A Tutorial on Support Vector Machines for Pattern Recognition</u>, Data Mining and Knowledge Discovery, 1998

## Questions

- What if the features are not 2d?
- What if the data is not linearly separable?
- What if we have more than just two categories?

## Questions

• What if the features are not 2d?

 Generalizes to d-dimensions – replace line with "hyperplane"

- What if the data is not linearly separable?
- What if we have more than just two categories?

#### Person detection with HoG's & linear SVM's



• Map each grid cell in the input window to a histogram counting the gradients per orientation.

• Train a linear SVM using training set of pedestrian vs. non-pedestrian windows.

Code available: http://pascal.inrialpes.fr/soft/olt/

Dalal & Triggs, CVPR 2005

# Person detection with HoG's & linear SVM's



- Histograms of Oriented Gradients for Human Detection, <u>Navneet Dalal</u>, <u>Bill Triggs</u>, International Conference on Computer Vision & Pattern Recognition - June 2005
- http://lear.inrialpes.fr/pubs/2005/DT05/

## Questions

- What if the features are not 2d?
- What if the data is not linearly separable?
- What if we have more than just two categories?

#### Non-linear SVMs

- Datasets that are linearly separable with some noise work out great:
- But what are we going to do if the dataset is just too hard?

х

х

 How about... mapping data to a higher-dimensional space:



#### Non-linear SVMs: feature spaces

 General idea: the original input space can be mapped to some higher-dimensional feature space where the training set is separable:



Slide from Andrew Moore's tutorial: http://www.autonlab.org/tutorials/svm.html

#### The "Kernel Trick"

- The linear classifier relies on dot product between vectors K(x<sub>i</sub>,x<sub>j</sub>)=x<sub>i</sub><sup>T</sup>x<sub>j</sub>
- If every data point is mapped into high-dimensional space via some transformation  $\Phi$ :  $x \rightarrow \phi(x)$ , the dot product becomes:

$$K(\mathbf{x}_i,\mathbf{x}_j) = \boldsymbol{\Phi}(\mathbf{x}_i)^{\mathsf{T}} \boldsymbol{\Phi}(\mathbf{x}_j)$$

• A *kernel function* is similarity function that corresponds to an inner product in some expanded feature space.

#### Example

2-dimensional vectors  $x=[x_1 \ x_2];$ let  $K(x_i,x_j)=(1 + x_i^T x_j)^2$ 

Need to show that  $K(x_i, x_i) = \varphi(x_i)^T \varphi(x_i)$ :  $K(x_i, x_i) = (1 + x_i^T x_i)^2$  $= 1 + x_{i1}^2 x_{i1}^2 + 2 x_{i1} x_{i1} x_{i2} x_{i2} + x_{i2}^2 x_{i2}^2 + 2 x_{i1} x_{i1} + 2 x_{i2} x_{i2}^2$  $= [1 \ x_{i1}^2 \ \sqrt{2} \ x_{i1} x_{i2} \ x_{i2}^2 \ \sqrt{2} x_{i1} \ \sqrt{2} x_{i2}]^T$  $\begin{bmatrix} 1 & x_{i1}^2 & \sqrt{2} & x_{i1} & x_{i2} & x_{i2}^2 & \sqrt{2} & x_{i1} & \sqrt{2} & x_{i2} \end{bmatrix}$  $= \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_i),$ where  $\varphi(\mathbf{x}) = \begin{bmatrix} 1 & x_1^2 & \sqrt{2} & x_1 x_2 & x_2^2 & \sqrt{2} & x_1 & \sqrt{2} & x_2 \end{bmatrix}$ 

from Andrew Moore's tutorial: http://www.autonlab.org/tutorials/svm.html

• The kernel trick: instead of explicitly computing the lifting transformation  $\varphi(\mathbf{x})$ , define a kernel function K such that

$$K(\mathbf{x}_i, \mathbf{x}_j) = \boldsymbol{\varphi}(\mathbf{x}_i) \cdot \boldsymbol{\varphi}(\mathbf{x}_j)$$

• This gives a nonlinear decision boundary in the original feature space:

$$\sum_{i} \alpha_{i} y_{i} K(\mathbf{x}_{i}, \mathbf{x}) + b$$

Examples of kernel functions  

$$K(x_i, x_j) = x_i^T x_j$$

Linear:

$$K(x_i, x_j) = \exp(-\frac{\|x_i - x_j\|^2}{2\sigma^2})$$

Gaussian RBF:

$$K(x_i, x_j) = \sum_k \min(x_i(k), x_j(k))$$
  
• Histogram intersection:

# SVMs for recognition

- 1. Define your representation for each example.
- 2. Select a kernel function.
- 3. Compute pairwise kernel values between labeled examples
- 4. Use this "kernel matrix" to solve for SVM support vectors & weights.
- 5. To classify a new example: compute kernel values between new input and support vectors, apply weights, check sign of output.



#### Example: learning gender with SVMs



Moghaddam and Yang, Learning Gender with Support Faces, TPAMI 2002.

Moghaddam and Yang, Face & Gesture 2000.





Processed faces

Moghaddam and Yang, Learning Gender with Support Faces, TPAMI 2002.

### Learning gender with SVMs

- Training examples:
  - 1044 males
  - 713 females
- Experiment with various kernels, select Gaussian RBF

$$K(\mathbf{x}_{i}, \mathbf{x}_{j}) = \exp(-\frac{\|\mathbf{x}_{i} - \mathbf{x}_{j}\|^{2}}{2\sigma^{2}})$$

#### **Support Faces**



Moghaddam and Yang, Learning Gender with Support Faces, TPAMI 2002.

# **Classifier Performance**

Classifier	Error Rate		
	Overall	Male	Female
SVM with RBF kernel	3.38%	2.05%	4.79%
SVM with cubic polynomial kernel	4.88%	4.21%	5.59%
Large Ensemble of RBF	5.54%	4.59%	6.55%
Classical RBF	7.79%	6.89%	8.75%
Quadratic classifier	10.63%	9.44%	11.88%
Fisher linear discriminant	13.03%	12.31%	13.78%
Nearest neighbor	27.16%	26.53%	28.04%
Linear classifier	58.95%	58.47%	59.45%

Moghaddam and Yang, Learning Gender with Support Faces, TPAMI 2002.

Gender perception experiment: How well can humans do?

- Subjects:
  - 30 people (22 male, 8 female)
  - Ages mid-20's to mid-40's
- Test data:
  - 254 face images (6 males, 4 females)
  - Low res and high res versions
- Task:
  - Classify as male or female, forced choice
  - No time limit

Moghaddam and Yang, Face & Gesture 2000.

#### Gender perception experiment: How well can humans do?



Results  $\rightarrow$  High-Res Low-Res 6.54% 30.7% Error Error

$$\sigma = 3.7\%$$

Moghaddam and Yang, Face & Gesture 2000.

#### Human vs. Machine

% Error Rates



 SVMs performed better than any single human test subject, at either resolution

Figure 6. SVM vs. Human performance

#### Hardest examples for humans



#### **Top five human misclassifications**

Moghaddam and Yang, Face & Gesture 2000.

## Questions

- What if the features are not 2d?
- What if the data is not linearly separable?
- What if we have more than just two categories?

# Multi-class SVMs

- Achieve multi-class classifier by combining a number of binary classifiers
- One vs. all
  - Training: learn an SVM for each class vs. the rest
  - Testing: apply each SVM to test example and assign to it the class of the SVM that returns the highest decision value

#### One vs. one

- Training: learn an SVM for each pair of classes
- Testing: each learned SVM "votes" for a class to assign to the test example

# SVMs: Pros and cons

- Pros
  - Many publicly available SVM packages: <u>http://www.kernel-machines.org/software</u>
  - <u>http://www.csie.ntu.edu.tw/~cjlin/libsvm/</u>
  - Kernel-based framework is very powerful, flexible
  - Often a sparse set of support vectors compact at test time
  - Work very well in practice, even with very small training sample sizes
- Cons
  - No "direct" multi-class SVM, must combine two-class SVMs
  - Can be tricky to select best kernel function for a problem
  - Computation, memory
    - During training time, must compute matrix of kernel values for every pair of examples
    - Learning can take a very long time for large-scale problems
# Summary

- Discriminative classifiers
  - Boosting
  - Nearest neighbors
  - Support vector machines
- Useful for object recognition when combined with "window-based" or holistic appearance descriptors

# Global window-based appearance representations



- These examples are truly global; each pixel in the window contributes to the representation.
- Classifier can account for relative relevance...
- When might this not be ideal?

 Part-based and local feature models for generic object recognition

# Part-based and local feature models for recognition



#### Main idea:

Rather than a representation based on holistic appearance, decompose the image into:

- local parts or patches, and
- their relative spatial relationships

# Part-based and local feature models for recognition



We'll look at three forms:

- 1. Bag of words (no geometry)
- 2. Implicit shape model (star graph for spatial model)
- Constellation model (fully connected graph for spatial model)

### Bag-of-words model

- Summarize entire image based on its distribution (histogram) of word occurrences.
  - Total freedom on spatial positions, relative geometry.
  - Vector representation easily usable by most classifiers.



#### **Bag-of-words model**



Our in-house database contains 1776 images in seven classes<sup>1</sup>: faces, buildings, trees, cars, phones, bikes and books. Fig. 2 shows some examples from this dataset.

Csurka et al. Visual Categorization with Bags of Keypoints, 2004

#### Words as parts



#### All local features

Local features from two selected clusters occurring in this image

Csurka et al. 2004

#### Naïve Bayes model for classification



What assumptions does the model make, and what are their significance?



#### **Confusion matrix**

True classes $\rightarrow$	faces	buildings	trees	cars	phones	bikes	books
faces	76	4	2	3	4	4	13
buildings	2	44	5	0	5	1	3
trees	3	2	80	0	0	5	0
cars	4	1	0	75	3	1	4
phones	9	15	1	16	70	14	11
bikes	2	15	12	0	8	73	0
books	4	19	0	6	7	2	69

Example bag of words + Naïve Bayes classification results for generic categorization of objects

### Clutter...or context?







Kristen Grauman

## Sampling strategies



#### Specific object



# Sampling strategies



Sparse, at interest points



Multiple interest operators





Dense, uniformly

Randomly

- To find specific, textured objects, sparse sampling from interest points more reliable.
- Multiple complementary interest operators offer more image coverage.
- For object categorization, dense sampling offers better coverage.

[See Nowak, Jurie & Triggs, ECCV 2006]

Kristen Grauman

Image credits: F-F. Li, E. Nowak, J. Sivic

#### Local feature correspondence for generic object categories



Kristen Grauman

# Local feature correspondence for generic object categories

- Comparing bags of words histograms coarsely reflects agreement between local "parts" (patches, words).
- *But* choice of quantization directly determines what we consider to be similar...



#### Pyramid match: main idea



Feature space partitions serve to "match" the local descriptors within successively wider regions.

[Grauman & Darrell, ICCV 2005]

#### Pyramid match: main idea





$$\mathcal{I}(H_X, H_Y) = \sum_j \min \left( H_X(j), H_Y(j) \right)$$
$$= 3$$

Histogram intersection counts number of possible matches at a given partitioning.

[Grauman & Darrell, ICCV 2005]

#### Pyramid match kernel





Optimal match: O(m<sup>3</sup>) Pyramid match: O(mL)



optimal partial matching

[Grauman & Darrell, ICCV 2005]

 $\mathbf{X} = \{ \vec{\mathbf{x}}_1, \dots, \vec{\mathbf{x}}_m \} \quad \mathbf{Y} = \{ \vec{\mathbf{y}}_1, \dots, \vec{\mathbf{y}}_n \}$ 

#### Unordered sets of local features: **No** spatial layout preserved!



Too much?

Too little?

# Spatial pyramid match

- Make a pyramid of bag-of-words histograms.
- Provides some loose (global) spatial layout information



Sum over PMKs computed in *image coordinate* space, one per word.

[Lazebnik, Schmid & Ponce, CVPR 2006]

Kristen Grauman

# Spatial pyramid match

Captures scene categories well---texture-like patterns but with some variability in the positions of all the local pieces.



**Confusion matrix** 

#### Kristen Grauman

# Spatial pyramid match

Captures scene categories well---texture-like patterns but with some variability in the positions of all the local pieces.

				Strong features	
				(vocabulary size: 200)	
			Level	Single-level	Pyramid
			$0(1 \times 1)$	$72.2 \pm 0.6$	
level 0	level 1	level 2	$1 (2 \times 2)$	$77.9 \pm 0.6$	$79.0 \pm 0.5$
			$2(4 \times 4)$	$79.4 \pm 0.3$	<b>81.1</b> ±0.3
			3 (8 × 8)	77.2 $\pm 0.4$	$80.7 \pm 0.3$

# Part-based and local feature models for recognition



We'll look at three forms:

- 1. Bag of words (no geometry)
- 2. Implicit shape model (star graph for spatial model)
- Constellation model (fully connected graph for spatial model)

# Shape representation in part-based models

"Star" shape model



> e.g. implicit shape model> Parts mutually independent

N image features, P parts in the model

# Implicit shape models

 Visual vocabulary is used to index votes for object position [a visual word = "part"]





visual codeword with displacement vectors

training image annotated with object localization info

B. Leibe, A. Leonardis, and B. Schiele, <u>Combined Object Categorization and Segmentation</u> <u>with an Implicit Shape Model</u>, ECCV Workshop on Statistical Learning in Computer Vision 2004

# Implicit shape models Visual vocabulary is used to index votes for object position [a visual word = "part"]



test image

B. Leibe, A. Leonardis, and B. Schiele, <u>Combined Object Categorization and Segmentation</u> <u>with an Implicit Shape Model</u>, ECCV Workshop on Statistical Learning in Computer Vision 2004

# Implicit shape models: Training

- 1. Build vocabulary of patches around extracted
  - interest mints using clustering  $\leftarrow$

# Implicit shape models: Training

- 1. Build vocabulary of patches around extracted interest points using clustering
- 2. Map the patch around each interest point to closest word



# Implicit shape models: Training

- 1. Build vocabulary of patches around extracted interest points using clustering
- 2. Map the patch around each interest point to closest word
- 3. For each word, store all positions it was found, relative to object center



# 1. Given new test image, extract patches, match to

- Given new test image, extract patches, match to vocabulary words
- 2. Cast votes for possible positions of object center
- 3. Search for maxima in voting space
- 4. (Extract weighted segmentation mask based on stored masks for the codebook occurrences)

# Implicit shape models: Testing






















## **Example: Results on Cows**





## Example: Results on Cows





#### K. Grauman, B. Leibe

## Detection Results Qualitative Performance

### - Recognizes different kinds of objects

- Robust to clutter, occlusion, noise, low contrast



# Shape representation in part-based models

"Star" shape model



> e.g. implicit shape model> Parts mutually independent

Fully connected constellation model



> e.g. Constellation Model> Parts fully connected

N image features, P parts in the model

## Probabilistic constellation model





#### Candidate parts

## Probabilistic constellation model

P(image | object) = P(appearance, shape | object)



Part 3

Part 2

## Probabilistic constellation model

P(image | object) = P(appearance, shape | object)

 $= \max_{h} P(appearance | h, object) p(shape | h, object) p(h | object)$ 

h: assignment of features to parts

Part 1



Part 3

## Example results from constellation model: data from four categories





Faces







**Motorbikes** 















### Spotted cats





















Slide from Li Fei-Fei http://www.vision.caltech.edu/feifeili/Resume.htm







Appearance: 10 patches closest to mean for each part

Fergus et al. CVPR 2003

#### Face model



Appearance: 10 patches closest to mean for each part

## Recognition results

Test images: size of circles indicates score of hypothesis

Fergus et al. CVPR 2003

Kristen Grauman



Appearance: 10 patches closest to mean for each part

Fergus et al. CVPR 2003



Appearance: 10 patches closest to mean for each part

### results

Kristen Grauman

Fergus et al. CVPR 2003

## Comparison



class	bag of features	bag of features	Part-based model
	Zhang et al. (2005)	Willamowski et al. (2004)	Fergus et al. (2003)
airplanes	98.8	97.1	90.2
cars (rear)	98.3	98.6	90.3
cars (side)	95.0	87.3	88.5
faces	100	99.3	96.4
motorbikes	98.5	98.0	92.5
spotted cats	97.0		90.0

# Shape representation in part-based models

"Star" shape model



- » e.g. implicit shape model
- > Parts mutually independent
- > Recognition complexity: O(NP)
- > Method: Gen. Hough Transform

Fully connected constellation model



- » e.g. Constellation Model
- » Parts fully connected
- > Recognition complexity: O(N<sup>P</sup>)
- > Method: Exhaustive search

N image features, P parts in the model

#### Slide credit: Rob Fergus

### Summary:

# part-based and local feature models for generic object recognition

- **Histograms of visual words** to capture global or local layout in the bag-of-words framework
  - Pyramid match kernels
  - Powerful in practice for image recognition
- **Part-based models** encode category's part appearance together with 2d layout and allow detection within cluttered image
  - "implicit shape model": shape based on layout of all parts relative to a reference part; Generalized Hough for detection
  - "constellation model": explicitly model mutual spatial layout between all pairs of parts; exhaustive search for best fit of features to parts

## **Recognition models**



recognition by alignment

Categories: Holistic appearance models (and sliding window detection)





Categories: Local feature and part-based models

Kristen Grauman