Recognizing and Learning Object Categories

Based on work by R. Fergus, P. Perona, A. Zisserman, J. Ponce, S. Lazebnik, C. Schmid, F. DiMaio, and others

Traditional Problem: Single Object Recognition



Most Objects Exhibit Considerable Intra-Class Variability



Task: Recognition of object categories











| Object categorization: the statistical viewpoint | | | | | |
|---|---------------------|-------------|--|--|--|
| p(zebra image) vs. p(no zebra/image) | | | | | |
| Sayes's rule: | | | | | |
| p(zebra image) | p(image zebra) | p(zebra) | | | |
| p(no zebra image) | p(image no zebra) | p(no zebra) | | | |
| posterior ratio | likelihood ratio | prior ratio | | | |





| Generative | | | | | |
|------------|--|--------------------|--|--|--|
| § Mode | § Model $p(image zebra)$ and $p(image no zebra)$ | | | | |
| | | | | | |
| | p(image zebra) | p(image no zebra)) | | | |
| 826 | Low | Middle | | | |
| | High | Middle Low | | | |

Three main issues

§ Representation

 $\ensuremath{\mathbb{S}}$ How to represent an object category

§ Learning

§ How to form the classifier, given training data

S Recognition

§ How the classifier is to be used on novel data

Approach 2: Generative Methods using Bag of Words Models

- S An image is represented by a collection of "visual words" and their corresponding counts given a universal dictionary
- S Object categories are modeled by the distributions of these visual words
- S Although "bag of words" models can use both generative and discriminative approaches, here we will focus on generative models



Approach 3: Generative Methods using Part-Based Models

- S An object in an image is represented by a collection of parts, characterized by both their visual appearances and locations
- S Object categories are modeled by the appearance and spatial distributions of these characteristic parts
- § Issues for such models include efficient methods for finding correspondences between the object and the scene







Model Structure

S Model shape using Gaussian distribution on image location between parts and scale of each part



- Model **appearance** as patches of pixel intensities
- S Represent object class as graph of P image patches with parameters θ



Representation of Occlusion

- § Explicit
 - § Additional match of each part to missing state
- § Implicit
 - § Truncated minimum probability of appearance











| Object categorization: the statistical viewpoint | | | | | |
|--|--|--|--|--|--|
| p(zebra image) vs. p(no zebra/image) | | | | | |
| § Bayes rule: | | | | | |
| $\frac{p(zebra image)}{p(no \ zebra image)} =$ posterior ratio | $\underbrace{\frac{p(image \mid zebra)}{p(image \mid no \ zebra)}}_{likelihood \ ratio}$ | $\frac{p(zebra)}{p(no \ zebra)}$ prior ratio | | | |

Model Structure • Assume prior ratio is known or learned • Find values for parameters θ that maximizes the likelihood ratio $p(X, S, A | \theta) = \sum_{h \in H} p(X, S, A, h | \theta)$ • H is the set of all valid correspondences of image features to model parts, so $|H| = O(N^P)$ • Factor the likelihood to simplify computation (using Chain Rule)







Recognition

- § For each of *P* parts, run template over all locations in image
- S Detect local maxima, giving possible locations of each part
- § Given learned model, find maximum likelihood ratio of $p(\mathbf{X}, \mathbf{S}, \mathbf{A}|\theta)/p(\mathbf{X}, \mathbf{S}, \mathbf{A}|\theta_{bg})$ for all possible correspondences $-O(N^2P)$ where N = number of locations of each part in image
- § If greater than a threshold, signify object detected























Probabilistic Parts and Structure Models Summary

- § Correspondence problem
- § Efficient methods for large # parts and # positions in image
- § Challenge to get representation with desired invariance
- § Minimal supervision
- § Future directions:
 - § Multiple views
 - § Approaches to learning
 - § Multiple category training

ROC equal error rates

Pre-scaled data (identical settings):

| | | | | Model | | |
|--------------|--------------------------|----------------------------|------------|-------|-----------|--------------|
| Dataset | Total size of dataset | ~ Object width (pixels) | Motorbikes | Faces | Airplanes | Spotted Cats |
| Motorbikes | 800 | 200 | 92.5 | 50 | 51 | 56 |
| Faces | 435 | 300 | 33 | 96.4 | 32 | 32 |
| Airplanes | 800 | 300 | 64 | 63 | 90.2 | 53 |
| Spotted Cats | 200 | 80 | 48 | 44 | 51 | 90.0 |

Scale-invariant learning and recognition:

| | Total size | Object size | Pre-scaled | Unscaled |
|-------------|------------|----------------|-------------|-------------|
| Dataset | of dataset | range (pixels) | performance | performance |
| Motorbikes | 800 | 200-480 | 95.0 | 93.3 |
| Airplanes | 800 | 200-500 | 94.0 | 93.0 |
| Cars (Rear) | 800 | 100-550 | 84.8 | 90.3 |