

Goal: Detect all instances of objects

EEE422/6082 Computational Vision

Object Category Detection

Ling Shao

Some slides from Derek Hoiem



Demo: face detection

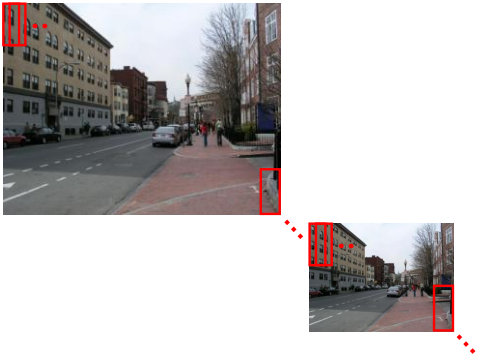


<http://demo.pittpatt.com/>

Influential Works in Detection

- Sung-Poggio (1994, 1998) : ~1450 citations
 - Basic idea of statistical template detection (I think), bootstrapping to get “face-like” negative examples, multiple whole-face prototypes (in 1994)
- Rowley-Baluja-Kanade (1996-1998) : ~2900
 - “Parts” at fixed position, non-maxima suppression, simple cascade, rotation, pretty good accuracy, fast
- Schneiderman-Kanade (1998-2000,2004) : ~1250
 - Careful feature engineering, excellent results, cascade
- Viola-Jones (2001, 2004) : ~6500
 - Haar-like features, Adaboost as feature selection, hyper-cascade, very fast, easy to implement
- Dalal-Triggs (2005) : 1025
 - Careful feature engineering, excellent results, HOG feature, online code
- Felzenszwalb-McAllester-Ramanan (2008)? 105 citations
 - Excellent template/parts-based blend

Sliding window detection

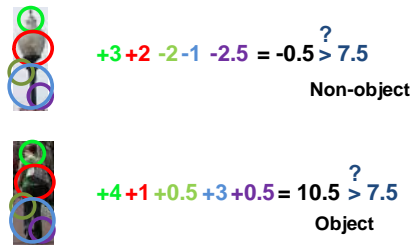


What the Detector Sees



Statistical Template

- Object model = log linear model of parts at fixed positions



Design challenges

- Part design
 - How to model appearance
 - Which “parts” to include
 - How to set part likelihoods
- How to make it fast
- How to deal with different viewpoints
- Implementation details
 - Window size
 - Aspect ratio
 - Translation/scale step size
 - Non-maxima suppression

Schneiderman and Kanade

Parts model

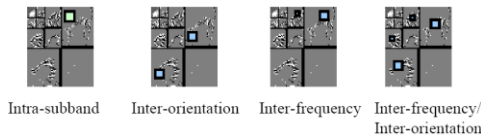
- Part = group of wavelet coefficients that are statistically dependent



Schneiderman and Kanade. [A Statistical Method for 3D Object Detection](#). (2000)

Parts: groups of wavelet coefficients

- Fixed parts within/across subbands



- 17 types of “parts” that can appear at each position
- Discretize wavelet coefficient to 3 values
- E.g., part with 8 coefficients has $3^8 = 6561$ values

Training

- Create training data
 - Get positive and negative patches
 - Pre-process (optional), compute wavelet coefficients, discretize
 - Compute parts values
- Learn statistics
 - Compute ratios of histograms by counting for positive and negative examples
 - Reweight examples using Adaboost, recount, etc.
- Get more negative examples (bootstrapping)

Training multiple viewpoints



Train new detector for each viewpoint.



Testing

- 1) Processing:
 - a) Lighting correction (optional)
 - b) Compute wavelet coefficients, quantize
- 2) Slide window over each position/scale (2 pixels, $2^{1/4}$ scale)
 - a) Compute part values
 - b) Lookup likelihood ratios
 - c) Sum over parts
 - d) Threshold
- 3) Use faster classifier to prune patches (cascade...more on this later)
- 4) Non-maximum suppression

Results: faces



Table 1. Face detection with out-of-plane rotation

γ	Detection (all faces)	Detection (profiles)	False Detections
0.0	92.7%	92.8%	700
1.5	85.5%	86.4%	91
2.5	75.2%	78.6%	12

208 images with 441 faces, 347 in profile

Results: cars



Table 3. Car detection

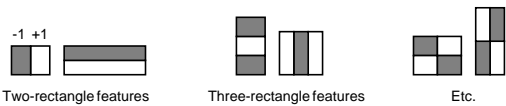
γ	Detections	False Detections
1.05	83%	7
1.0	86%	10
0.9	92%	71

Viola and Jones

Fast detection through two mechanisms

Integral Images

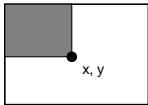
- “Haar-like features”
 - Differences of sums of intensity
 - Thousands, computed at various positions and scales within detection window



Viola and Jones. [Rapid Object Detection using a Boosted Cascade of Simple Features](#) (2001).

Integral Images

- `ii = cumsum(cumsum(lm, 1), 2)`



`ii(x,y)` = Sum of the values in the grey region

A	B	
	1	2
C	D	4
	3	

How to compute B-A?

How to compute A+D-B-C?

Adaboost as feature selection

- Create a large pool of parts (180K)
- “Weak learner” = feature + threshold + parity

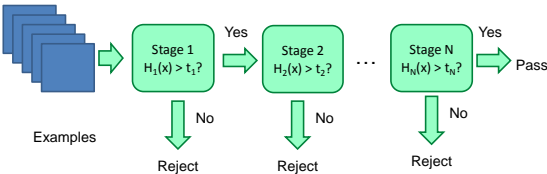
$$h_j(x) = \begin{cases} 1 & \text{if } p_j f_j(x) < p_j \theta_j \\ 0 & \text{otherwise} \end{cases}$$

- Choose weak learner that minimizes error on the weighted training set
- Reweight

Adaboost

- Given example images $(x_1, y_1), \dots, (x_m, y_m)$ where $y_i = 0, 1$ for negative and positive examples respectively.
- Initialize weights $w_{1,j} = \frac{1}{m} \frac{1}{2}$ for $y_i = 0, 1$ respectively, where m and l are the number of negatives and positives respectively.
- For $t = 1, \dots, T$:
 1. Normalize the weights.
$$w_{t,j} \leftarrow \frac{w_{t-1,j}}{\sum_{j=1}^l w_{t-1,j}}$$
so that $w_{t,j}$ is a probability distribution.
 2. For each feature, j , train a classifier $h_{t,j}$ which is restricted to using a single feature. The error is evaluated with respect to $w_{t,j}$, $e_{t,j} = \sum_i w_{t,j} |h_{t,j}(x_i) - y_i|$.
 3. Choose the classifier, $h_{t,j}$, with the lowest error $e_{t,j}$.
 4. Update the weights:
$$w_{t+1,j} := w_{t,j} \beta_j^{1-e_{t,j}}$$
where $e_{t,j} = 0$ if example x_i is classified correctly, $e_{t,j} = 1$ otherwise, and $\beta_j = \frac{e_{t,j}}{1-e_{t,j}}$.
- The final strong classifier is:
$$h(x) = \begin{cases} 1 & \sum_{t=1}^T \alpha_t h_t(x) \geq \frac{1}{2} \sum_{t=1}^T \alpha_t \\ 0 & \text{otherwise} \end{cases}$$
where $\alpha_t = \log \frac{1}{\beta_t}$.

Cascade for Fast Detection



- Choose threshold for low false negative rate
- Fast classifiers early in cascade
- Slow classifiers later, but most examples don't get there

Viola-Jones details

- 38 stages with 1, 10, 25, 50 ... features
 - 6061 total used out of 180K candidates
 - 10 features evaluated on average
- Examples
 - 4916 positive examples
 - 10000 negative examples collected after each stage
- Scanning
 - Scale detector rather than image
 - Scale steps = 1.25, Translation 1.0*s to 1.5*s
- Non-max suppression: average coordinates of overlapping boxes
- Train 3 classifiers and take vote

Viola Jones Results



Detector	False detections						
	10	31	50	65	78	95	167
Viola-Jones	76.1%	88.4%	91.4%	92.0%	92.1%	92.9%	93.9%
Viola-Jones (voting)	81.1%	89.7%	92.1%	93.1%	93.1%	93.2%	93.7%
Rowley-Baluja-Kanade	-	-	86.0%	-	-	89.2%	90.1%
Schneidersman-Kanade	-	-	-	94.4%	-	-	-
Roth-Yang-Aluja	-	-	-	-	(94.8%)	-	-

MIT + CMU face dataset

Schneiderman later results

Schneiderman 2004

Viola-Jones 2001

Roth et al. 1999

Schneiderman-Kanade 2000

	89.7%	93.1%	94.4%	94.8%	95.7%
Bayesian Network *	1	8	19	36	56
Semi-Naive Bayes*	6	19	29	35	46
[6]	31	65	--	--	--
[7]*	--	--	--	78	--
[16]*	--	--	65	--	--

Table 2. False alarms as a function of recognition rate on the MIT-CMU Test Set for Frontal Face Detection. * indicates exclusion of the 5 images of hand-drawn faces.

Speed: frontal face detector

- Schneiderman-Kanade (2000): 5 seconds
- Viola-Jones (2001): 15 fps

Strengths and Weaknesses of Statistical Template Approach

- Strengths
- Works very well for non-deformable objects: faces, cars, upright pedestrians
 - Fast detection
- Weaknesses
- Not so well for highly deformable objects
 - Not robust to occlusion
 - Requires lots of training data

SK vs. VJ

- Schneiderman-Kanade
- Wavelet features
 - Log linear model via boosted histogram ratios
 - Bootstrap training
 - Two-stage cascade
 - NMS: Remove overlapping weak boxes
 - Slow but very accurate
- Viola-Jones
- Similar to Haar wavelets
 - Log linear model via boosted stubs
 - Bootstrap training
 - Multistage cascade, integrated into training
 - NMS: average coordinates of overlapping boxes
 - Less accurate but very fast

Things to remember

- Excellent results require careful feature engineering
- Sliding window for search
- Features based on differences of intensity (gradient, wavelet, etc.)
- Boosting for feature selection (also L1-logistic regression)
- Integral images, cascade for speed
- Bootstrapping to deal with many, many negative examples

