



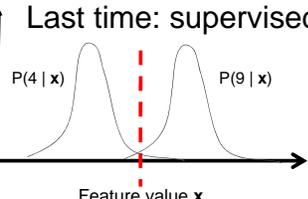
Window-based models for generic object detection

Monday, April 11
Kristen Grauman
UT-Austin

Previously

- Instance recognition
 - Local features: detection and description
 - Local feature matching, scalable indexing
 - Spatial verification
- Intro to generic object recognition
- Supervised classification
 - Main idea
 - Skin color detection example

Last time: 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

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

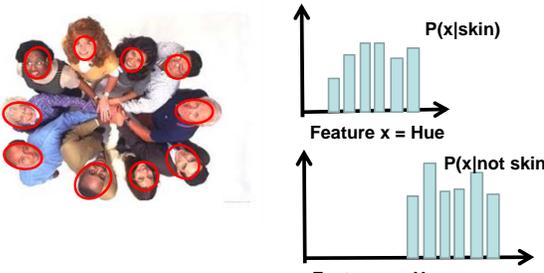
To classify a new point, choose class with lowest expected loss; i.e., choose "four" if

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

Kristen Grauman

Last time: Example: skin color classification

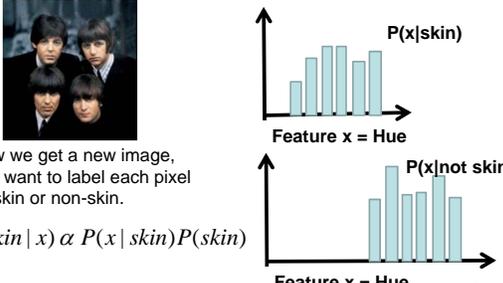
- We can represent a class-conditional density using a histogram (a "non-parametric" distribution)



Kristen Grauman

Last time: Example: skin color classification

- 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.

$$P(\text{skin} | x) \propto P(x | \text{skin}) P(\text{skin})$$

Kristen Grauman

Last time: Example: skin color classification

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



Brighter pixels \rightarrow higher probability of being skin

Classify pixels based on these probabilities

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

Kristen Grauman

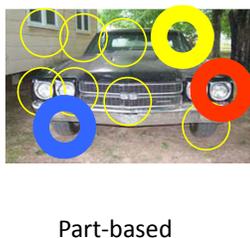
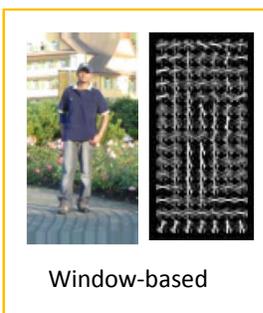
Today

- Window-based generic object detection
 - basic pipeline
 - boosting classifiers
 - face detection as case study

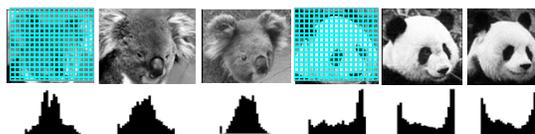
Generic category recognition: basic framework

- Build/train object model
 - Choose a representation
 - Learn or fit parameters of model / classifier
- Generate candidates in new image
- Score the candidates

Generic category recognition: representation choice



Window-based models Building an object model



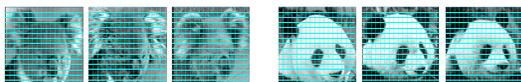
Simple holistic descriptions of image content

- grayscale / color histogram
- vector of pixel intensities

Kristen Grauman

Window-based models Building an object model

- Pixel-based representations sensitive to small shifts

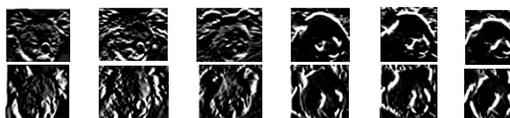


- Color or grayscale-based appearance description can be sensitive to illumination and intra-class appearance variation

Kristen Grauman

Window-based models Building an object model

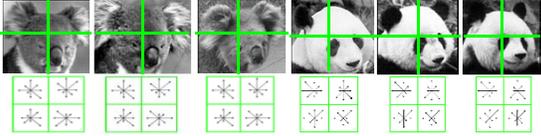
- Consider edges, contours, and (oriented) intensity gradients



Kristen Grauman

Window-based models Building an object model

- Consider edges, contours, and (oriented) intensity gradients

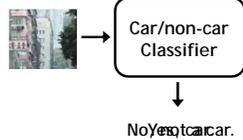


- Summarize local distribution of gradients with histogram
 - Locally orderless: offers invariance to small shifts and rotations
 - Contrast-normalization: try to correct for variable illumination

Kristen Grauman

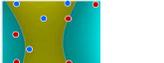
Window-based models Building an object model

Given the representation, train a binary classifier



Kristen Grauman

Discriminative classifier construction

<p>Nearest neighbor</p>  <p>Shakhnarovich, Viola, Darrell 2003 Berg, Berg, Malik 2005...</p>	<p>Neural networks</p>  <p>LeCun, Bottou, Bengio, Haffner 1998 Rowley, Baluja, Kanade 1998 ...</p>
<p>Support Vector Machines</p>  <p>Guyon, Vapnik Heisele, Serre, Poggio, 2001,...</p>	<p>Boosting</p>  <p>Viola, Jones 2001, Torralba et al. 2004, Opelt et al. 2006,...</p>
<p>Conditional Random Fields</p>  <p>McCallum, Freitag, Pereira 2000; Kumar, Hebert 2003 ...</p>	

Slide adapted from Antonio Torralba

Generic category recognition: basic framework

- Build/train object model
 - Choose a representation
 - Learn or fit parameters of model / classifier
- Generate candidates in new image
- Score the candidates

Window-based models Generating and scoring candidates



Kristen Grauman

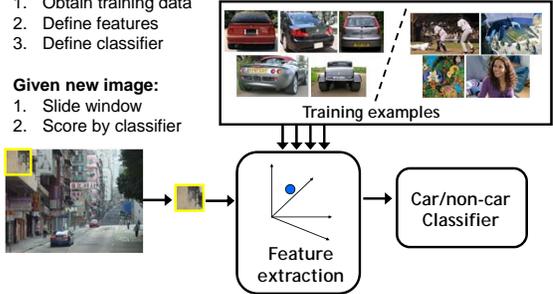
Window-based object detection: recap

Training:

- Obtain training data
- Define features
- Define classifier

Given new image:

- Slide window
- Score by classifier



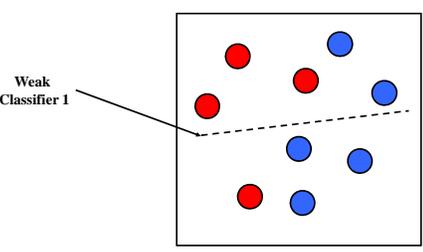
Kristen Grauman

Discriminative classifier construction

<p>Nearest neighbor</p>  <p>10⁶ examples</p> <p>Shakhnarovich, Viola, Darrell 2003 Berg, Berg, Malik 2005...</p>	<p>Neural networks</p>  <p>LeCun, Bottou, Bengio, Haffner 1998 Rowley, Baluja, Kanade 1998 ...</p>	
<p>Support Vector Machines</p>  <p>Guyon, Vapnik Heisele, Serre, Poggio, 2001,....</p>	<p>Boosting</p>  <p>Viola, Jones 2001, Torralba et al. 2004, Opelt et al. 2006,....</p>	<p>Conditional Random Fields</p>  <p>McCallum, Freitag, Pereira 2000; Kumar, Hebert 2003 ...</p>

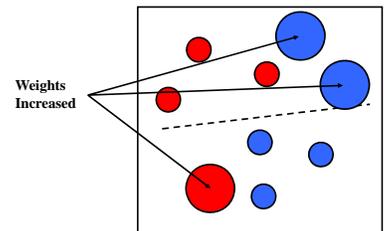
Slide adapted from Antonio Torralba

Boosting intuition

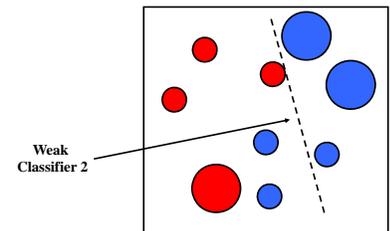


Slide credit: Paul Viola

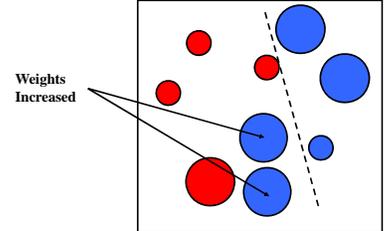
Boosting illustration



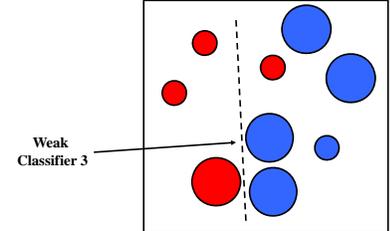
Boosting illustration



Boosting illustration

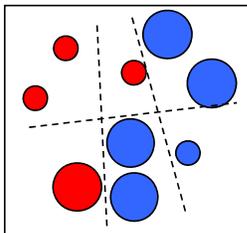


Boosting illustration



Boosting illustration

Final classifier is a combination of weak classifiers



Boosting: training

- Initially, weight each training example equally
- In each boosting round:
 - Find the weak learner that achieves the lowest *weighted* training error
 - Raise weights of training examples misclassified by current weak learner
- Compute final classifier as linear combination of all weak learners (weight of each learner is directly proportional to its accuracy)
- Exact formulas for re-weighting and combining weak learners depend on the particular boosting scheme (e.g., AdaBoost)

Slide credit: Lana Lazebnik

Boosting: pros and cons

- Advantages of boosting
 - Integrates classification with feature selection
 - Complexity of training is linear in the number of training examples
 - Flexibility in the choice of weak learners, boosting scheme
 - Testing is fast
 - Easy to implement
- Disadvantages
 - Needs many training examples
 - Often found not to work as well as an alternative discriminative classifier, support vector machine (SVM)
 - especially for many-class problems

Slide credit: Lana Lazebnik

Viola-Jones face detector

ACCEPTED CONFERENCE ON COMPUTER VISION AND PATTERN RECOGNITION 2001

Rapid Object Detection using a Boosted Cascade of Simple Features

Paul Viola
viola@merl.com
Mitsubishi Electric Research Labs
201 Broadway, 8th FL
Cambridge, MA 02139

Michael Jones
mjones@crl.dec.com
Compaq CRL
One Cambridge Center
Cambridge, MA 02142

Abstract

This paper describes a machine learning approach for vi-

ected at 15 frames per second on a conventional 700 MHz Intel Pentium III. In other face detection systems, auxiliary information, such as image differences in video sequences,

Viola-Jones face detector

Main idea:

- Represent local texture with efficiently computable "rectangular" features within window of interest
- Select discriminative features to be weak classifiers
- Use boosted combination of them as final classifier
- Form a cascade of such classifiers, rejecting clear negatives quickly

Kristen Grauman

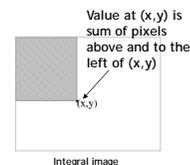
Viola-Jones detector: features



"Rectangular" filters

Feature output is difference between adjacent regions

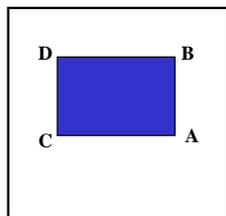
Efficiently computable with integral image: any sum can be computed in constant time.



Kristen Grauman

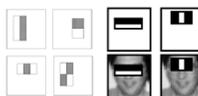
Computing sum within a rectangle

- Let A,B,C,D be the values of the integral image at the corners of a rectangle
- Then the sum of original image values within the rectangle can be computed as:
 $sum = A - B - C + D$
- Only 3 additions are required for any size of rectangle!



Lana Lazebnik

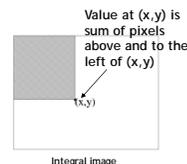
Viola-Jones detector: features



"Rectangular" filters
 Feature output is difference between adjacent regions

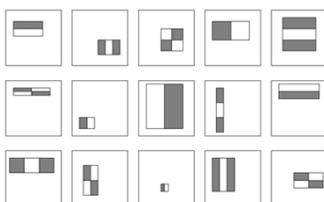
Efficiently computable with integral image: any sum can be computed in constant time

Avoid scaling images → scale features directly for same cost



Kristen Grauman

Viola-Jones detector: features



Considering all possible filter parameters: position, scale, and type:

180,000+ possible features associated with each 24 x 24 window

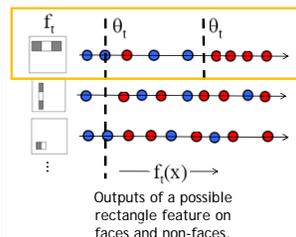
Which subset of these features should we use to determine if a window has a face?

Use AdaBoost both to select the informative features and to form the classifier

Kristen Grauman

Viola-Jones detector: AdaBoost

- Want to select the single rectangle feature and threshold that best separates positive (faces) and negative (non-faces) training examples, in terms of weighted error.



Resulting weak classifier:

$$h_t(x) = \begin{cases} +1 & \text{if } f_t(x) > \theta_t \\ -1 & \text{otherwise} \end{cases}$$

For next round, reweight the examples according to errors, choose another filter/threshold combo.

Kristen Grauman

- Given example images $(x_1, y_1), \dots, (x_n, y_n)$ where $y_i = 0, 1$ for negative and positive examples respectively.
- Initialize weights $w_{1,j} = \frac{1}{2m}, \frac{1}{2l}$ for $y_i = 0, 1$ respectively, where m and l are the number of negatives and positives respectively.
- For $t = 1, \dots, T$:
 - Normalize the weights,

$$w_{t,j} \leftarrow \frac{w_{t,j}}{\sum_{j=1}^n w_{t,j}}$$
 so that w_t is a probability distribution.
 - 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 , $\epsilon_j = \sum_i w_{t,j} |h_{t,j}(x_i) - y_i|$.
 - Choose the classifier, $h_{t,i}$, with the lowest error $\epsilon_{t,i}$.
 - Update the weights:

$$w_{t+1,i} = w_{t,i} \beta_i^{\pm 1}$$
 where $\epsilon_i = 0$ if example x_i is classified correctly, $\epsilon_i = 1$ otherwise, and $\beta_i = \frac{\epsilon_i}{1 - \epsilon_i}$.
- The final strong classifier is:

$$h(x) = \begin{cases} 1 & \sum_{t=1}^T \alpha_t h_t(x) \geq \frac{1}{2} \sum_{t=1}^T \alpha_t \\ \text{otherwise} \end{cases}$$
 where $\alpha_t = \log \frac{1}{\beta_t}$

AdaBoost Algorithm

Start with uniform weights on training examples



For T rounds

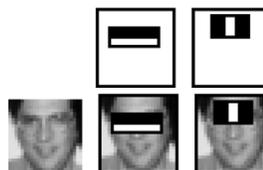
Evaluate weighted error for each feature, pick best.

Re-weight the examples:
 Incorrectly classified → more weight
 Correctly classified → less weight

Final classifier is combination of the weak ones, weighted according to error they had.

Freund & Schapire 1995

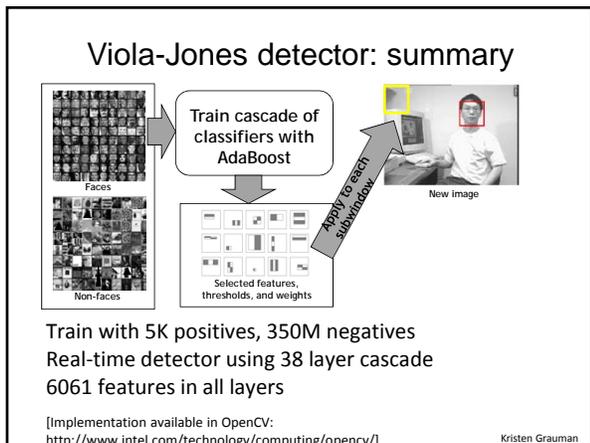
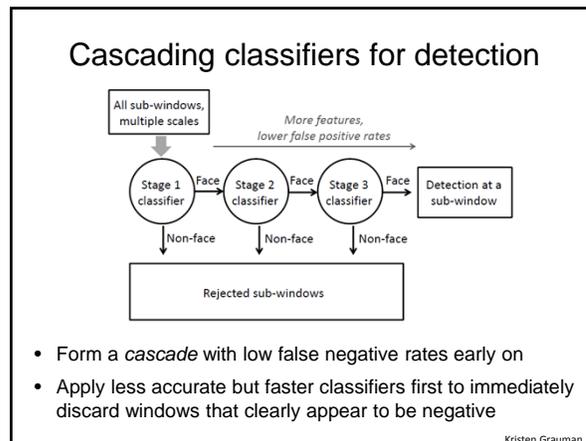
Viola-Jones Face Detector: Results



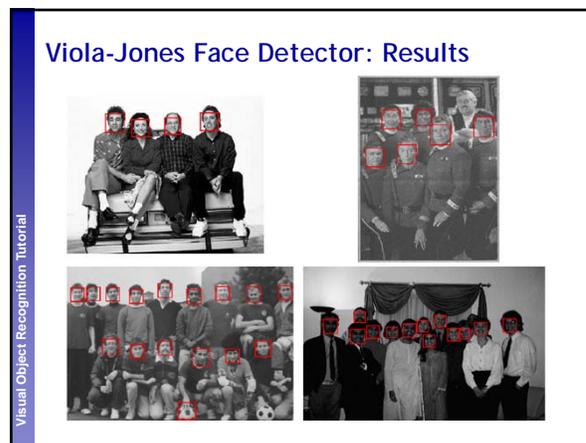
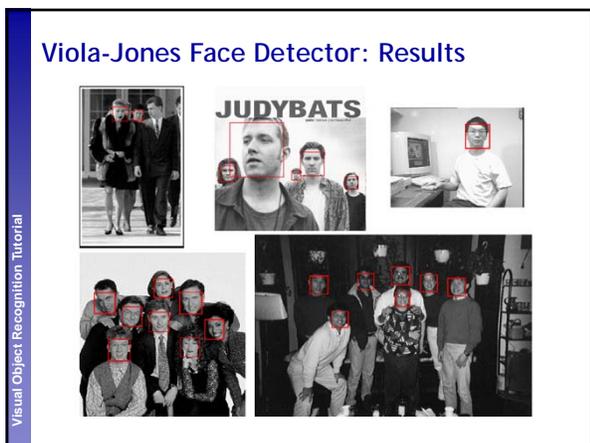
First two features selected

Visual Object Recognition Tutorial

- Even if the filters are fast to compute, each new image has a lot of possible windows to search.
- How to make the detection more efficient?



- ### Viola-Jones detector: summary
- A seminal approach to real-time object detection
 - Training is slow, but detection is very fast
 - Key ideas
 - *Integral images* for fast feature evaluation
 - *Boosting* for feature selection
 - *Attentional cascade* of classifiers for fast rejection of non-face windows
- P. Viola and M. Jones. [Rapid object detection using a boosted cascade of simple features](#). CVPR 2001.
- P. Viola and M. Jones. [Robust real-time face detection](#). IJCV 57(2), 2004.



Viola-Jones Face Detector: Results

Visual Object Recognition Tutorial

Detecting profile faces?

Can we use the same detector?

Visual Object Recognition Tutorial

Viola-Jones Face Detector: Results

Visual Object Recognition Tutorial

Example using Viola-Jones detector

Frontal faces detected and then tracked, character names inferred with alignment of script and subtitles.

Everingham, M., Sivic, J. and Zisserman, A. "Hello! My name is... Buffy" - Automatic naming of characters in TV video, BMVC 2006. <http://www.robots.ox.ac.uk/~vgg/research/nface/index.html>

ZDNet Asia
Where Technology Meets Business

TECH SHOWCASE
See how he stays with Cisco Collaboration Solutions

Home News Insight Reviews TechGuides Jobs Blogs Videos Community Downloads IT Library

Software Hardware Security Communications Business Internet Photo

Search ZDNet Asia

News > Internet

News from Countries/Region

- Singapore
- India
- China/HK
- Malaysia
- Philippines
- ASEAN
- Thailand
- Indonesia
- Apa Pacific

What's Hot | Latest News

- Is eBay being safer now?
- Report: Amazon may again be making Netflix
- Mobile maps out Jetpack and on transition plan
- Google begins search for Florida East lobbyist
- Google still thinks it can change China

advertisement

ZDNet Asia
TECH SHOWCASE
Cisco Collaboration Solution

Brought to you by CIS

Google now erases faces, license plates on Map Street View

By Elnor Ima, CNET News.com
Friday, August 24, 2007 9:32 PM

Google has gotten a lot of flack from privacy advocates for photographing faces and license plate numbers and displaying them on the Street View in Google Maps. Originally, the company said only people who identified themselves could ask the company to remove their image.

But Google has quietly changed that policy, partly in response to criticism, and now anyone can alert the company and have an image of a license plate or a recognizable face removed, not just the owner of the face or car, says Marissa Mayer, vice president of search products and user experience at Google.

"It's a good policy for us and also clarifies the intent of the product," she said in an interview following her keynote at the Search Engine Strategies conference in San Jose, Calif., Wednesday.

The policy change was made about 10 days after the launch of the product in late May, but was not publicly announced, according to Mayer. The company is removing images only when someone notifies them and not proactively, she said. "It was definitely a big policy change inside."

Consumer application: iPhoto 2009

<http://www.apple.com/ilife/iphoto/>

Slide credit: Lana Lazebni

Consumer application: iPhoto 2009

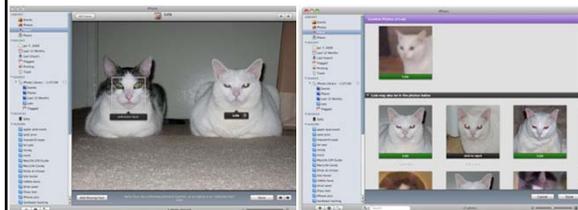
Things iPhoto thinks are faces



Slide credit: Lana Lazebnik

Consumer application: iPhoto 2009

Can be trained to recognize pets!



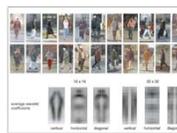
http://www.maclife.com/article/news/iphotos_faces_recognizes_cats

Slide credit: Lana Lazebnik

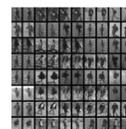
What other categories are amenable to *window-based representation*?

Pedestrian detection

- Detecting upright, walking humans also possible using sliding window's appearance/texture; e.g.,



SVM with Haar wavelets [Papageorgiou & Poggio, IJCV 2000]



Space-time rectangle features [Viola, Jones & Snow, ICCV 2003]



SVM with HoGs [Dalal & Triggs, CVPR 2005]

Visual Object Recognition Tutorial

Kristen Grauman

Window-based detection: strengths

- Sliding window detection and global appearance descriptors:
 - Simple detection protocol to implement
 - Good feature choices critical
 - Past successes for certain classes

Visual Object Recognition Tutorial

Kristen Grauman

Window-based detection: Limitations

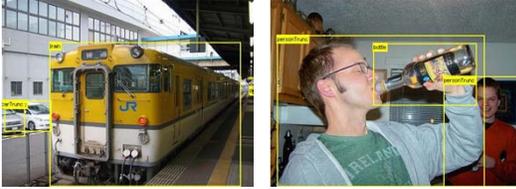
- High computational complexity
 - For example: 250,000 locations x 30 orientations x 4 scales = 30,000,000 evaluations!
 - If training binary detectors independently, means cost increases linearly with number of classes
- With so many windows, false positive rate better be low

Visual Object Recognition Tutorial

Kristen Grauman

Limitations (continued)

- Not all objects are "box" shaped



Kristen Grauman

Limitations (continued)

- Non-rigid, deformable objects not captured well with representations assuming a fixed 2d structure; or must assume fixed viewpoint
- Objects with less-regular textures not captured well with holistic appearance-based descriptions



Kristen Grauman

Limitations (continued)

- If considering windows in isolation, context is lost

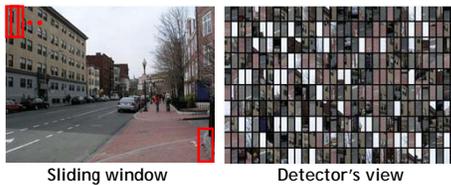


Figure credit: Derek Hoiem

Kristen Grauman

Limitations (continued)

- In practice, often entails large, cropped training set (expensive)
- Requiring good match to a global appearance description can lead to sensitivity to partial occlusions



Image credit: Adam, Rivlin, & Shimshoni

Kristen Grauman

Summary

- Basic pipeline for window-based detection
 - Model/representation/classifier choice
 - Sliding window and classifier scoring
- Boosting classifiers: general idea
- Viola-Jones face detector
 - Exemplar of basic paradigm
 - Plus key ideas: rectangular features, Adaboost for feature selection, cascade
- Pros and cons of window-based detection

- Given example images $(x_1, y_1), \dots, (x_n, y_n)$ where $y_i = 0, 1$ for negative and positive examples respectively.
- Initialize weights $w_{1,i} = \frac{1}{2m}, \frac{1}{2l}$ 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,

$$w_{t,i} \leftarrow \frac{w_{t,i}}{\sum_{j=1}^n w_{t,j}}$$

so that w_t is a probability distribution.

2. For each feature, j , train a classifier h_j which is restricted to using a single feature. The error is evaluated with respect to w_t , $\epsilon_j = \sum_i w_i |h_j(x_i) - y_i|$.
3. Choose the classifier, h_t , with the lowest error ϵ_t .
4. Update the weights:

$$w_{t+1,i} = w_{t,i} \beta_t^{1-e_i}$$

where $e_i = 0$ if example x_i is classified correctly, $e_i = 1$ otherwise, and $\beta_t = \frac{e_t}{1-e_t}$.

- 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 $\alpha_t = \log \frac{1}{\beta_t}$

Table 1: The AdaBoost algorithm for classifier learning. Each round of boosting selects one feature from the 180,000 potential features.

number of features are retained (perhaps a few hundred or thousand).

3.2. Learning Results

While details on the training and performance of the final system are presented in Section 5, several simple results merit discussion. Initial experiments demonstrated that a frontal face classifier constructed from 200 features yields a detection rate of 95% with a false positive rate of 1 in 14084. These results are compelling, but not sufficient for many real-world tasks. In terms of computation, this classifier is probably faster than any other published system, requiring 0.7 seconds to scan an 384 by 288 pixel image. Unfortunately, the most straightforward technique for improving detection performance, adding features to the classifier, directly increases computation time.

For the task of face detection, the initial rectangle features selected by AdaBoost are meaningful and easily interpreted. The first feature selected seems to focus on the property that the region of the eyes is often darker than the region

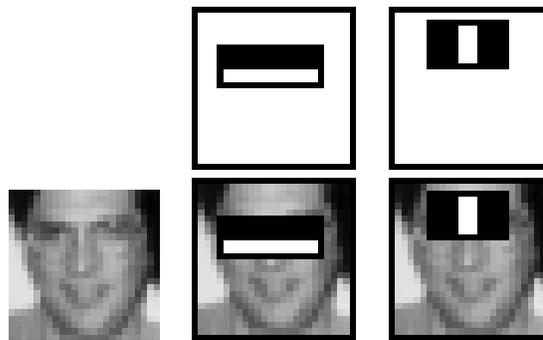


Figure 3: The first and second features selected by AdaBoost. The two features are shown in the top row and then overlaid on a typical training face in the bottom row. The first feature measures the difference in intensity between the region of the eyes and a region across the upper cheeks. The feature capitalizes on the observation that the eye region is often darker than the cheeks. The second feature compares the intensities in the eye regions to the intensity across the bridge of the nose.

of the nose and cheeks (see Figure 3). This feature is relatively large in comparison with the detection sub-window, and should be somewhat insensitive to size and location of the face. The second feature selected relies on the property that the eyes are darker than the bridge of the nose.

4. The Attentional Cascade

This section describes an algorithm for constructing a cascade of classifiers which achieves increased detection performance while radically reducing computation time. The key insight is that smaller, and therefore more efficient, boosted classifiers can be constructed which reject many of the negative sub-windows while detecting almost all positive instances (i.e. the threshold of a boosted classifier can be adjusted so that the false negative rate is close to zero). Simpler classifiers are used to reject the majority of sub-windows before more complex classifiers are called upon to achieve low false positive rates.

The overall form of the detection process is that of a degenerate decision tree, what we call a “cascade” (see Figure 4). A positive result from the first classifier triggers the evaluation of a second classifier which has also been adjusted to achieve very high detection rates. A positive result from the second classifier triggers a third classifier, and so on. A negative outcome at any point leads to the immediate rejection of the sub-window.

Stages in the cascade are constructed by training classifiers using AdaBoost and then adjusting the threshold to minimize false negatives. Note that the default AdaBoost threshold is designed to yield a low error rate on the training data. In general a lower threshold yields higher detec-