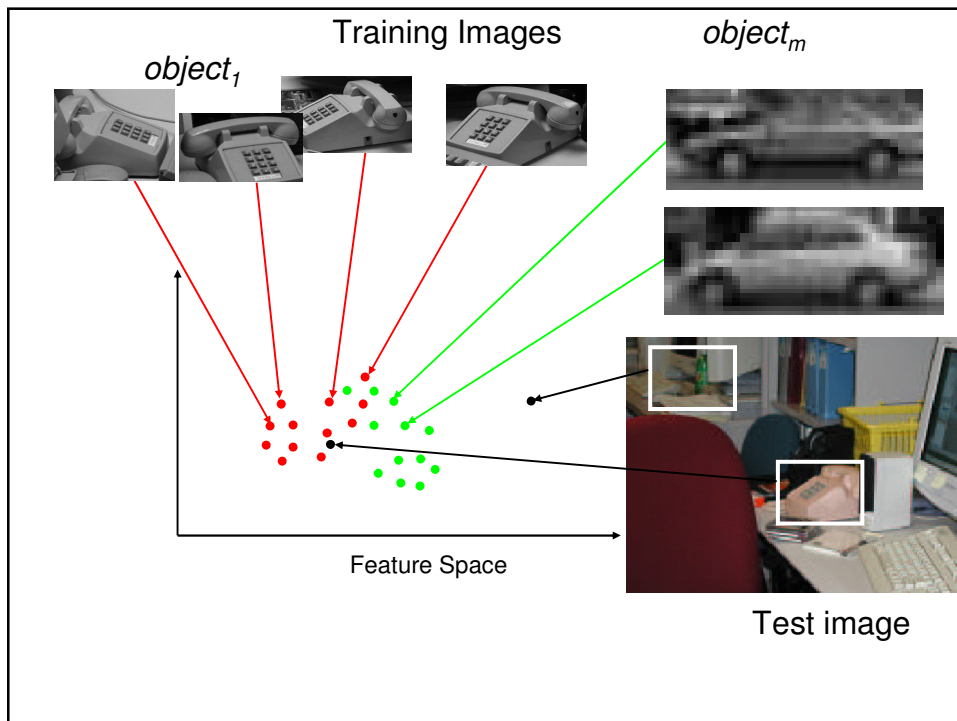
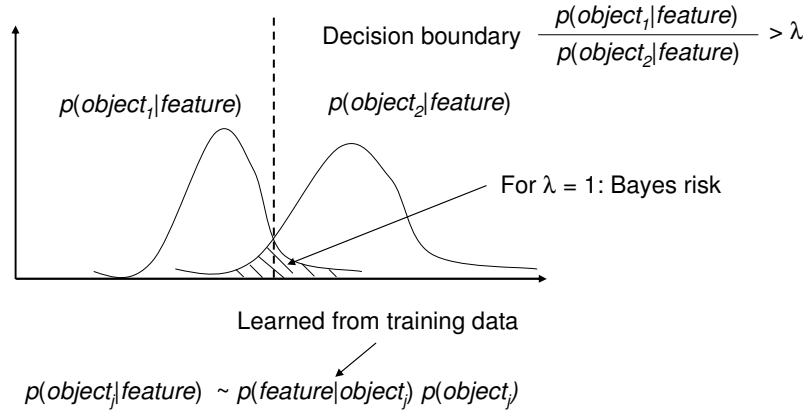


Recognition

Recognition by Templates Classifiers



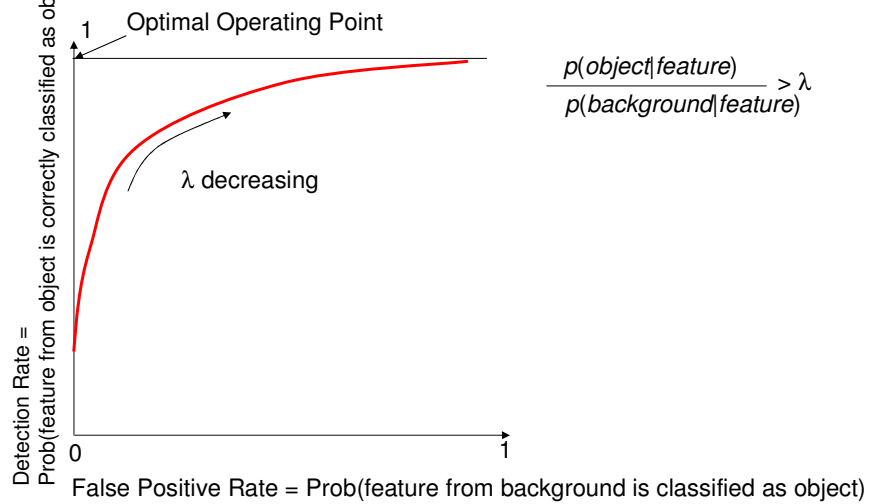
Probabilistic Formulation



- How to represent and learn $p(feature|object_j)$ or decision boundary?
- How to approach Bayes risk given small number of samples?
- What features to use?
- How to reduce the feature space?

How to evaluate classifier performance?

Receiver Operating Characteristic (ROC)



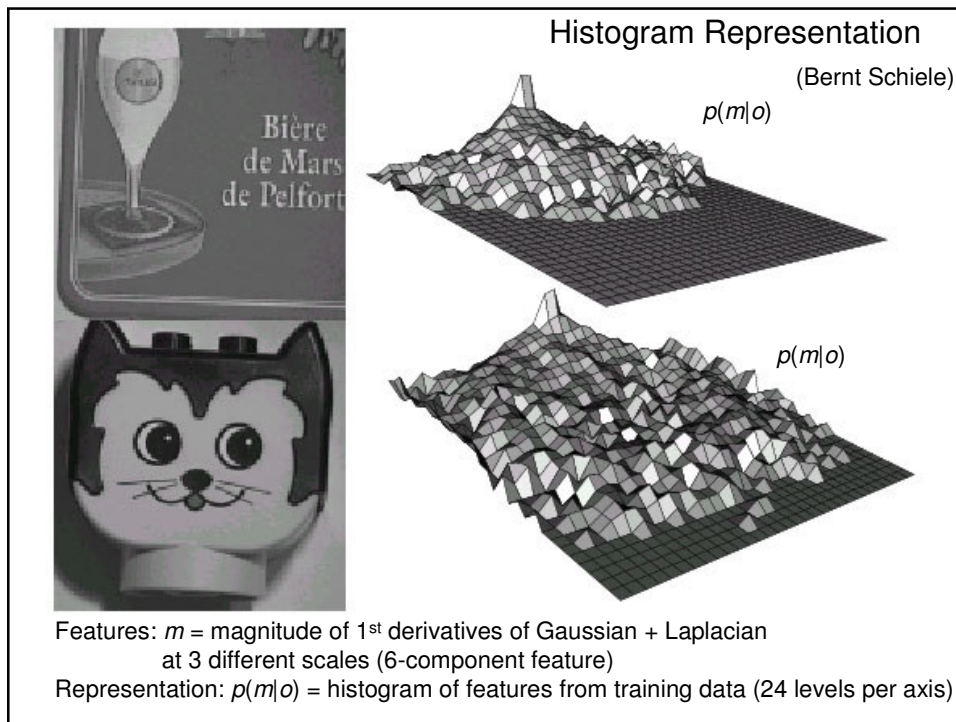
Approaches

- Every single pattern classification/learning approach has been applied to this problem
- Pick your favorite:
 - Naïve Bayes
 - Boosting
 - Neural networks
 - SVMs
 - NNs
 - PCA/LDA/ICA dimensionality reduction

.....
Your favorite buzzword goes here.

Approaches

- Every single pattern classification/learning approach has been applied to this problem
- Pick your favorite:
 - Naïve Bayes ←
 - Boosting
 - Neural networks
 - SVMs
 - NNs
 - PCA/LDA/ICA dimensionality reduction



m_i = 6-vector at location i
Scan over all possible locations in the window

$P(m_i | object_n)$

Computed by looking up
the histogram tables
computed at training time

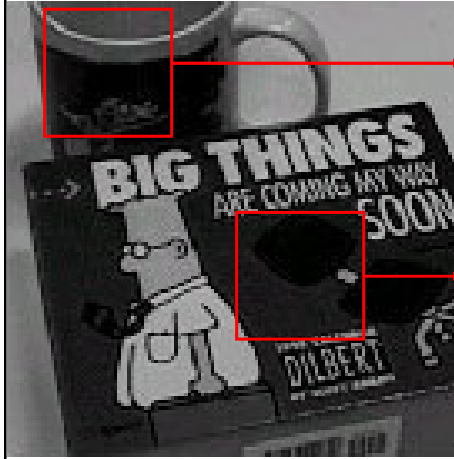
$P(m_j | object_n)$

Return object that maximizes:

$$P(object_n | image) \propto \prod_i P(m_i | object_n) P(object_n)$$

[Example from Bernt Schiele]

For complex scenes: Scan the image and evaluate “probability” at scanned window locations



$$P(\text{object}_n | W_k) \propto \prod_{i \in W_k} P(m_i | \text{object}_n) P(\text{object}_n)$$

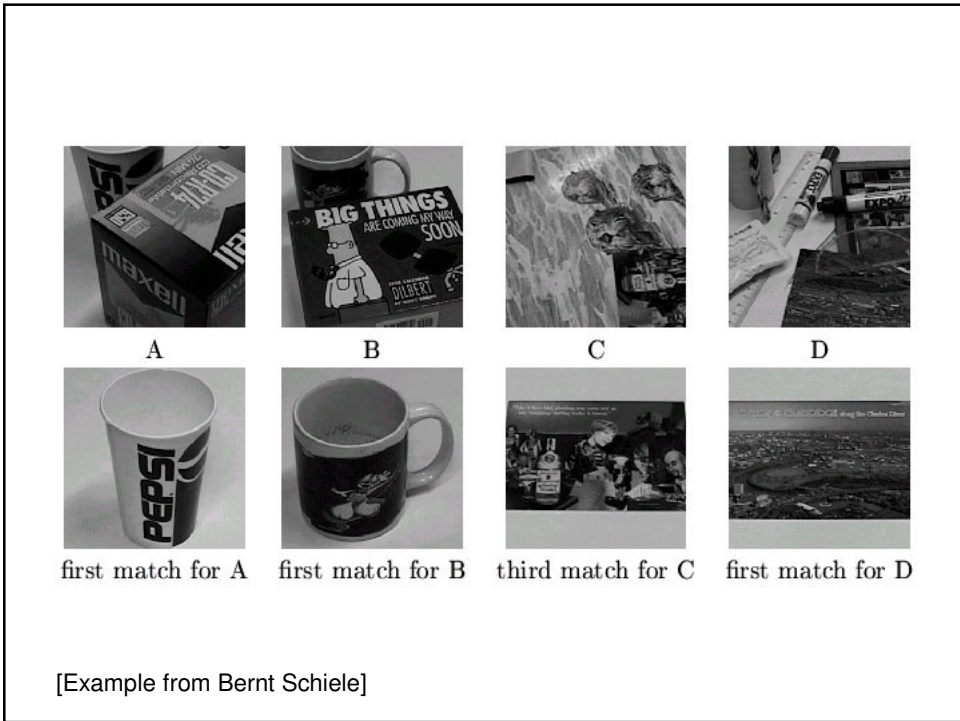
$$P(\text{object}_n | W_l) \propto \prod_{i \in W_l} P(m_i | \text{object}_n) P(\text{object}_n)$$

Object n is in the image if many windows “vote” for the image, e.g., by evaluating $\text{vote}(\text{object}_n) = \sum_k P(\text{object}_n | W_k)$

Note: This is not necessarily the best way to do this....see later



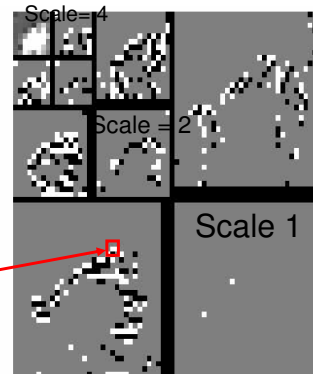
[Example from Bernt Schiele]



More Complicated Features



$C(x,y,s)$ =
Wavelet
coefficient at
position x,y at
scale s



- Feature = Set of coefficients $S = (C_1, \dots, C_N)$

Example from Henry Schneiderman

- Given features S_1, \dots, S_r computed from a window, threshold the likelihood ratio

$$\log \frac{P(S_1 \cdots S_r | \omega_1)}{P(S_1 \cdots S_r | \omega_2)} =$$

Assume independence
(Naïve Bayes)

$$\log \frac{P(S_1 | \omega_1)}{P(S_1 | \omega_2)} + \log \frac{P(S_2 | \omega_1)}{P(S_2 | \omega_2)} + \dots + \log \frac{P(S_r | \omega_1)}{P(S_r | \omega_2)} > \lambda ?$$

How can we compute these probabilities?

Example from Henry Schneiderman

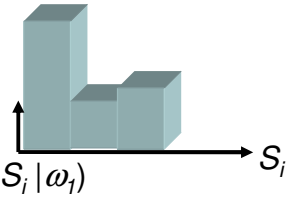
Estimating the Probabilities

- Collect the values of the features for training data in histograms that approximate the probabilities

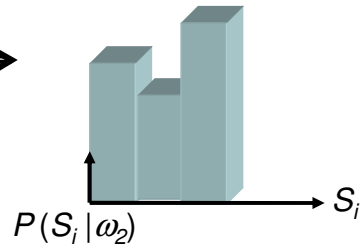


50-2,000 original images

~1,000 synthetic variations per original image



~10,000,000 examples



Example from Henry Schneiderman

Compute the values of all the features in the window
 For each feature, compute the probabilities of coming from the object or non-object class
 Aggregate into likelihood ratio



$$f_1(0, 0) = \#5710$$

$$P_1(\#5710, 0, 0 | \text{obj}) = 0.53$$

$$P_1(\#5710, 0, 0 | \text{non-obj}) = 0.56$$



$$f_1(0, 1) = \#3214$$

$$P_1(\#3214, 0, 1 | \text{obj}) = 0.57$$

$$P_1(\#3214, 0, 1 | \text{non-obj}) = 0.48$$



$$f_N(n, m) = \#723$$

$$P_N(\#723, n, m | \text{obj}) = 0.83$$

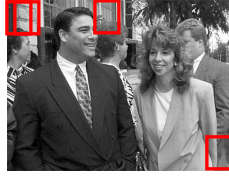
$$P_N(\#723, n, m | \text{non-obj}) = 0.19$$

$$0.53 * 0.57 * \dots * 0.83$$

$$0.56 * 0.48 * \dots * 0.19$$

$> \lambda$

From Windows to Images



Search in position

- Move a window to all possible positions and all possible scales
- At each (position, scale) evaluate the classifier
- Return detection if above threshold



Search in scale

Example from Henry Schneiderman

Feature Selection Problem

- Each feature is a set of variables (wavelet coefficients) $S = \{C_1, \dots, C_N\}$
- Problem:
 - If N is large, the feature is very discriminative (S is equivalent to the entire window if N is the total number of variables) but representing the corresponding distribution is very expensive
 - If N is small, the feature is not discriminative but classification is very fast

Example from Henry Schneiderman

Solution: Classifier Cascade

- Standard problem:
 - We can have either discriminative or efficient features but not both!
 - Cannot do classification in one shot
- Standard solution: *Classifier Cascade*
 - Apply first a classifier with simple features → Fast and will eliminate the most obvious non-object locations
 - Then apply a classifier with more complex features → More expensive but applied *only* to these locations that survived the previous stage

Example from Henry Schneiderman

Cascade Example

Cascade Stage 1



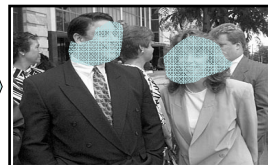
Apply classifier with very simple (and fast) features → Eliminates most of the image

Cascade Stage 2



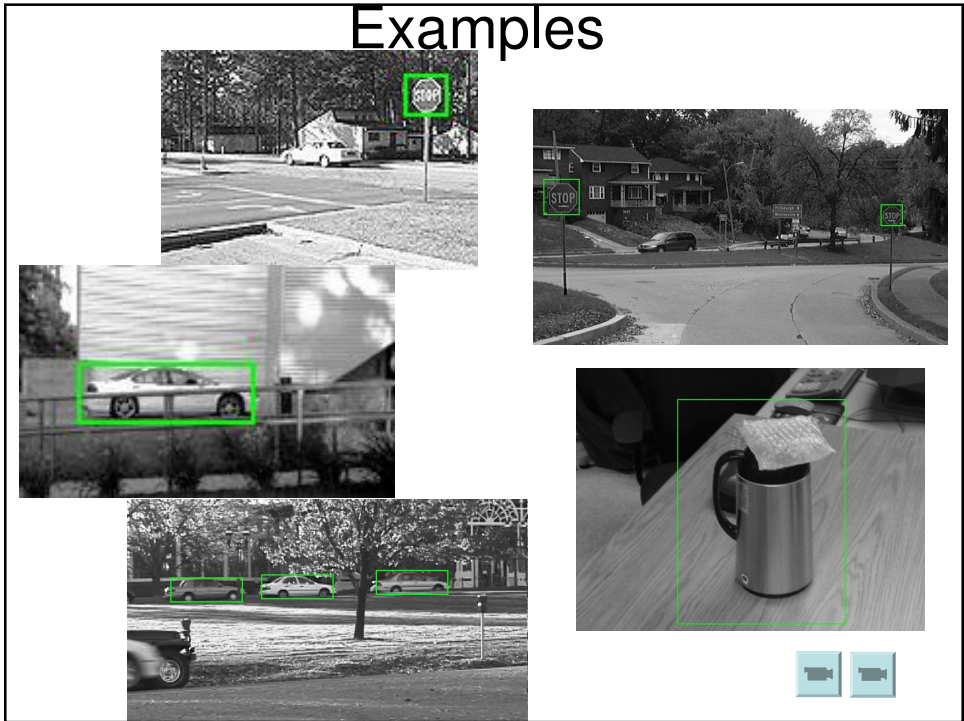
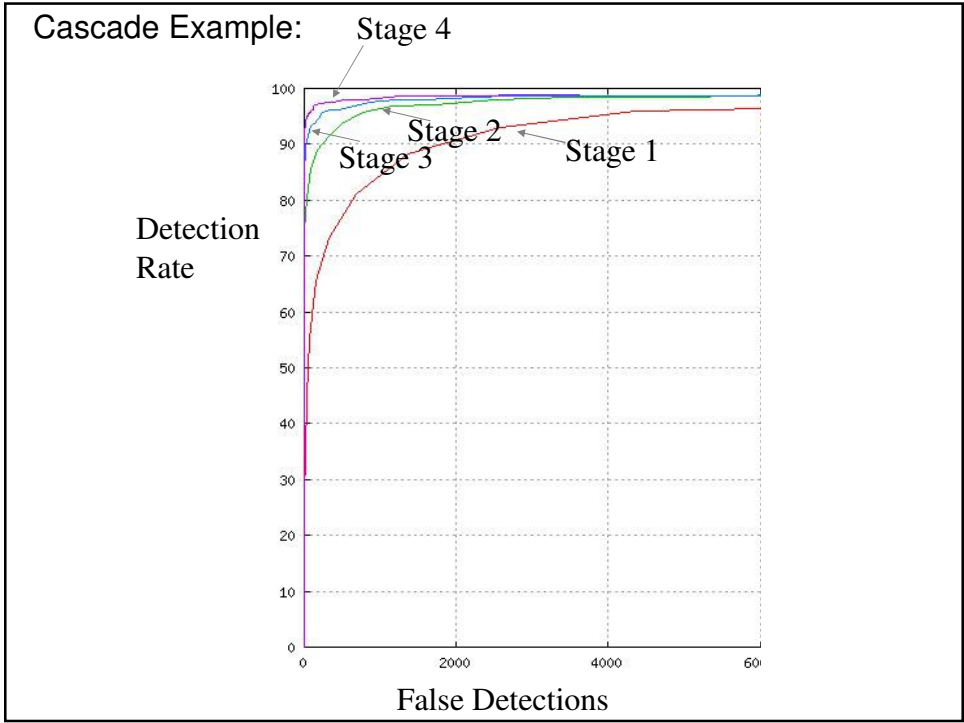
Apply classifier with more complex features on what is left

Cascade Stage 3




Apply classifier with more complex features on what is left

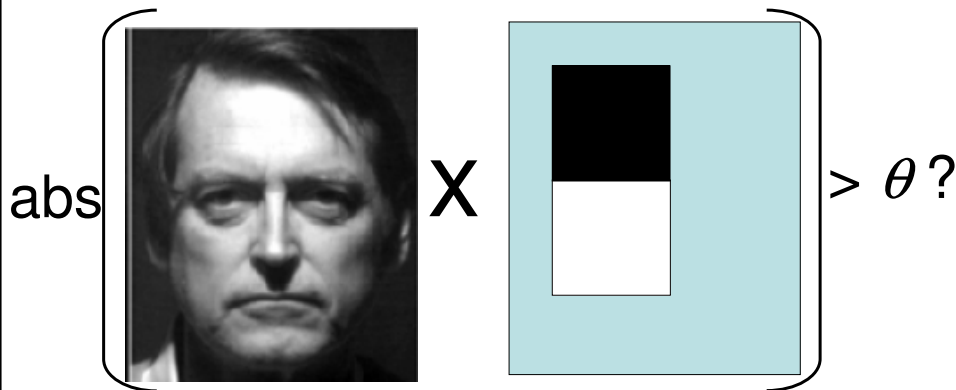
Example from Henry Schneiderman



Approaches

- Every single pattern classification/learning approach has been applied to this problem
- Pick your favorite:
 - Naïve Bayes
 - Boosting 
 - Neural networks
 - SVMs
 - NNs
 - PCA/LDA/ICA dimensionality reduction
 -

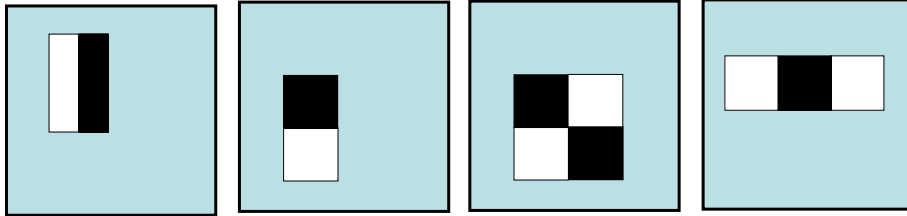
Using Weak Features



- Don't try to design strong features from the beginning, just use really stupid but really fast features (and a lot of them)
- *Weak learner* = Very fast (but very inaccurate) classifier
- Example: Multiply input window by a very simple box operator and threshold output

(Example from Paul Viola, Distributed by Intel as part of the OpenCV library)

Feature Selection



- Operators defined over all possible shapes and positions within the window
- For a 24x24 window \rightarrow 45,396 combinations!!
- How to select the “useful” features?
- How to combine them into classifiers?

(Example from Paul Viola)

- Input: Training examples $\{x_i\}$ with labels (“face” or “non-face” = +/-1) $\{y_i\}$ + weights w_i (initially $w_i = 1$)

- Choose the feature (weak classifier h_t) with minimum error: $\varepsilon_t = \sum_i w_i [h_t(x_i) \neq y_i]$

- Update the weights such that
 - w_i is increased if x_i is misclassified
 - w_i is decreased if x_i is correctly classified

- Compute a weight α_t for classifier h_t
 - α_t large if ε_t is small

- Final classifier: $H(x) = \text{sgn}\left(\sum_t \alpha_t h_t(x)\right)$

Repeat T times

This is a general description of a *boosting* algorithm. Well-defined rules for updating w and for computing α guarantee convergence and “good” classification performance.

Repeat T times

Choose the feature (weak classifier h_t) with minimum error: ϵ_t

The training examples that are not correctly classified contribute more through higher weights

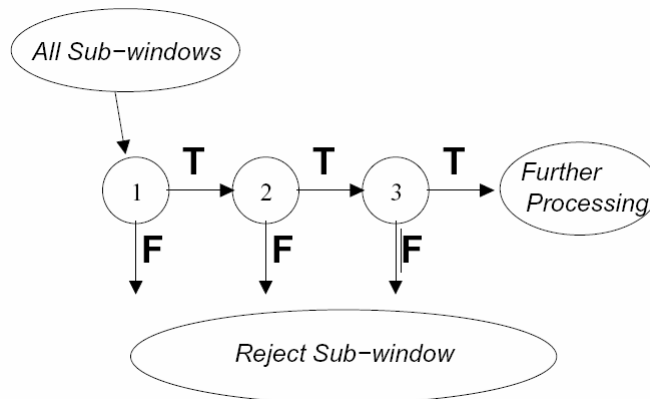
- Update the weights such that
 - w_i is increased if x_i is misclassified
 - w_i is decreased if x_i is correctly classified
- Compute a weight α_t for classifier h_t
 - α_t large if ϵ_t is small
- Final classifier:
$$H(x) = \text{sgn}\left(\sum_t \alpha_t h_t(x)\right)$$

Features that yield good classification performance receive higher weights

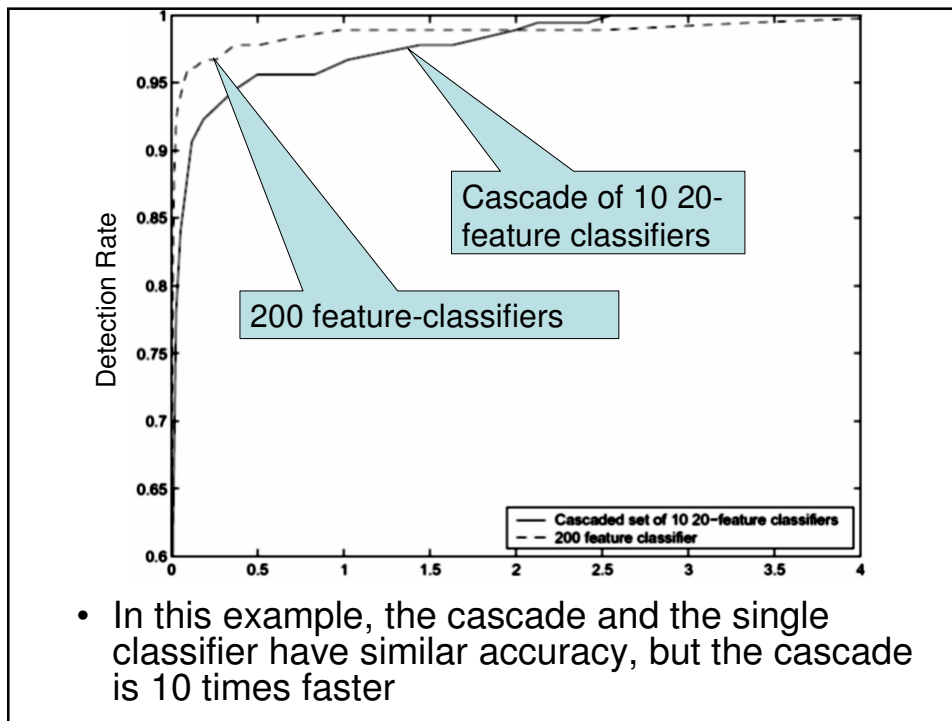
1st feature selected 2nd feature selected

The automatic selection process selects “natural” features
(Example from Paul Viola)

Using a Cascade (Again)



- Same problem as before: It is too hard (or impossible) to build a single accurate classifier
- Key reason: An image containing one face may have 10^5 possible locations but only 1 "correct" location \rightarrow *rare event detection* \rightarrow Would require an enormous number of features
- Solution: Use a cascade of classifiers. Each classifier eliminates more of the non-object locations while retaining the "object" locations.

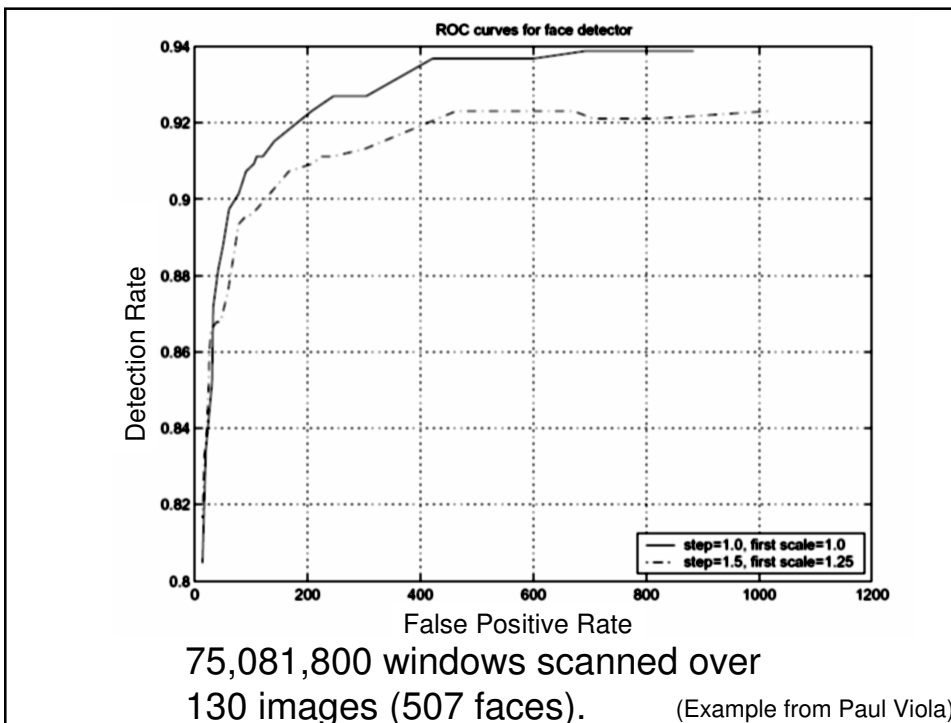


- In this example, the cascade and the single classifier have similar accuracy, but the cascade is 10 times faster



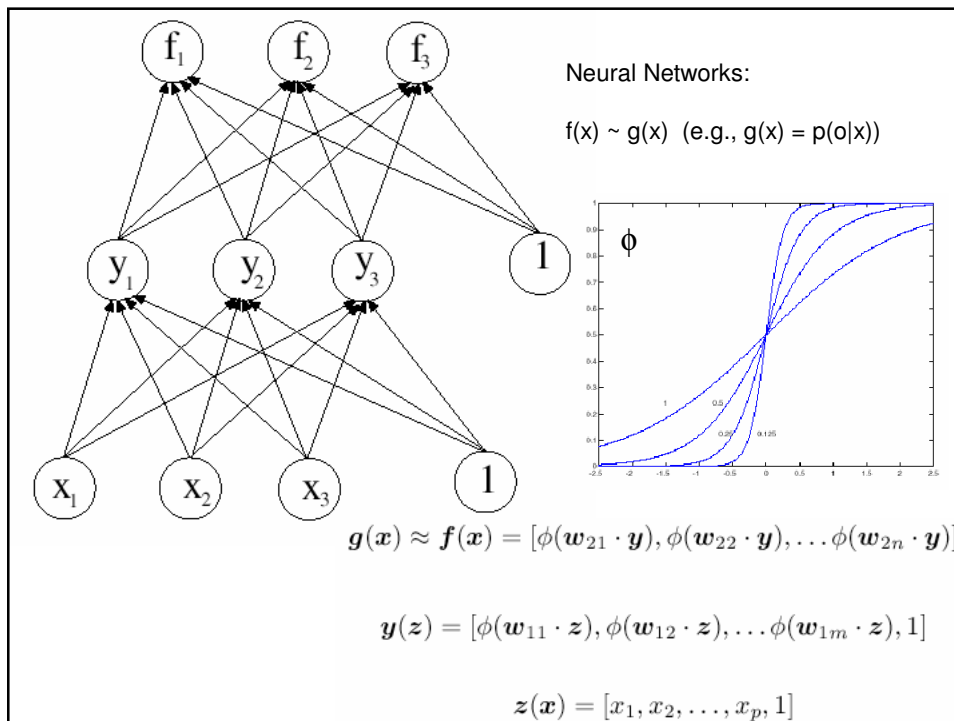
Training: ~5000 face images
+ ~10000 non-face windows

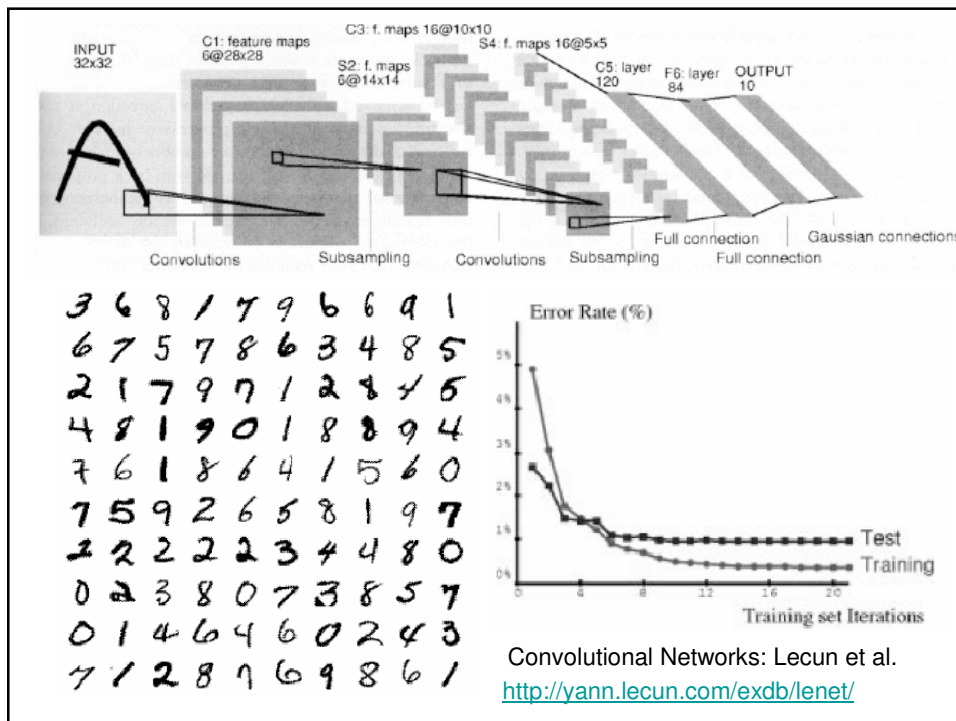
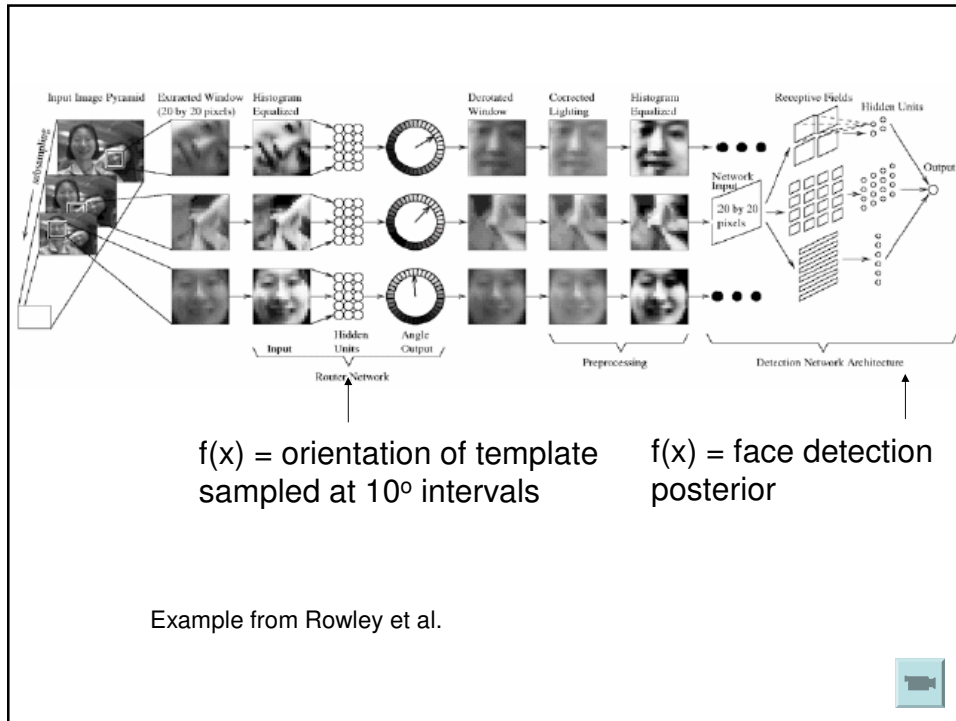
Run-Time: Apply the cascade of classifiers over all positions and scales and return those positions/scales that survive the sequence of classifiers



Approaches

- Every single pattern classification/learning approach has been applied to this problem
- Pick your favorite:
 - Naïve Bayes
 - Boosting
 - Neural networks ←
 - SVMs
 - NNs
 - PCA/LDA/ICA dimensionality reduction
 -



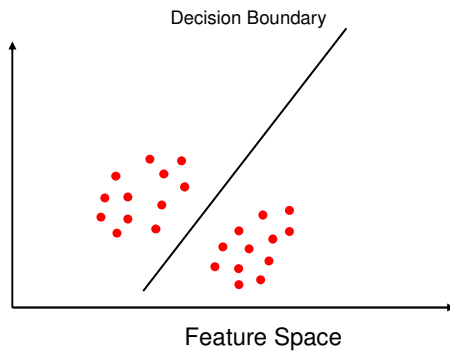


Approaches

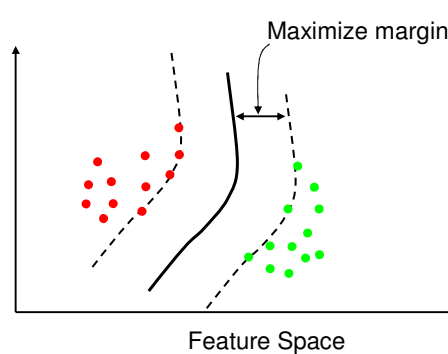
- Every single pattern classification/learning approach has been applied to this problem
- Pick your favorite:
 - Naïve Bayes
 - Boosting
 - Neural networks
 - SVMs ←
 - NNs
 - PCA/LDA/ICA dimensionality reduction
 -

Discriminative Approaches

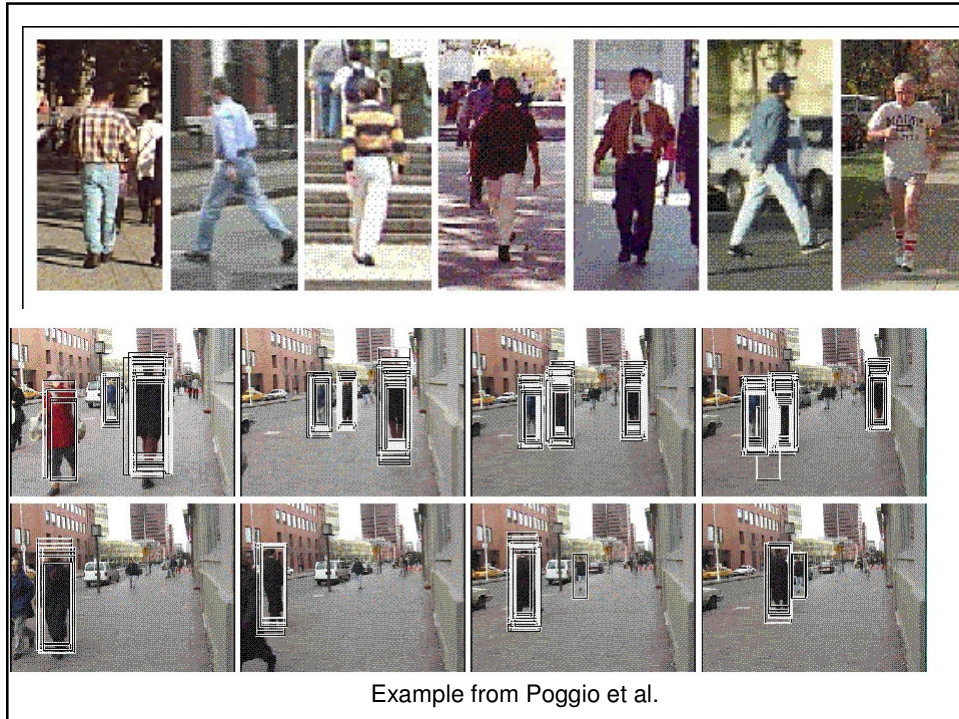
Linear Discriminant



General: Large-Margin Classifiers, SVMs



- Difficult to represent the distribution in high-dimensional feature spaces → Find decision boundary directly
- General idea: Much less training data is needed to construct the decision boundary than the distributions
- Maximize separation between the classes for better generalization
- Fewer parameters: 2 Gaussians with equal covariance = 7 params, line = 2 params



Approaches

- Every single pattern classification/learning approach has been applied to this problem
- Pick your favorite:
 - Naïve Bayes
 - Boosting
 - Neural networks
 - SVMs
 - NNs ←
 - PCA/LDA/ICA dimensionality reduction ←
 -

Nearest Neighbors

Feature Space

- Does not require recovery of distributions or decision surfaces
- Asymptotically twice Bayes risk at most
- Choice of distance metric critical
- Indexing may be difficult

Large Feature Spaces: PCA

\mathbf{X} = feature vector of high dimension
 → Difficult indexing in high-dimensional space
 → Most of the dimensions are probably not useful

Principal Component: Dominant eigenvectors of scatter matrix
 $\tilde{\mathbf{X}} = \mathbf{X} - \bar{\mathbf{X}}$

$$\sum_i \tilde{\mathbf{X}}_i \tilde{\mathbf{X}}_i^T = \begin{bmatrix} \sum_i x_{i1}^2 & \cdots & \sum_i x_{i1} x_{in} \\ \vdots & \ddots & \vdots \\ \sum_i x_{i1} x_{in} & \cdots & \sum_i x_{in}^2 \end{bmatrix}$$

Most of the information is contained in the Space spanned by $(\mathbf{V}_1, \dots, \mathbf{V}_k)$

$$\tilde{\mathbf{X}} \approx \lambda_1 \mathbf{V}_1 \cdot \tilde{\mathbf{X}} + \cdots + \lambda_k \mathbf{V}_k \cdot \tilde{\mathbf{X}}$$

PCA: Project first in the lower-dimensional space spanned by the principal component
 → Indexing in much lower dimensional space
 → Feature selection

PCA for Recognition

(assume centered features)

Training features: $\mathbf{X}_1, \dots, \mathbf{X}_m$

Compute principal directions: $\mathbf{V}_1, \dots, \mathbf{V}_k$

Project training features
onto principal directions

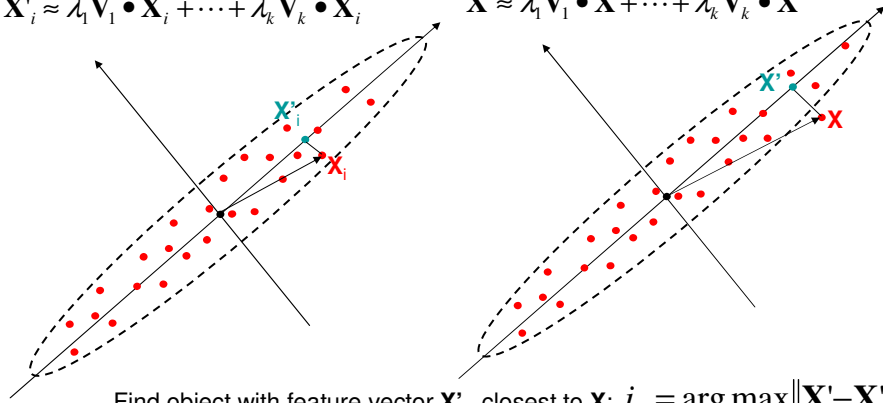
$$\mathbf{X}'_i \approx \lambda_1 \mathbf{V}_1 \cdot \mathbf{X}_i + \dots + \lambda_k \mathbf{V}_k \cdot \mathbf{X}_i$$

Input: Feature vector \mathbf{X}

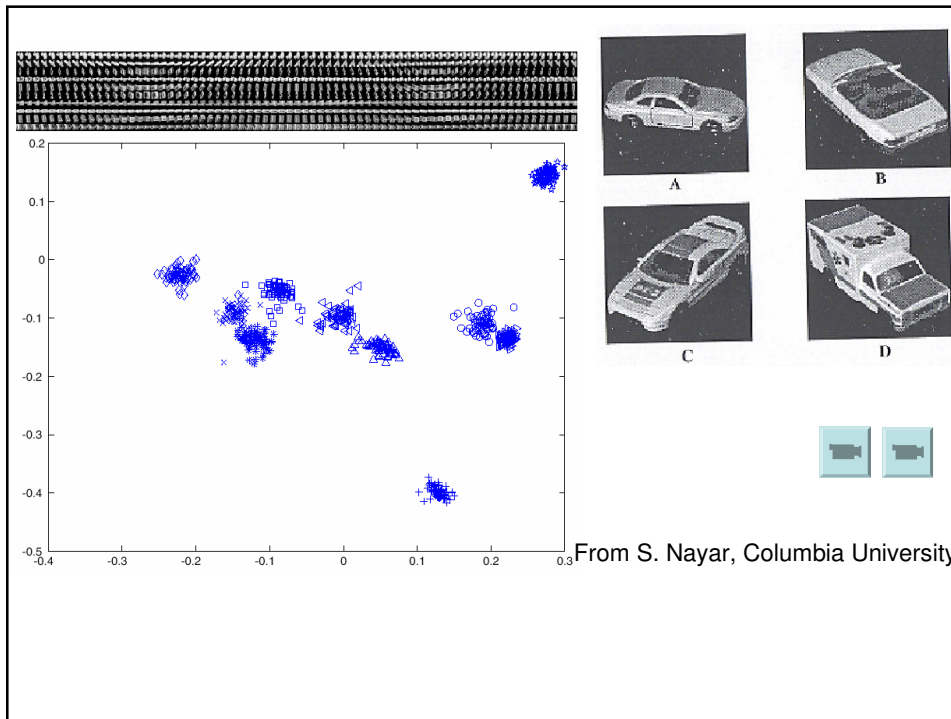


Project \mathbf{X} onto principal component space:

$$\mathbf{X}' \approx \lambda_1 \mathbf{V}_1 \cdot \mathbf{X} + \dots + \lambda_k \mathbf{V}_k \cdot \mathbf{X}$$

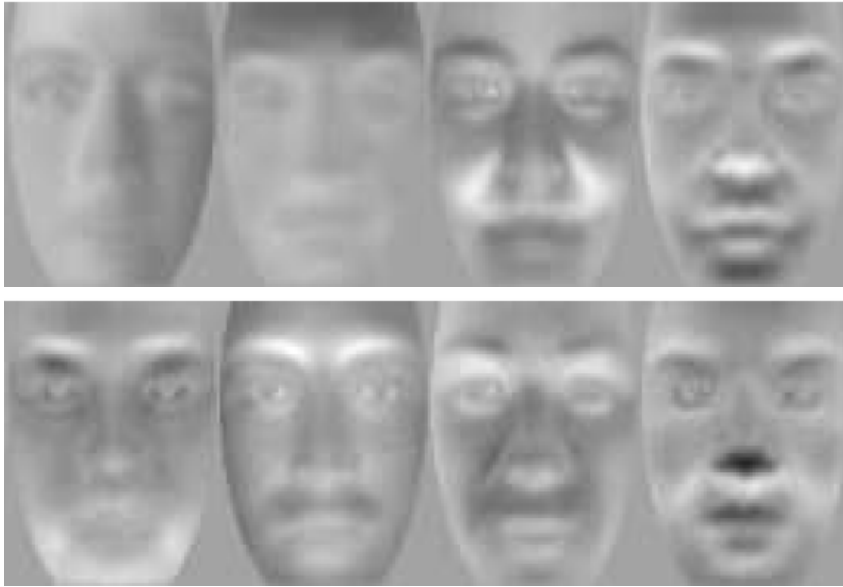


Find object with feature vector \mathbf{X}'_{i_o} closest to \mathbf{X} : $i_o = \arg \max_i \|\mathbf{X}' - \mathbf{X}'_i\|$



From S. Nayar, Columbia University

Example Eigenvectors:



[Example from Draper et al.]

Extreme case: $X = \text{image itself}$
Example: EigenFaces

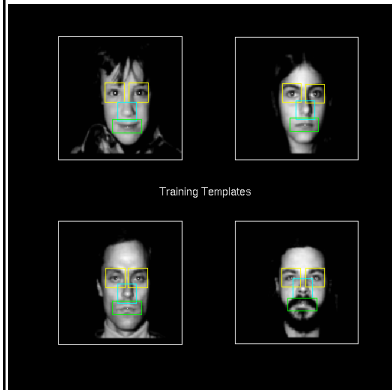
The screenshot shows the EigenFaces web interface. On the left, there is a grid of 20 grayscale face images, which are the first 20 eigenvectors. On the right, there is a search interface with a grid of 15 grayscale face images, which are the 15 most similar faces among 7,562 faces (3,000 subjects). The search interface includes buttons for 'Initialize', 'Shuffle', 'Load Query', 'Save Query', 'Text...', 'Symbols...', 'Label...', 'Hooks...', 'G Label...', 'Resize', 'Refresh Cache', 'Page Up/Down', 'Page 1 of 473', 'Jump to page', and 'Jump to item'. There is also a search input field.

Eigenfaces: Projection on the first 20 eigenvectors from 128 face images

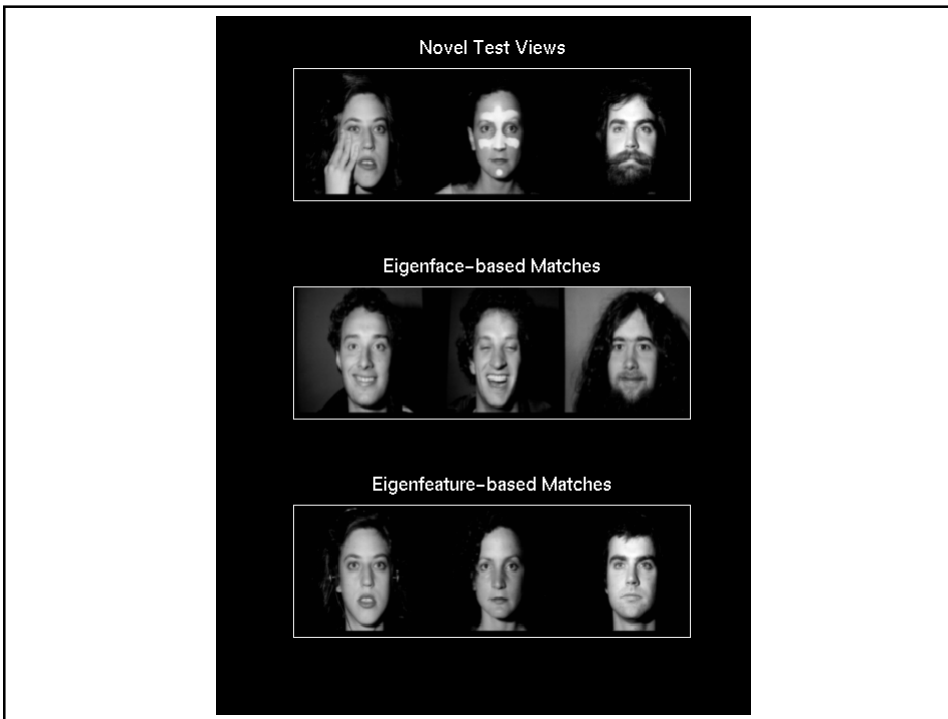
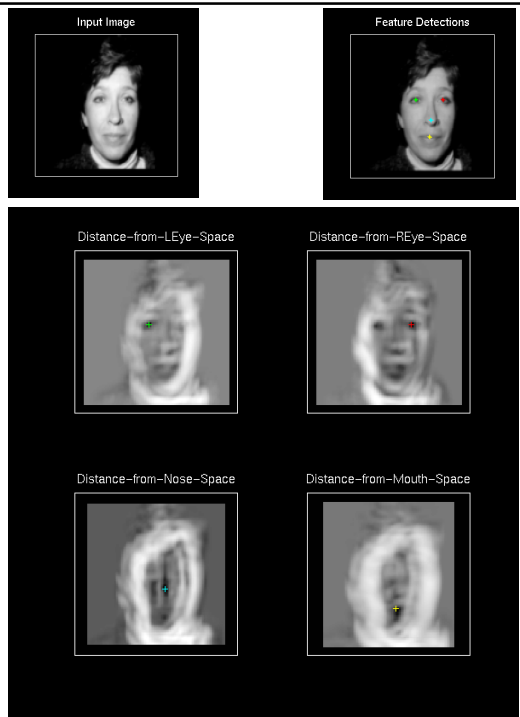
15 most similar faces among 7,562 faces (3000 subjects)

<http://www-white.media.mit.edu/vismod/demos/facerec/basic.html>

EigenFeatures



Features trained on 128 face images, retaining the first 10 eigenvectors



Problems

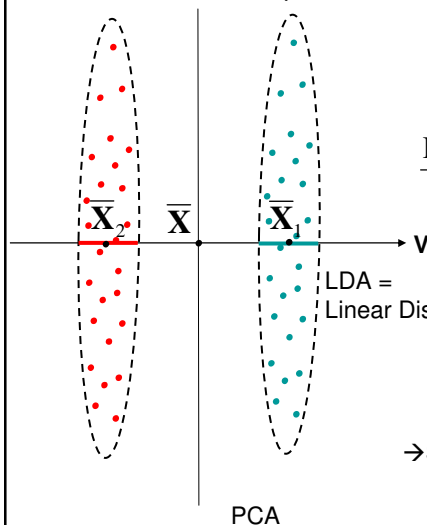


Variation in appearance due to illumination and expression >>> Variation in appearance due to identity

Example from the Yale Face DataBase

Problems

- Assume "linear" distribution of features
- Best choice for compression may not be the best choice for discrimination



LDA =
Linear Discriminant

LDA → Find projection direction \mathbf{V} that separates the 2 classes best

Maximize:

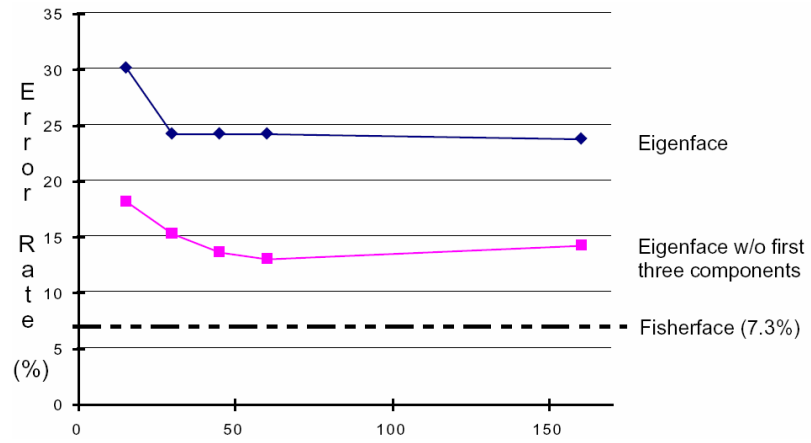
$$\frac{\text{Distance between classes after projection}}{\text{Scatter of classes after projection}} =$$

$$\frac{(\mathbf{V} \cdot (\bar{\mathbf{X}}_1 - \bar{\mathbf{X}}_2))^2}{\mathbf{V}(\mathbf{C}_1 + \mathbf{C}_2)\mathbf{V}^T}$$

→ Generalized eigenvalue problem

→ Similar application of LDA to faces: FisherFaces (Belhumeur, Yale/Columbia)

Example



Example from Belhumeur et al.