

Lecture 13

Visual recognition

- Announcements



Lecture 13

Visual recognition



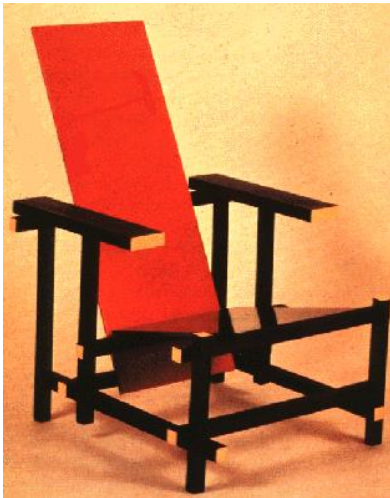
- Object classification bag of words models
 - Discriminative methods
 - Generative methods
- Object classification by PCA and FLD

Challenges

Variability due to:

- View point
- Illumination
- Occlusions
- Intra-class variability

Challenges: intra-class variation

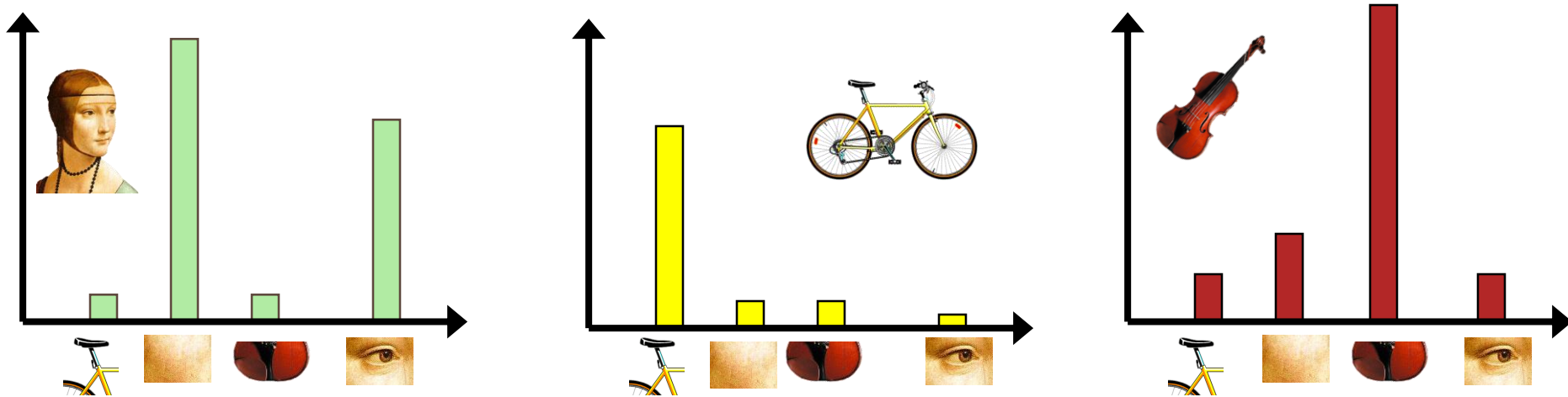


Basic properties

- Representation
 - How to represent an object category; which classification scheme?
- Learning
 - How to learn the classifier, given training data
- Recognition
 - How the classifier is to be used on novel data

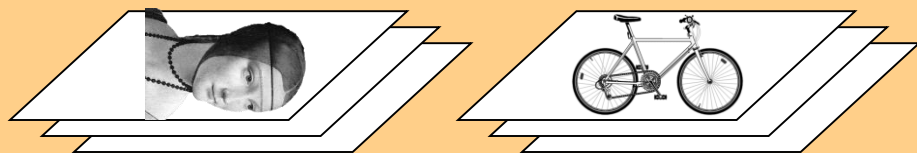
definition of “BoW”

– Histogram of visual words (codewords)



codewords dictionary

Representation



feature detection
& representation

codewords dictionary

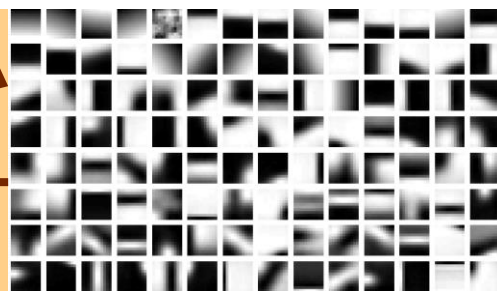
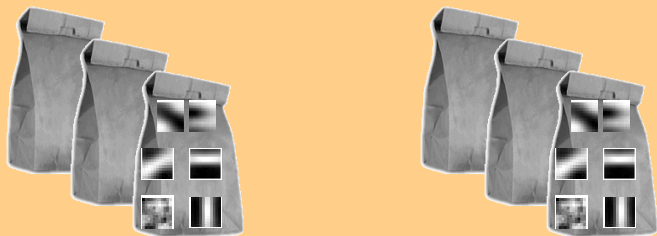


image representation



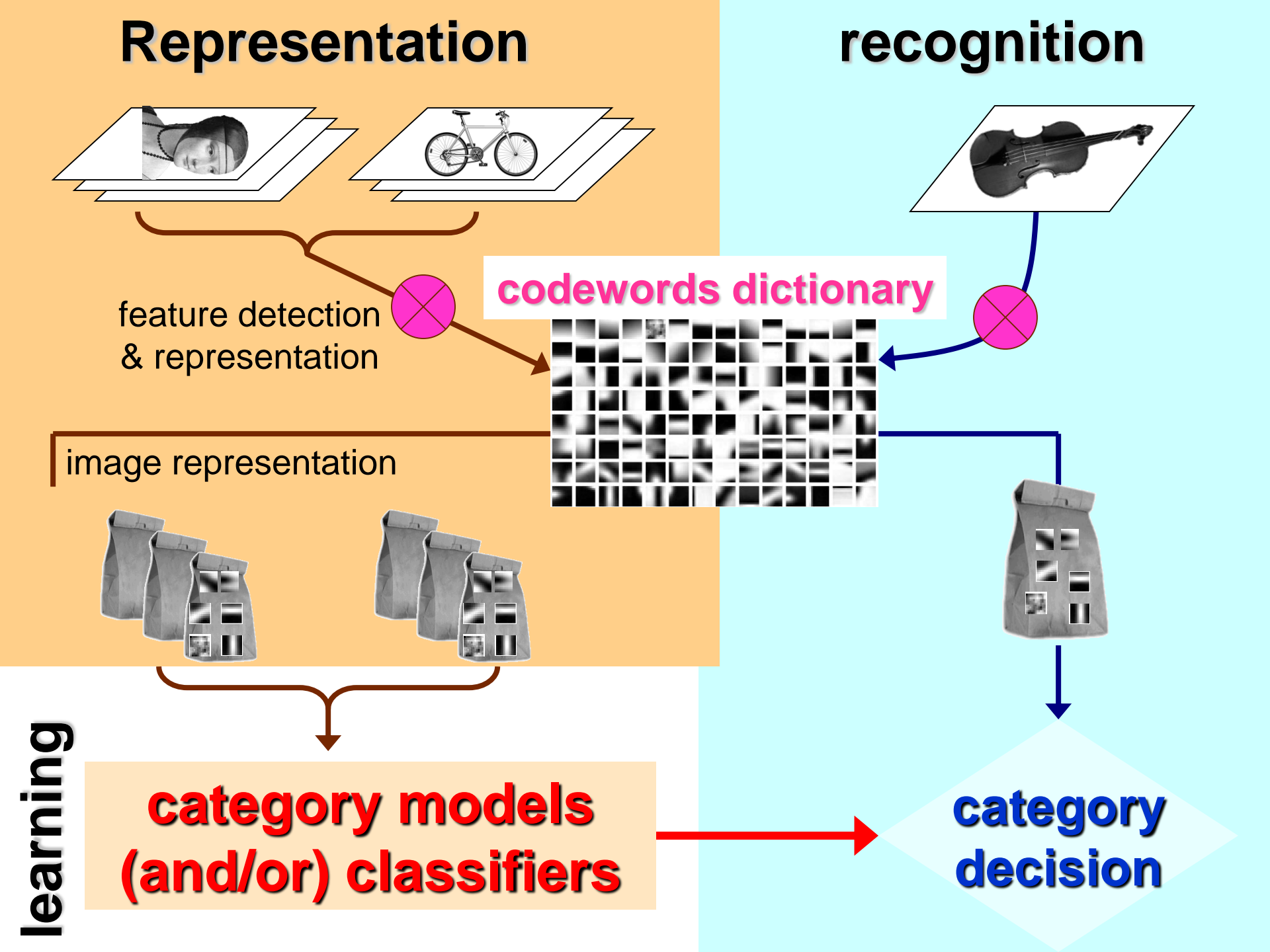
**category models
(and/or) classifiers**

recognition



**category
decision**

learning

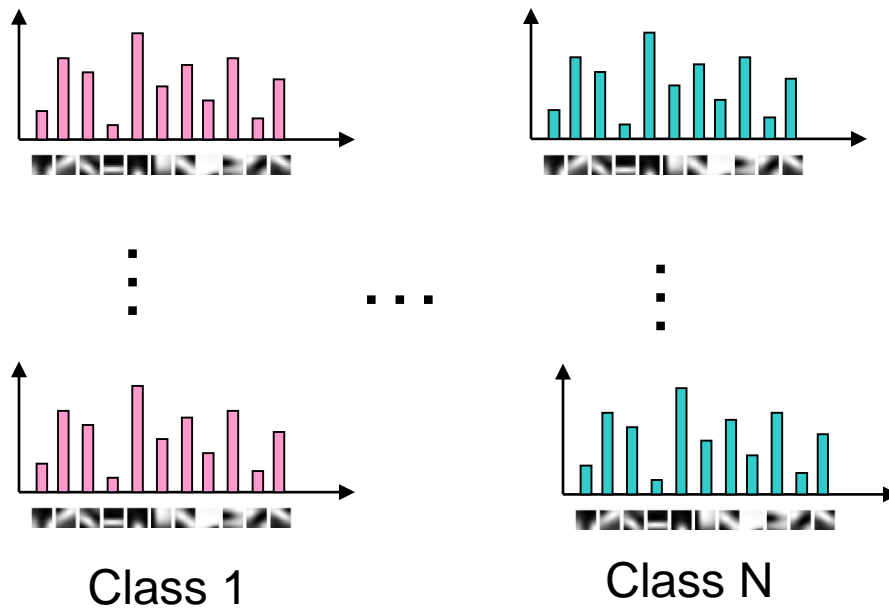


Classification

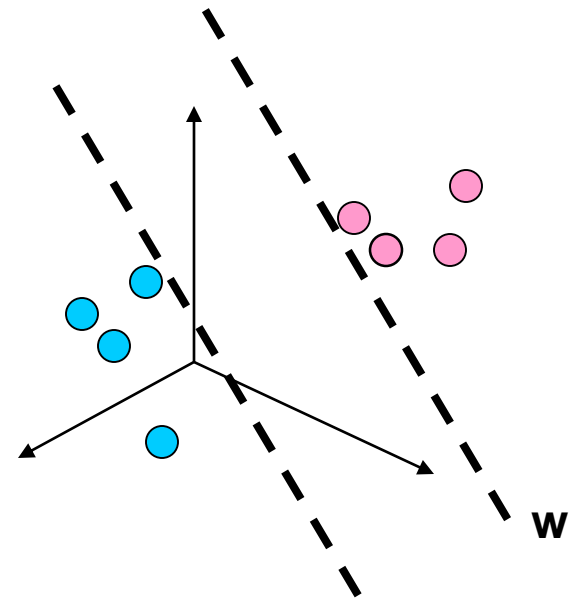
- Discriminative methods
 - Nearest neighbors
 - Linear classifier
 - SVM
- Generative methods

SVM classification

category models

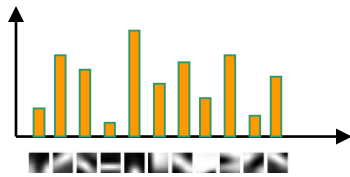


Model space



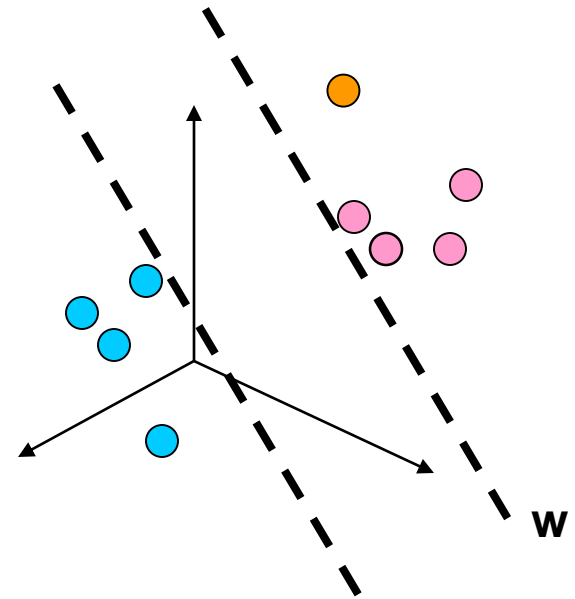
SVM classification

Query image



Winning class: pink

Model space



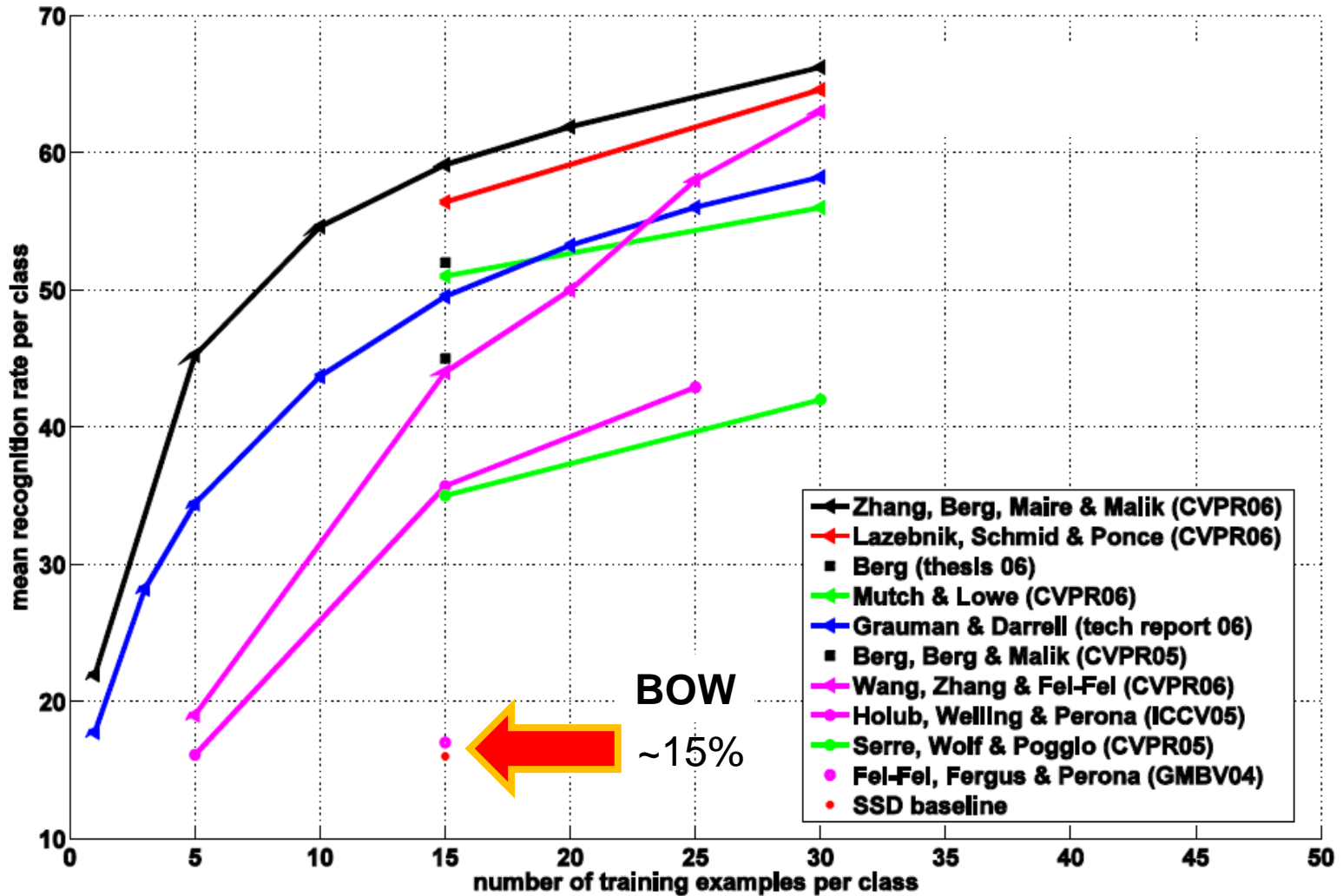
Caltech 101

Fei-Fei et al. (2004)

http://www.vision.caltech.edu/Image_Datasets/Caltech101/Caltech101.html



Caltech 101

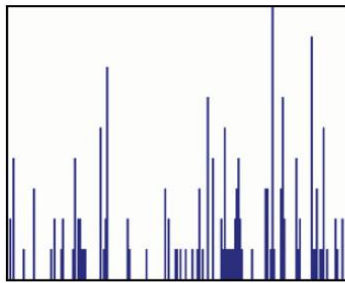
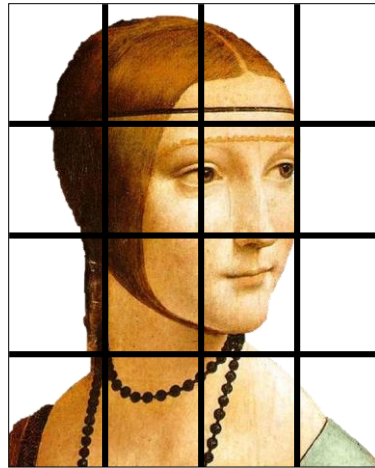


Major drawback of BOW models

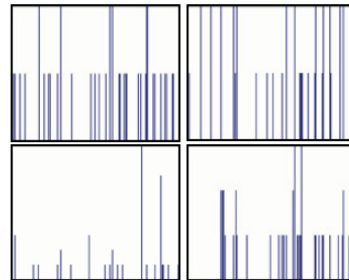
Don't capture spatial information!

Spatial Pyramid Matching

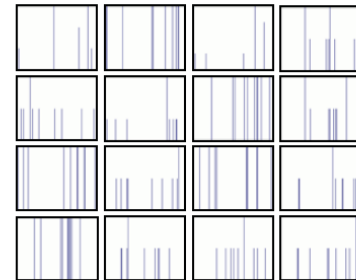
Beyond Bags of Features: Spatial Pyramid Matching for Recognizing Natural Scene Categories. S. Lazebnik, C. Schmid, and J. Ponce.. 2006



level 0



level 1



level 2

$$SPM(x_i, x) = \frac{1}{2^L} HIK_0(x_i, x) + \dots + \frac{1}{L - l + 1} HIK_l(x_i, x) + \dots + HIK_L(x_i, x)$$

$$I(h_1, h_2) = I(h_1, h_2) + \frac{1}{2} I(h_1, h_2)$$

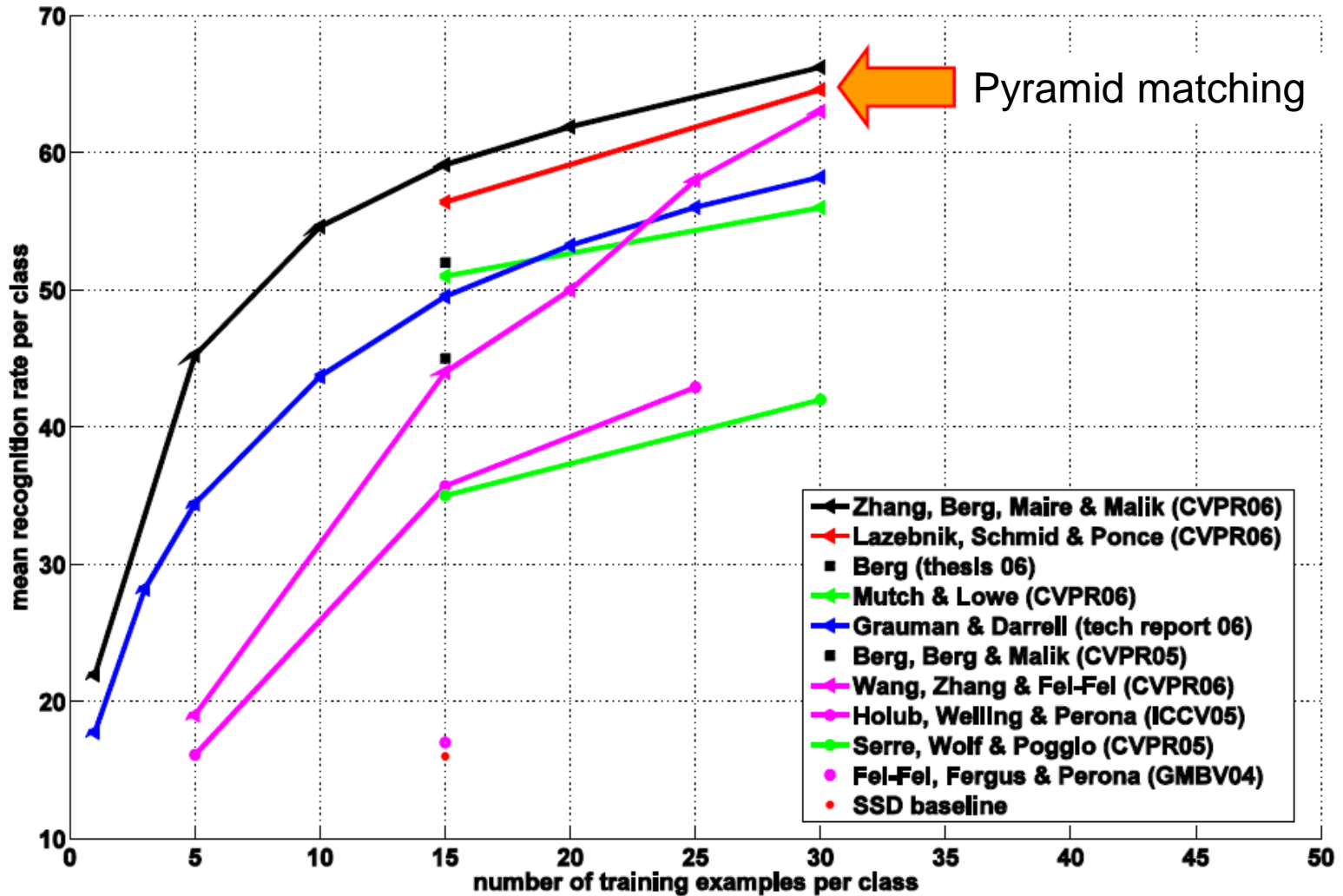
$$I(h_1, h_2) = \sum_{i=1}^N \min(h_1(i), h_2(i))$$

Caltech 101

Multi-class classification results (30 training images per class)

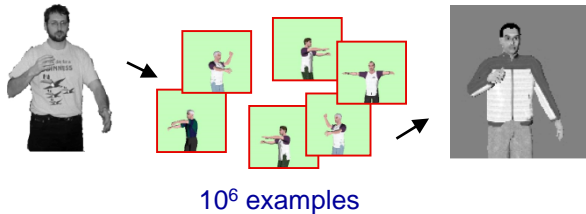
	Weak features (16)		Strong features (200)	
Level	Single-level	Pyramid	Single-level	Pyramid
0	15.5 \pm 0.9		41.2 \pm 1.2	
1	31.4 \pm 1.2	32.8 \pm 1.3	55.9 \pm 0.9	57.0 \pm 0.8
2	47.2 \pm 1.1	49.3 \pm 1.4	63.6 \pm 0.9	64.6 \pm 0.8
3	52.2 \pm 0.8	54.0 \pm 1.1	60.3 \pm 0.9	64.6 \pm 0.7

Caltech 101



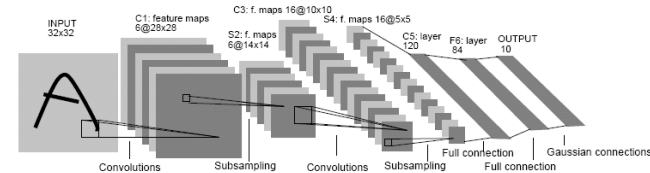
Discriminative models

Nearest neighbor



Shakhnarovich, Viola, Darrell 2003
Berg, Berg, Malik 2005...

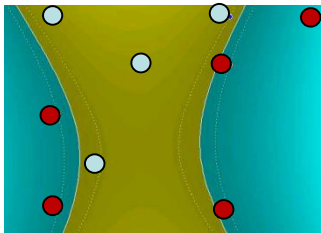
Neural networks



LeCun, Bottou, Bengio, Haffner 1998
Rowley, Baluja, Kanade 1998

...

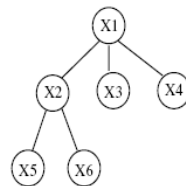
Support Vector Machines



Guyon, Vapnik, Heisele,
Serre, Poggio...

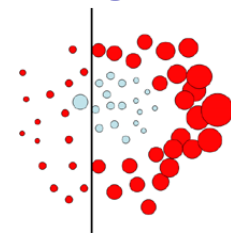
Latent SVM

Structural SVM



Felzenszwalb 00
Ramanan 03...

Boosting

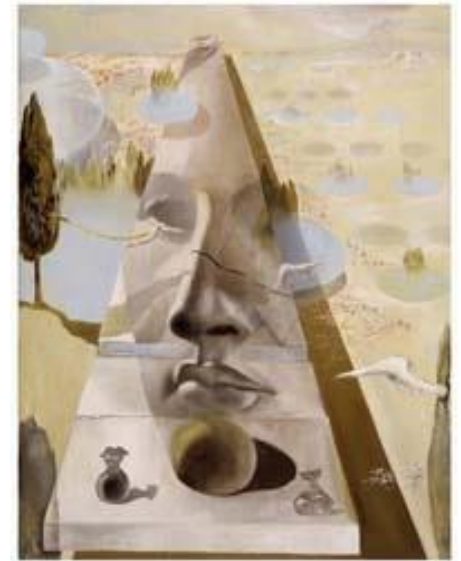


Viola, Jones 2001,
Torralba et al. 2004,
Opelt et al. 2006,...

Random forests

Lecture 13

Visual recognition



- Object classification bag of words models
 - Discriminative methods
 - Generative methods
- Object classification by PCA and FLD

Image classification



$$p(\textit{zebra} \mid \textit{image})$$

vs.

$$p(\textit{no zebra} \mid \textit{image})$$

- Bayes rule:

$$\frac{p(\textit{zebra} \mid \textit{image})}{p(\textit{no zebra} \mid \textit{image})}$$

posterior ratio

$$\frac{p(\textit{image} \mid \textit{zebra})}{p(\textit{image} \mid \textit{no zebra})} \cdot \frac{p(\textit{zebra})}{p(\textit{no zebra})}$$

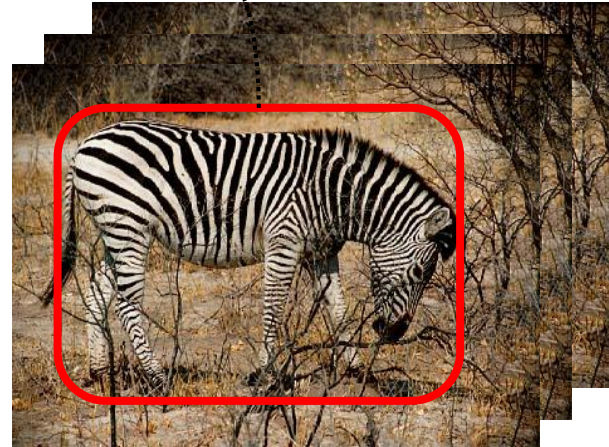
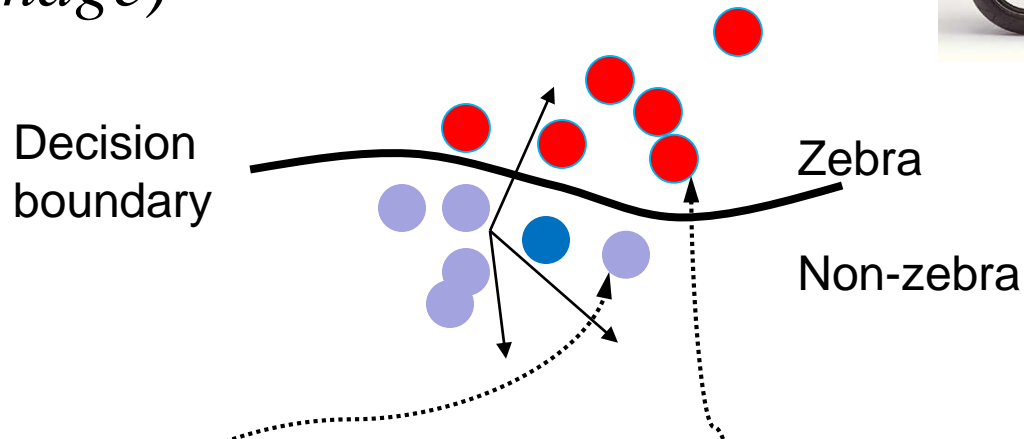
likelihood ratio

prior ratio

Discriminative methods

- Modeling the posterior ratio:

$$\frac{p(\text{zebra} | \text{image})}{p(\text{no zebra} | \text{image})}$$



Generative methods



$$p(\textit{zebra} \mid \textit{image})$$

vs.

$$p(\textit{no zebra} \mid \textit{image})$$

- Bayes rule:

$$\underbrace{\frac{p(\textit{zebra} \mid \textit{image})}{p(\textit{no zebra} \mid \textit{image})}}_{\text{posterior ratio}} = \underbrace{\frac{p(\textit{image} \mid \textit{zebra})}{p(\textit{image} \mid \textit{no zebra})}}_{\text{likelihood ratio}} \cdot \underbrace{\frac{p(\textit{zebra})}{p(\textit{no zebra})}}_{\text{prior ratio}}$$

Generative models

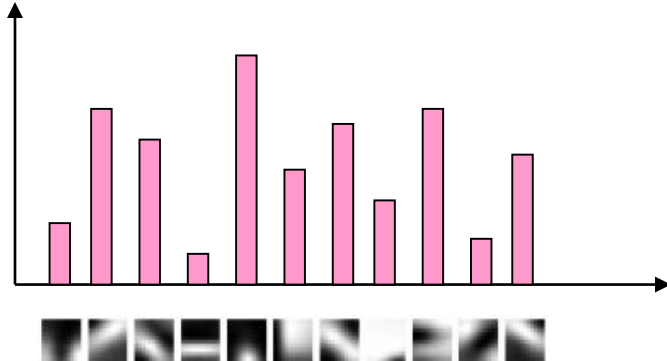
1. Naïve Bayes classifier

- Csurka Bray, Dance & Fan, 2004

2. Hierarchical Bayesian text models (pLSA and LDA)

- Background: Hoffman 2001, Blei, Ng & Jordan, 2004
- Object categorization: Sivic et al. 2005, Sudderth et al. 2005
- Natural scene categorization: Fei-Fei et al. 2005

Some notations



- **w**: a collection of all N codewords in the image

$$\mathbf{w} = [w_1, w_2, \dots, w_N]$$

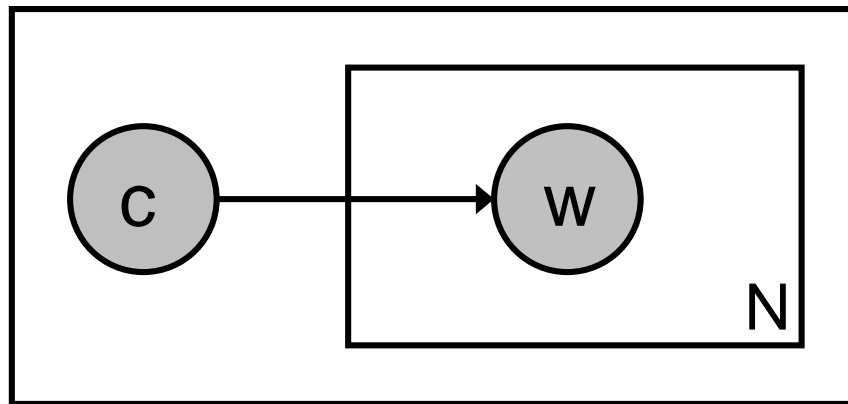
- **c**: category of the image

the Naïve Bayes model

$$p(c | w) \sim p(c)p(w | c) = p(c) p(w_1, \dots, w_N | c)$$

Prior prob. of
the object classes

Image likelihood
given the class



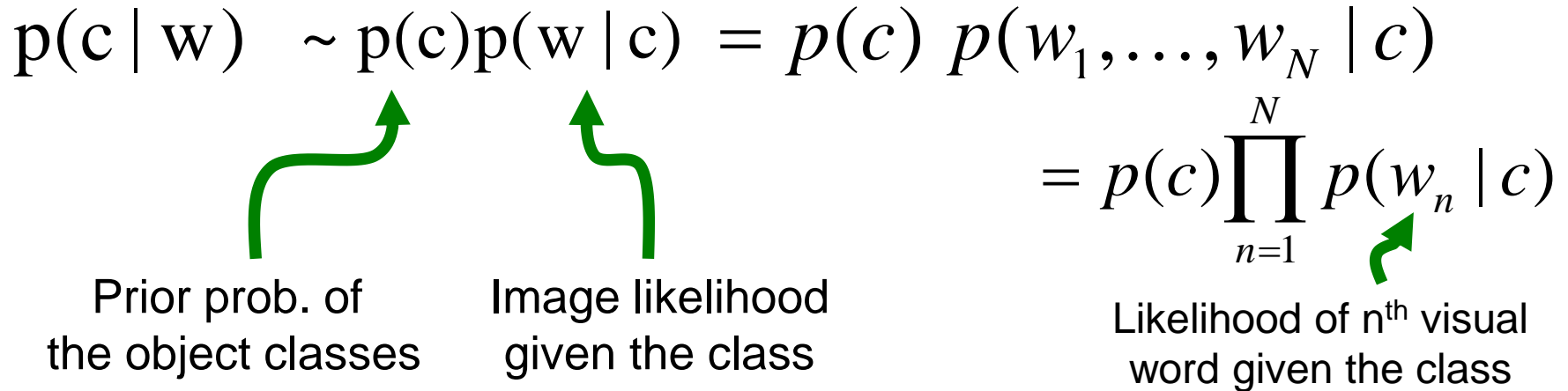
the Naïve Bayes model

$$p(c | w) \sim p(c)p(w | c) = p(c) p(w_1, \dots, w_N | c)$$
$$= p(c) \prod_{n=1}^N p(w_n | c)$$

Prior prob. of the object classes

Image likelihood given the class

Likelihood of n^{th} visual word given the class



- Assume that each feature (codewords) is conditionally independent *given the class*

$$p(w_1, \dots, w_N | c) = \prod_{i=1}^N p(w_i | c)$$

the Naïve Bayes model

$$p(c | w) \sim p(c)p(w | c) = p(c) p(w_1, \dots, w_N | c)$$
$$= p(c) \prod_{n=1}^N p(w_n | c)$$

Prior prob. of the object classes

Image likelihood given the class

Likelihood of n^{th} visual word given the class

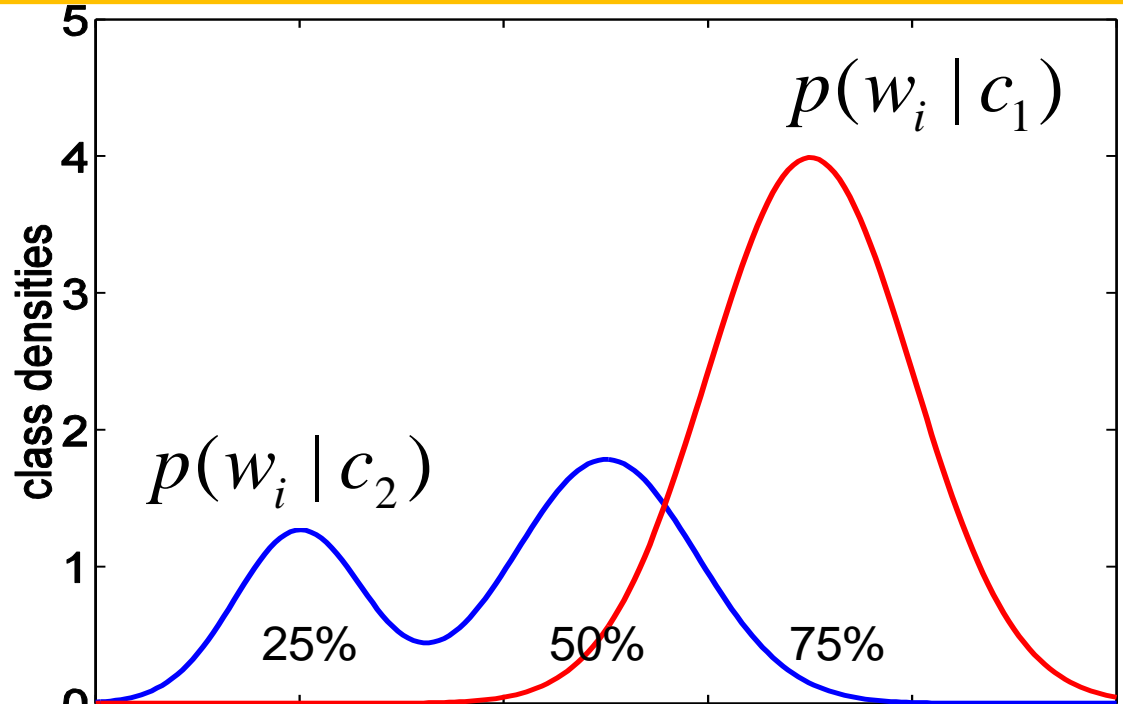
Example:

2 classes:
bananas vs oranges

Histogram of colors

w_i = number of pixels colored in yellow in the image

x-axis: percentage of pixel that are colored in yellow in the image



the Naïve Bayes model

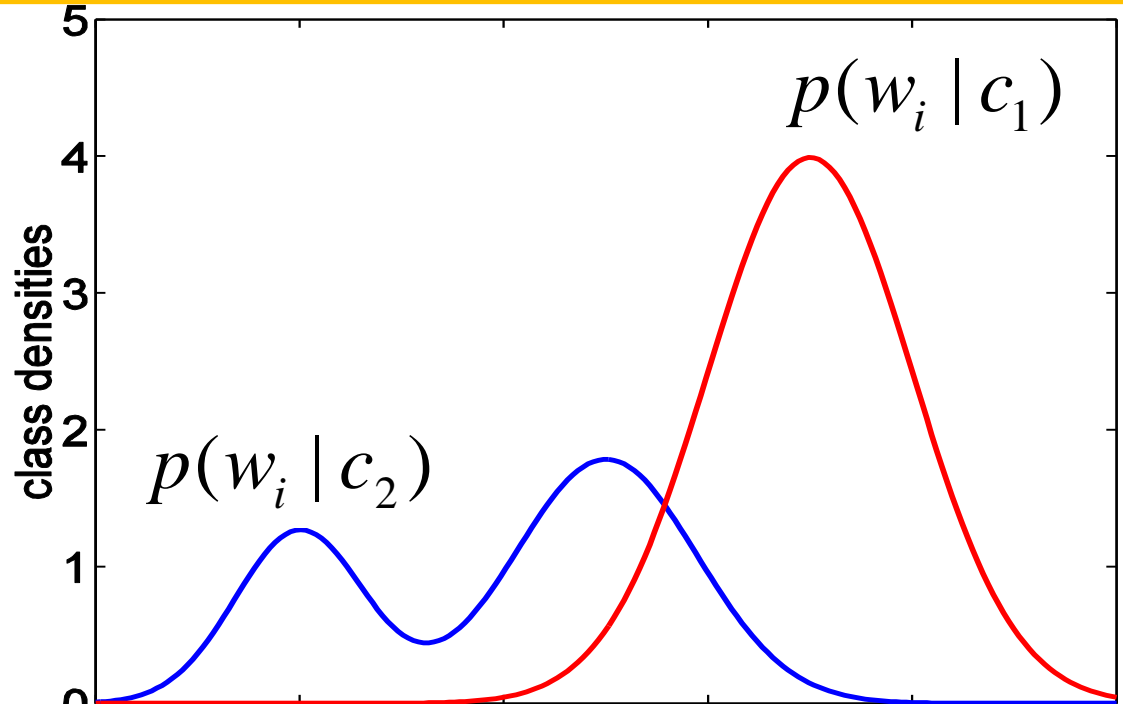
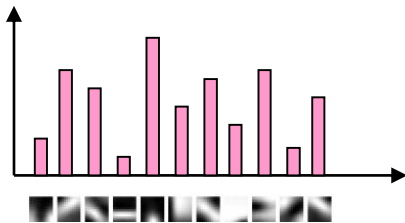
$$p(c | w) \sim p(c)p(w | c) = p(c) p(w_1, \dots, w_N | c)$$
$$= p(c) \prod_{n=1}^N p(w_n | c)$$

Prior prob. of the object classes

Image likelihood given the class

Likelihood of n^{th} visual word given the class

- How do we learn $P(w_i|c_j)$?
- From empirical frequencies of code words in images from a given class



Classification/Recognition

$$c^* = \arg \max_c p(c | w) \propto p(c) \prod_{n=1}^N p(w_n | c)$$

Object class
decision

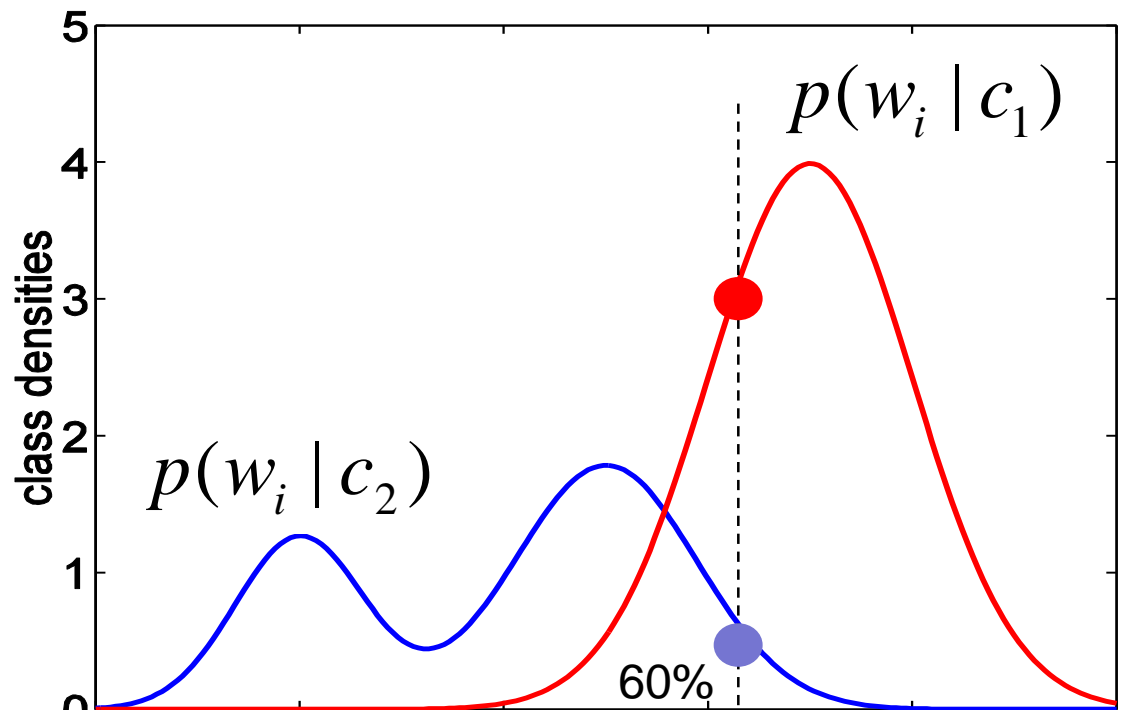
Example:

2 classes:
bananas vs oranges

Query image contains a banana

Look at how many pixels are
yellow: say 60%

Look at corresponding likelihood
values given the two class
hypotheses → banana!



Summary: Generative models

- Naïve Bayes
 - *Unigram models* in document analysis
 - Assumes conditional independence of words given class
 - Parameter estimation: frequency counting

Csurka's dataset – 7 classes

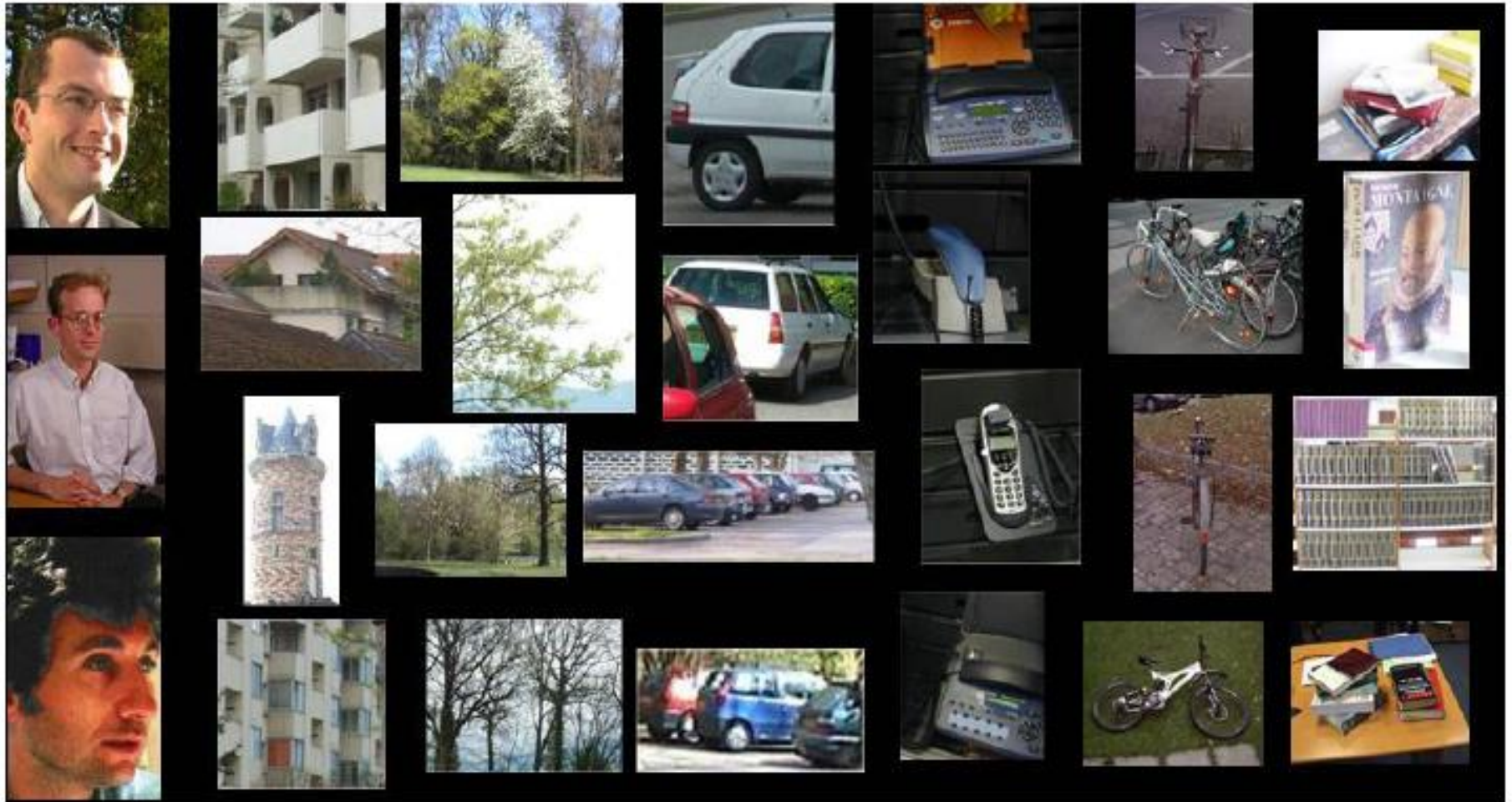


Table 1. Confusion matrix and the mean rank for the best vocabulary ($k=1000$).

True classes →	<i>faces</i>	<i>buildings</i>	<i>trees</i>	<i>cars</i>	<i>phones</i>	<i>bikes</i>	<i>books</i>
<i>faces</i>	76	4	2	3	4	4	13
<i>buildings</i>	2	44	5	0	5	1	3
<i>trees</i>	3	2	80	0	0	5	0
<i>cars</i>	4	1	0	75	3	1	4
<i>phones</i>	9	15	1	16	70	14	11
<i>bikes</i>	2	15	12	0	8	73	0
<i>books</i>	4	19	0	6	7	2	69
<i>Mean ranks</i>	1.49	1.88	1.33	1.33	1.63	1.57	1.57

E = 28%

Table 2. Confusion matrix and mean rank for SVM ($k=1000$, linear kernel).

True classes →	<i>faces</i>	<i>buildings</i>	<i>trees</i>	<i>cars</i>	<i>phones</i>	<i>bikes</i>	<i>books</i>
<i>faces</i>	98	14	10	10	34	0	13
<i>buildings</i>	1	63	3	0	3	1	6
<i>trees</i>	1	10	81	1	0	6	0
<i>cars</i>	0	1	1	85	5	0	5
<i>phones</i>	0	5	4	3	55	2	3
<i>bikes</i>	0	4	1	0	1	91	0
<i>books</i>	0	3	0	1	2	0	73
<i>Mean ranks</i>	1.04	1.77	1.28	1.30	1.83	1.09	1.39

E = 15%

Generative vs discriminative

- Discriminative methods
 - Computationally efficient & fast
- Generative models
 - Convenient for weakly- or un-supervised, incremental training
 - Prior information
 - Flexibility in modeling parameters

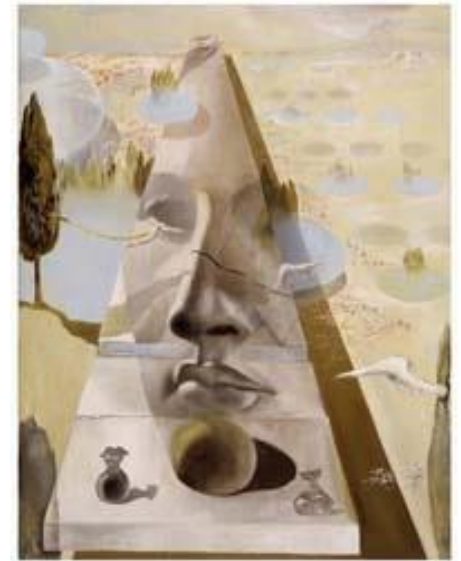
Weakness of BoW the models



- All have equal probability for bag-of-words methods
- Location information is important
- No rigorous geometric information of the object components
- Segmentation and localization unclear

Lecture 13

Visual recognition



- Object classification bag of words models
 - Discriminative methods
 - Generative methods
- Object classification by PCA and FLD

Object classification by...

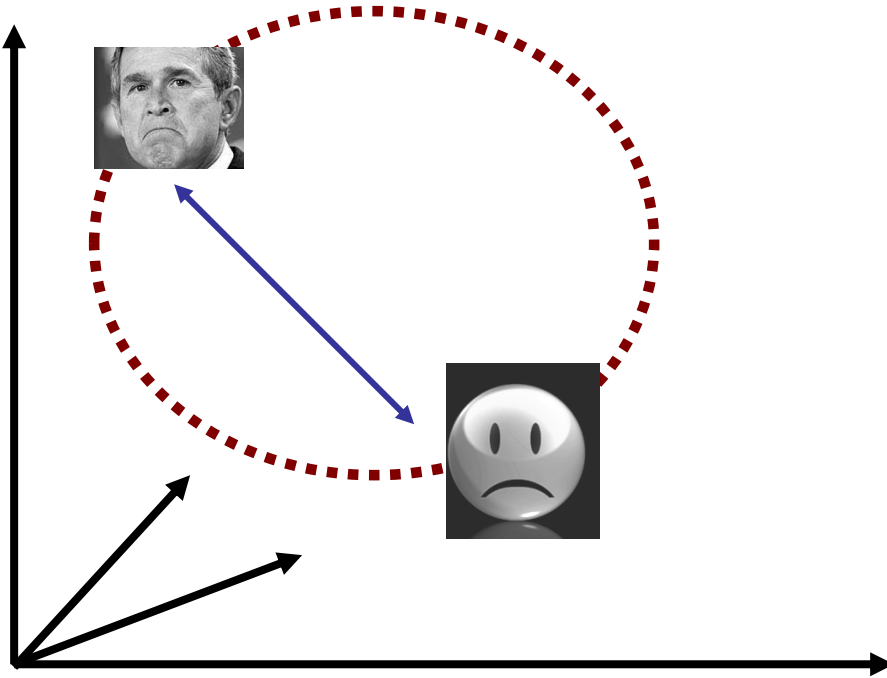
- Principle Component Analysis (PCA)
- Linear Discriminant Analysis (LDA)

Originally introduced for faces:

- Eigenfaces and Fisherfaces

Turk & Penland, 91
Belhumeur et al.,

The Space of images or histograms



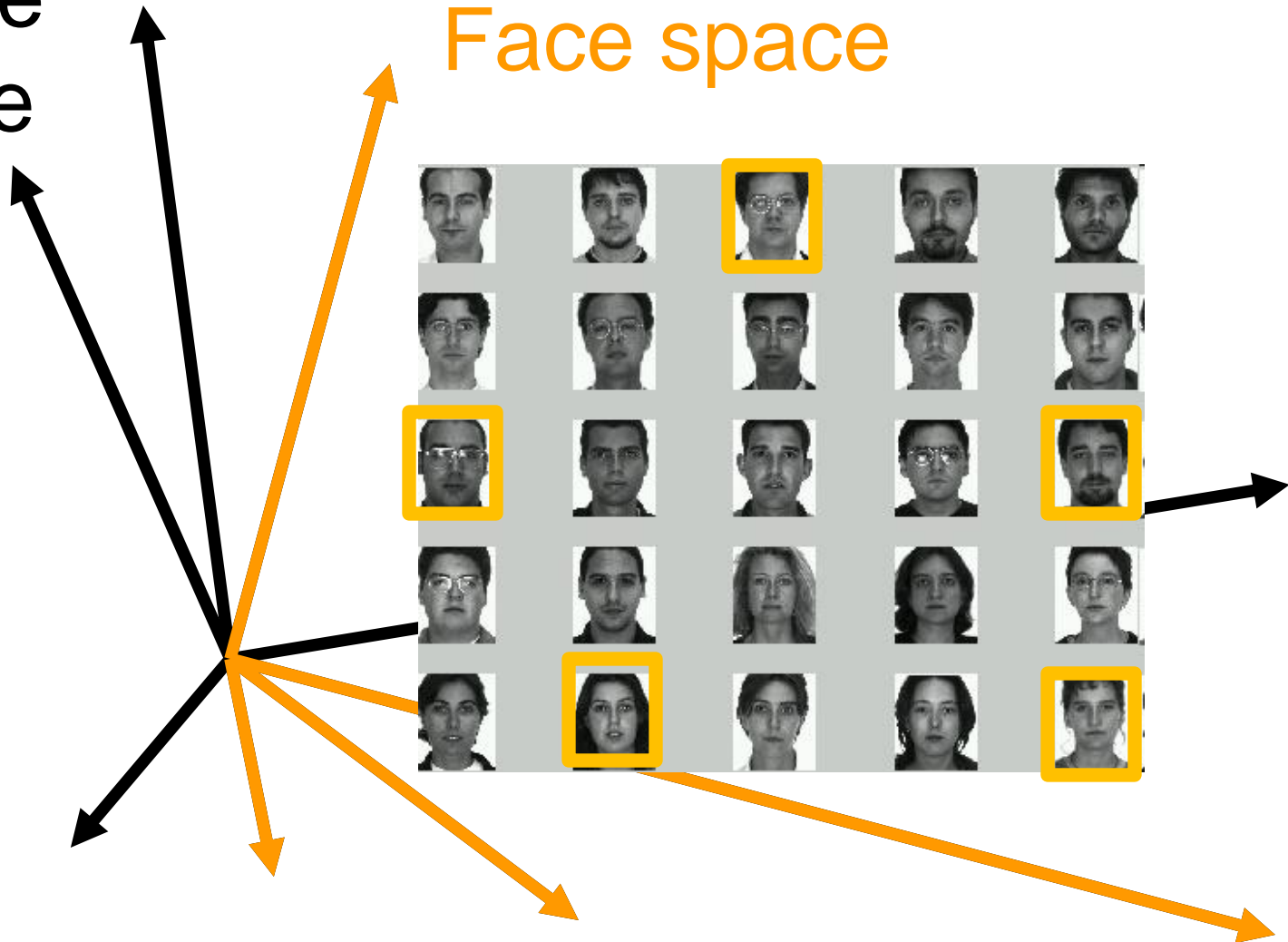
- An image (or histogram) H is a point in a high dimensional space
 - An $N \times M$ image is a point in \mathbb{R}^{NM}

Key Idea

- H in the possible set $\mathcal{X} = \{\hat{\mathbf{x}}\}$ are highly correlated.
- So, compress them to a low-dimensional subspace that captures key appearance characteristics of the visual DOFs.
- **USE PCA for estimating the sub-space**
(dimensionality reduction)
- Compare two objects by projecting the images into the subspace and measuring the EUCLIDEAN distance between them.

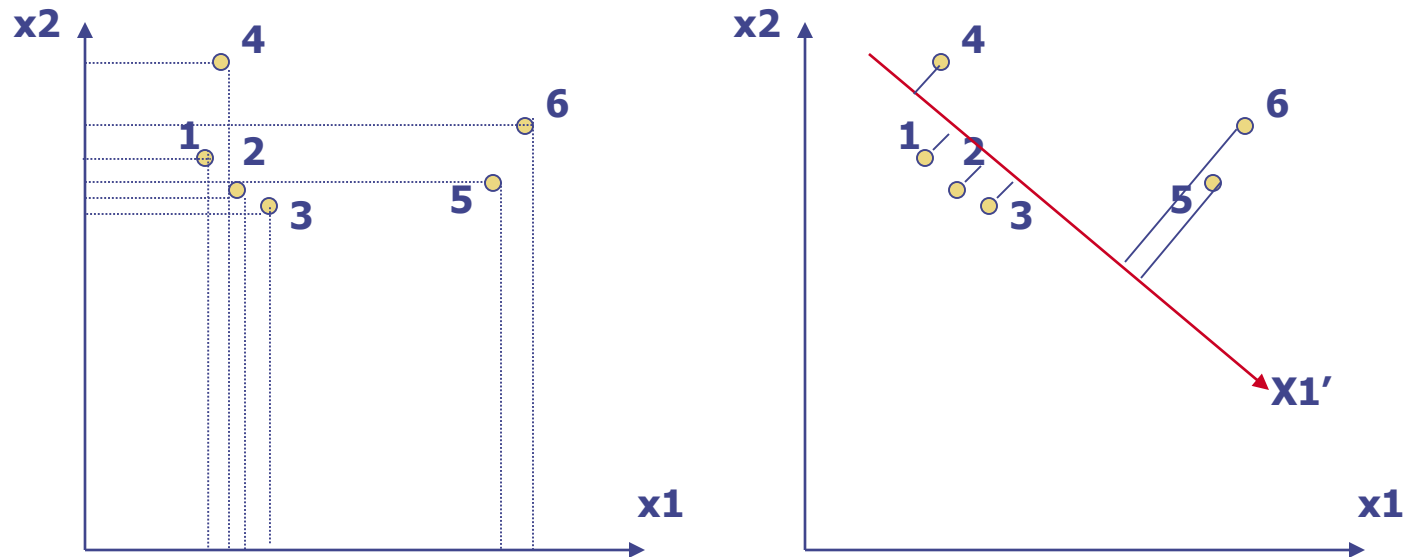
Image space

Face space



- Computes n-dim subspace such that the projection of the data points onto the subspace has **the largest variance** among all n-dim subspaces.
- **Maximize the scatter** of the training images in face space

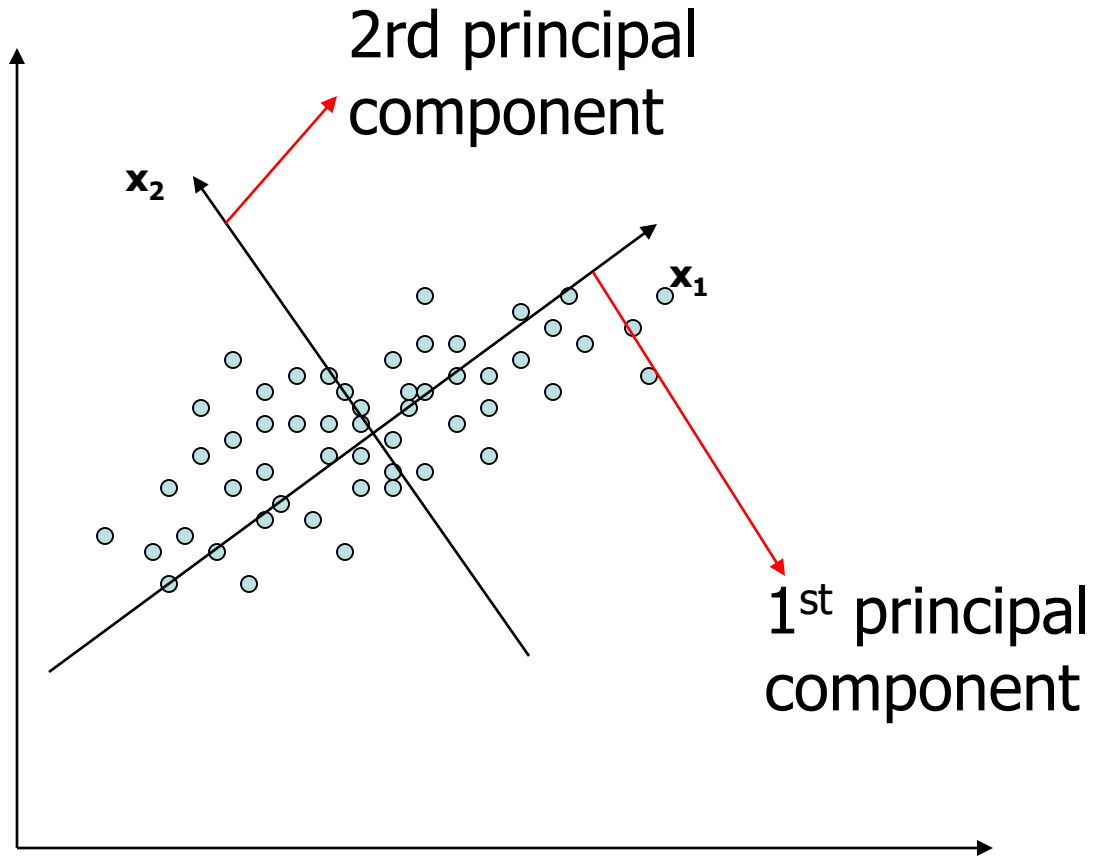
USE PCA for estimating the sub-space



PCA projection

- Computes n -dim subspace such that the projection of the data points onto the subspace has **the largest variance** among all n -dim subspaces.

USE PCA for estimating the sub-space



PCA Mathematical Formulation

PCA = eigenvalue decomposition of a data covariance matrix

Define a transformation, W ,

$$y_j = W^T x_j \quad j = 1, 2 \dots N$$



$$S_T = \sum_{j=1}^N (x_j - \bar{x})(x_j - \bar{x})^T = \text{Data Scatter matrix}$$

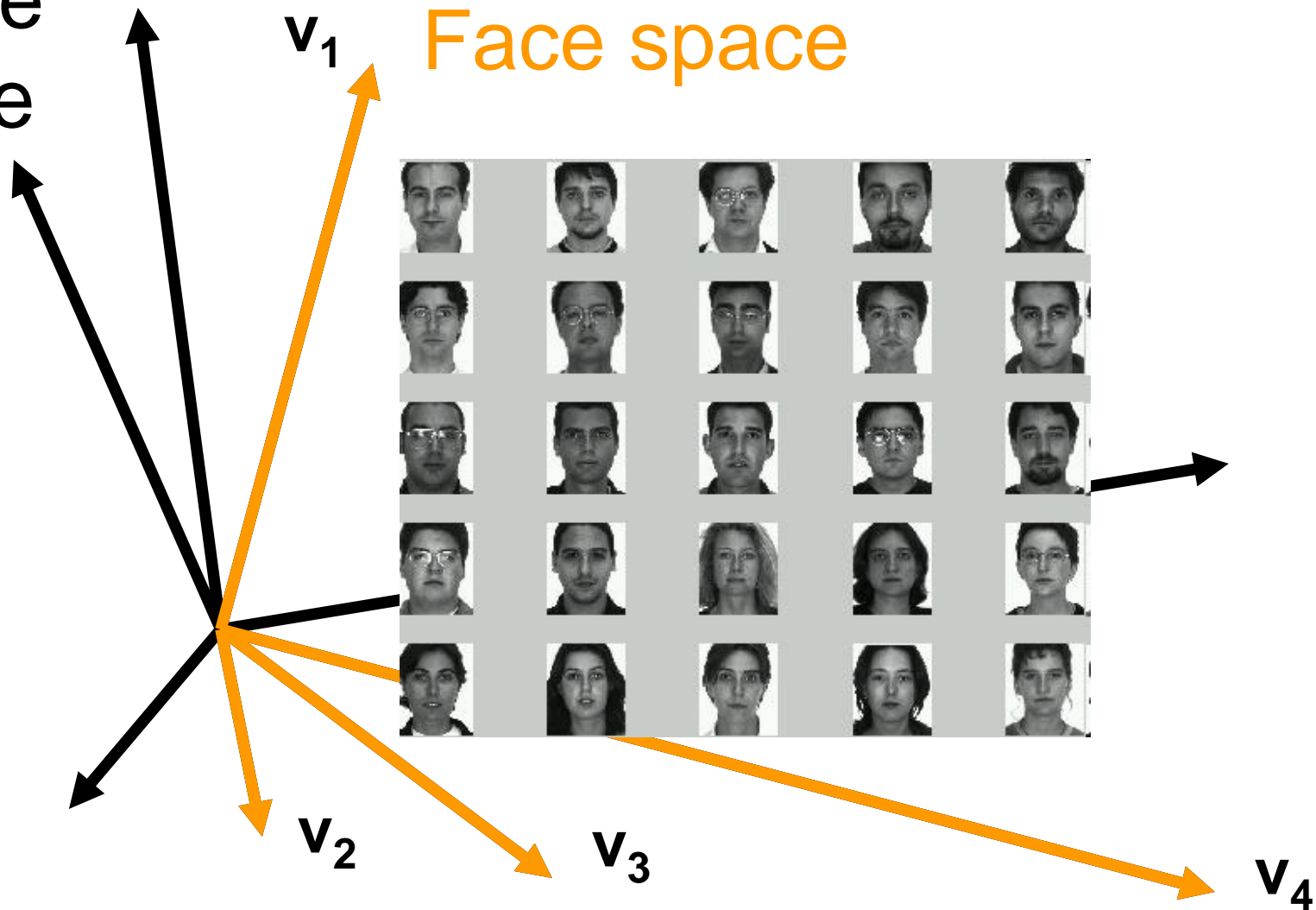
$$\tilde{S}_T = \sum_{j=1}^N (y_j - \bar{y})(y_j - \bar{y})^T = W^T S_T W = \text{Transf. data scatter matrix}$$

Measure data scatter **Eigenvectors of S_T**

$$W_{opt} = \arg \max_W |W^T S_T W| = [\mathbf{v}_1 \quad \mathbf{v}_2 \quad \dots \quad \mathbf{v}_m]$$

Image space

Face space



Projecting onto the Eigenfaces

- The eigenfaces $\mathbf{v}_1, \dots, \mathbf{v}_K$ span the space of faces
 - A face is converted to eigenface coordinates by

$$\mathbf{x} \rightarrow \left(\underbrace{(\mathbf{x} - \bar{\mathbf{x}}) \cdot \mathbf{v}_1}_{a_1}, \underbrace{(\mathbf{x} - \bar{\mathbf{x}}) \cdot \mathbf{v}_2}_{a_2}, \dots, \underbrace{(\mathbf{x} - \bar{\mathbf{x}}) \cdot \mathbf{v}_K}_{a_K} \right)$$

$$\mathbf{x} \approx \bar{\mathbf{x}} + a_1 \mathbf{v}_1 + a_2 \mathbf{v}_2 + \dots + a_K \mathbf{v}_K$$



\mathbf{x}



$a_1 \mathbf{v}_1$ $a_2 \mathbf{v}_2$ $a_3 \mathbf{v}_3$ $a_4 \mathbf{v}_4$ $a_5 \mathbf{v}_5$ $a_6 \mathbf{v}_6$ $a_7 \mathbf{v}_7$ $a_8 \mathbf{v}_8$



Algorithm

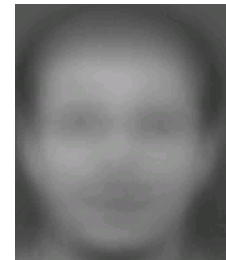
Training

1. Align training images x_1, x_2, \dots, x_N



Note that each image is formulated into a long vector!

2. Compute average face $\bar{x} = 1/N \sum x_i$



3. Compute the difference image $x_i - \bar{x}$

Algorithm

4. Compute the covariance matrix (total scatter matrix)

$$S_T = \sum_{j=1}^N (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})^T$$

5. Compute the eigenvectors of the covariance matrix S_T
6. Compute training projections a_1, a_2, \dots, a_N

Testing

1. Take query image X
2. Project X into Eigenface space ($W = \{\text{eigenfaces}\}$) and compute projection ω_i
3. Compare projection ω_i with all training N projections a_i

Illustration of Eigenfaces

- The visualization of eigenvectors:



These are the first 4 eigenvectors from a training set of 400 images (ORL Face Database).



Eigenfaces look somewhat like generic faces.

Reconstruction and Errors

$P = 4$



$P = 200$



$P = 400$



- Only selecting the top P eigenfaces \rightarrow reduces the dimensionality.
- Fewer eigenfaces result in more information loss, and hence less discrimination between faces.

Summary for Eigenface

Pros

- Non-iterative, globally optimal solution

Limitations

- PCA projection is **optimal for reconstruction** from a low dimensional basis, but **may NOT be optimal for discrimination...**

Extensions

- Generalized PCA:

R. Vidal, Y. Ma, and S. Sastry. [Generalized Principal Component Analysis \(GPCA\)](#). IEEE Transactions on Pattern Analysis and Machine Intelligence, volume 27, number 12, pages 1 - 15,

2005.

- Tensor Faces:

"[Multilinear Analysis of Image Ensembles: TensorFaces](#)," M.A.O. Vasilescu, D. Terzopoulos, *Proc. 7th European Conference on Computer Vision (ECCV'02)*, Copenhagen, Denmark, May, 2002

- PCA-SIFT

[PCA-SIFT: A More Distinctive Representation for Local Image Descriptors](#) -
Y Ke, R Sukthankar - IEEE CVPR 04

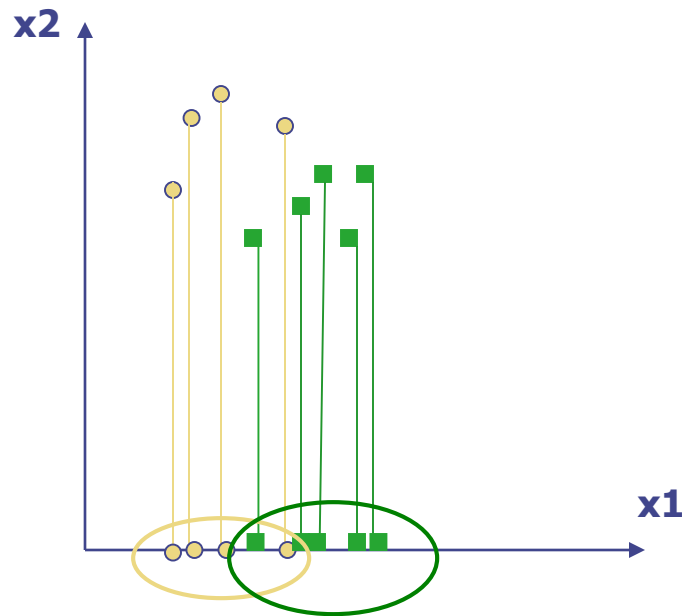
Linear Discriminant Analysis (LDA)

Fisher's Linear Discriminant (FLD)

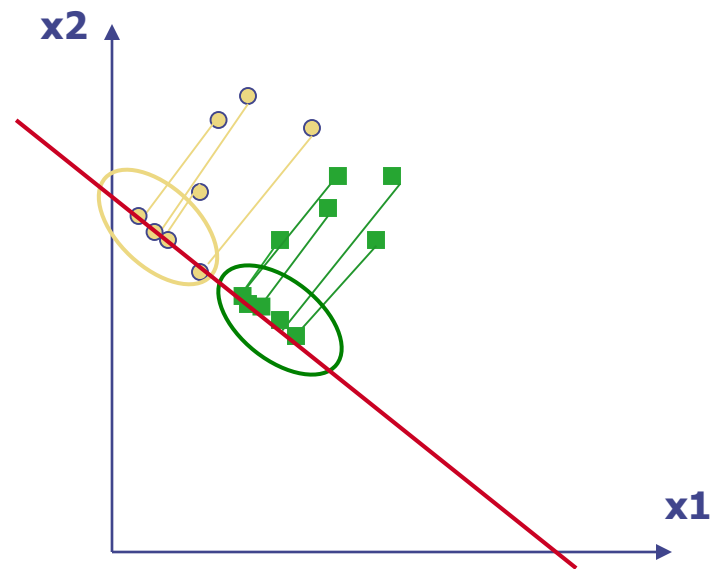
- Eigenfaces exploit the max scatter of the training images in face space
- Fisherfaces attempt to maximise the **between class scatter**, while minimising the **within class scatter**.

Illustration of the Projection

- ◆ Using two classes as example:



Poor Projection



Good

Variables

- N Sample images: $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$
- c classes: $\{\mathcal{X}_1, \dots, \mathcal{X}_c\}$
- Average of each class: $\mu_i = \frac{1}{N_i} \sum_{\mathbf{x}_k \in \mathcal{X}_i} \mathbf{x}_k$
- Total average: $\mu = \frac{1}{N} \sum_{k=1}^N \mathbf{x}_k$

Scatters

- Scatter of class i :

$$S_i = \sum_{x_k \in \mathcal{X}_i} (x_k - \mu_i)(x_k - \mu_i)^T$$

- Within class scatter:

$$S_W = \sum_{i=1}^c S_i$$

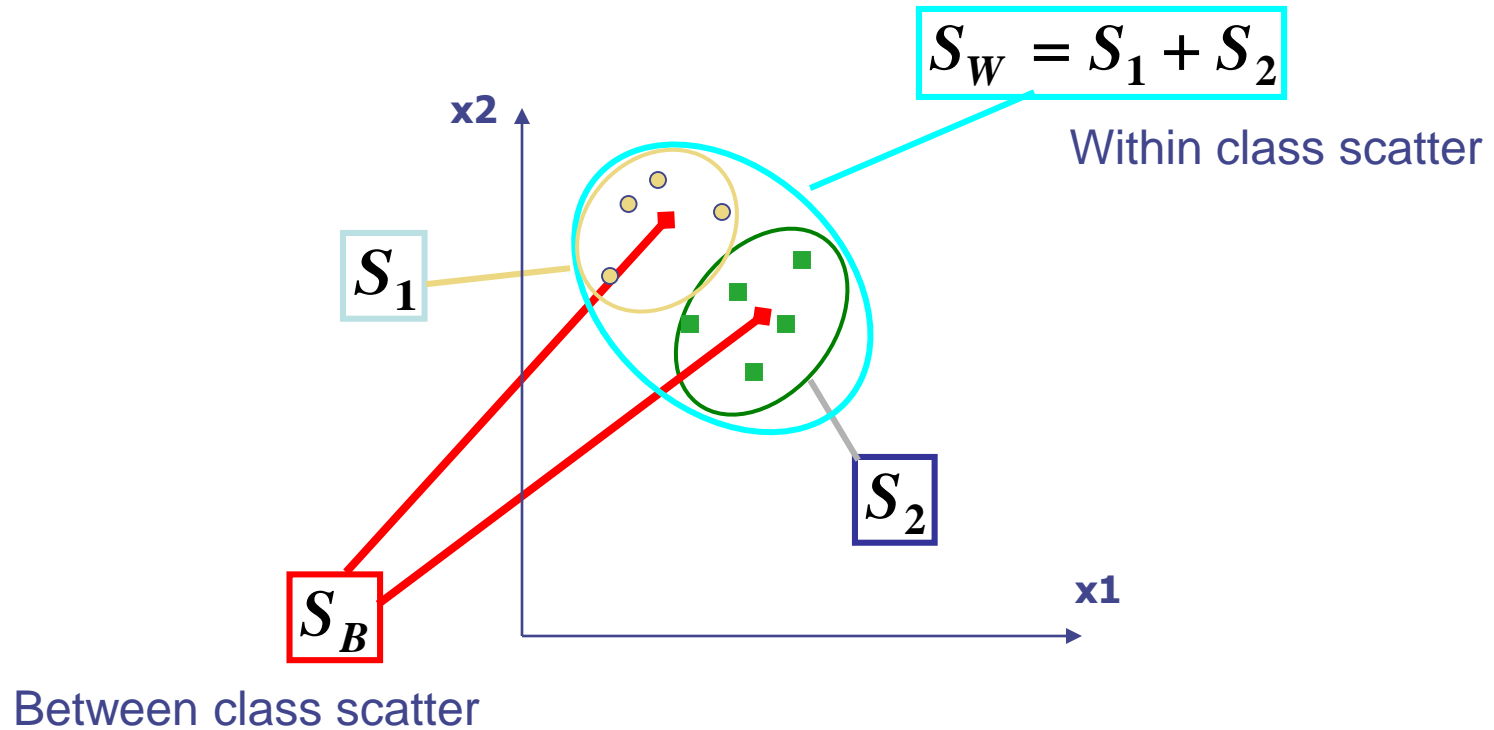
- Between class scatter:

$$S_B = \sum_{i=1}^c |\mathcal{X}_i| (\mu_i - \mu)(\mu_i - \mu)^T$$

- Total scatter:

$$S_T = S_W + S_B$$

Illustration



$$S_i = \sum_{x_k \in \mathcal{X}_i} (x_k - \mu_i)(x_k - \mu_i)^T$$

$$S_W = \sum_{i=1}^c S_i$$

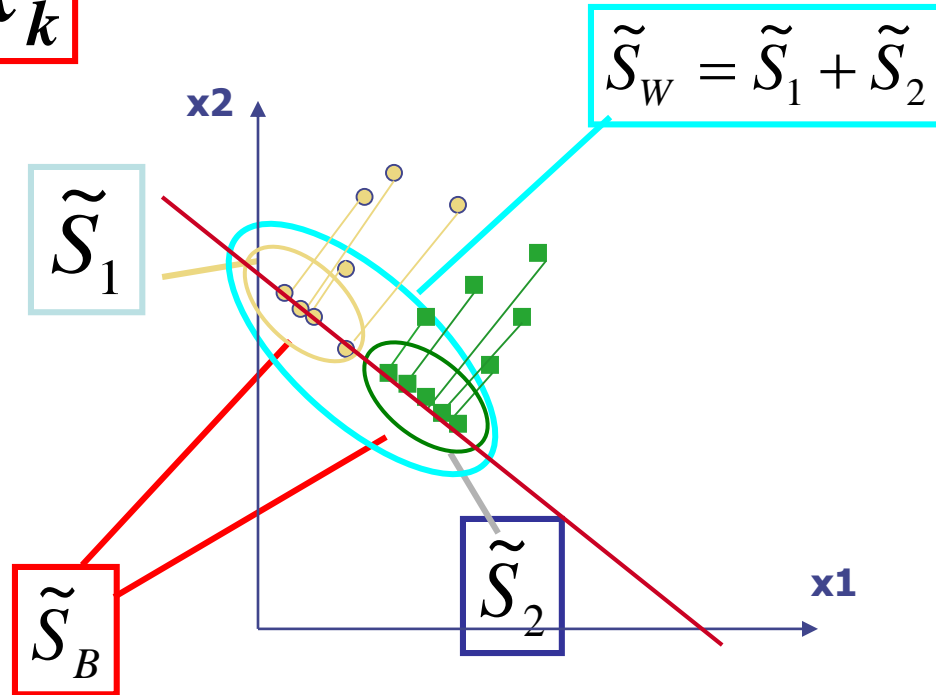
$$S_B = \sum_{i=1}^c |\mathcal{X}_i| (\mu_i - \mu)(\mu_i - \mu)^T$$

Mathematical Formulation (1)

- After projection: $y_k = W^T x_k$
- Between class scatter (of y's): $\tilde{S}_B = W^T S_B W$
- Within class scatter (of y's): $\tilde{S}_W = W^T S_W W$

Illustration

$$\mathbf{y}_k = \mathbf{W}^T \mathbf{x}_k$$



$$\mathbf{S}_W = \sum_{i=1}^c \mathbf{S}_i$$

$$\mathbf{S}_B = \sum_{i=1}^c |\chi_i| (\mu_i - \mu)(\mu_i - \mu)^T$$

$$\tilde{\mathbf{S}}_W = \mathbf{W}^T \mathbf{S}_W \mathbf{W}$$

$$\tilde{\mathbf{S}}_B = \mathbf{W}^T \mathbf{S}_B \mathbf{W}$$

Mathematical Formulation

- The desired projection:

$$\mathbf{W}_{opt} = \arg \max_{\mathbf{W}} \frac{|\tilde{\mathbf{S}}_B|}{|\tilde{\mathbf{S}}_W|} = \arg \max_{\mathbf{W}} \frac{|\mathbf{W}^T \mathbf{S}_B \mathbf{W}|}{|\mathbf{W}^T \mathbf{S}_W \mathbf{W}|}$$

- How is it found ? \rightarrow Generalized Eigenvectors
- If \mathbf{S}_W has full rank, the generalized eigenvectors are eigenvectors of $\mathbf{S}_W^{-1} \mathbf{S}_B$ with largest eigen-values

$$\mathbf{S}_B \mathbf{w}_i = \lambda_i \mathbf{S}_W \mathbf{w}_i \quad i = 1, \dots, m$$

Results: Eigenface vs. Fisherface

- Input: 160 images of 16 people
- Train: 159 images
- Test: 1 image
- Variation in Facial Expression, Eyewear, and Lighting

With
glasses

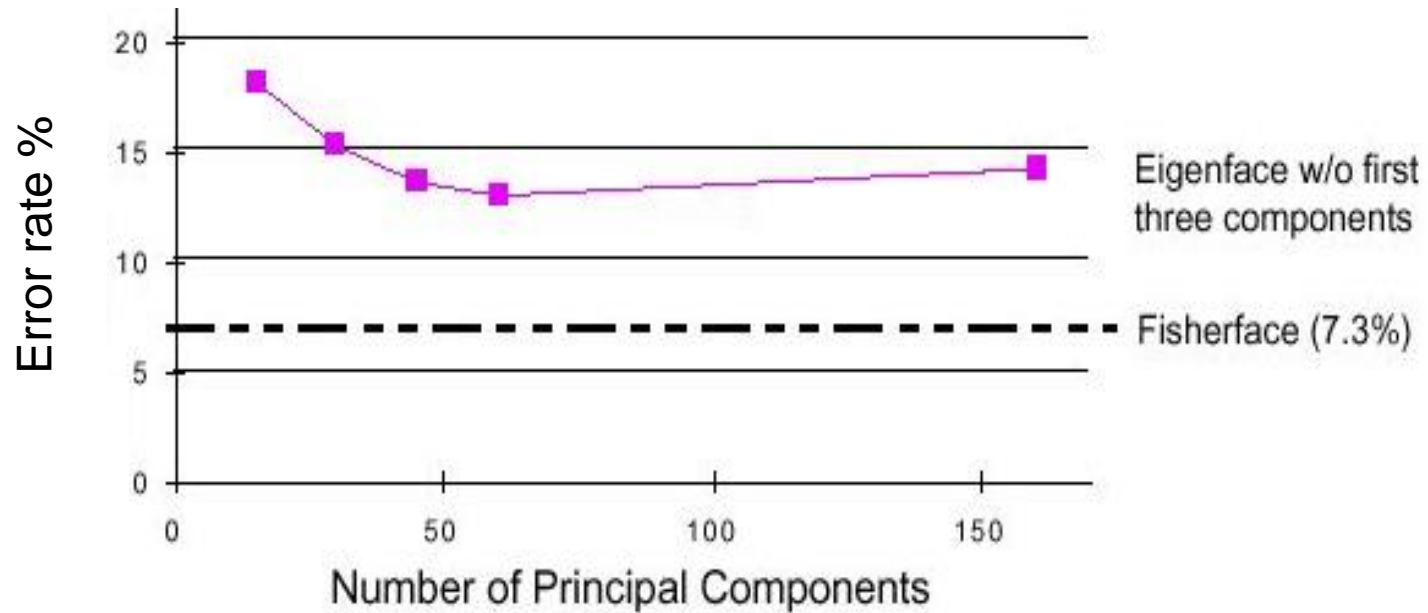
Without
glasses

3 Lighting
conditions

5 expressions



Results: Eigenface vs. Fisherface



Next lecture

- Object detection