

# CEE598 - Visual Sensing for Civil Infrastructure Eng. & Mgmt.

## Session 19 – Object Recognition II

***Mani Golparvar-Fard***

*Department of Civil and Environmental Engineering*

*3129D, Newmark Civil Engineering Lab*

*e-mail: [mgolpar@illinois.edu](mailto:mgolpar@illinois.edu)*

# Fun Time



Professor

Hard working students

Me?

# Outline

- Object Recognition
  - Introduction
  - Recognition of single 3D objects
    - Bag of word models
    - Part based models
    - Models for 3D objects categorization

# Challenges: object intra-class variation



# Usual Challenges

- Variability due to:
  - View point
  - Illumination
  - Occlusions

# Object categorization: the statistical viewpoint



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

vs.

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

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

- Bayes rule:

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

# Object categorization: the statistical viewpoint



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

vs.

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

- **Bayes rule:**

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

posterior ratio

likelihood ratio

prior ratio

# Object categorization: the statistical viewpoint

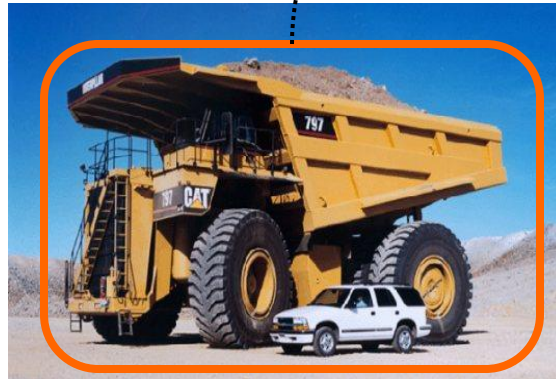
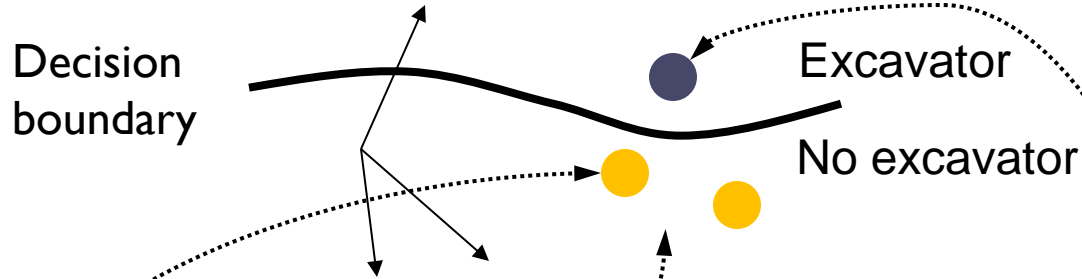
- **Discriminative methods** model posterior
- **Generative methods** model likelihood and prior

- **Bayes rule:**

$$\underbrace{\frac{p(\text{excavator} | \text{image})}{p(\text{no excavator} | \text{image})}}_{\text{posterior ratio}} = \underbrace{\frac{p(\text{image} | \text{excavator})}{p(\text{image} | \text{no excavator})}}_{\text{likelihood ratio}} \cdot \underbrace{\frac{p(\text{excavator})}{p(\text{no excavator})}}_{\text{prior ratio}}$$

# Discriminative

- Direct modeling of  $\frac{p(\text{excavator} \mid \text{image})}{p(\text{no excavator} \mid \text{image})}$



# Generative

$p(\text{image} \mid \text{excavator})$



$p(\text{image} \mid \text{no excavator})$



$p(\text{image} \mid \text{excavator})$	$p(\text{image} \mid \text{no excavator})$
Low	high
High	Low



# Learning

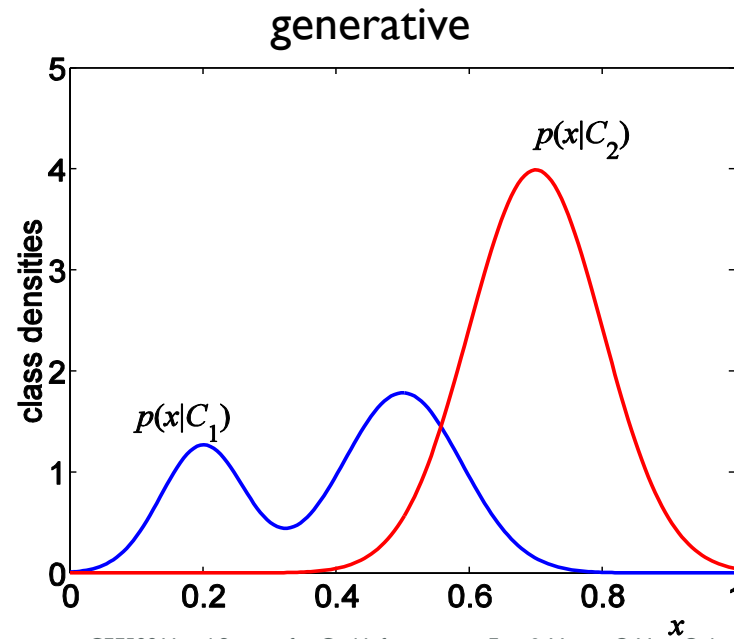
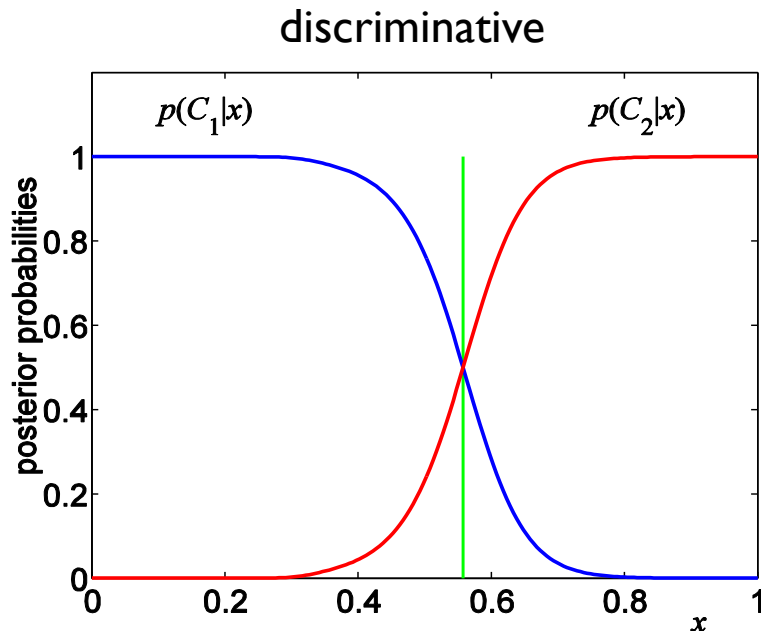
- machine learning useful to model intraclass variability



# Learning

- Learning parameters: What are you **maximizing**?  
**Likelihood** (Gen.) or **performances** on train/validation set (Disc.)

$$\underbrace{\frac{p(\text{excavator} | \text{image})}{p(\text{no excavator} | \text{image})}}_{\text{posterior ratio}} = \underbrace{\frac{p(\text{image} | \text{excavator})}{p(\text{image} | \text{no excavator})}}_{\text{likelihood ratio}} \cdot \underbrace{\frac{p(\text{excavator})}{p(\text{no excavator})}}_{\text{prior ratio}}$$



# Learning

- Learning parameters: What are you maximizing? Likelihood (Gen.) or performances on train/validation set (Disc.)
- Level of supervision
  - Manual segmentation; bounding box; image labels; noisy labels
- Batch/incremental (on category and image level; user-feedback )
- Priors
- Training images:
  - Issue of overfitting
  - Negative images for discriminative methods

Contains a excavator





# Bag of Words



Part of this segment is based on the tutorial “*Recognizing and Learning Object Categories: Year 2007*”

by Prof A. Torralba, R. Fergus and F. Li

# Related works

- Early “bag of words” models: mostly texture recognition
  - Cula & Dana, 2001; Leung & Malik 2001; Mori, Belongie & Malik, 2001; Schmid 2001; Varma & Zisserman, 2002, 2003; Lazebnik, Schmid & Ponce, 2003;
- Hierarchical Bayesian models for documents (pLSA, LDA, etc.)
  - Hoffman 1999; Blei, Ng & Jordan, 2004; Teh, Jordan, Beal & Blei, 2004
- Object categorization
  - Csurka, Bray, Dance & Fan, 2004; Sivic, Russell, Efros, Freeman & Zisserman, 2005; Sudderth, Torralba, Freeman & Willsky, 2005;
- Natural scene categorization
  - Vogel & Schiele, 2004; Fei-Fei & Perona, 2005; Bosch, Zisserman & Munoz, 2006

Object

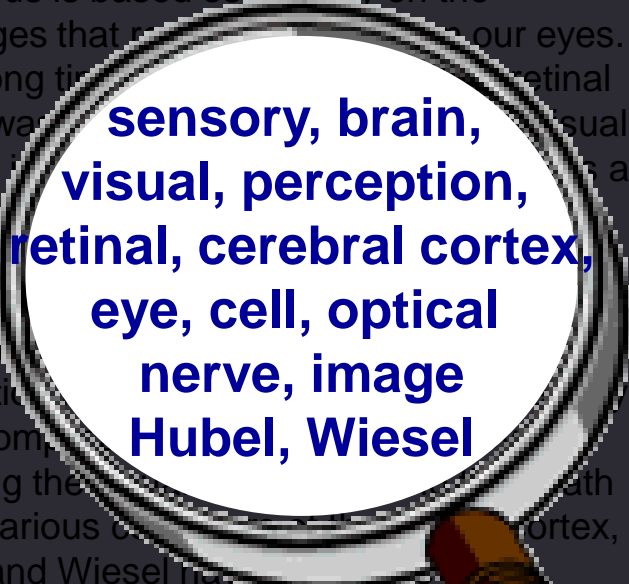


Bag of 'words'



# Analogy to documents

Of all the sensory impressions proceeding to the brain, the visual experiences are the dominant ones. Our perception of the world around us is based essentially on the messages that reach our eyes. For a long time, the retinal image was considered as a movie screen. The image is discovered and we know that perception is more complex following the path to the various cortex, Hubel and Wiesel have demonstrated that the *message about the image falling on the retina undergoes a wise analysis in a system of nerve cells stored in columns. In this system each cell has its specific function and is responsible for a specific detail in the pattern of the retinal image.*



**sensory, brain,  
visual, perception,  
retinal, cerebral cortex,  
eye, cell, optical  
nerve, image  
Hubel, Wiesel**

China is forecasting a trade surplus of \$90bn (£51bn) to \$100bn this year, a threefold increase on 2004's \$32bn. The Commerce Ministry said the surplus would be created by a predicted 30% increase in exports to \$750bn, compared with \$560bn in 2004. The government also needs to increase demand so that the country. China has permitted it to trade within a narrow band but the US wants the yuan to be allowed to move freely. However, Beijing has made it clear that it will take its time and tread carefully before allowing the yuan to rise further in value.

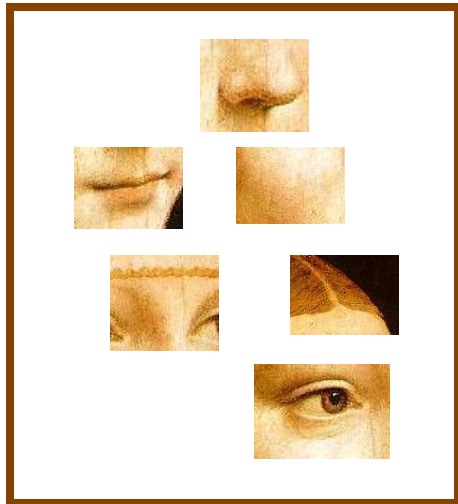


**China, trade,  
surplus, commerce,  
exports, imports, US,  
yuan, bank, domestic,  
foreign, increase,  
trade, value**

# Definition of “BoW”

- Independent features

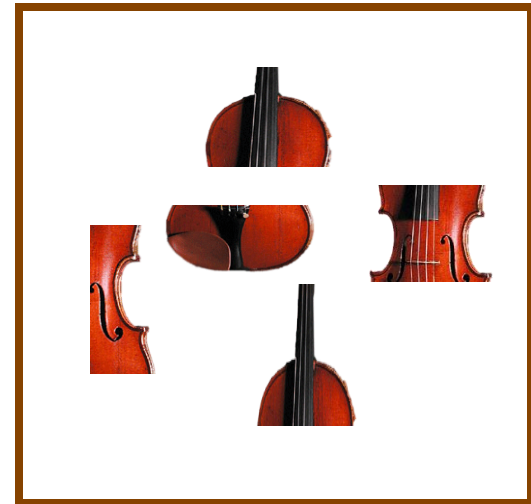
face



bike

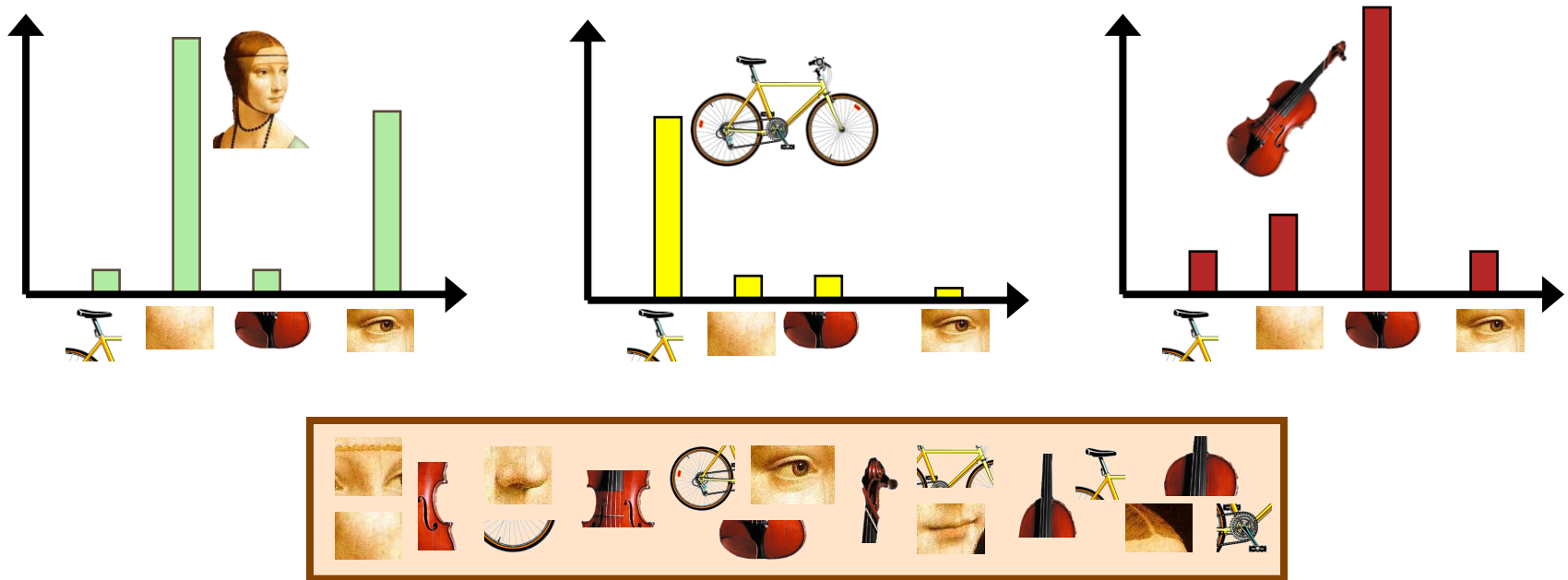


violin



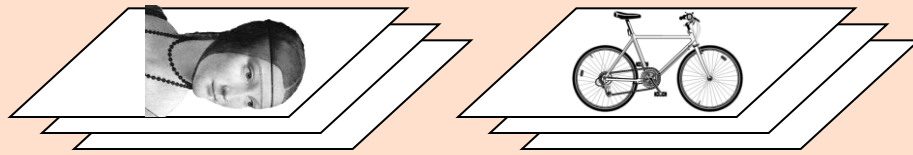
# definition of “BoW”

- Independent features
- histogram representation

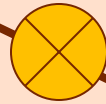


# Representation

# Learning and Recognition<sup>21</sup>



1. feature detection & representation



codewords dictionary

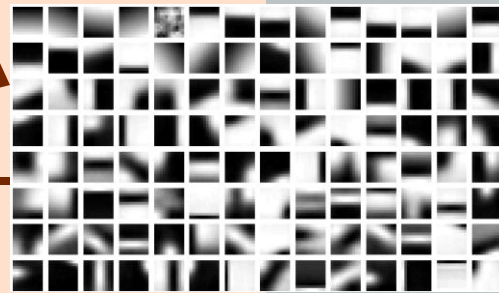
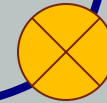


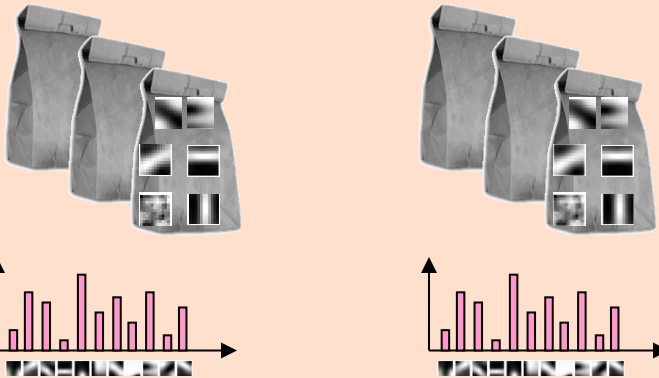
image representation



2.



3.

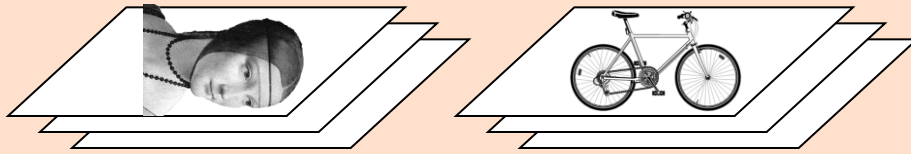


category models

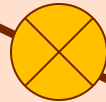


category decision

# Representation



1. feature detection & representation



2. **codewords dictionary**

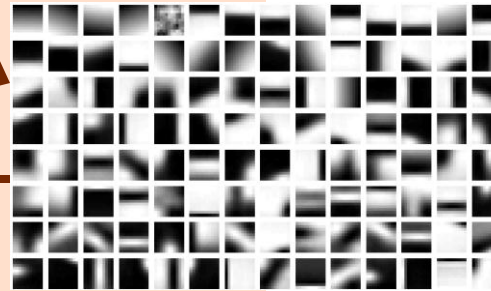
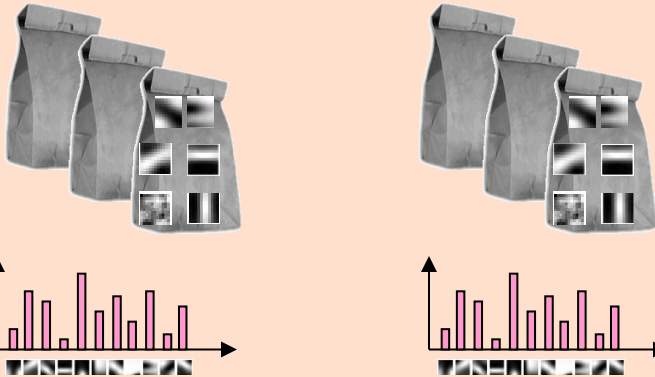


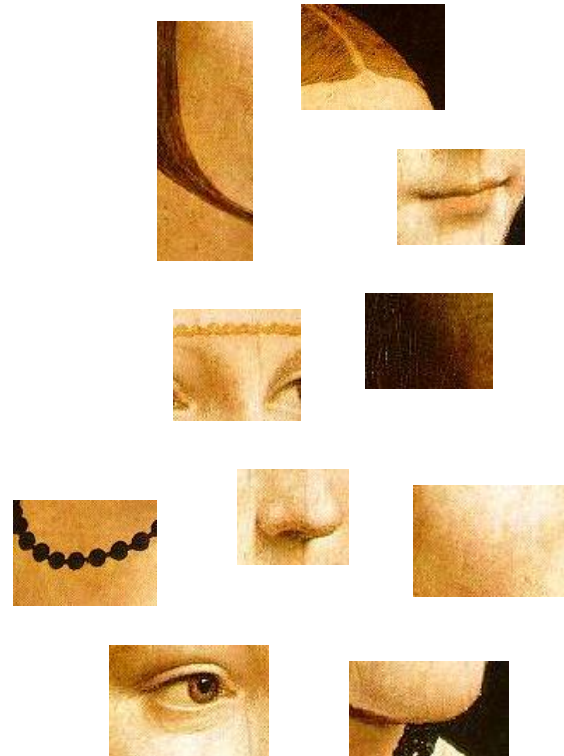
image representation

3.



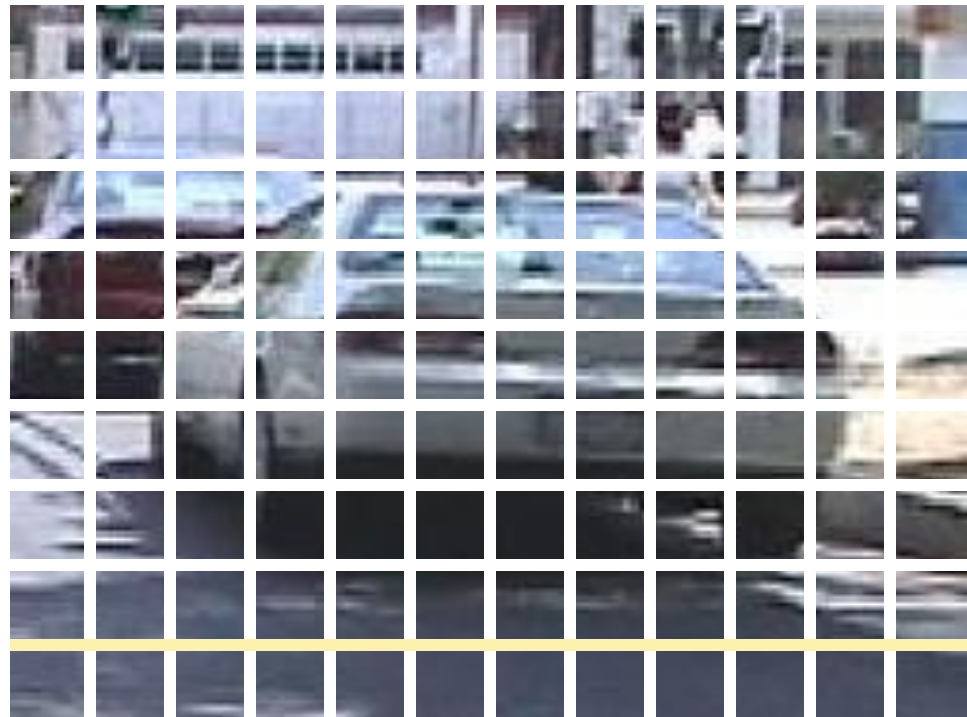
**category models**

# I. Feature detection and representation



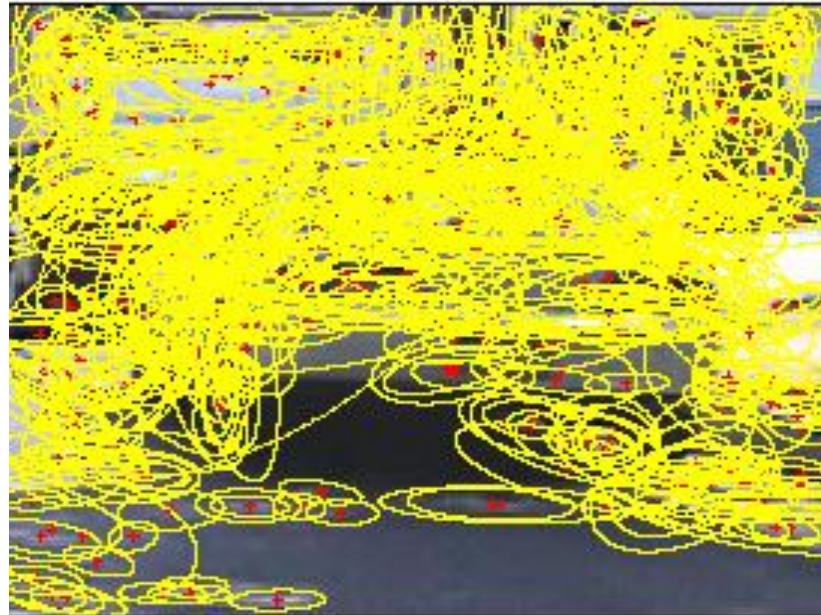
# I. Feature detection and representation

- Regular grid
  - Vogel & Schiele, 2003
  - Fei-Fei & Perona, 2005



# I. Feature detection and representation


- Regular grid
  - Vogel & Schiele, 2003
  - Fei-Fei & Perona, 2005
- Interest point detector
  - Csurka, et al. 2004
  - Fei-Fei & Perona, 2005
  - Sivic, et al. 2005



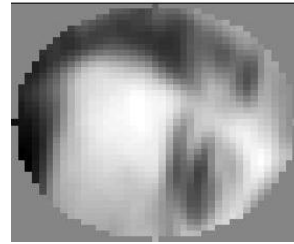
# I. Feature detection and representation

- Regular grid
  - Vogel & Schiele, 2003
  - Fei-Fei & Perona, 2005
  
- Interest point detector
  - Csurka, Bray, Dance & Fan, 2004
  - Fei-Fei & Perona, 2005
  - Sivic, Russell, Efros, Freeman & Zisserman, 2005
  
- Other methods
  - Random sampling (Vidal-Naquet & Ullman, 2002)
  - Segmentation based patches (Barnard, Duygulu, Forsyth, de Freitas, Blei, Jordan, 2003)

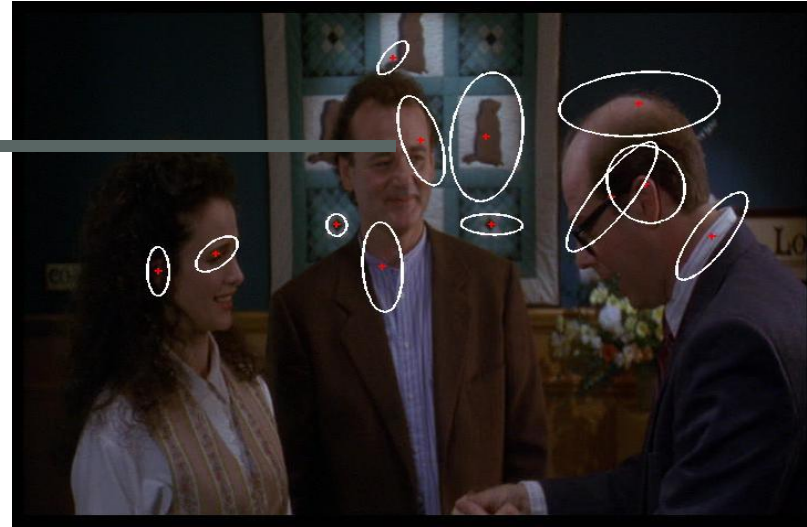
# I. Feature detection and representation



**Compute  
SIFT  
descriptor**  
[Lowe'99]



**Normalize  
patch**



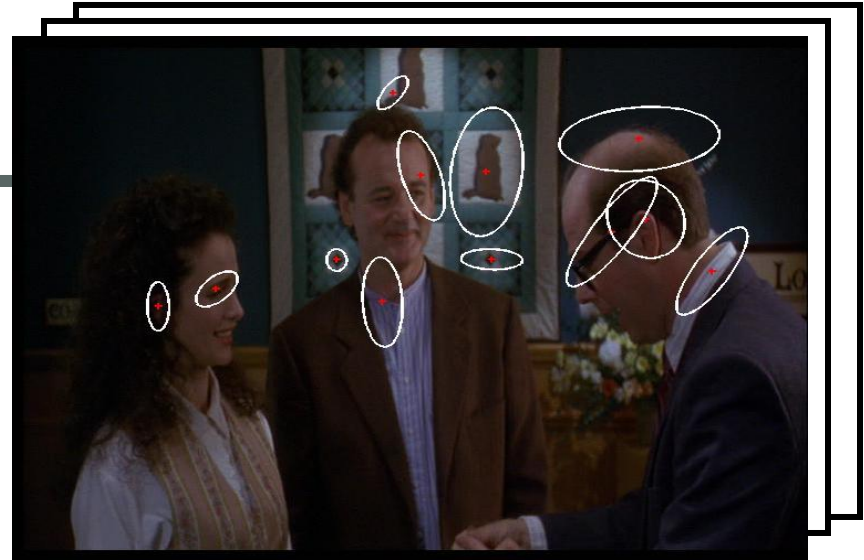
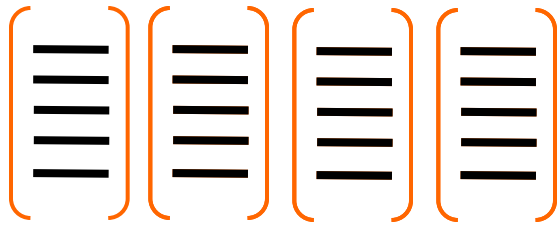
**Detect patches**

[Mikojaczyk and Schmid '02]

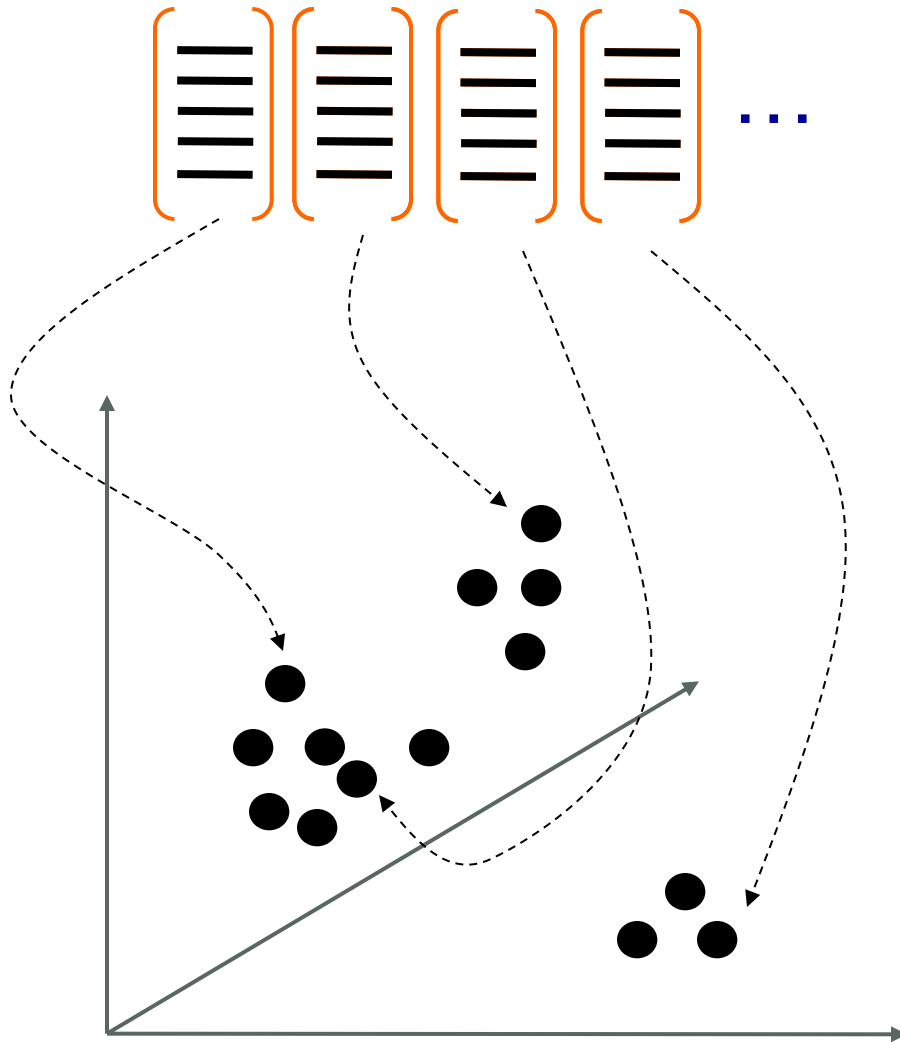
[Mata, Chum, Urban & Pajdla, '02]

[Sivic & Zisserman, '03]

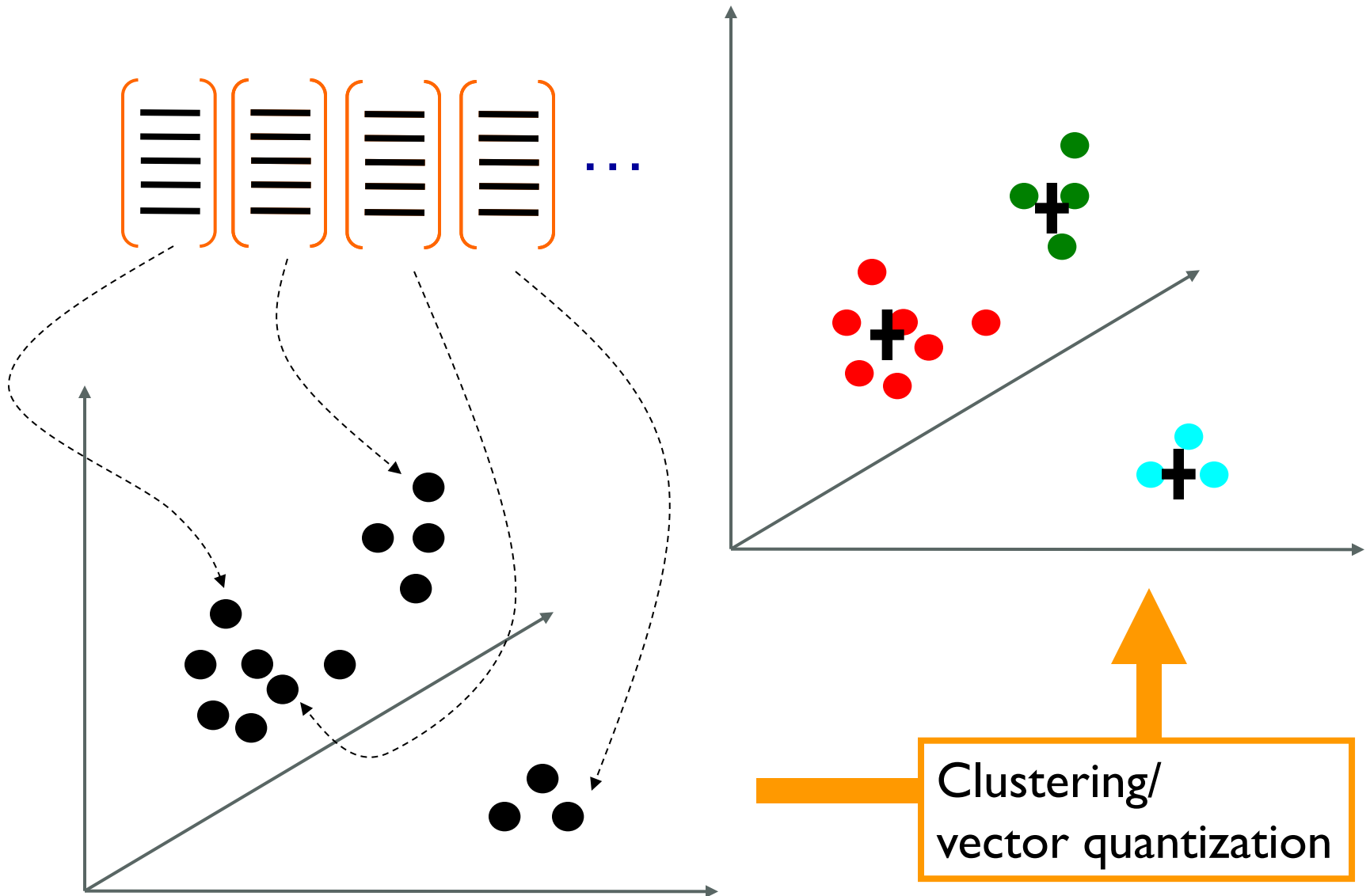
# I. Feature detection and representation



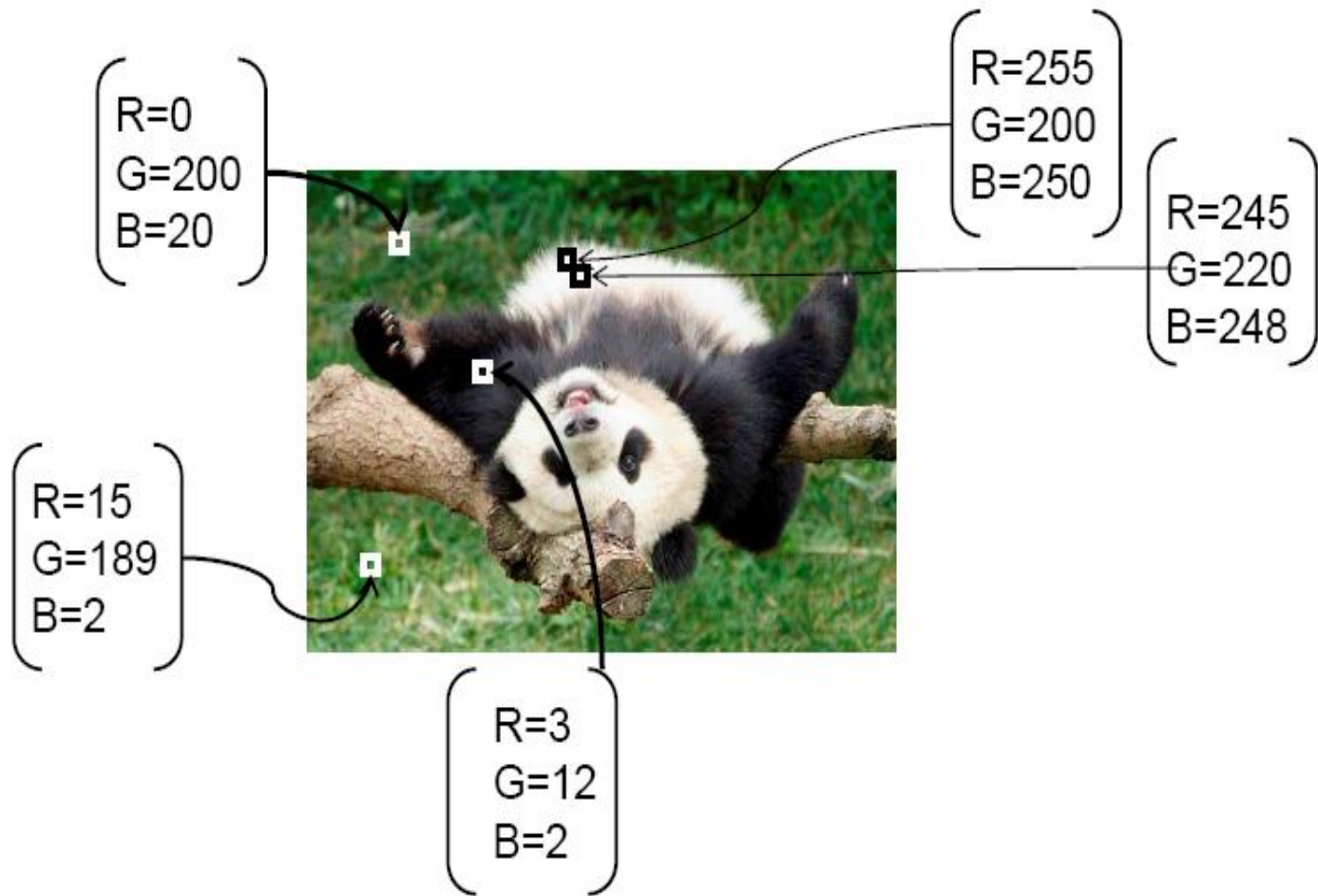
## 2. Codewords dictionary formation



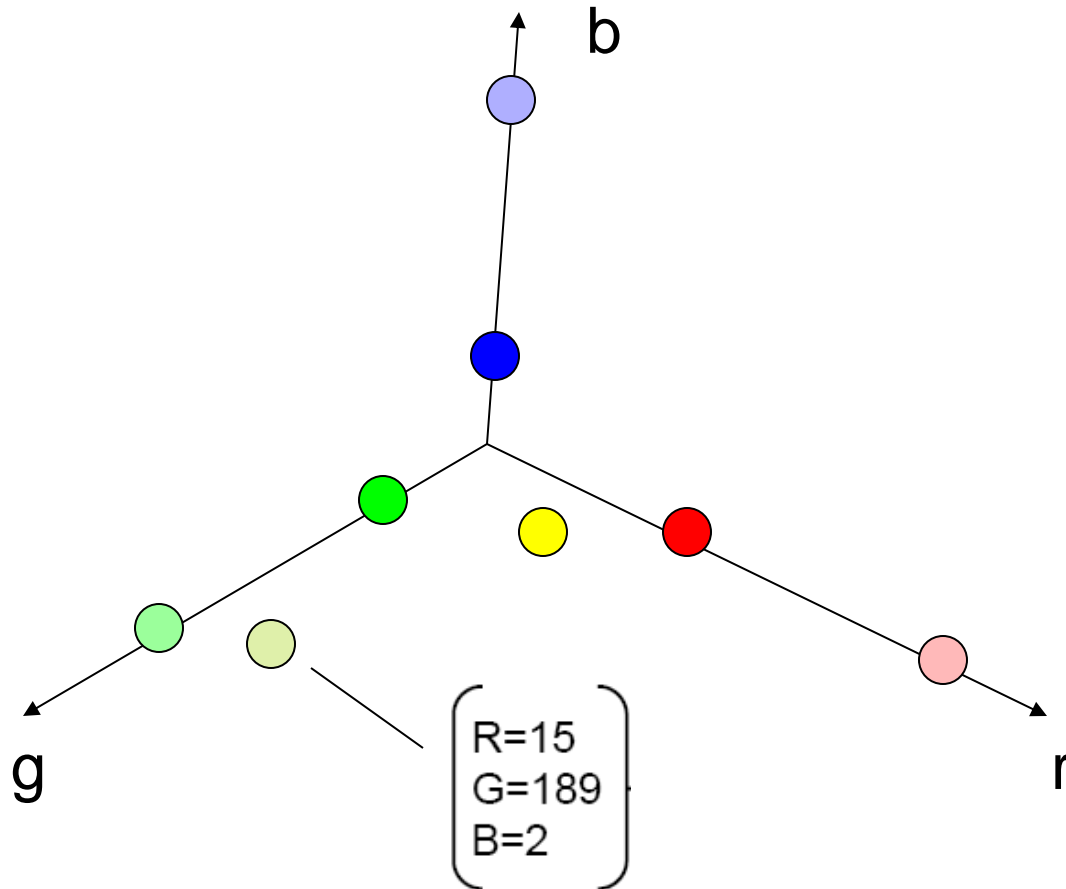
## 2. Codewords dictionary formation



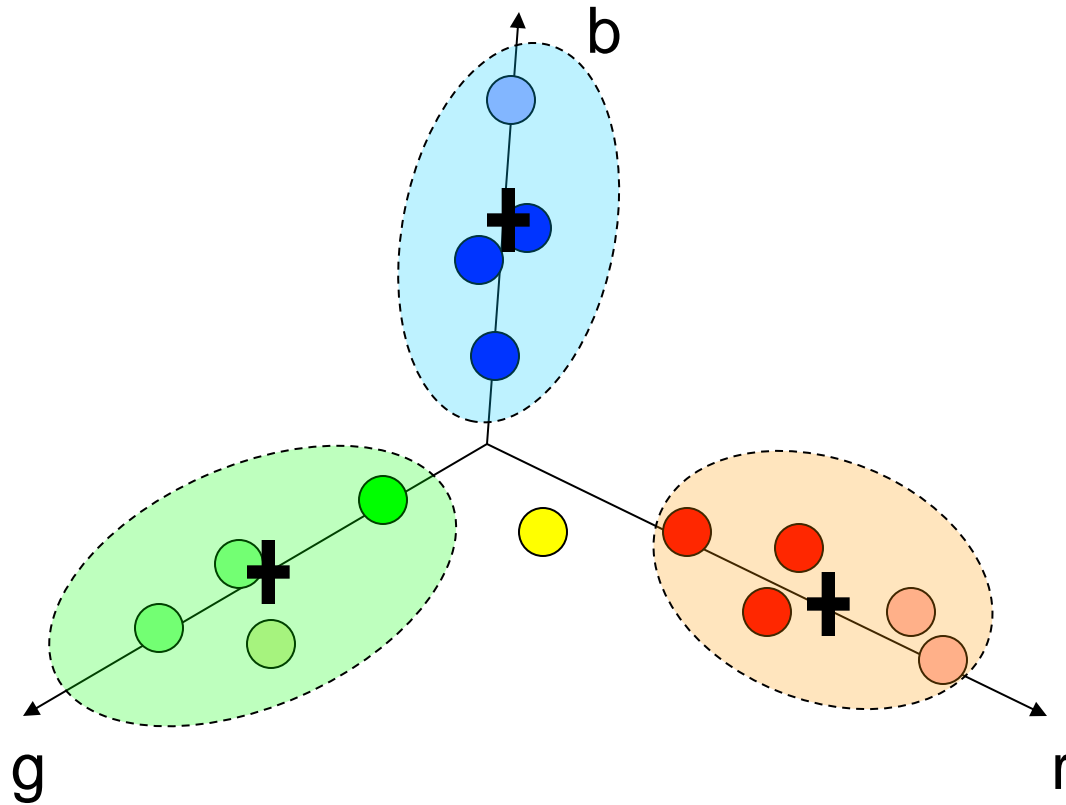
# Example: color feature



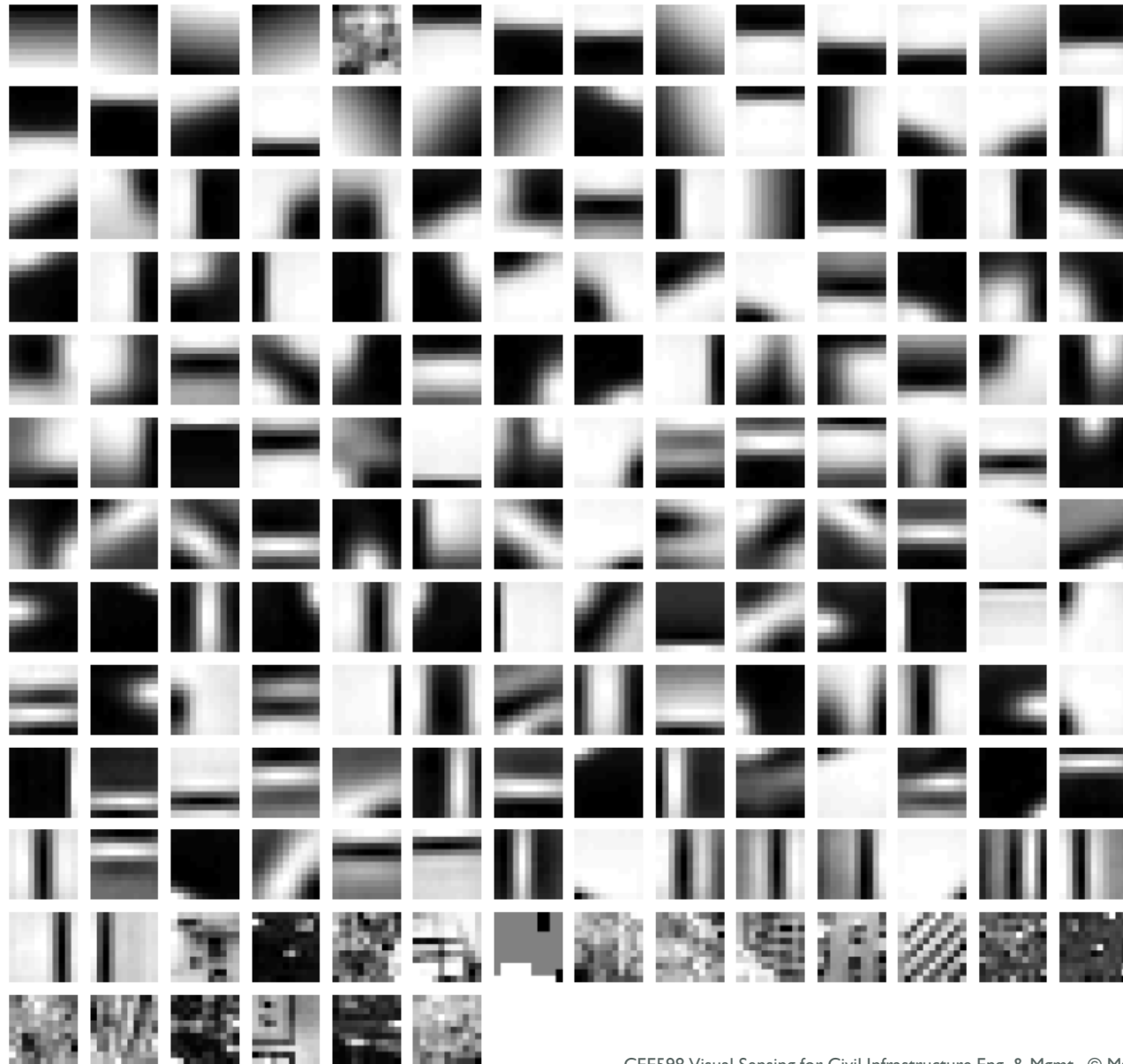
# Example: color feature



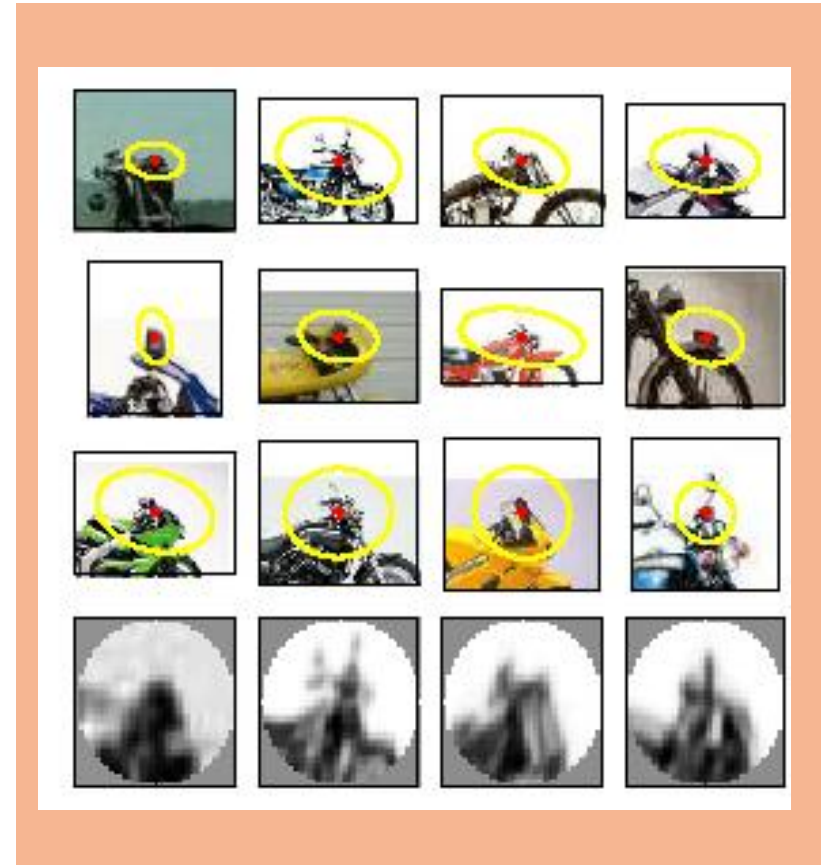
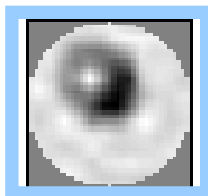
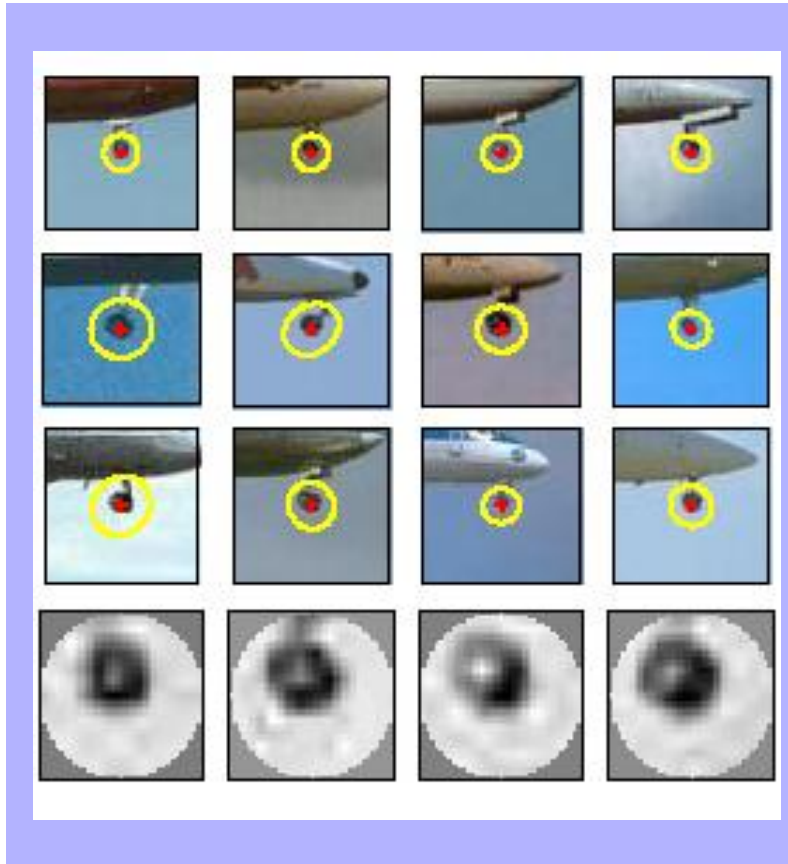
# Example: color feature



## 2. Codewords dictionary formation

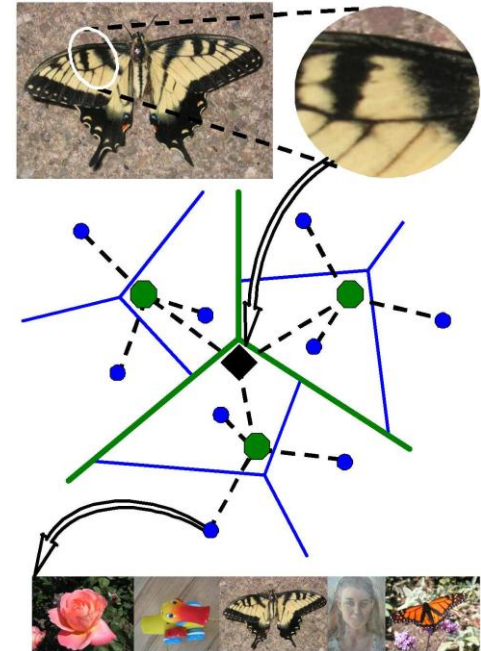


# Image patch examples of codewords

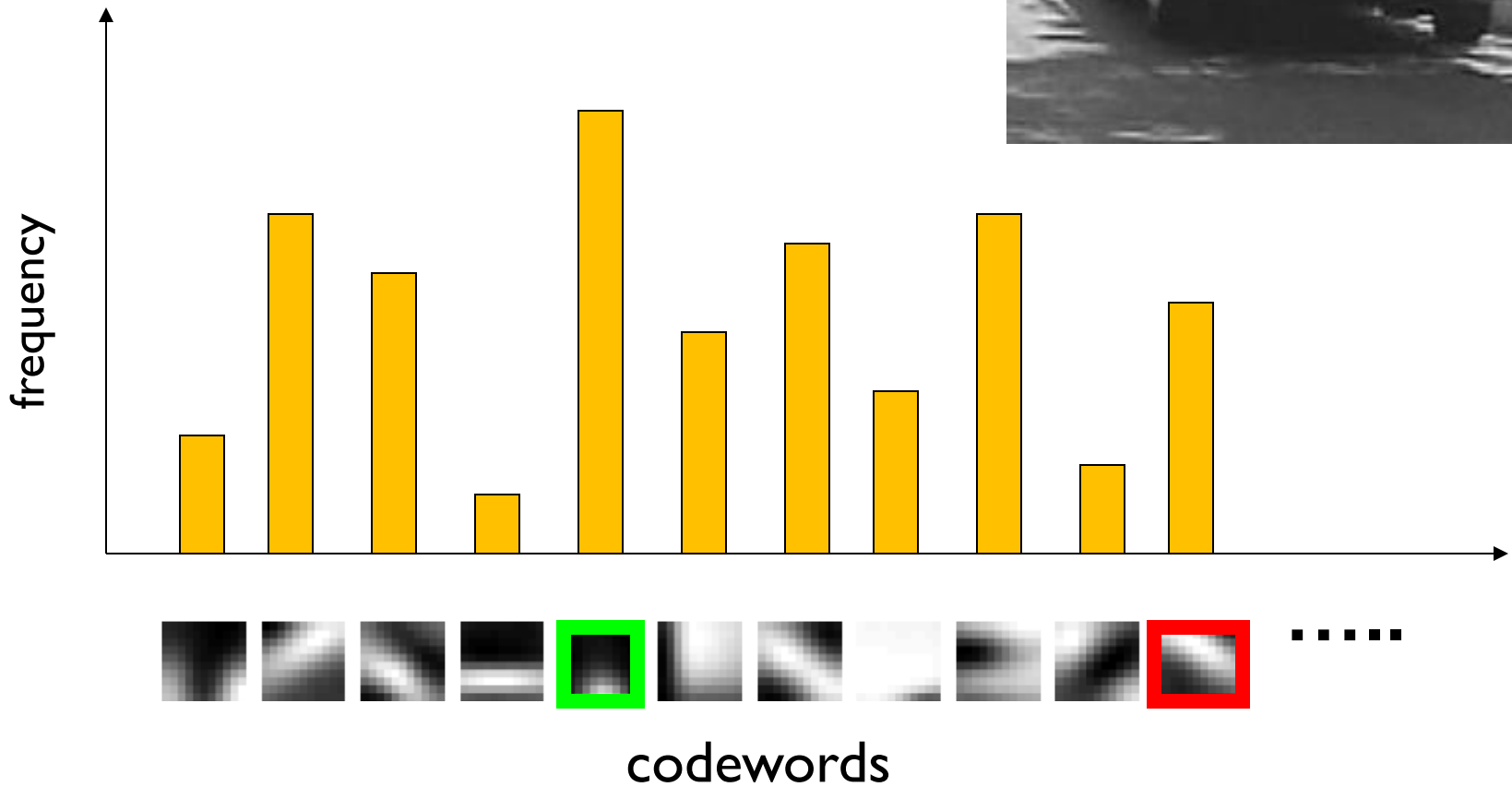
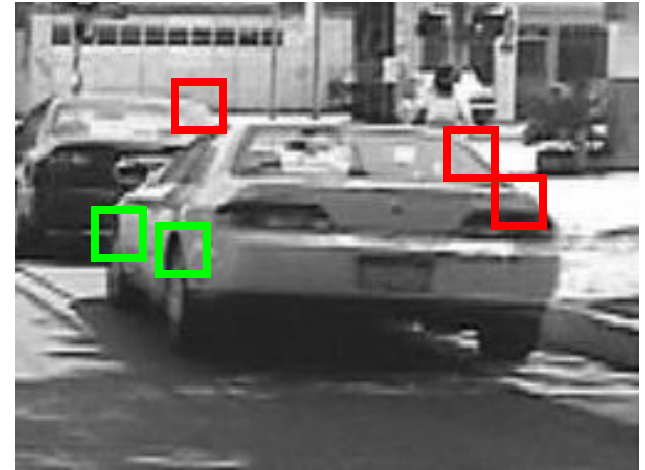


# Visual vocabularies: Issues

- How to choose vocabulary size?
  - Too small: visual words not representative of all patches
  - Too large: quantization artifacts, overfitting
- Computational efficiency
  - Vocabulary trees  
(Nister & Stewenius, 2006)

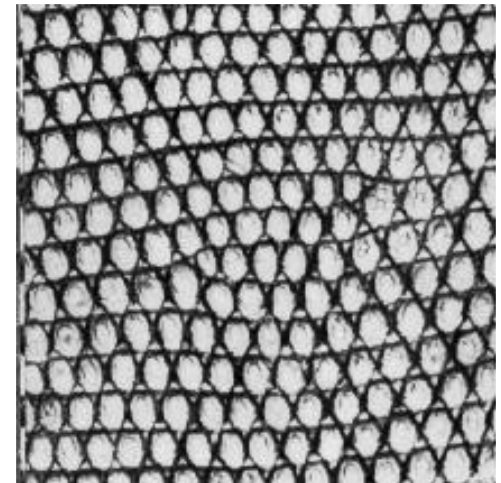
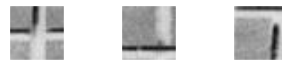
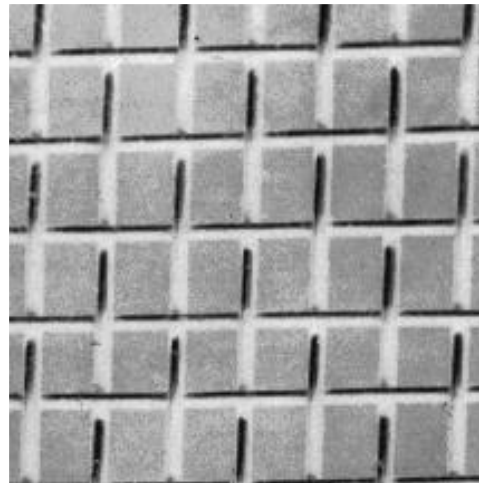
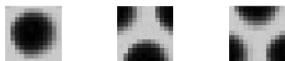
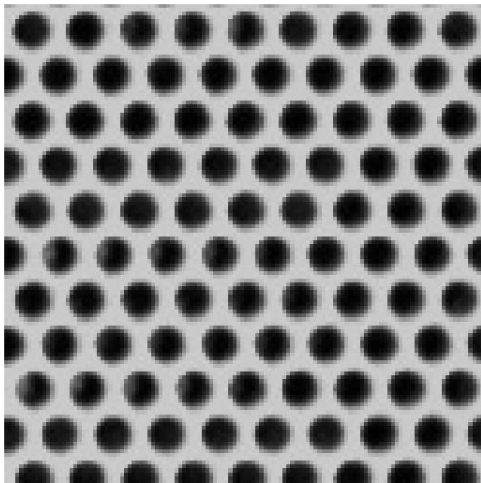


# 3. Image representation



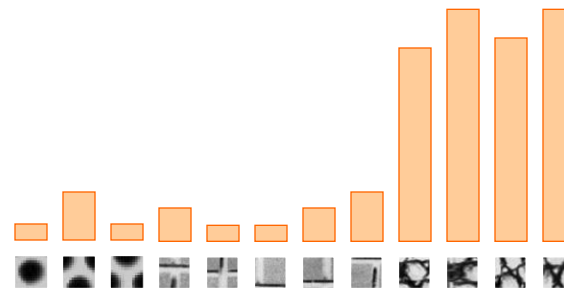
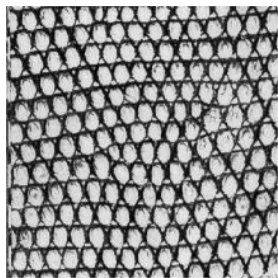
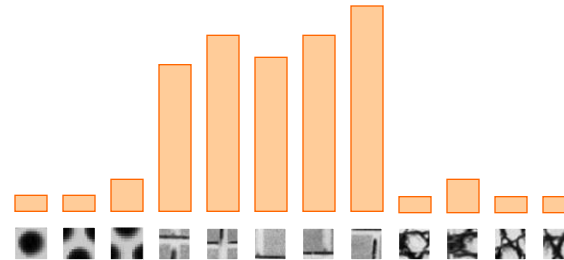
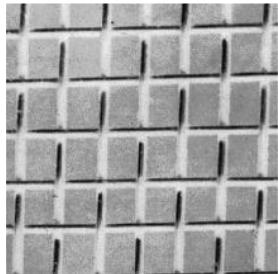
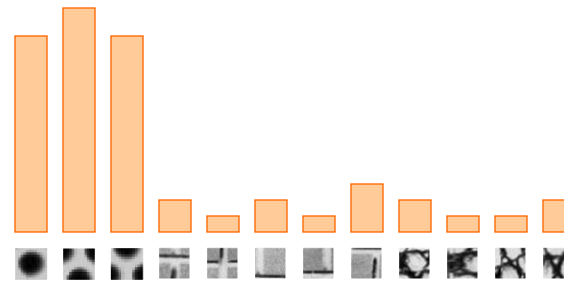
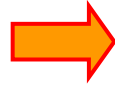
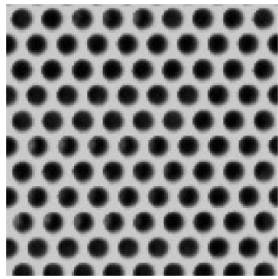
# Representing textures

- Texture is characterized by the repetition of basic elements or *textons*
- For stochastic textures, it is the identity of the textons, not their spatial arrangement, that matters



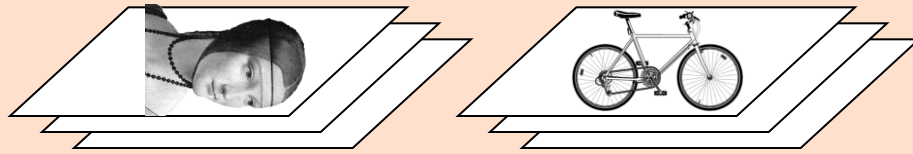
Julesz, 1981; Cula & Dana, 2001; Leung & Malik 2001; Mori, Belongie & Malik, 2001; Schmid 2001; Varma & Zisserman, 2002, 2003; Lazebnik, Schmid & Ponce, 2003

# Representing textures

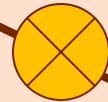


Julesz, 1981; Cula & Dana, 2001; Leung & Malik 2001; Mori, Belongie & Malik, 2001; Schmid 2001; Varma & Zisserman, 2002, 2003; Lazebnik, Schmid & Ponce, 2003

# Representation



1. feature detection & representation



2. **codewords dictionary**

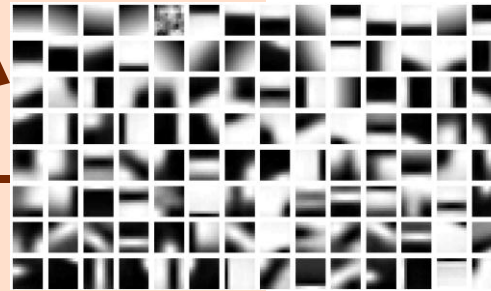
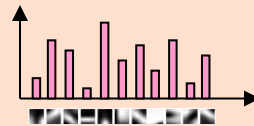
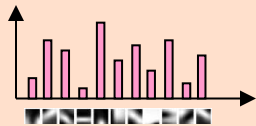
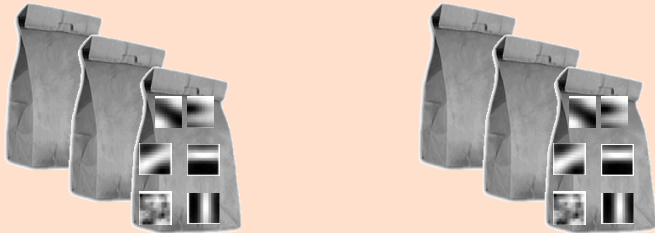


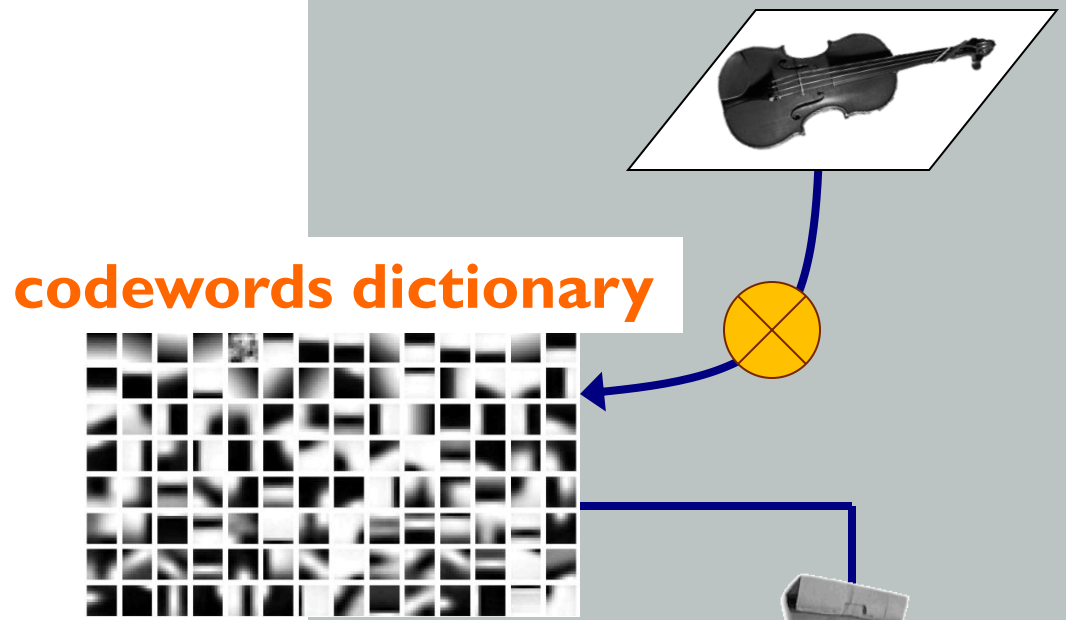
image representation

3.



**category models**

# Learning and Recognition



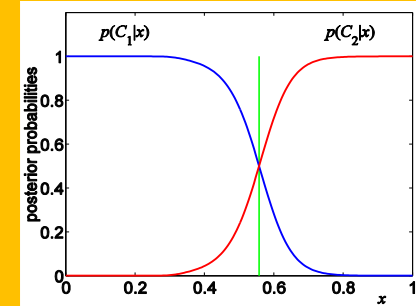
**codewords dictionary**

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

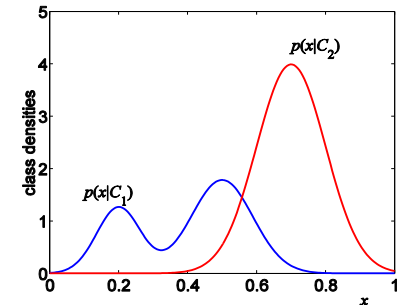
**category  
decision**

# Learning and Recognition

1. Discriminative method:
  - NN (Nearest Neighbor)
  - SVM (Support Vector Machine)



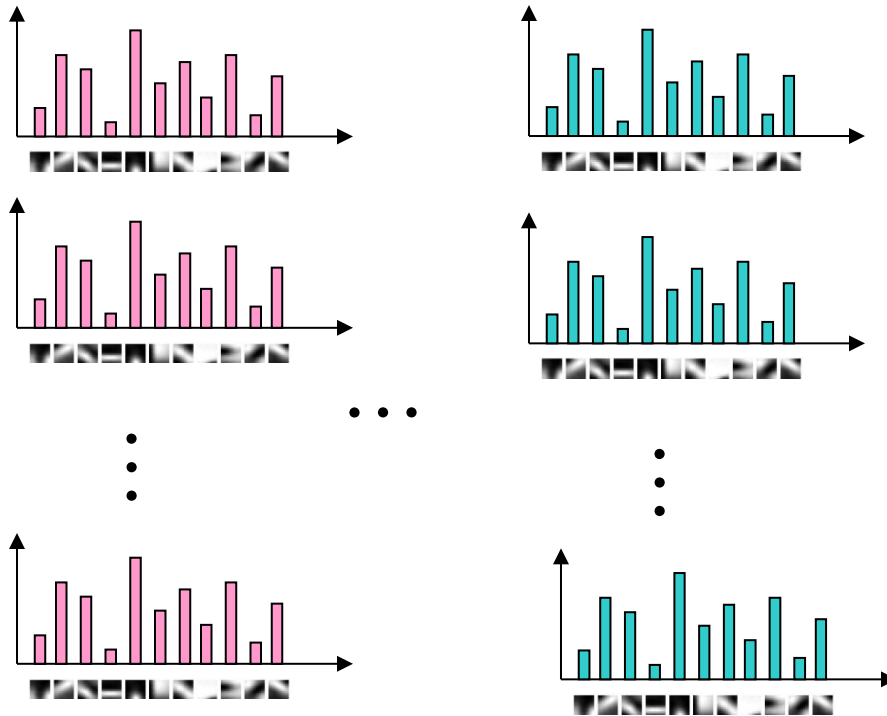
2. Generative method:
  - Graphical models



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

# Discriminative classifiers

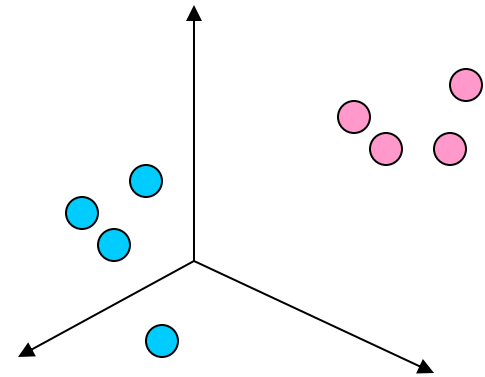
## category models



Class I

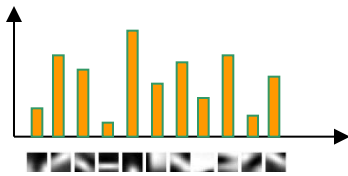
Class N

## Model space



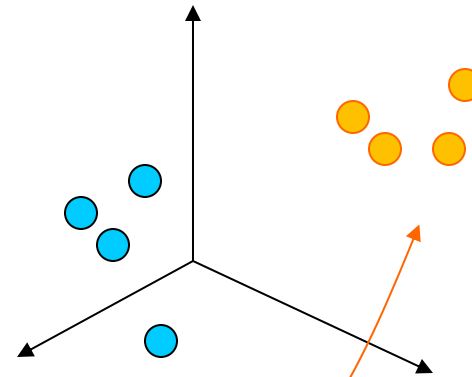
# Discriminative classifiers

Query image



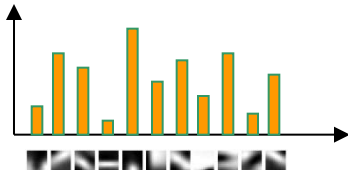
Winning class: Orange

Model space



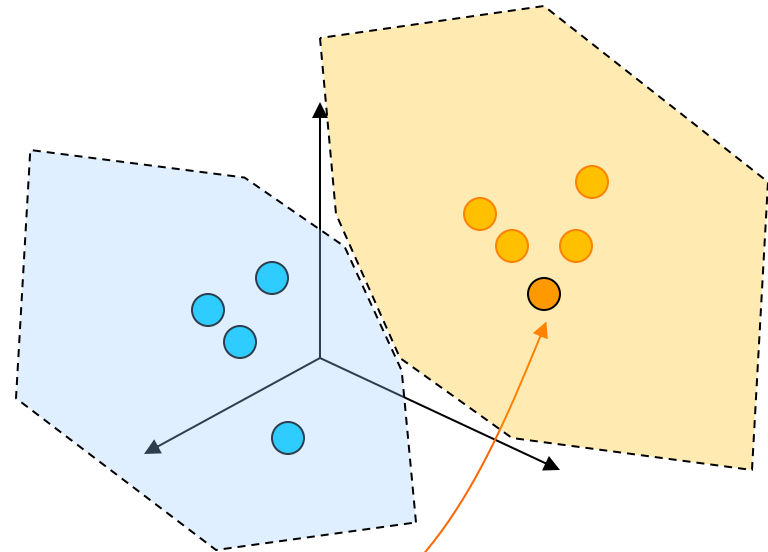
# Discriminative classifiers

Query image



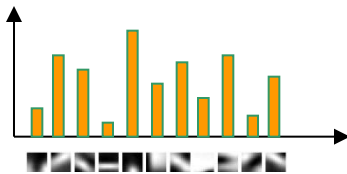
Winning class: Orange

Model space

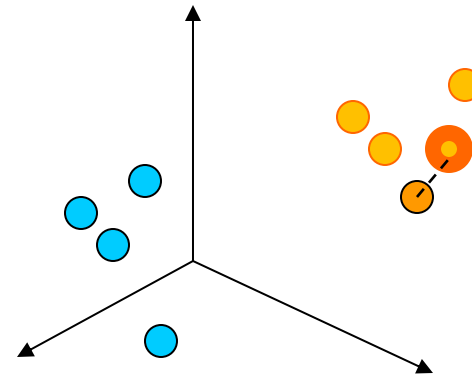


# Discriminative classifiers

Query image



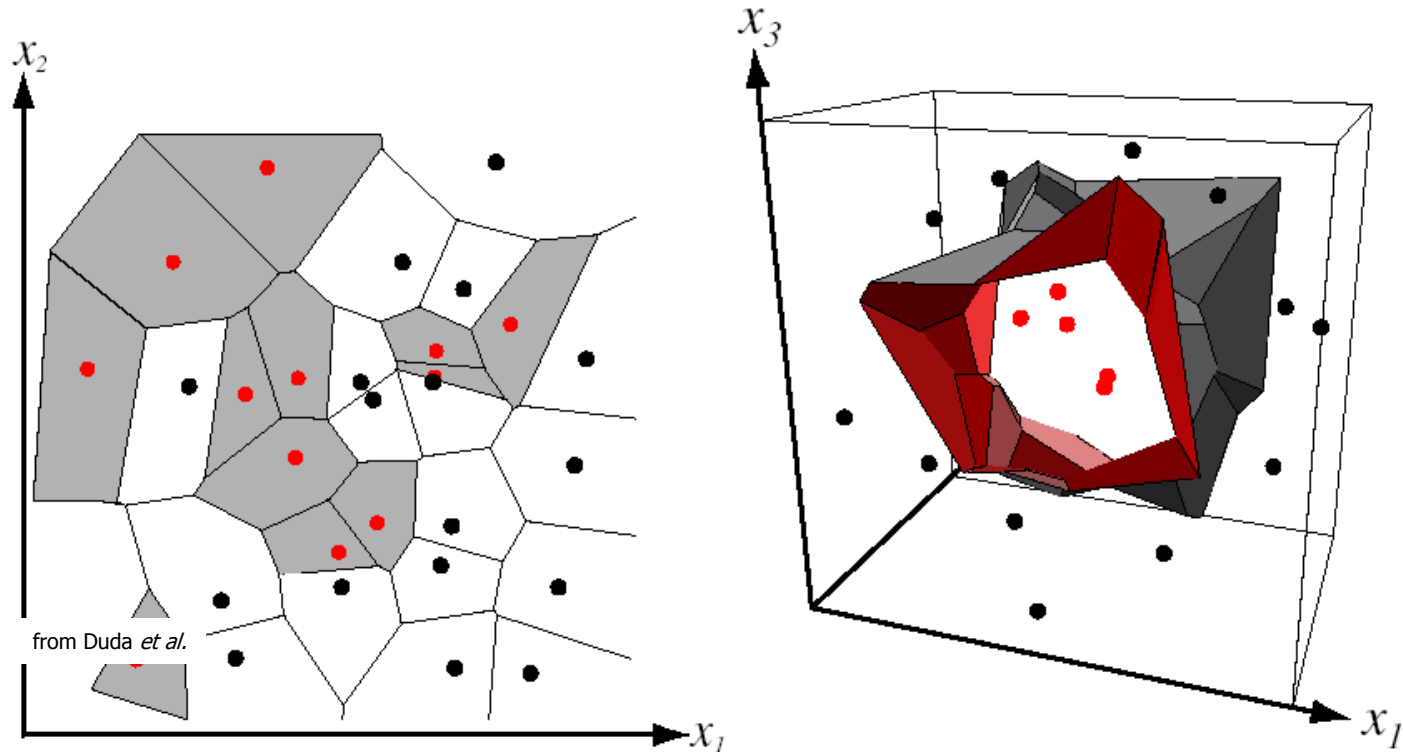
Model space



Winning class: **Orange**

- Assign label of nearest training data point to each test data point

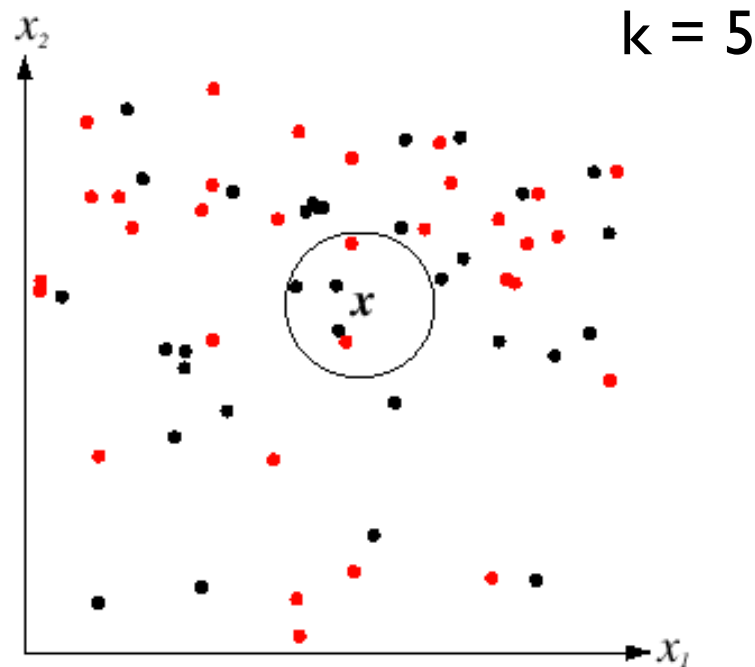
# Nearest Neighbors classifier



Voronoi partitioning of feature space  
for 2-category 2-D and 3-D data

# K-Nearest Neighbors

- For a new point, find the  $k$  closest points from training data
- Labels of the  $k$  points “vote” to classify
- Works well provided there is lots of data and the distance function is good



# Functions for comparing histograms

- L1 distance

$$D(h_1, h_2) = \sum_{i=1}^N |h_1(i) - h_2(i)|$$

- $\chi^2$  distance

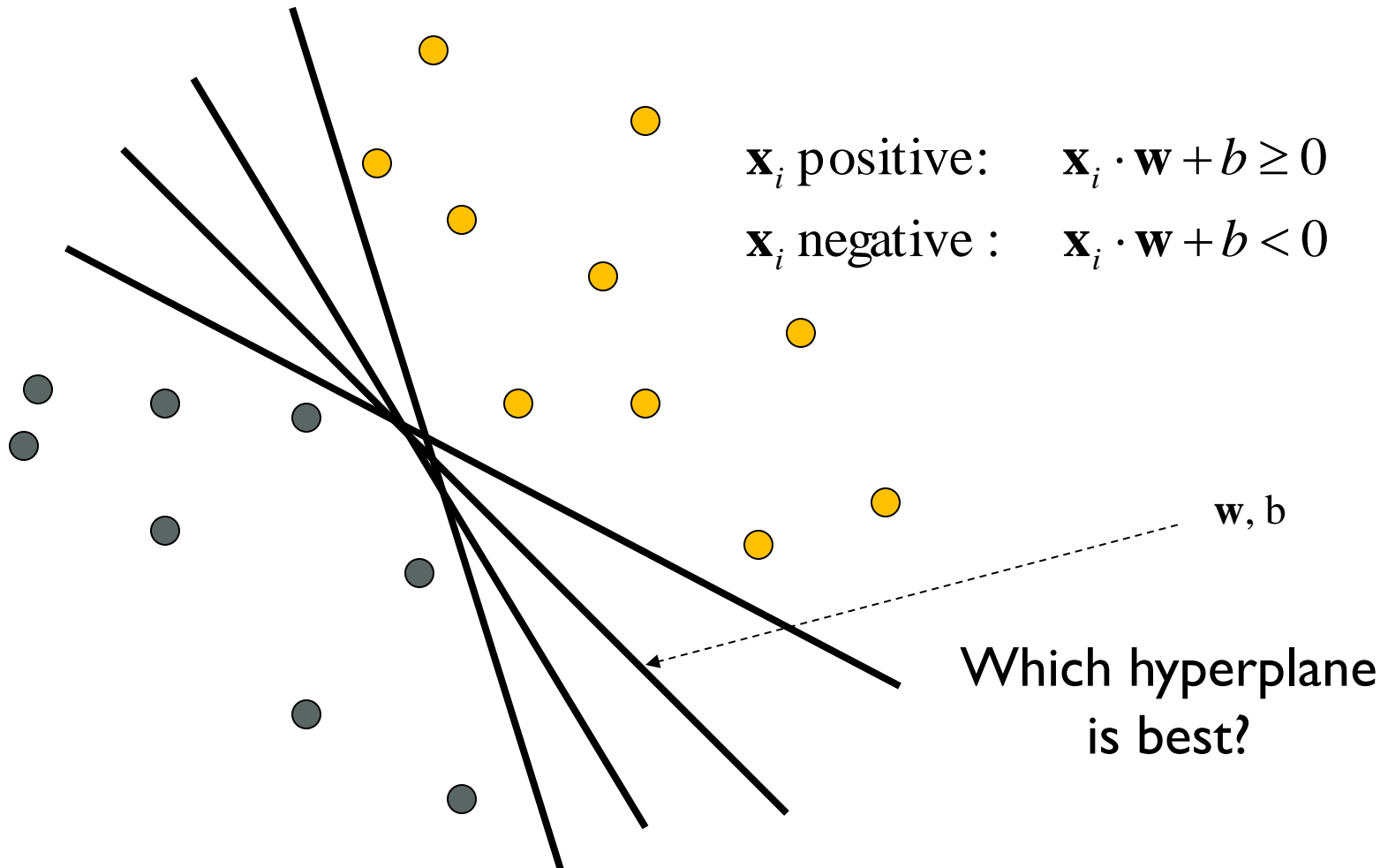
$$D(h_1, h_2) = \sum_{i=1}^N \frac{(h_1(i) - h_2(i))^2}{h_1(i) + h_2(i)}$$

- Quadratic distance (*cross-bin*)

$$D(h_1, h_2) = \sum_{i,j} A_{ij} (h_1(i) - h_2(j))^2$$

# Linear classifiers

- Find linear function (*hyperplane*) to separate positive and negative examples



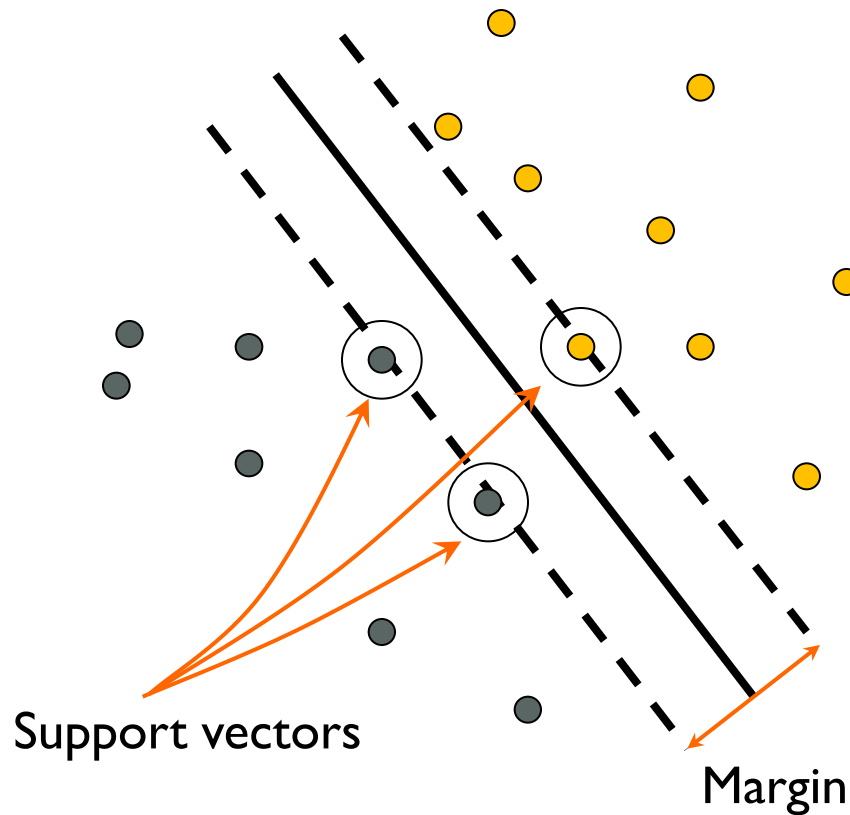
# Support vector machines

C. Burges

A Tutorial on Support Vector Machines for  
Pattern Recognition,  
Data Mining and Knowledge Discovery, 1998

# Support vector machines

- Find **hyperplane** that maximizes the *margin* between the positive and negative examples



$$\mathbf{x}_i \text{ positive } (y_i = 1): \quad \mathbf{x}_i \cdot \mathbf{w} + b \geq 1$$

$$\mathbf{x}_i \text{ negative } (y_i = -1): \quad \mathbf{x}_i \cdot \mathbf{w} + b \leq -1$$

For support, vectors,  $\mathbf{x}_i \cdot \mathbf{w} + b = \pm 1$

Distance between point and hyperplane:  $\frac{|\mathbf{x}_i \cdot \mathbf{w} + b|}{\|\mathbf{w}\|}$

Therefore, the margin is  $2 / \|\mathbf{w}\|$

# Finding the maximum margin hyperplane

1. Maximize margin  $2/\|\mathbf{w}\|$
2. Correctly classify all training data:

$$\mathbf{x}_i \text{ positive } (y_i = 1): \quad \mathbf{x}_i \cdot \mathbf{w} + b \geq 1$$

$$\mathbf{x}_i \text{ negative } (y_i = -1): \quad \mathbf{x}_i \cdot \mathbf{w} + b \leq -1$$

- *Quadratic optimization problem:*

$$\text{Minimize } \min_{\mathbf{w}, b} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=0}^n \alpha_i [c_i (\mathbf{w} \cdot \mathbf{x}_i - b) - 1] \right\}$$

$$\text{Subject to } y_i (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$$

# Finding the maximum margin hyperplane

- Solution:

$$\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$$

learned weight      Support vector

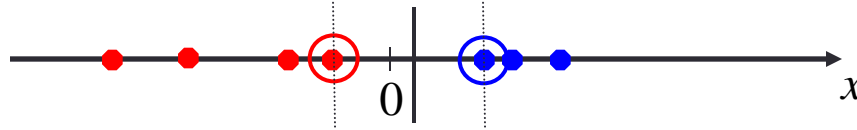
- Classification function (decision boundary):

$$\mathbf{w} \cdot \mathbf{x} + b = \sum_i \alpha_i y_i \mathbf{x}_i \cdot \mathbf{x} + b$$

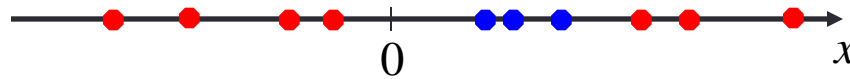
Test point

# Nonlinear SVMs

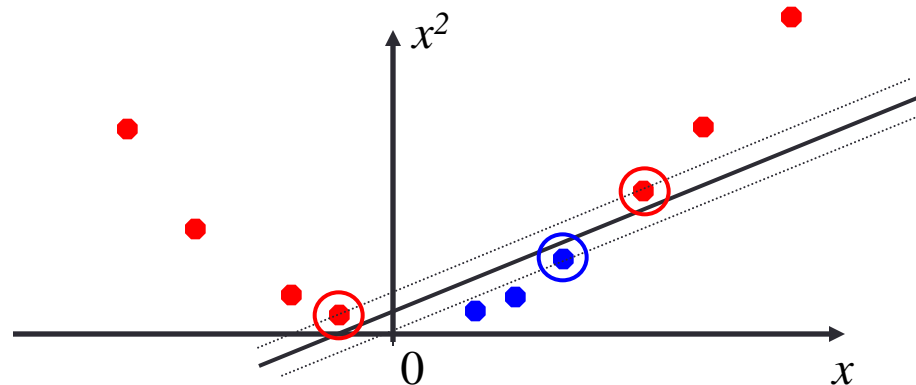
- Datasets that are linearly separable work out great:



- But what if the dataset is just too hard?

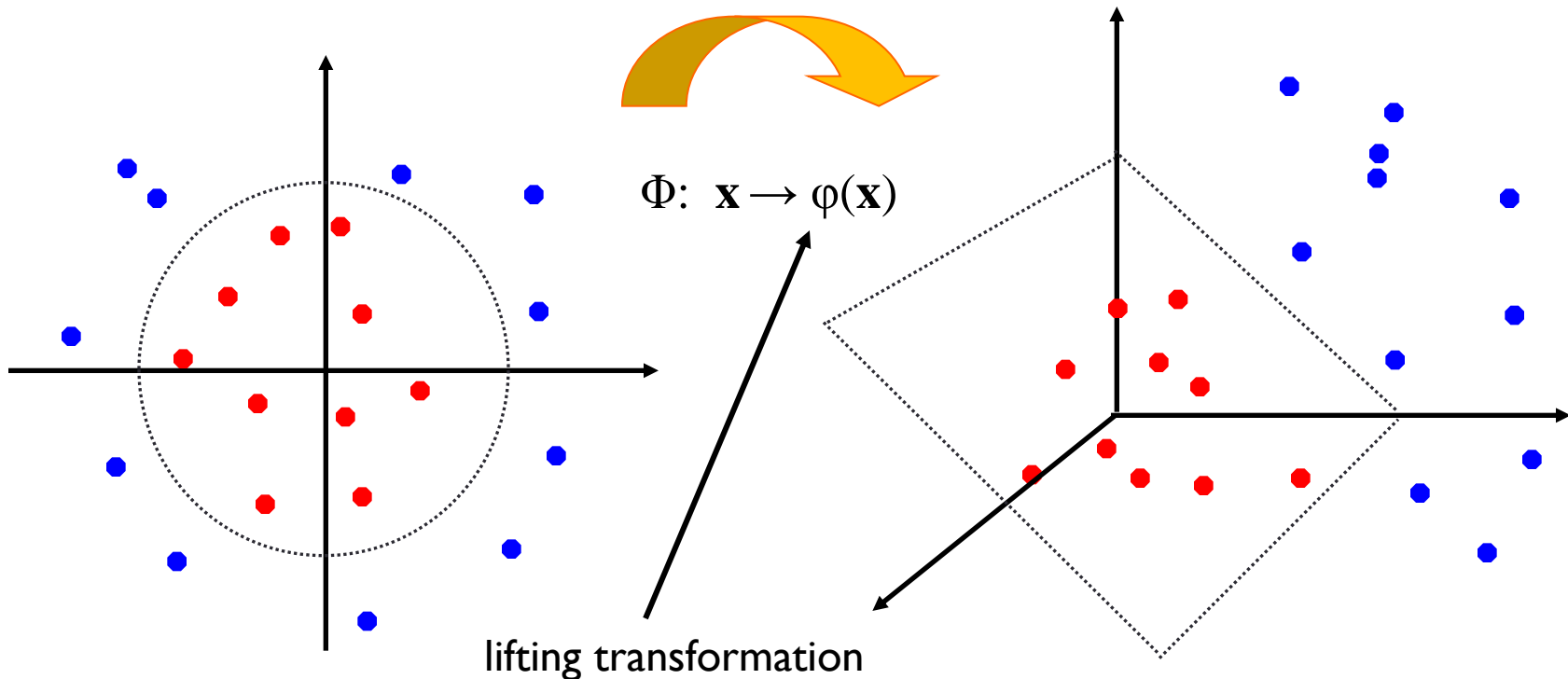


- We can map it to a higher-dimensional space:



# Nonlinear SVMs

- General idea: the original input space can always be mapped to some higher-dimensional feature space where the training set is separable:



# Nonlinear SVMs

- Nonlinear decision boundary in the original feature space:

$$\sum_i \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + b$$

- *The kernel  $K$*  = product of the lifting transformation  $\varphi(\mathbf{x})$ :

$$K(\mathbf{x}_i, \mathbf{x}_j) = \varphi(\mathbf{x}_i) \cdot \varphi(\mathbf{x}_j)$$

NOTE:

- It is not required to compute  $\varphi(\mathbf{x})$  explicitly:

C. Burges, [A Tutorial on Support Vector Machines for Pattern Recognition](#),  
Data Mining and Knowledge Discovery, 1998

# Kernels for bags of features

- Histogram intersection kernel:

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

- Generalized Gaussian kernel:

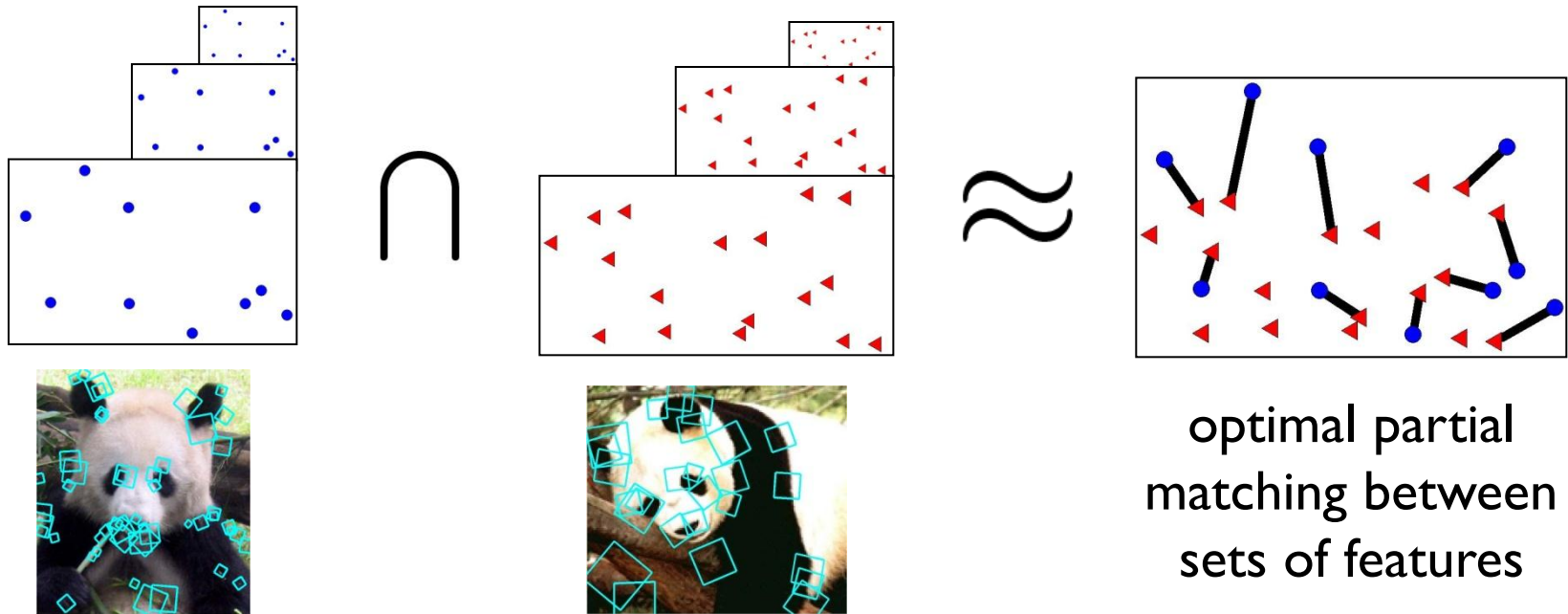
$$K(h_1, h_2) = \exp\left(-\frac{1}{A} D(h_1, h_2)^2\right)$$

- $D$  can be Euclidean distance,  $\chi^2$  distance etc...

J. Zhang, M. Marszalek, S. Lazebnik, and C. Schmid,  
Local Features and Kernels for Classification of Texture and Object Categories: A  
Comprehensive Study,  
IJCV 2007

# Pyramid match kernel

- Fast approximation of Earth Mover's Distance
- Weighted sum of histogram intersections at multiple resolutions (linear in the number of features instead of cubic)



K. Grauman and T. Darrell. The Pyramid Match Kernel: Discriminative Classification with Sets of Image Features, ICCV 2005.

# Summary: SVMs for image classification

1. Pick an image representation (in our case, bag of features)
2. Compute the matrix of kernel values between every pair of training examples
3. Feed the kernel matrix into your favorite SVM solver to obtain support vectors and weights
4. At test time: compute kernel values for your test example and each support vector, and combine them with the learned weights to get the value of the decision function

# What about multi-class SVMs?

- No “definitive” multi-class SVM formulation
- In practice, we have to obtain a multi-class SVM by combining multiple two-class SVMs
  
- One vs. others
  - Training: learn an SVM for each class vs. the others
  - Testing: apply each SVM to test example and assign to it the class of the SVM that returns the highest decision value
  
- One vs. one
  - Training: learn an SVM for each pair of classes
  - Testing: each learned SVM “votes” for a class to assign to the test example

# Object recognition results

- ETH-80 database

(Eichhorn and Chapelle 2004)

- Features:

- Harris detector
- PCA-SIFT descriptor,  $d=10$

## 8 object classes



Kernel	Complexity	Recognition rate
Match [Wallraven et al.]	$O(dm^2)$	84%
Bhattacharyya affinity [Kondor & Jebara]	$O(dm^3)$	85%
Pyramid match	$O(dmL)$	84%

# SVMs: Pros and cons

## ■ Pros

- Many publicly available SVM packages:  
<http://www.kernel-machines.org/software>
- Kernel-based framework is very powerful, flexible
- SVMs work very well in practice, even with very small training sample sizes

## ■ Cons

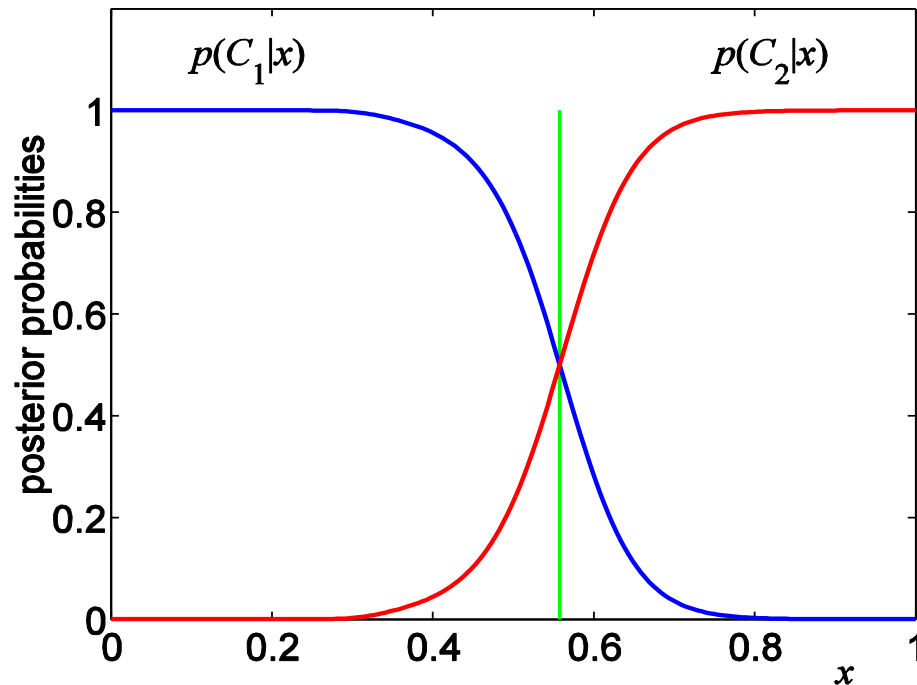
- No “direct” multi-class SVM, must combine two-class SVMs
- Computation, memory
  - During training time, must compute matrix of kernel values for every pair of examples
  - Learning can take a very long time for large-scale problems

# Of course, there are many other classifiers out there

- Neural networks
- Boosting
- Decision trees, ...

# Summary

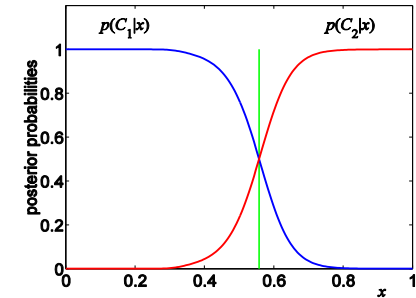
$$\frac{p(\text{class} \mid \text{image})}{p(\text{no class} \mid \text{image})} = \frac{p(\text{image} \mid \text{class})}{p(\text{image} \mid \text{no class})} \cdot \frac{p(\text{class})}{p(\text{no class})}$$



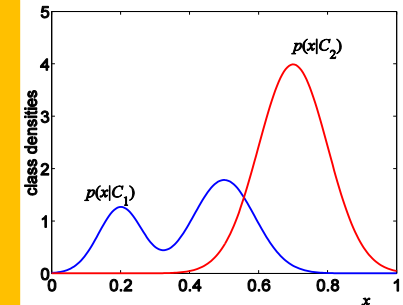
- Discriminative methods model posterior

# Learning and Recognition

1. Discriminative method:
  - NN (Nearest Neighbor)
  - SVM (Support Vector Machine)



2. Generative method:
  - Graphical models



→ Model the probability distribution that produces a given bag of features

# 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

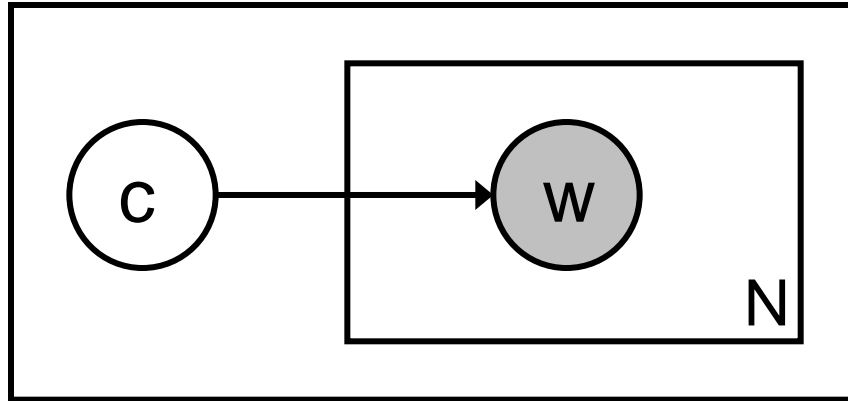
# Some Notation

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

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

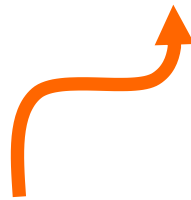
- **C**: category of the image

# The Naïve Bayes model



Graphical model

$$p(c | w) \sim p(c)p(w | c)$$



Prior prob. of  
the object classes



Image likelihood  
given the class

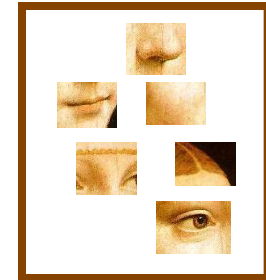
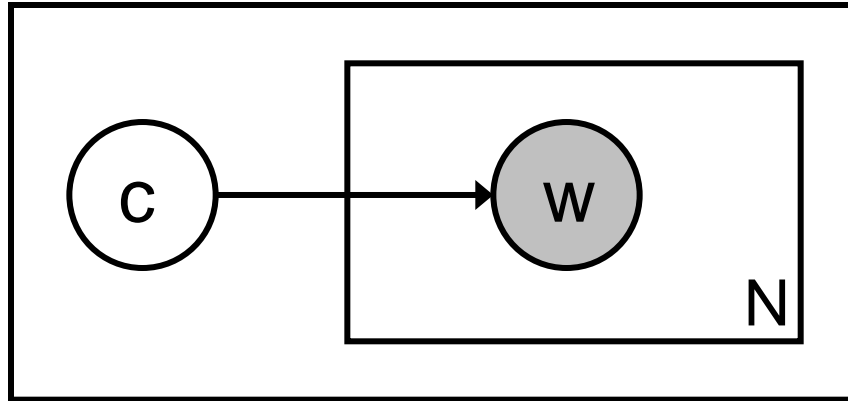
# The Naïve Bayes model

$$p(w_1, \dots, w_N | c)$$

- Assume that each feature 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



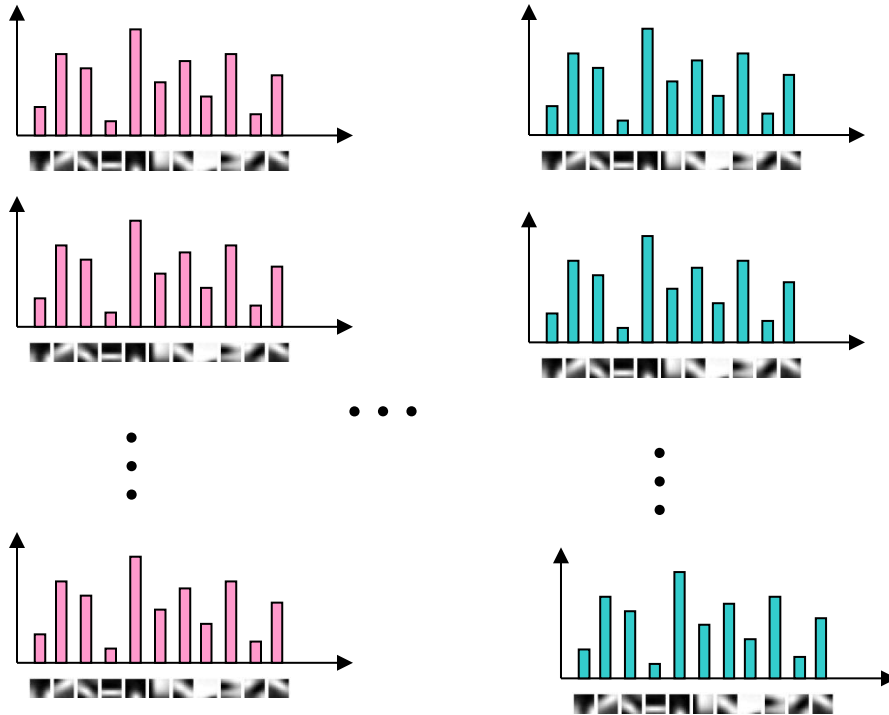
$$c^* = \arg \max_c p(c | w) \propto p(c) p(w | c) = p(c) \prod_{n=1}^N p(w_n | c)$$

Object class  
decision

Likelihood of *n*th visual word  
given the class

Estimated by empirical frequencies of code  
words in images from a given class

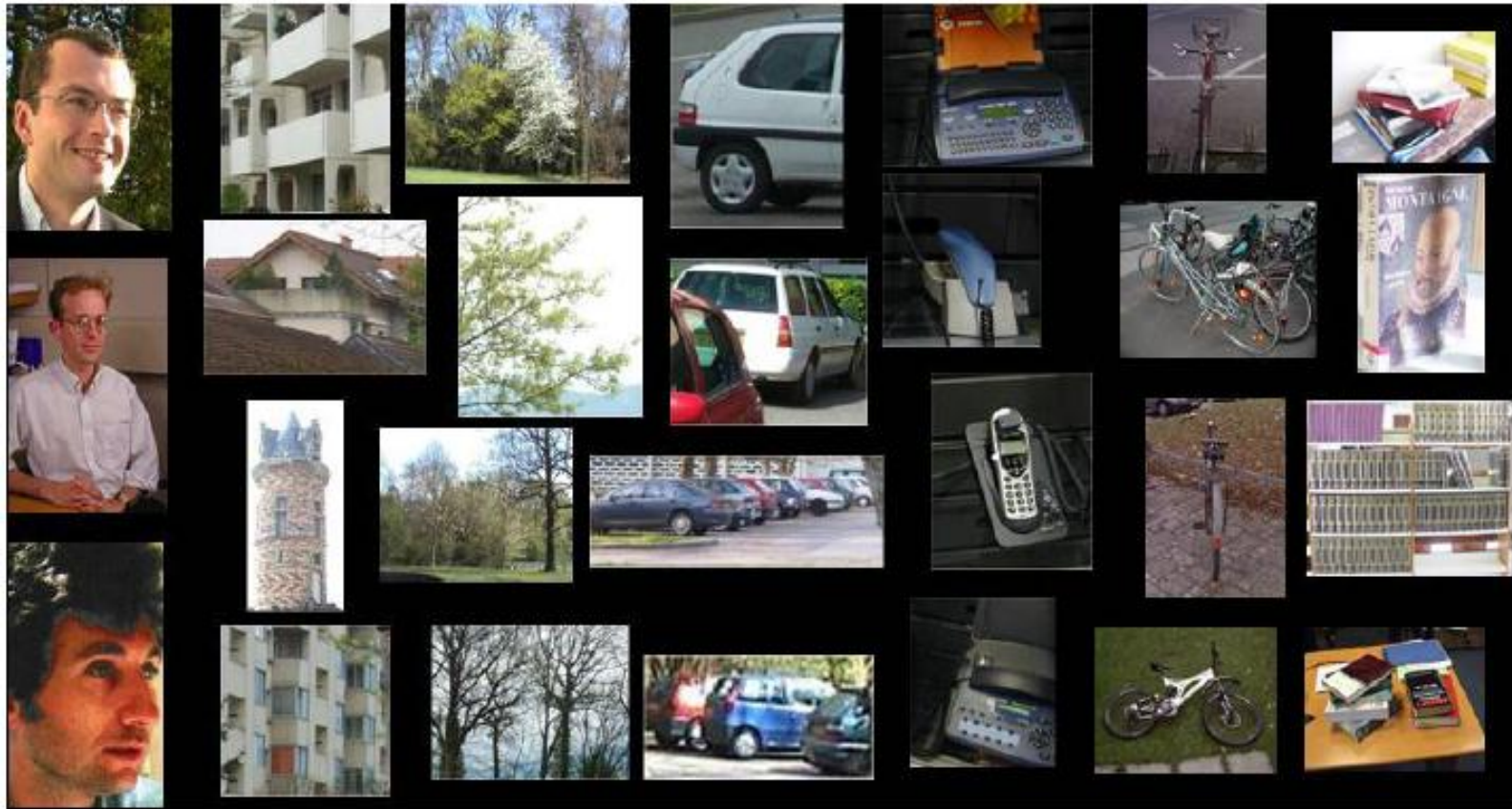
# Category Models



Class I

Class N

Our in-house database contains 1776 images in seven classes<sup>1</sup>: faces, buildings, trees, cars, phones, bikes and books. Fig. 2 shows some examples from this dataset.



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

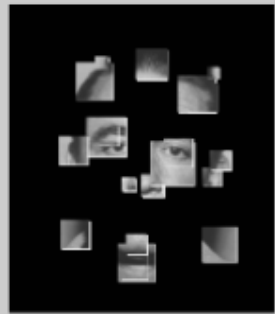
# Summary: Generative models

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

# Invariance issues



- Scale and rotation
  - Implicit
  - Detectors and descriptors



# Invariance issues



- Scale and rotation
- Occlusion
  - Implicit in the models
  - Codeword distribution: small variations
  - (In theory) Theme ( $z$ ) distribution: different occlusion patterns

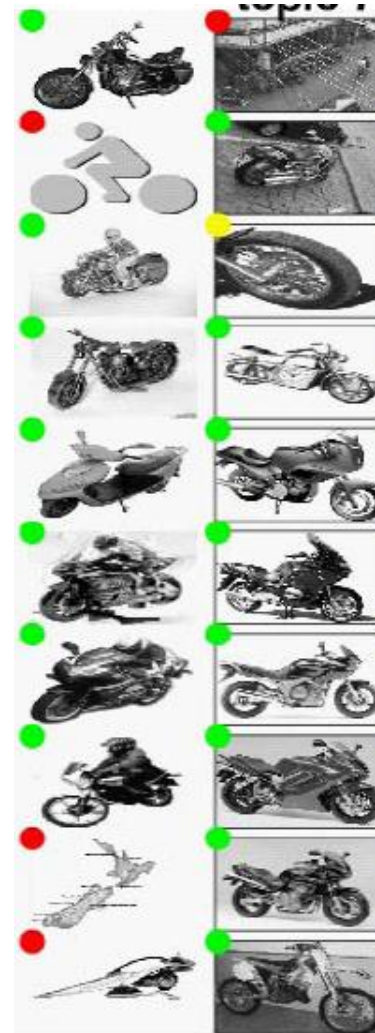
# Invariance issues



- Scale and rotation
- Occlusion
- Translation
  - Encode (relative) location information
    - Sudderth, Torralba, Freeman & Willsky, 2005, 2006
    - Niebles & Fei-Fei, 2007

# Invariance issues

- Scale and rotation
- Occlusion
- Translation
- View point (in theory)
  - Codewords: detector and descriptor
  - Theme distributions: different view points





# Model properties

- Intuitive
  - Analogy to documents

Of all the sensory impressions proceeding to the brain, the visual experiences are the dominant ones. Our perception of the world around us is based essentially on the messages that reach our eyes. For a long time, the retinal image was considered as a visual centers in the brain. In a movie image, the retinal image is discovered. We know that perception is more complex following the path to the various cortical areas. Hubel and Wiesel have demonstrated that the message about the image falling on the retina undergoes a fine-grained analysis in a system of nerve cells stored in columns. In this system each cell has its specific function and is responsible for a specific detail in the pattern of the retinal image.

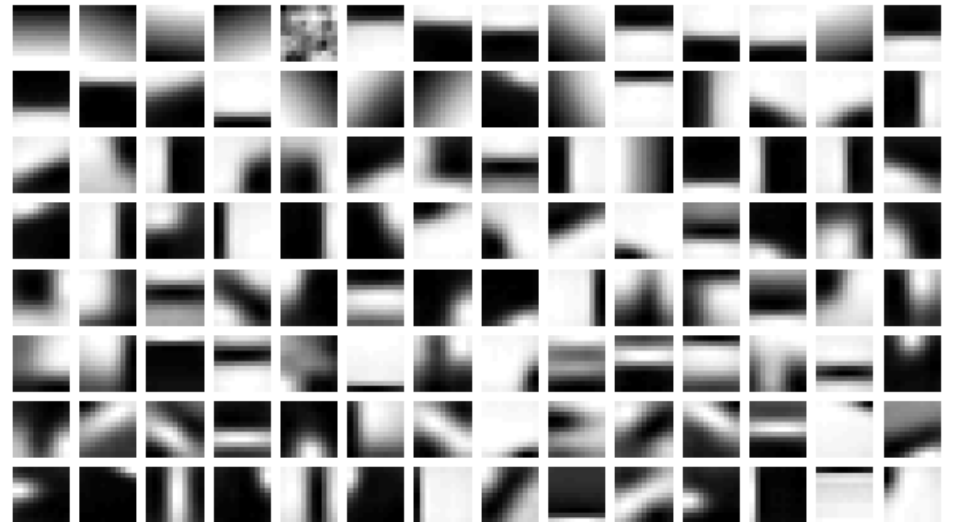
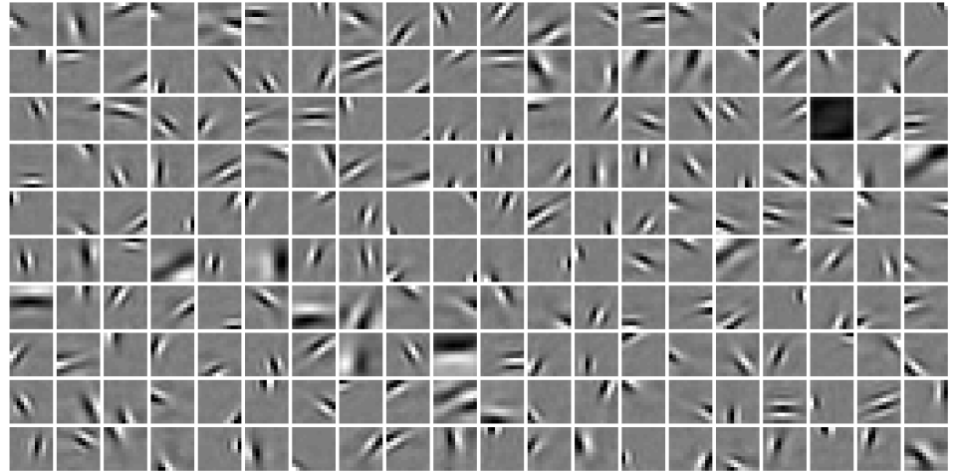
**sensory, brain,  
visual, perception,  
retinal, cerebral cortex,  
eye, cell, optical  
nerve, image  
Hubel, Wiesel**

# Model properties



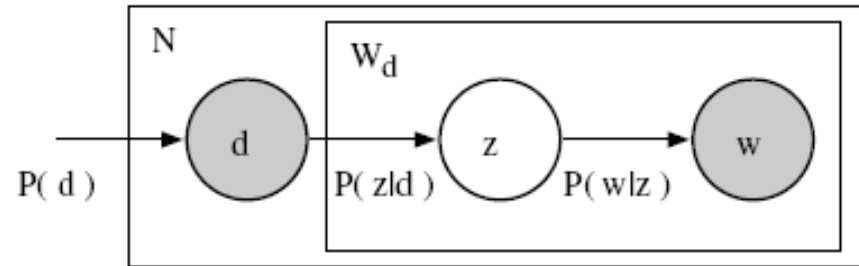
- Intuitive

- Analogy to documents
- Analogy to human vision



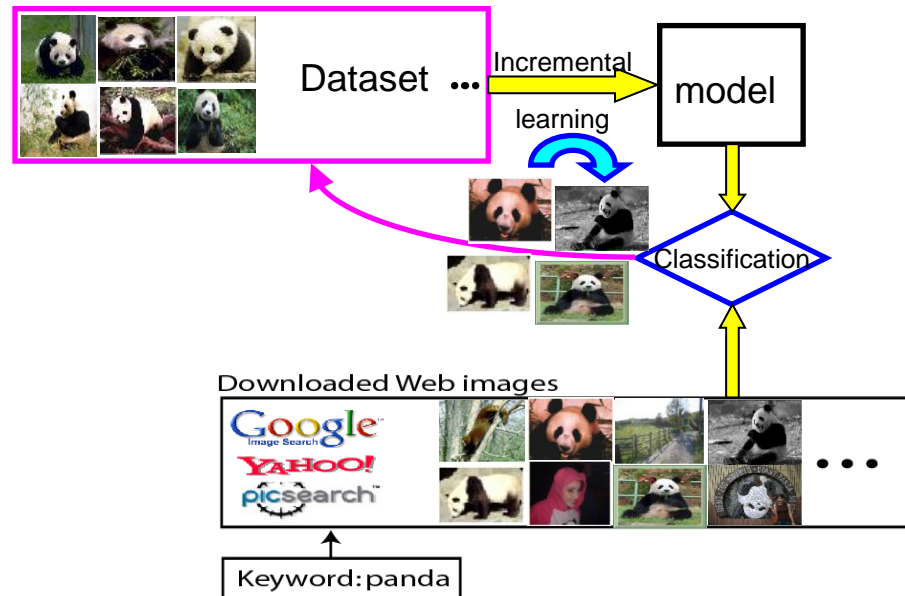


# Model properties



Sivic, Russell, Efros, Freeman, Zisserman, 2005

- Intuitive
- generative models
  - Convenient for weakly- or un-supervised, incremental training
  - Prior information
  - Flexibility (e.g. HDP)

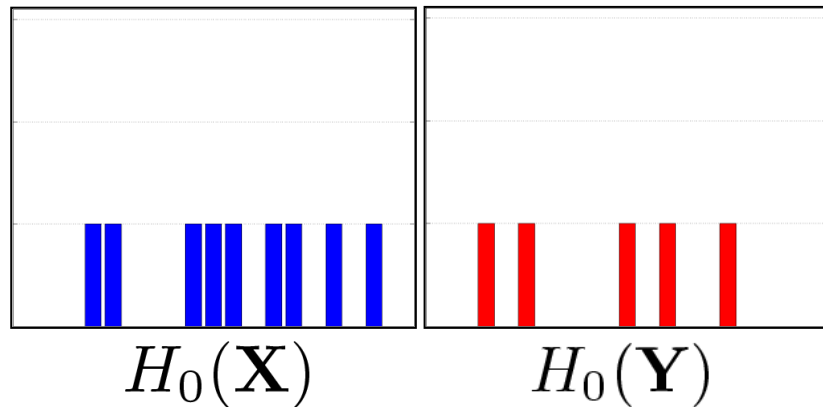
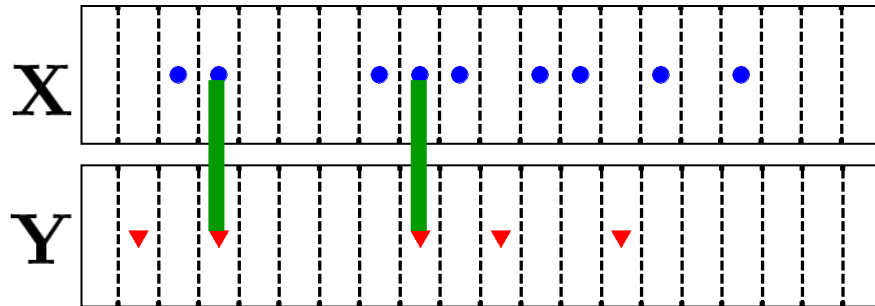


Li, Wang & Fei-Fei, CVPR 2007



# Model properties

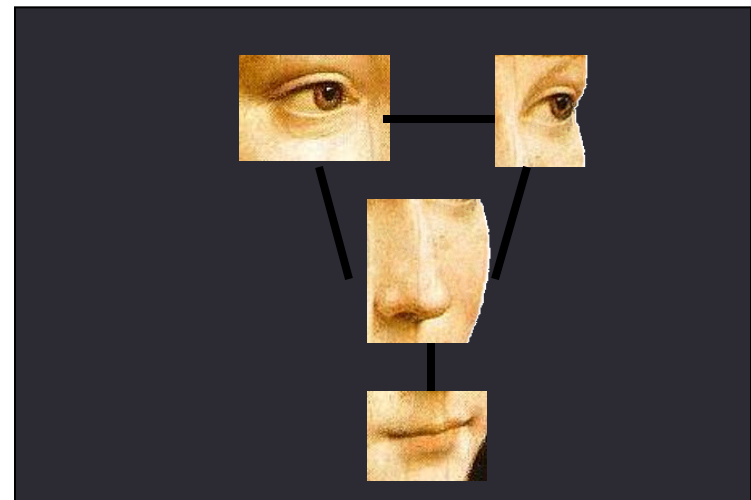
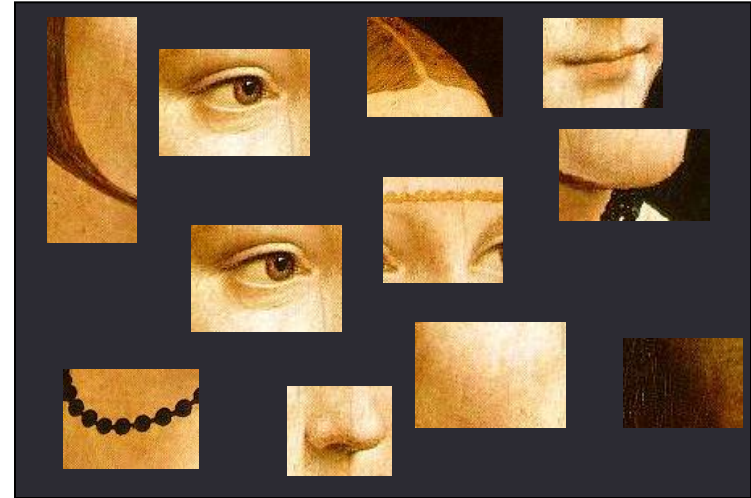
- Intuitive
- generative models
- **Discriminative method**
  - Computationally efficient
  - fast



# Model properties



- Intuitive
- generative models
- Discriminative method
- Learning and recognition relatively fast
  - Compare to other methods



# Weakness of the model



- No rigorous geometric information of the object components
- It's intuitive to most of us that objects are made of parts – no such information
- Not extensively tested yet for
  - View point invariance
  - Scale invariance
- Segmentation and localization unclear

# Outline

- **Object Recognition**
  - Introduction
  - Recognition of single 3D objects
    - Bag of world models
    - Part based models
    - Models for 3D objects categorization