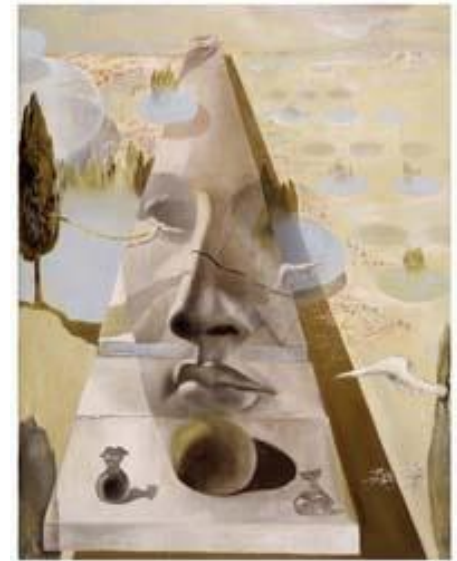


Lecture 12

Visual recognition



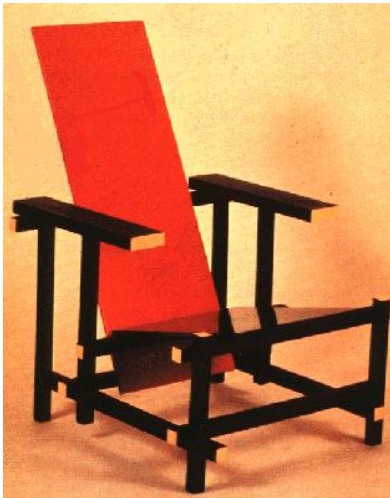
- Bag of words models for object recognition and classification
 - Discriminative methods
 - Generative methods

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



Part 1: Bag-of-words models

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. It is now discovered that the image is analyzed in a more complex manner following the path to the various centers of the cortex, Hubel and Wiesel have demonstrated that the *message about the image falling on the retina undergoes a point-by-point 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**

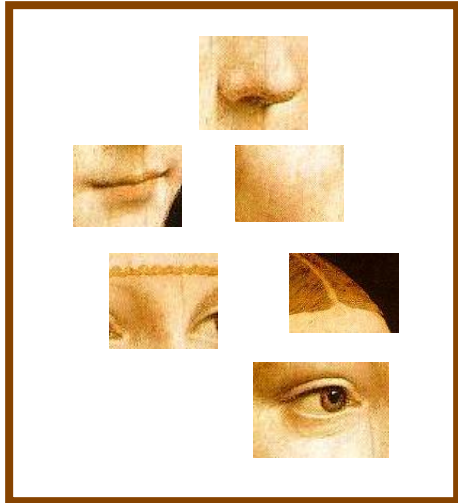
China is forecasting a trade surplus of \$90bn (£51bn) to \$100bn this year, a threefold increase on 2004's \$32bn. The Commerce Ministry said the surplus would be created by a predicted 30% increase in exports to \$750bn, compared with \$575bn in 2004. The surplus of \$660bn. The surplus will annoy the US. China's government has deliberately agreed to keep the yuan is pegged to the US dollar. The government also needs to keep the demand so high for the country. China has kept the yuan against the dollar and permitted it to trade within a narrow band but the US wants the yuan to be allowed to trade 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

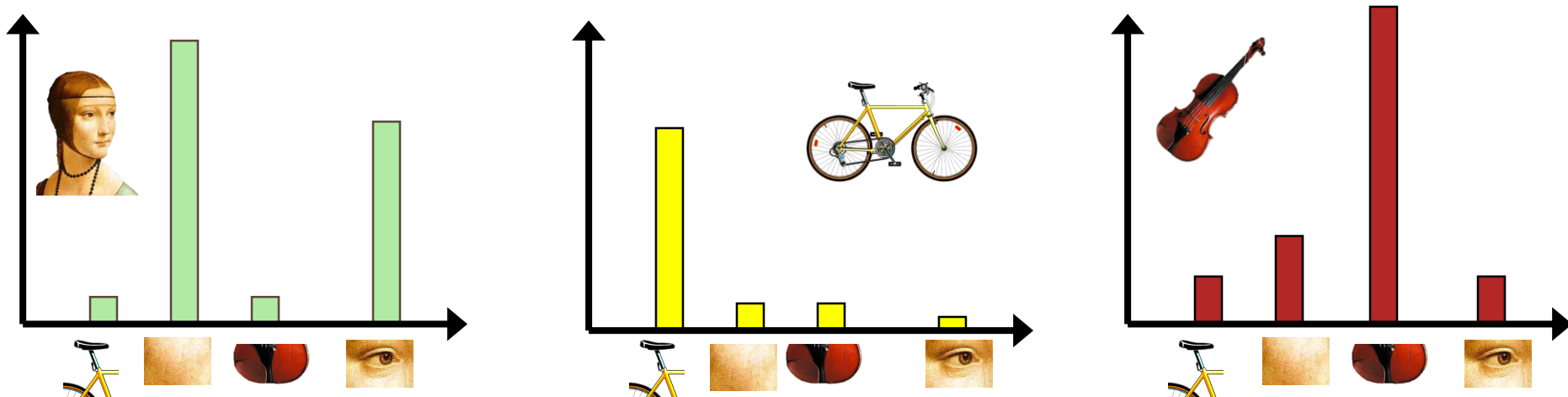


violin



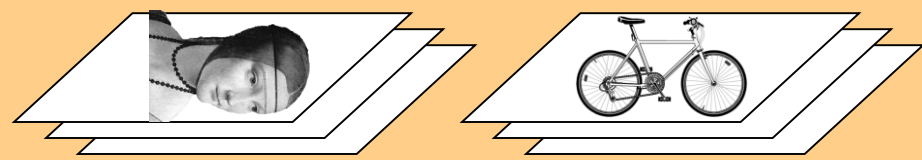
definition of “BoW”

- Independent features
- histogram representation



codewords dictionary

Representation



feature detection & representation



codewords dictionary

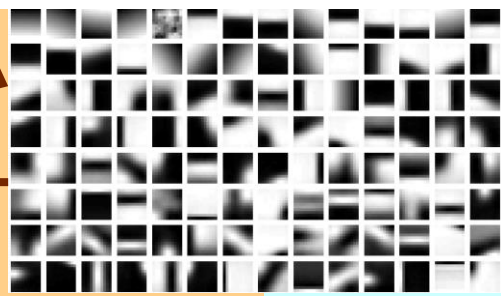
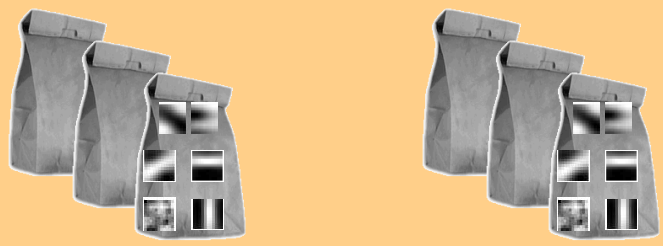


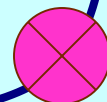
image representation



learning

category models (and/or) classifiers

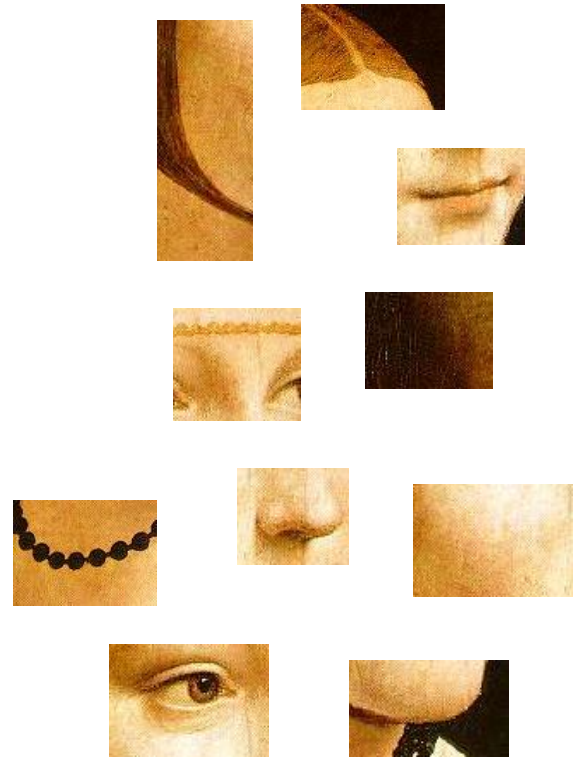
recognition



category decision

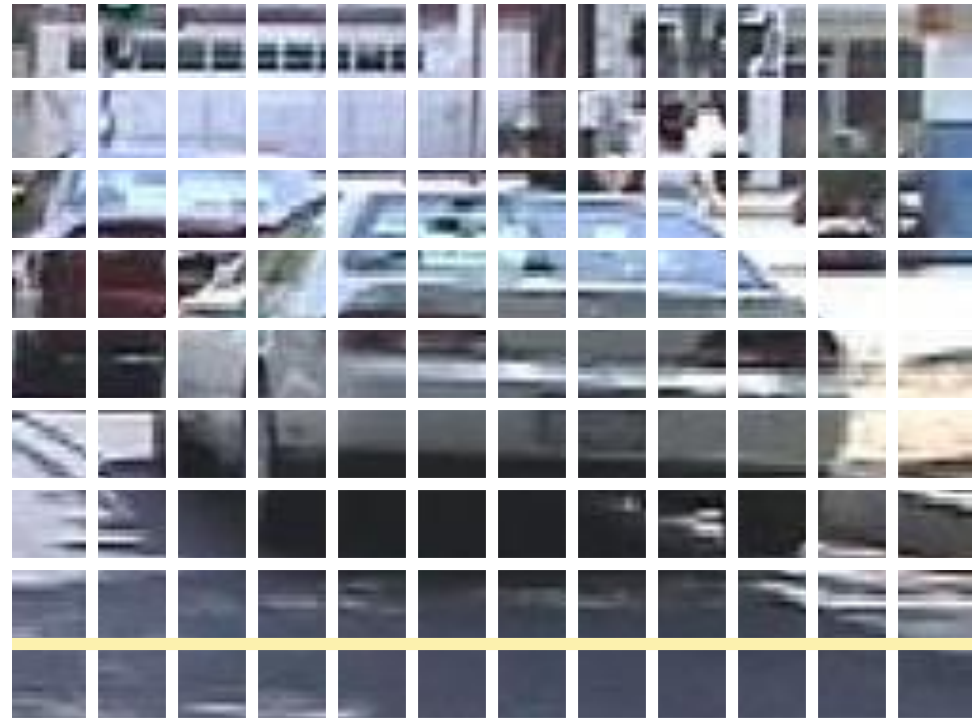


1.Feature detection and description



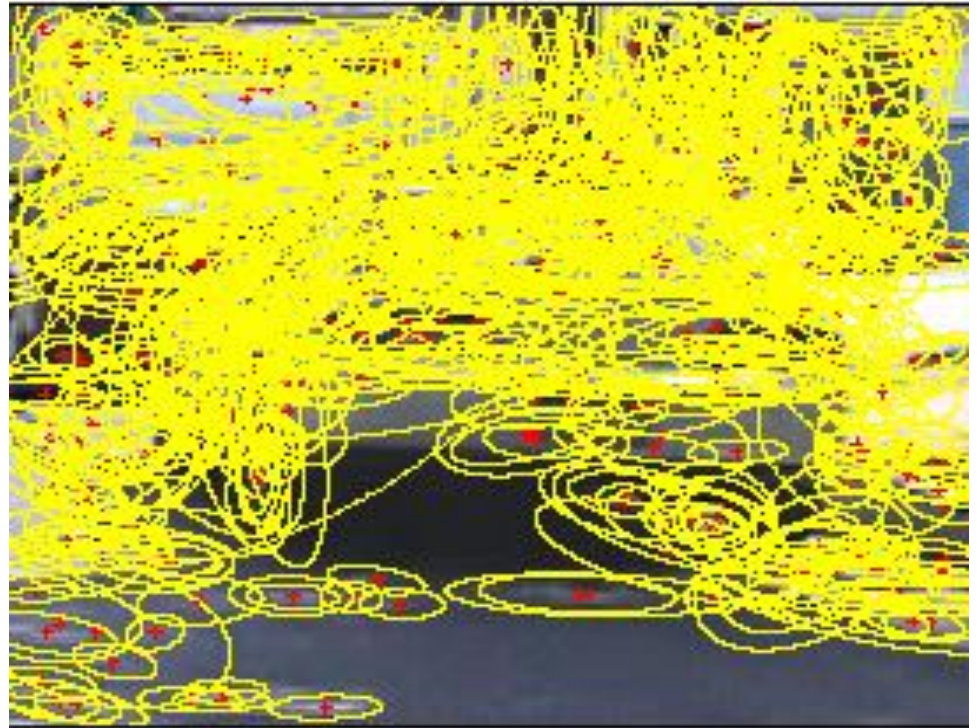
1. Feature detection and description

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



1. Feature detection and description

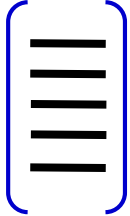
- Regular grid
 - Vogel & Schiele, 2003
 - Fei-Fei & Perona, 2005
- Interest point detector
 - Csurka, et al. 2004
 - Fei-Fei & Perona, 2005
 - Sivic, et al. 2005



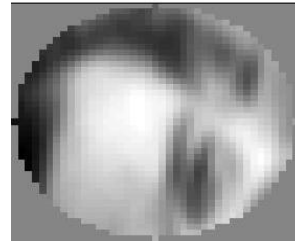
1.Feature detection and description

- 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)

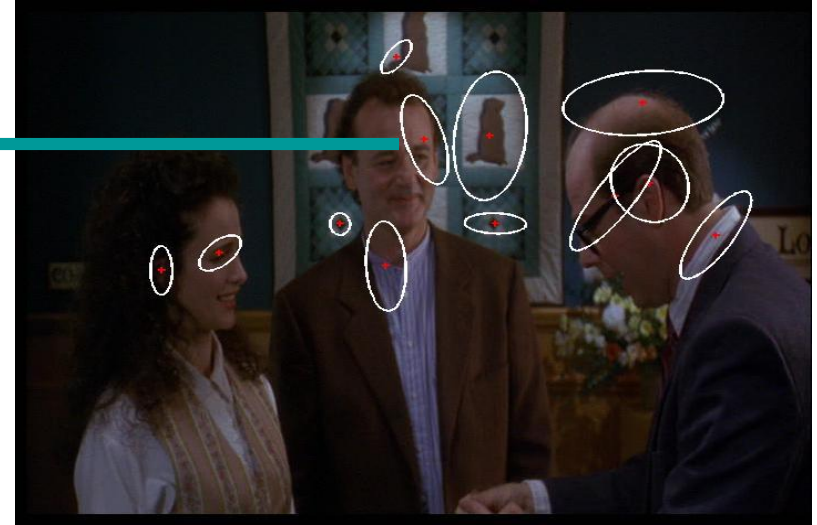
1. Feature detection and description



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



**Normalize
patch**



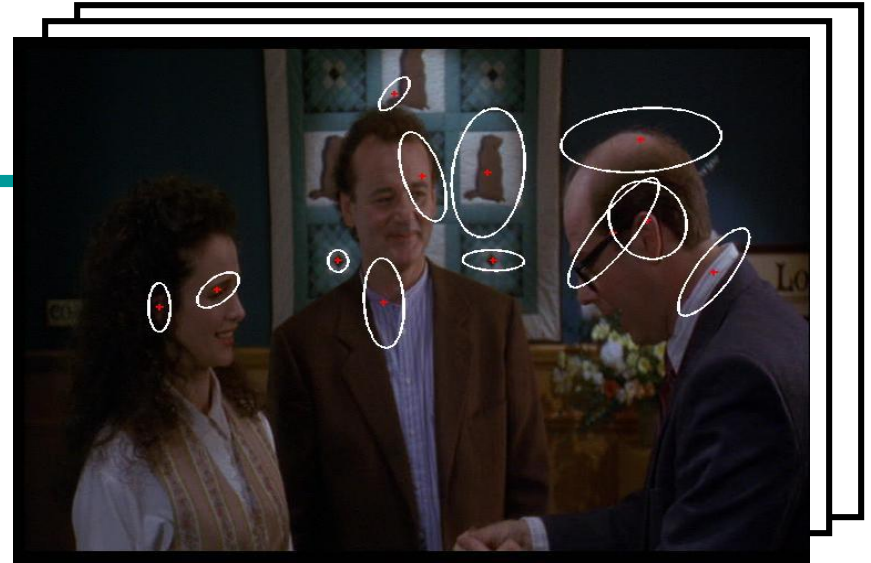
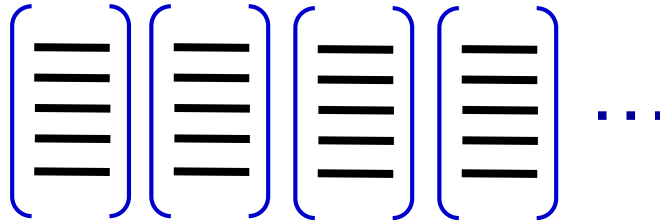
Detect patches

[Mikojczyk and Schmid '02]

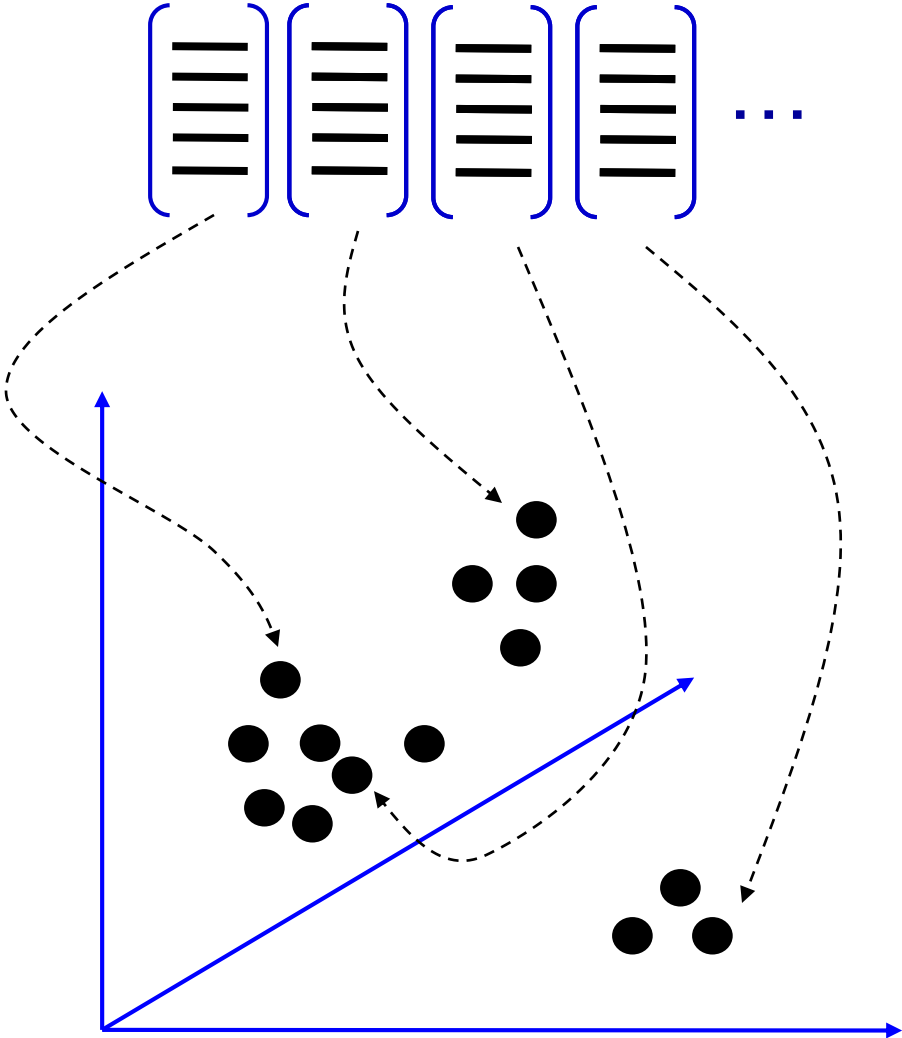
[Mata, Chum, Urban & Pajdla, '02]

[Sivic & Zisserman, '03]

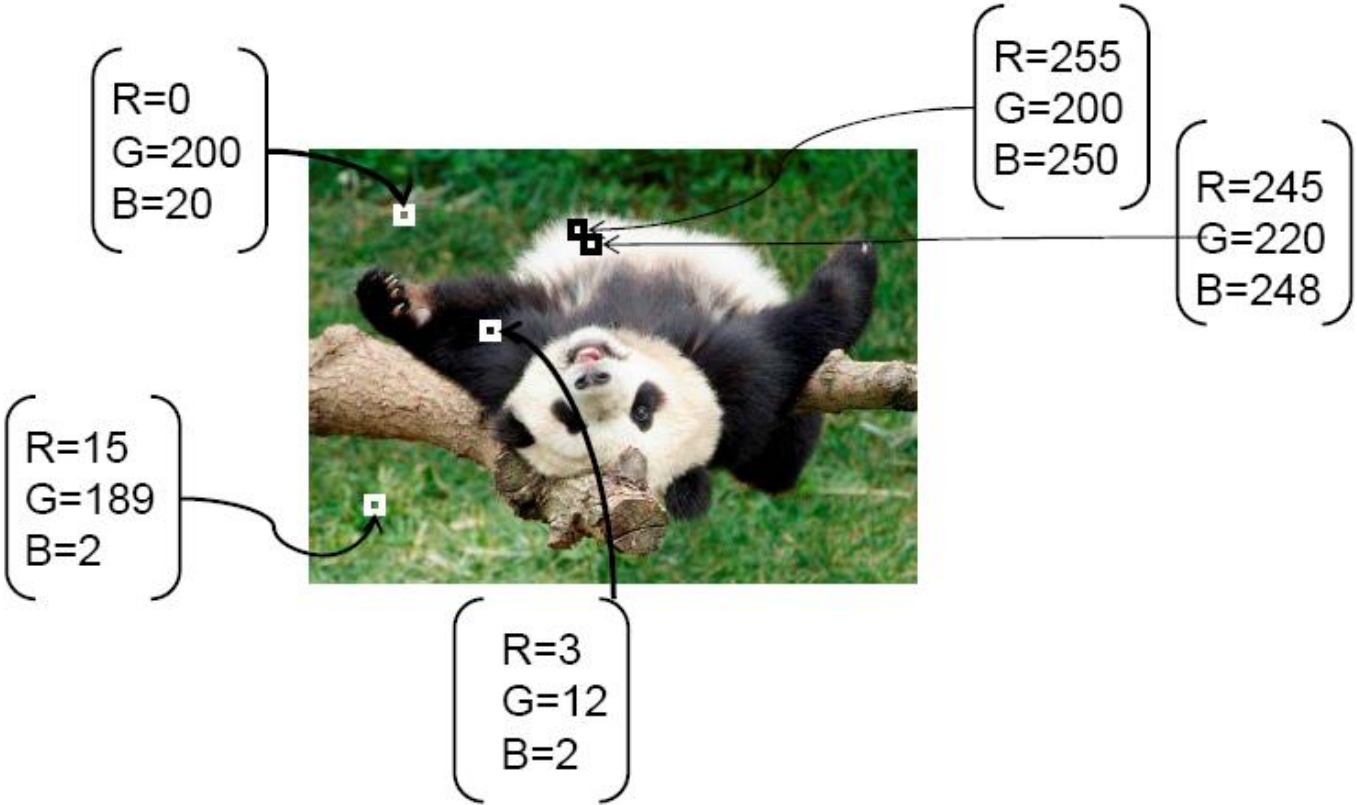
2. Codewords dictionary formation



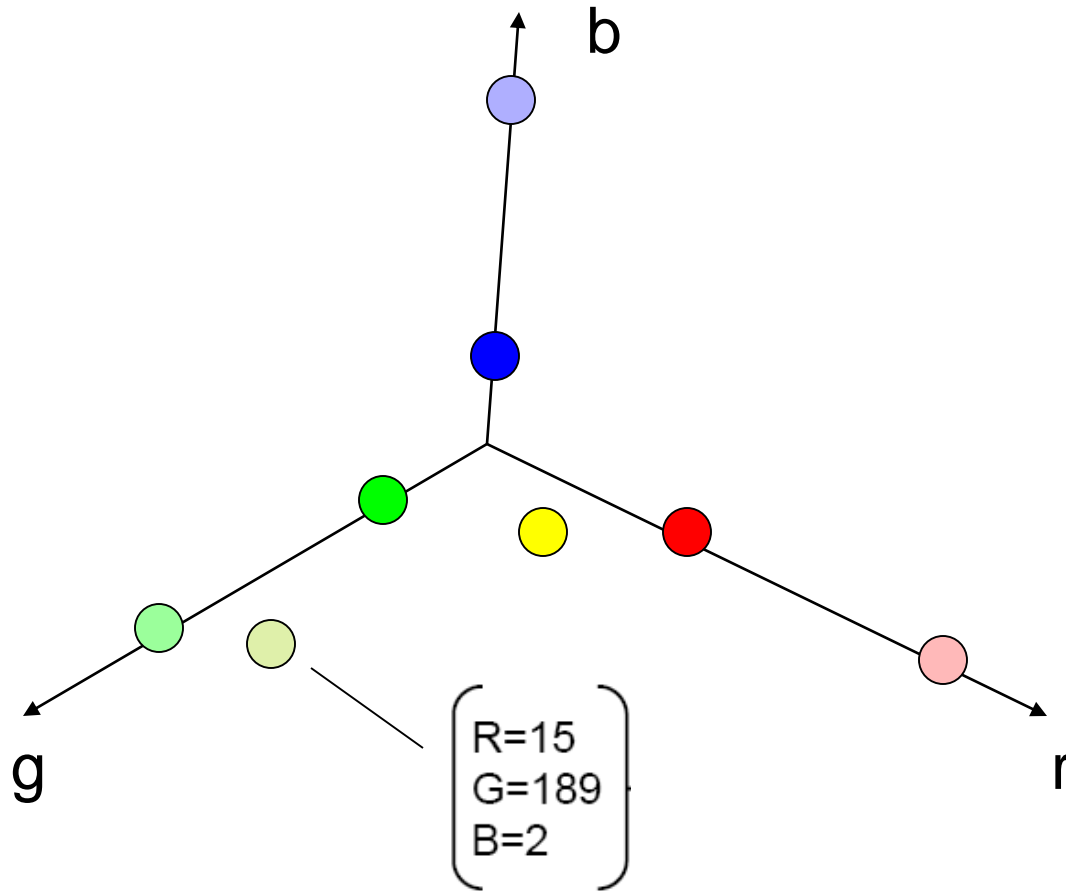
2. Codewords dictionary formation



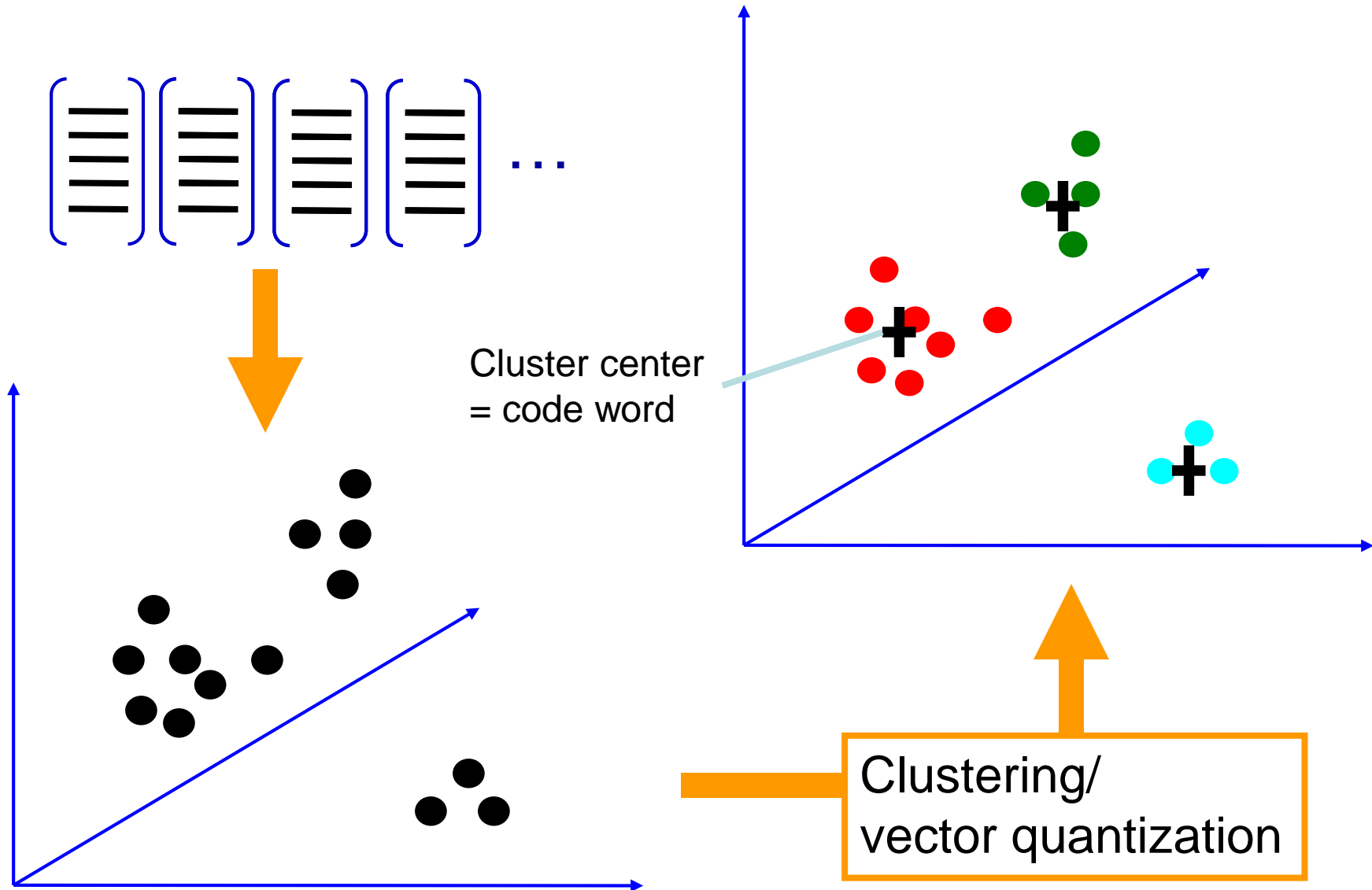
Example: color feature



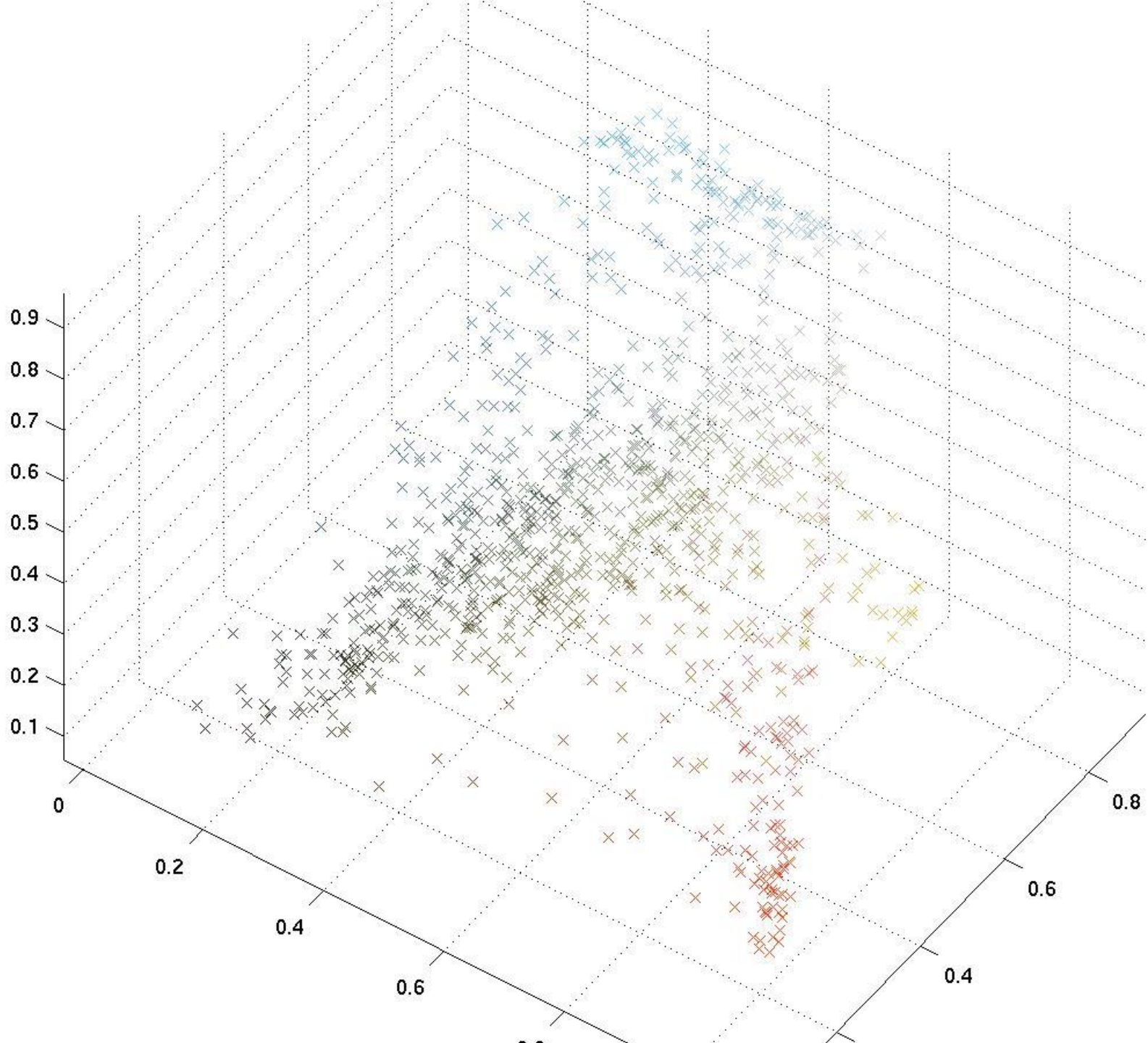
Example: color feature

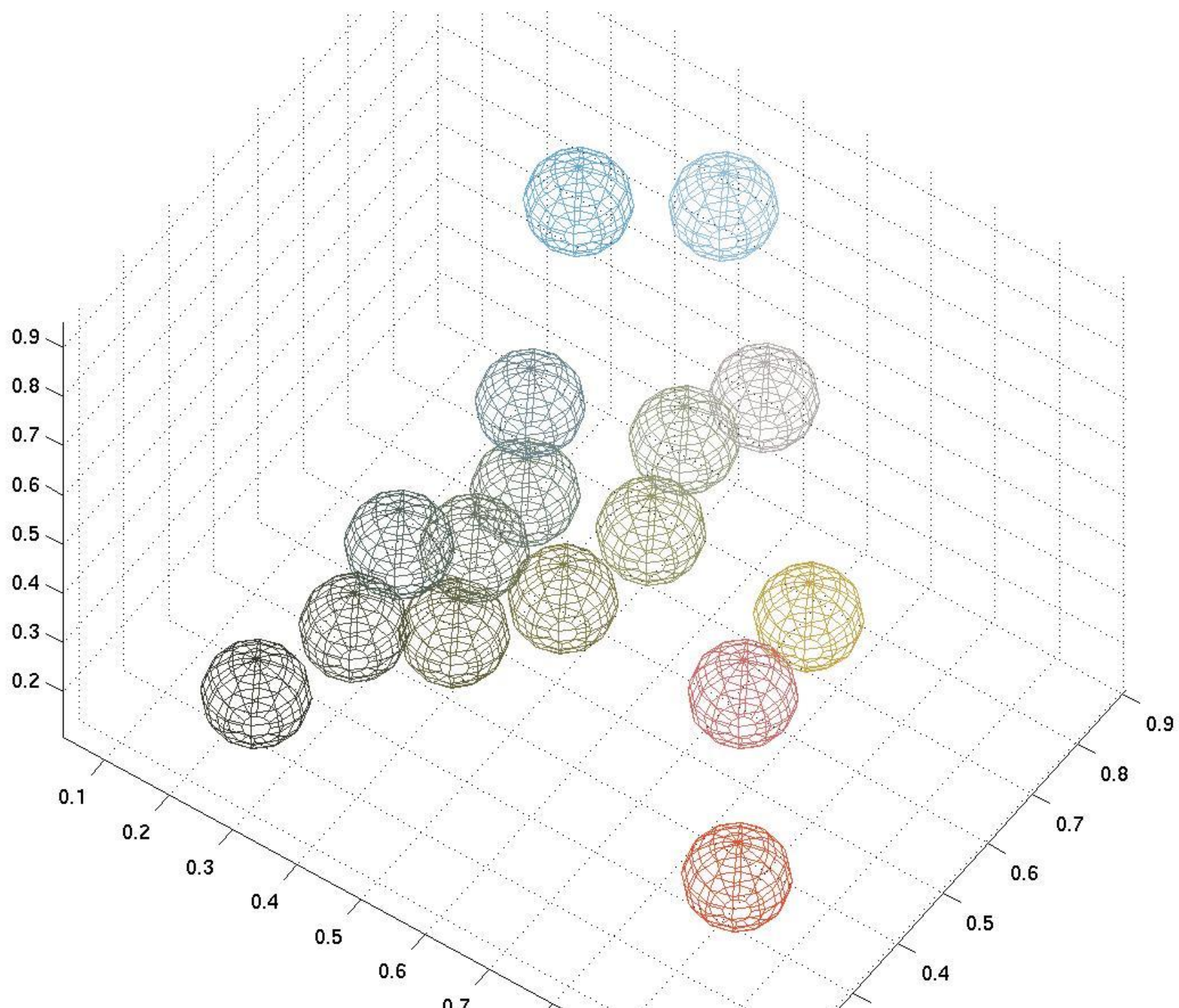


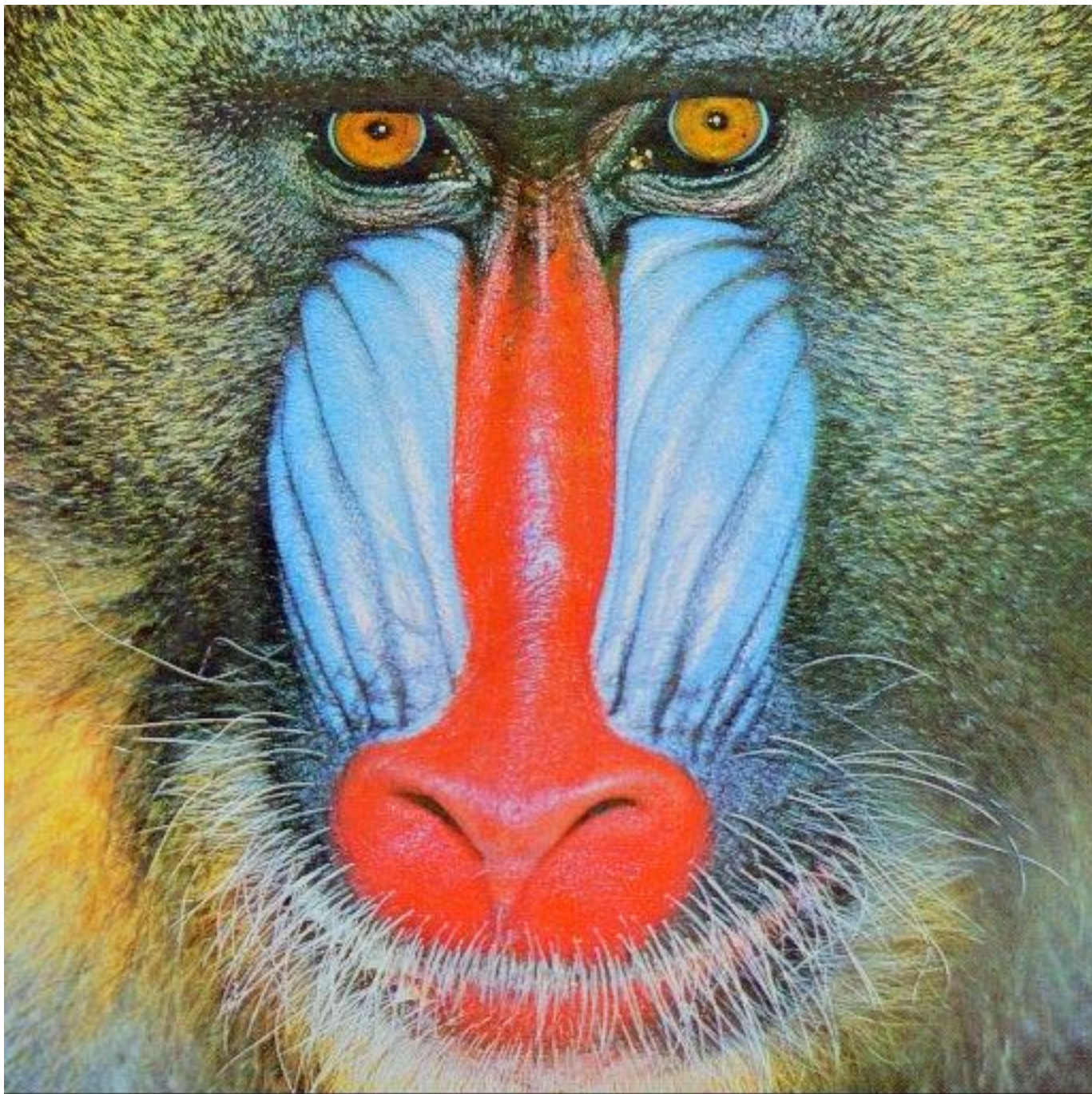
2. Codewords dictionary formation

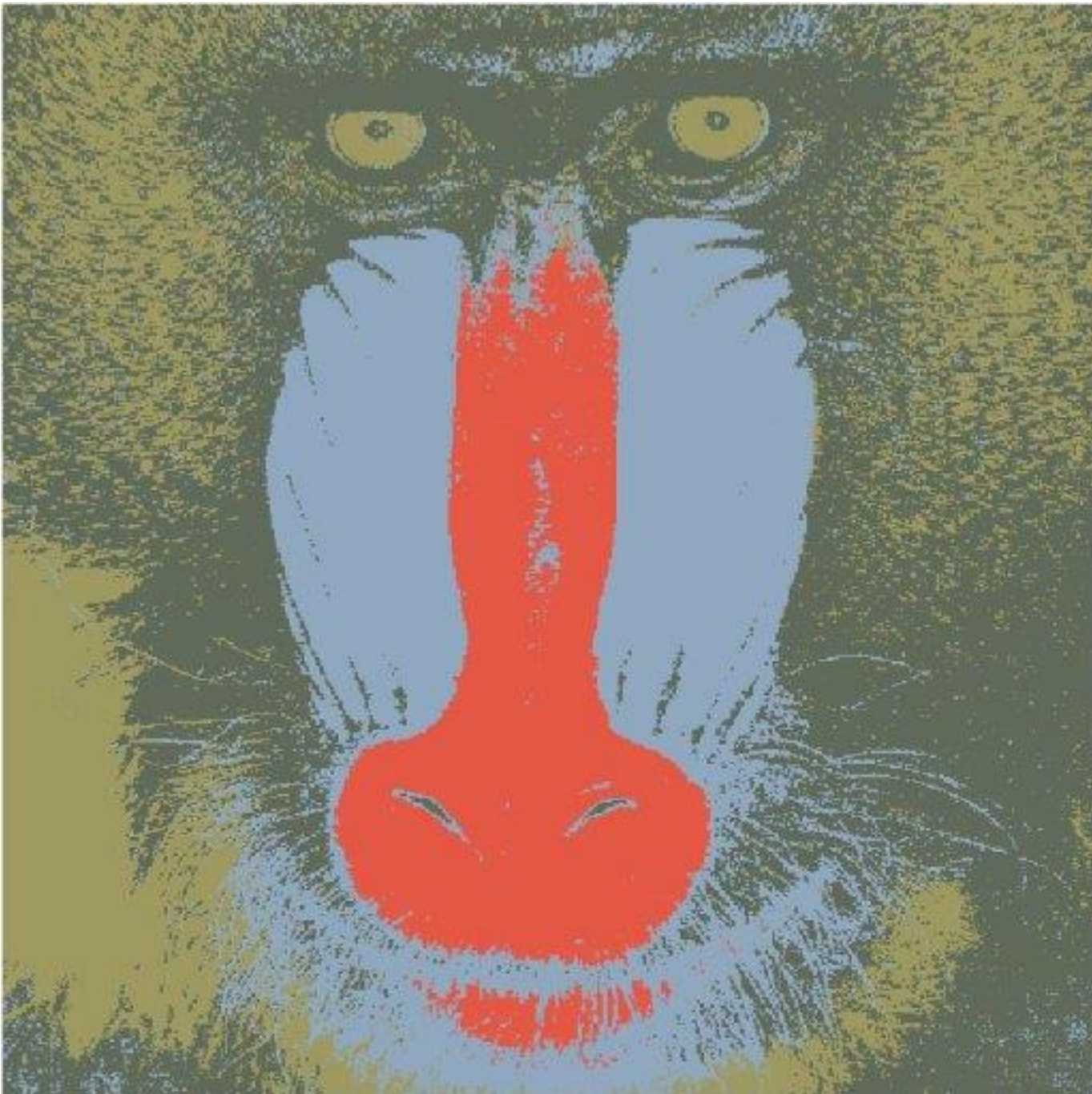


E.g., Kmeans, see CS131A



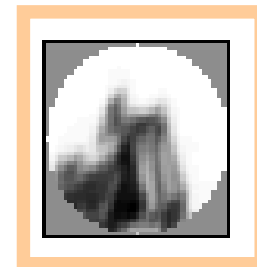
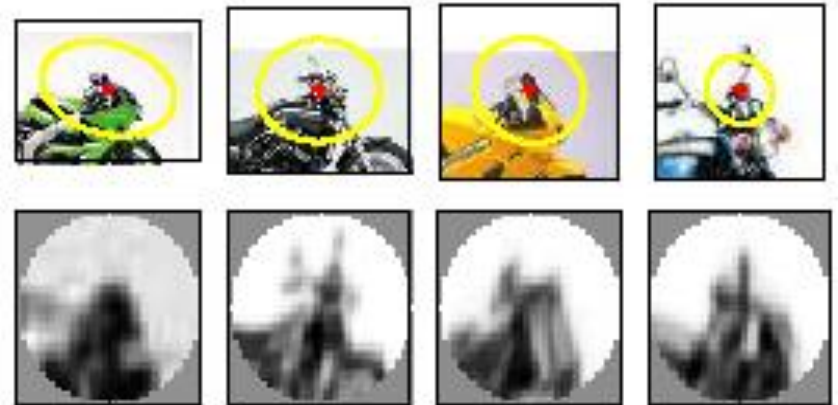
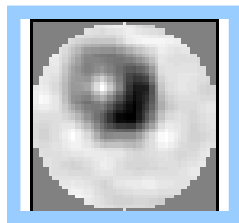
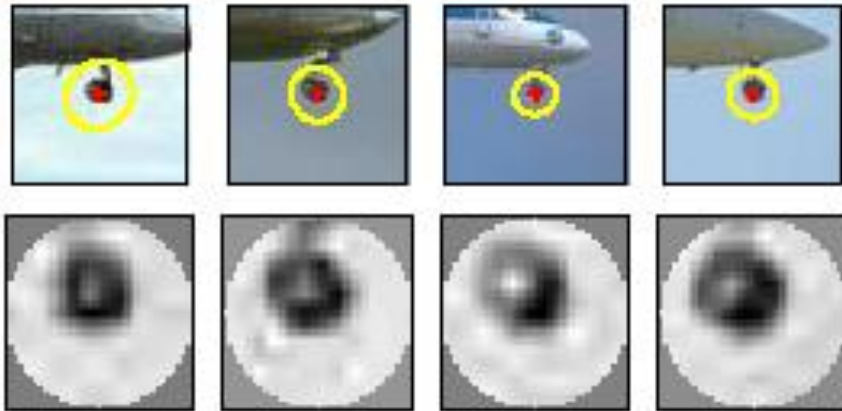




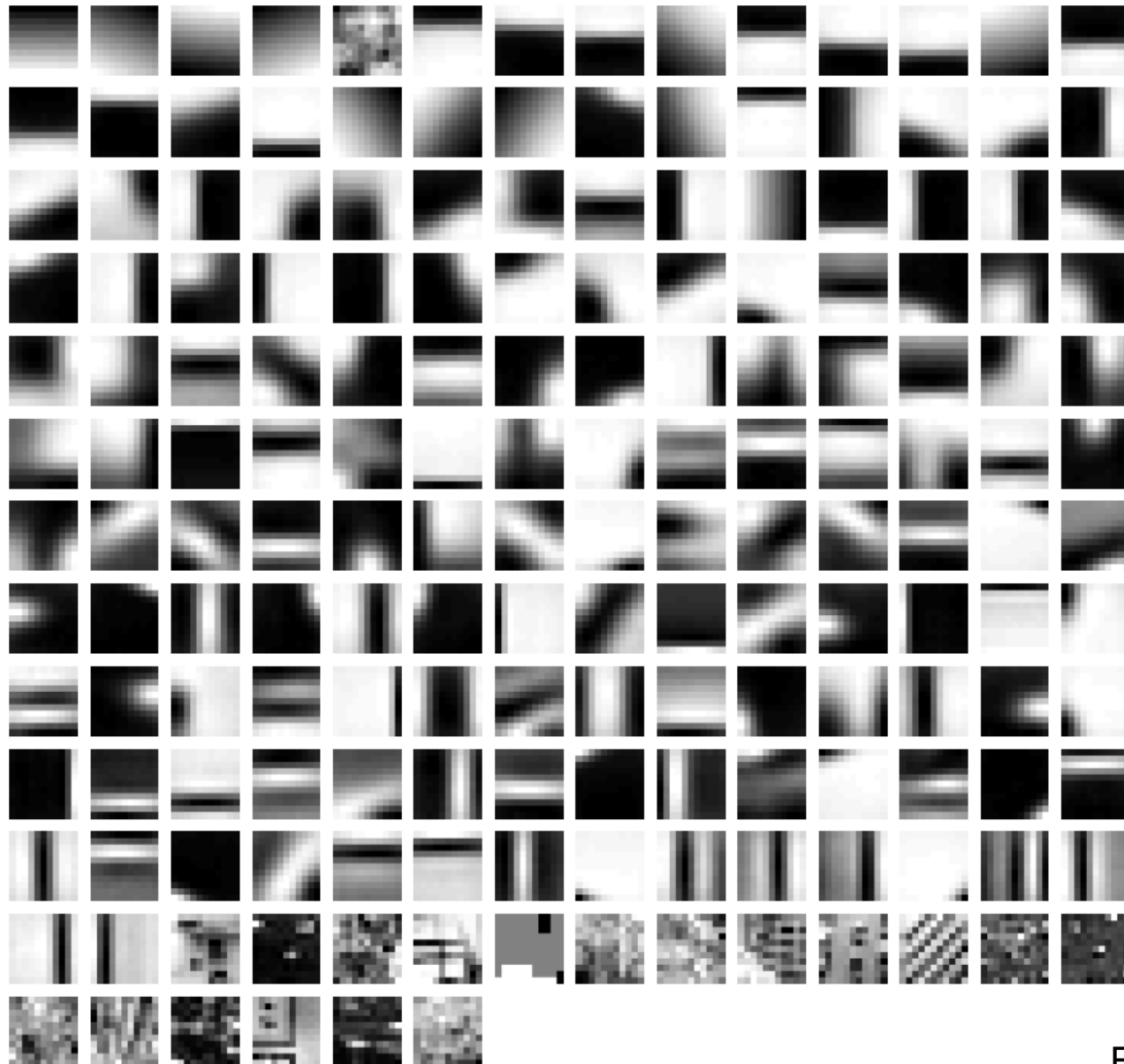


2. Codewords dictionary formation

- Image patch examples of codewords



2. Codewords dictionary formation



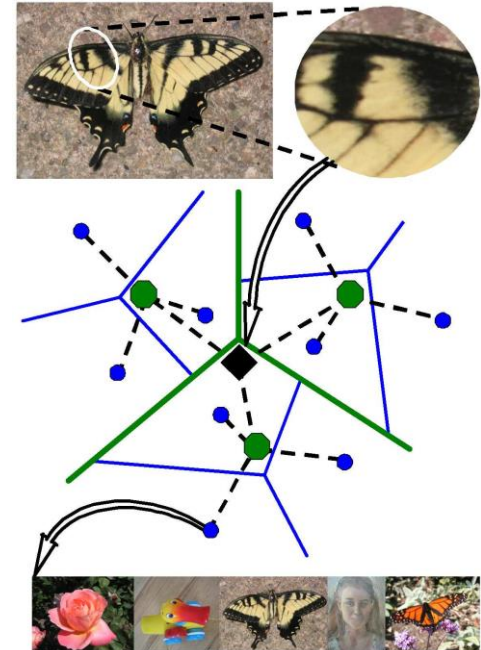
2. Codewords dictionary formation

- Typically a codeword dictionary is obtained from a training set comprising all the object classes of interests

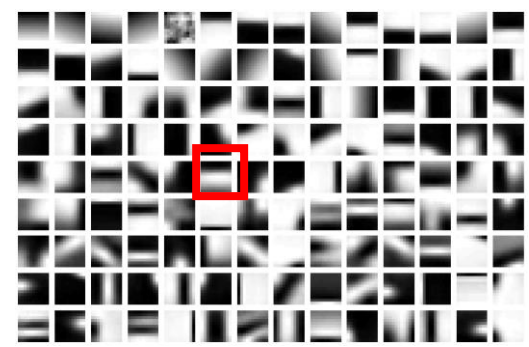
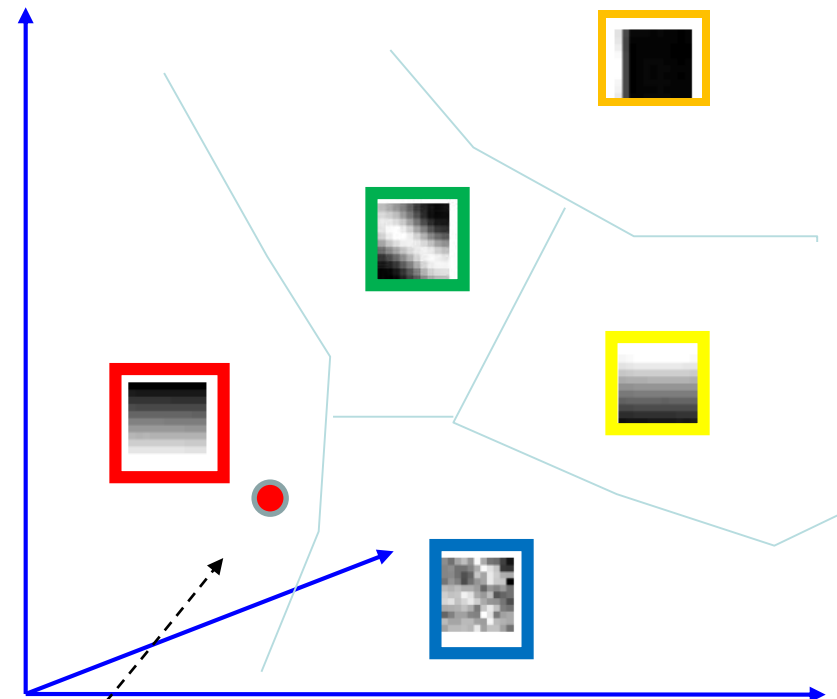
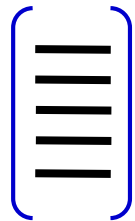
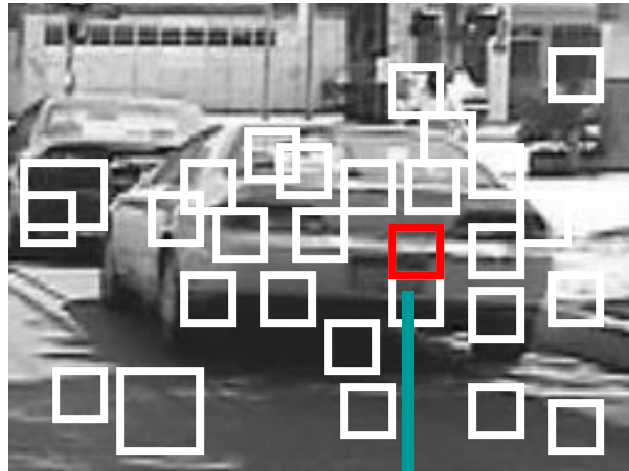


Visual vocabularies: Issues

- 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)



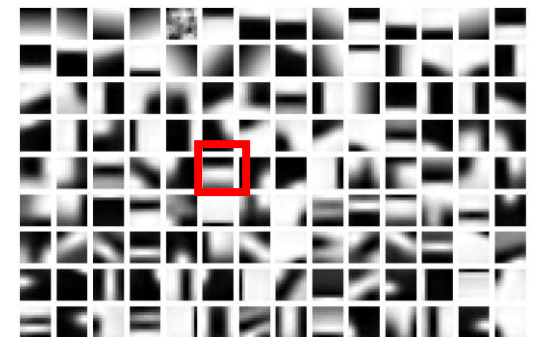
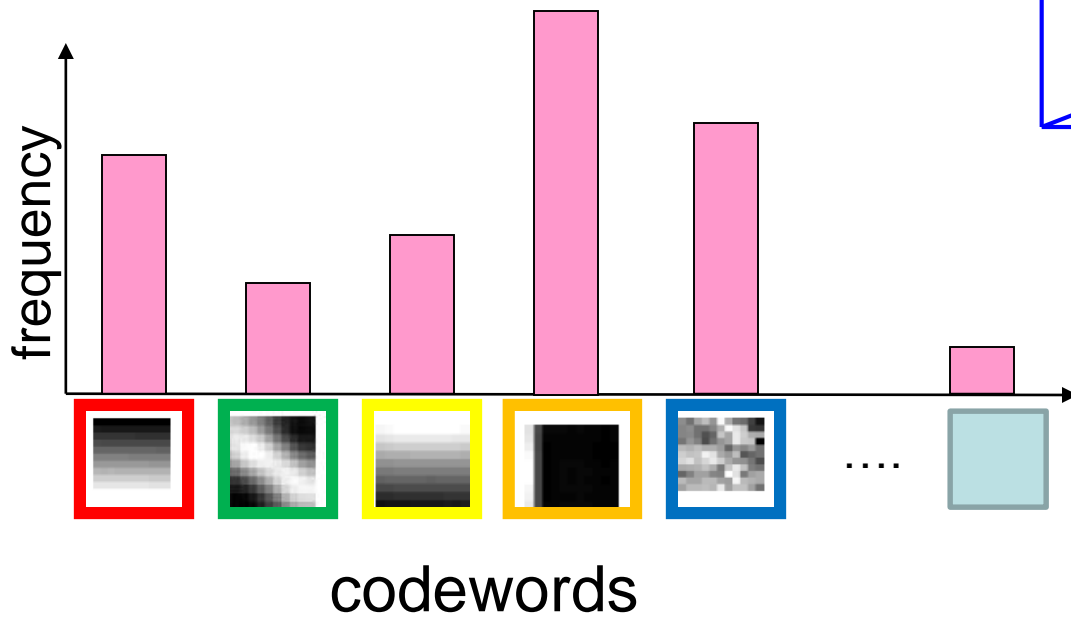
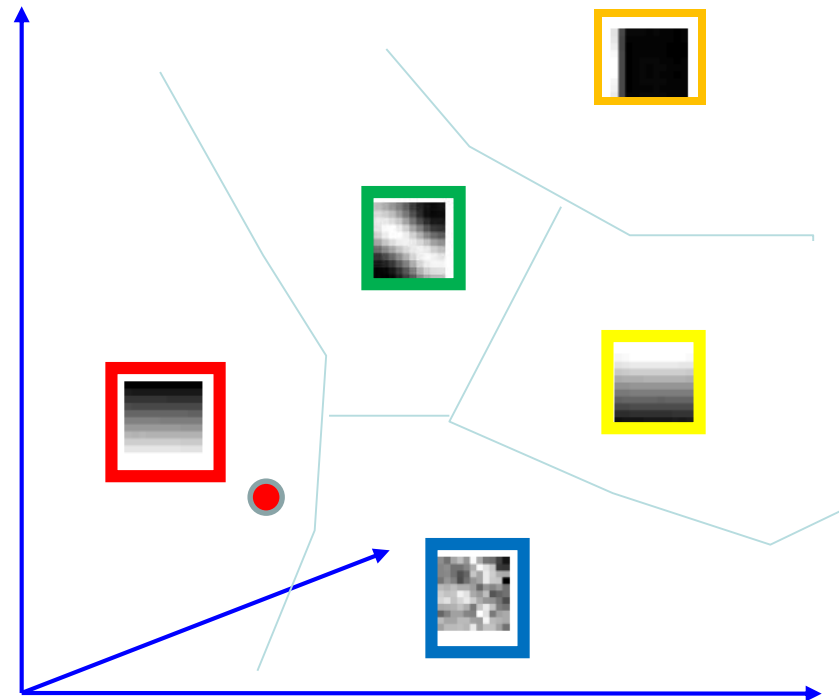
3. Bag of word representation



Codewords dictionary

- Nearest neighbors assignment
- K-D tree search strategy

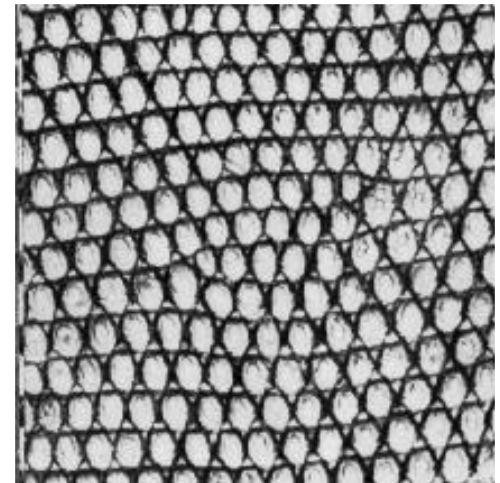
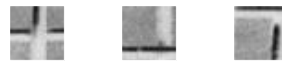
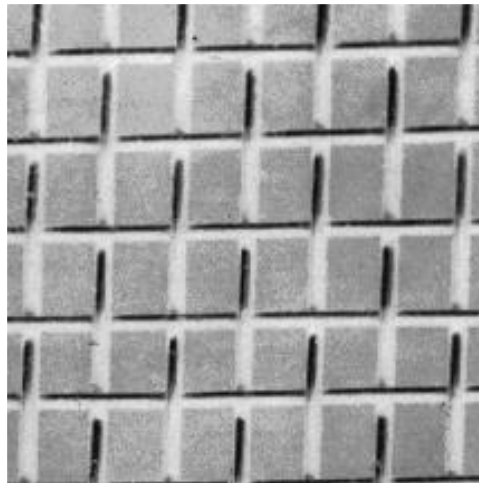
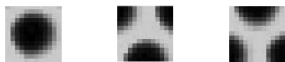
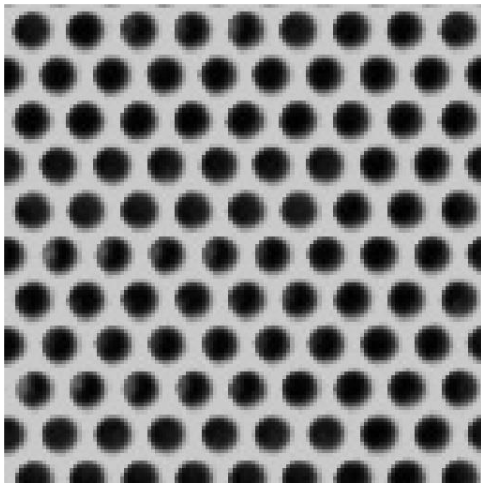
3. Bag of word representation



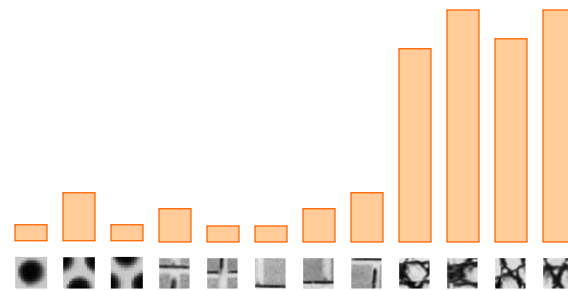
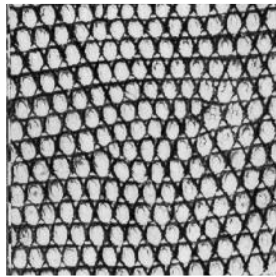
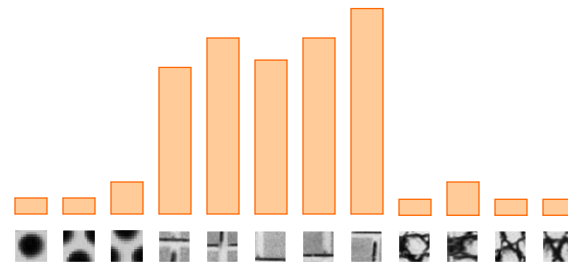
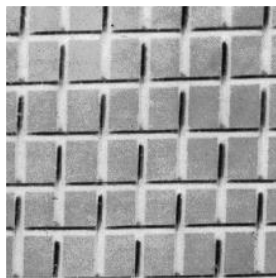
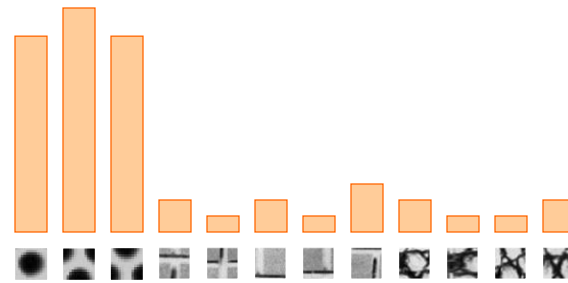
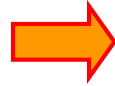
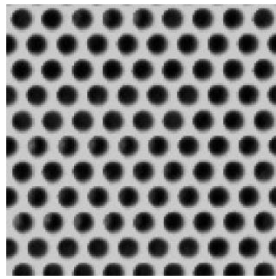
Codewords dictionary

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



Representing textures

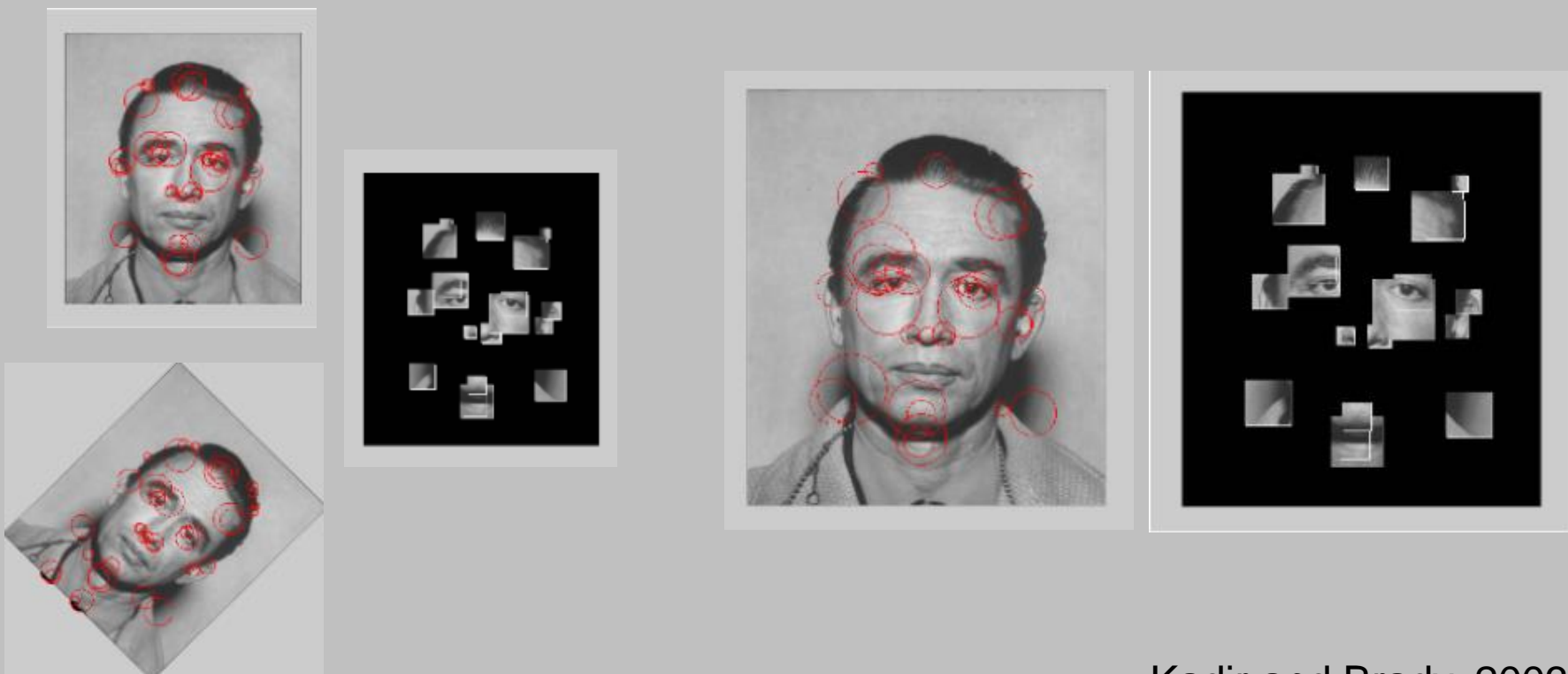


Julesz, 1981; Cula & Dana, 2001; Leung & Malik 2001; Mori, Belongie & Malik, 2001; Schmid 2001; Varma & Zisserman, 2002, 2003; Lazebnik, Schmid & Ponce, 2003

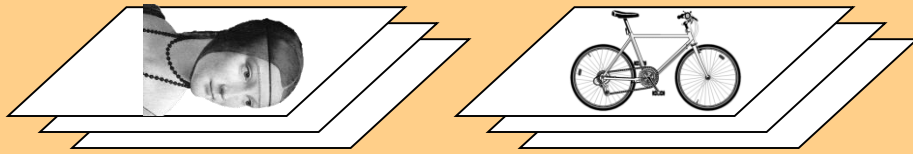
Credit slide: S. Lazebnik

Invariance issues

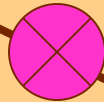
- Scale? Rotation? View point? Occlusions?
 - Implicit;
 - depends on detectors and descriptors



Representation



1. feature detection & representation



2. codewords dictionary

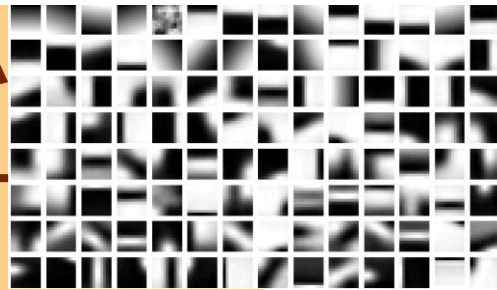
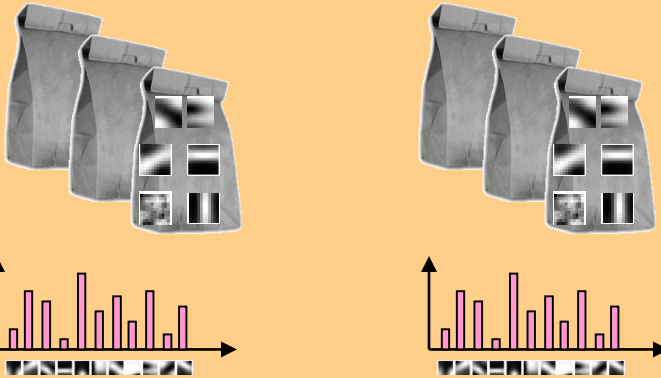


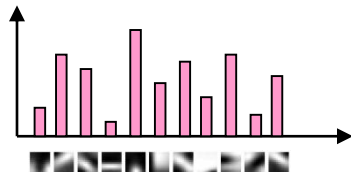
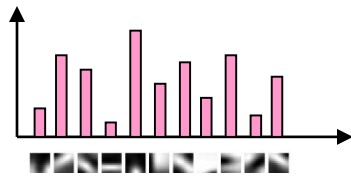
image representation

3.

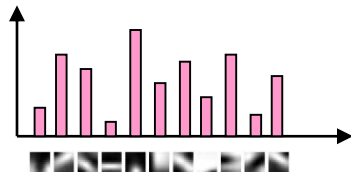


category models

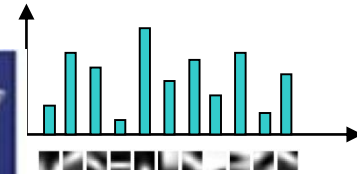
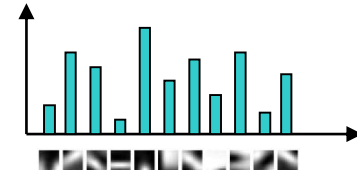
Category models



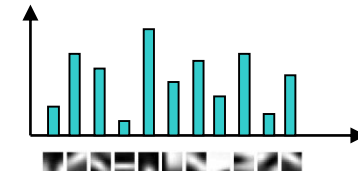
⋮



Class 1



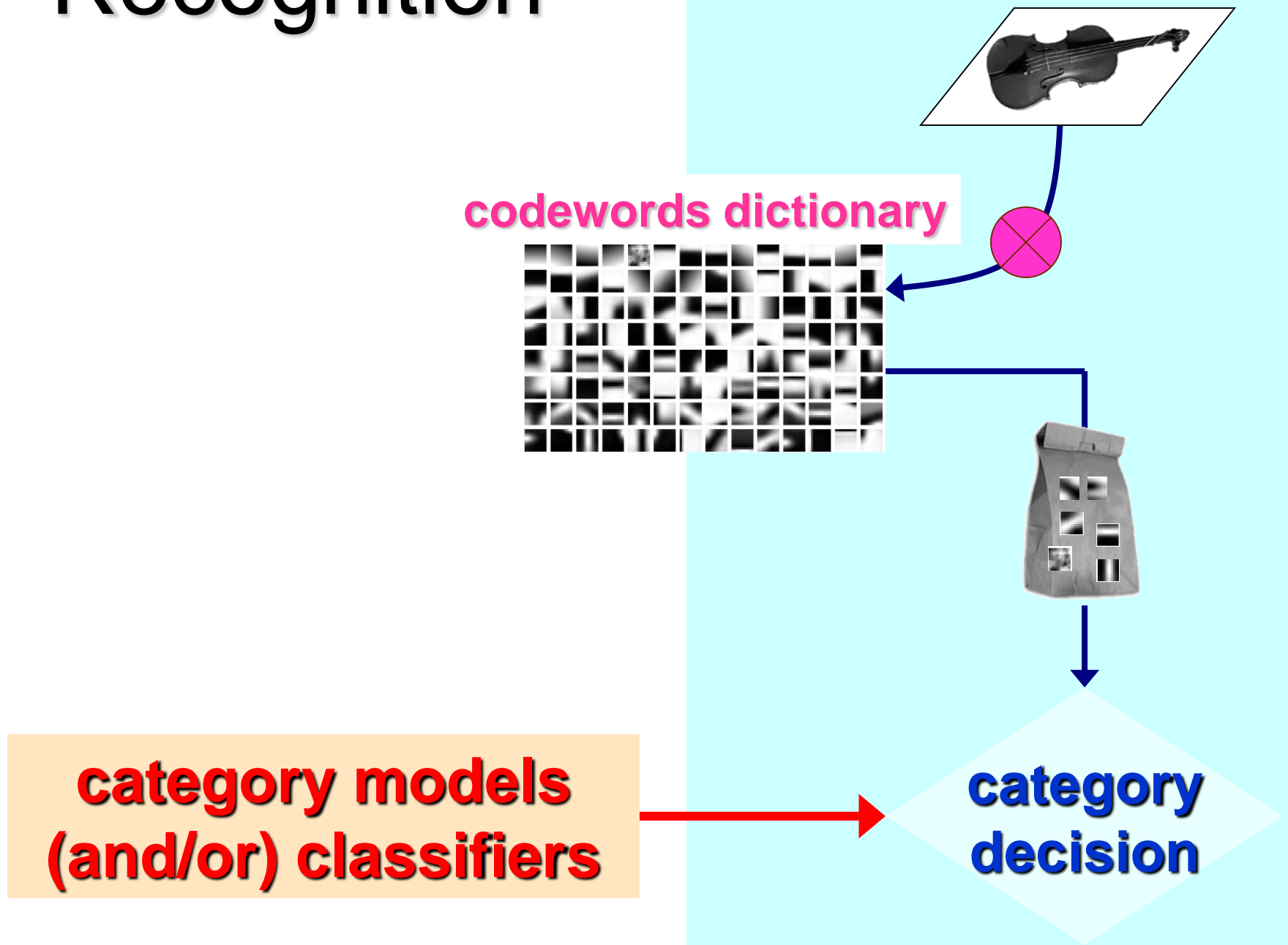
⋮



Class N

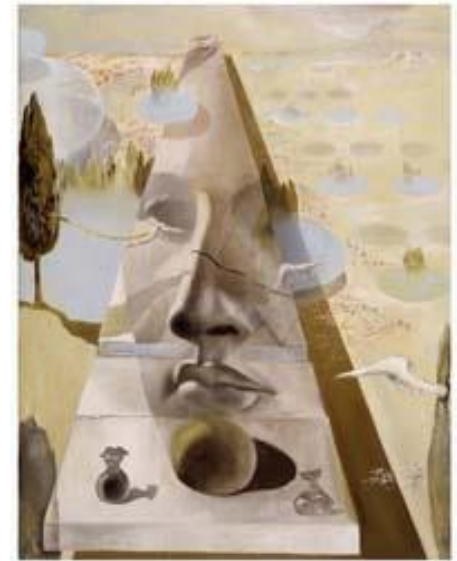
...

Recognition



Lecture 12

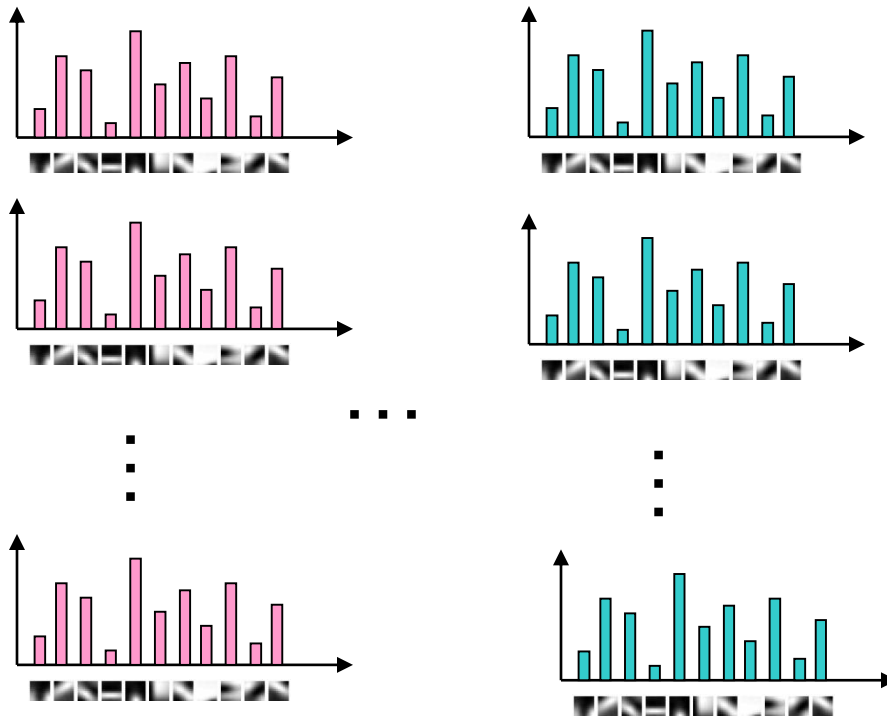
Visual recognition



- Bag of words models for object recognition and classification
 - Discriminative methods
 - Nearest neighborhood
 - Linear classifier
 - SVM
 - Generative methods

Discriminative classifiers

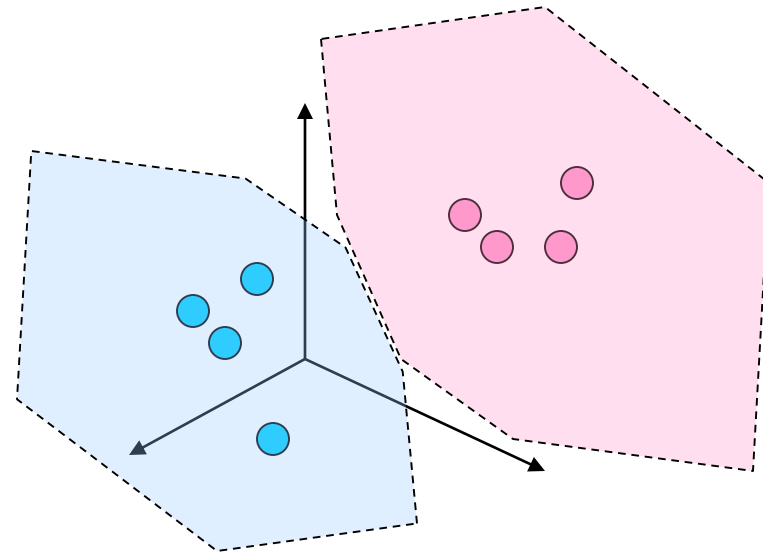
category models



Class 1

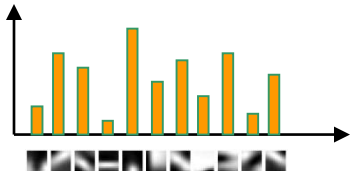
Class N

Model space



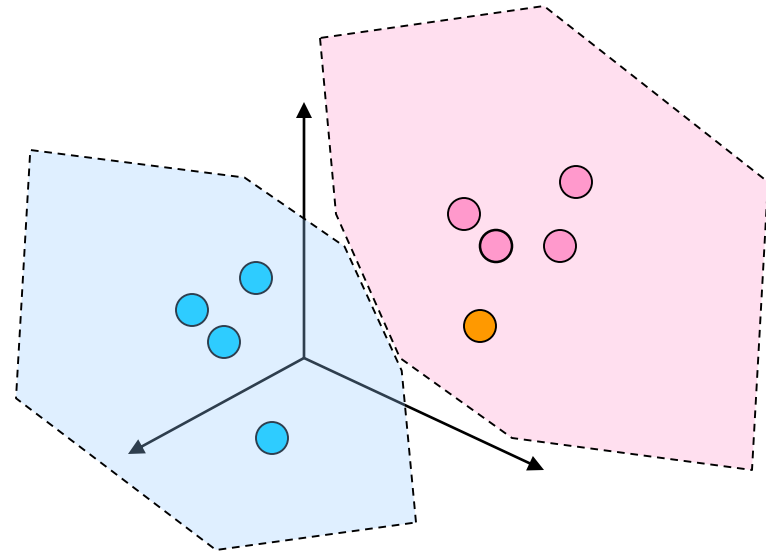
Discriminative classifiers

Query image



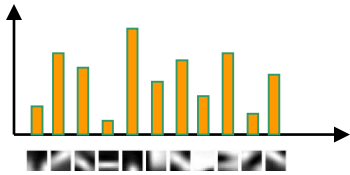
Winning class: pink

Model space

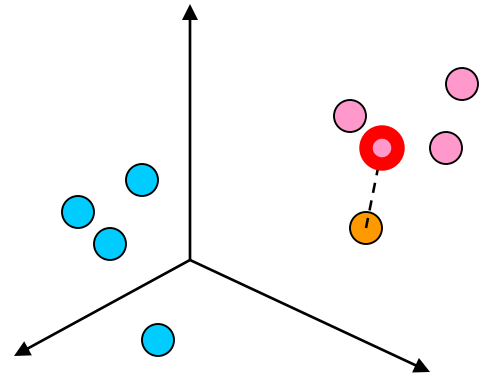


Nearest Neighbors classifier

Query image



Model space

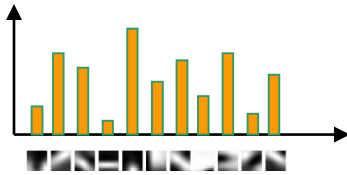


Winning class: pink

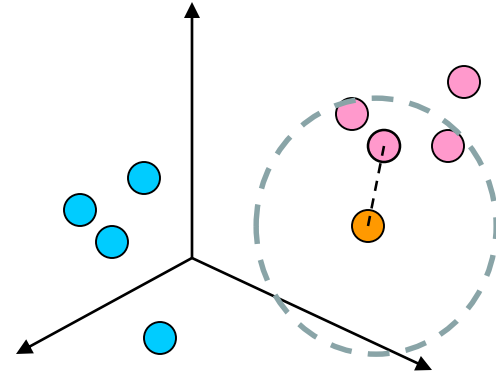
- Assign label of nearest training data point to each test data point

K- Nearest Neighbors classifier

Query image



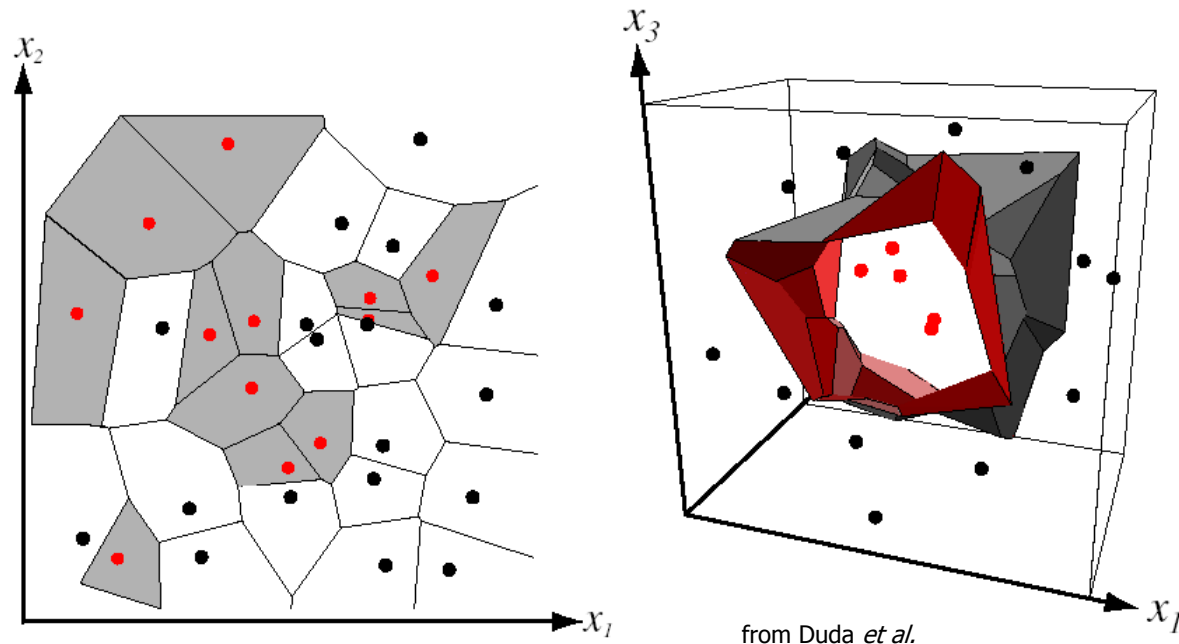
Model space



Winning class: pink

- 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

K- Nearest Neighbors classifier



- Voronoi partitioning of feature space for 2-category 2-D and 3-D data
- For k dimensions: k -D tree = space-partitioning data structure for organizing points in a k -dimensional space
- Enable efficient search
- Nice tutorial: <http://www.cs.umd.edu/class/spring2002/cmsc420-0401/pbasic.pdf>

Functions for comparing histograms

Jan Puzicha, Yossi Rubner, Carlo Tomasi, Joachim M. Buhmann: [Empirical Evaluation of Dissimilarity Measures for Color and Texture](#). ICCV 1999

- L1 distance

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

- χ^2 distance

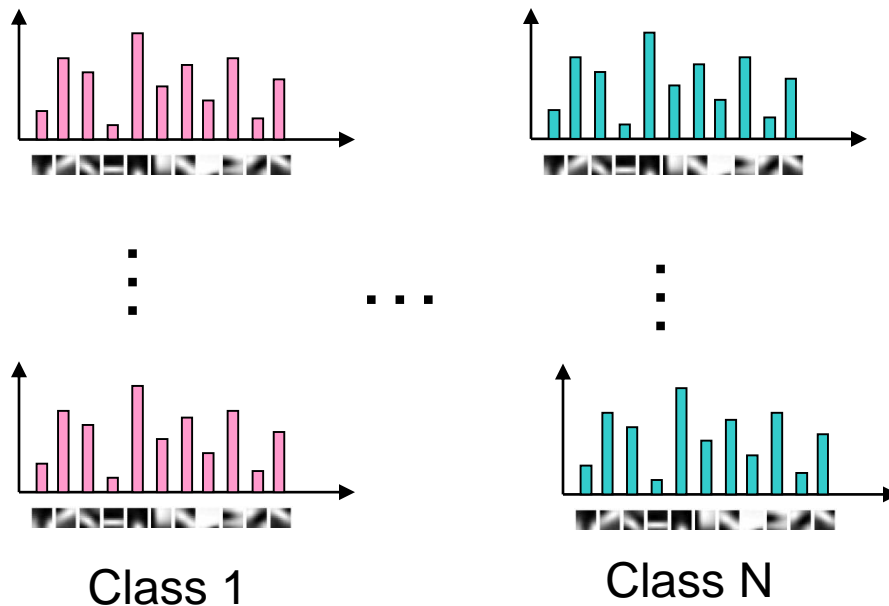
$$D(h_1, h_2) = \sum_{i=1}^N \frac{(h_1(i) - h_2(i))^2}{h_1(i) + h_2(i)}$$

- Quadratic distance (*cross-bin*)

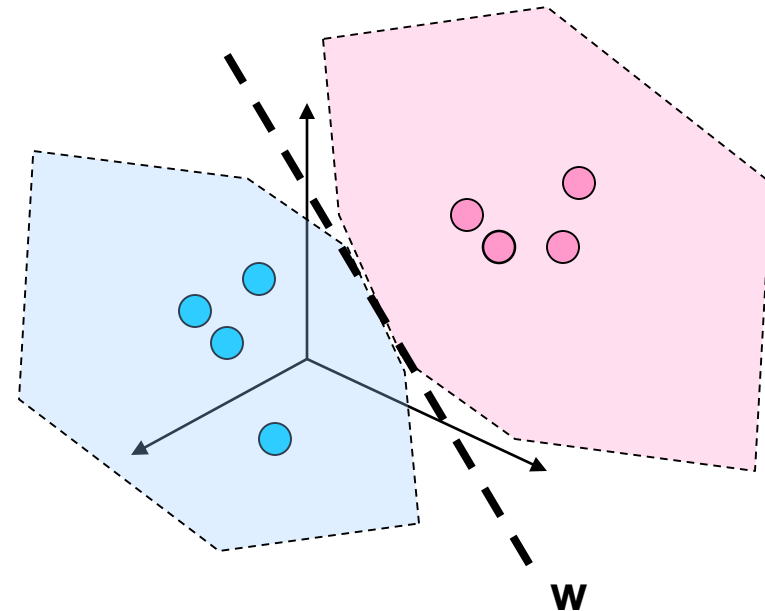
$$D(h_1, h_2) = \sum_{i,j} A_{ij} (h_1(i) - h_2(j))^2$$

Discriminative classifiers (linear classifier)

category models



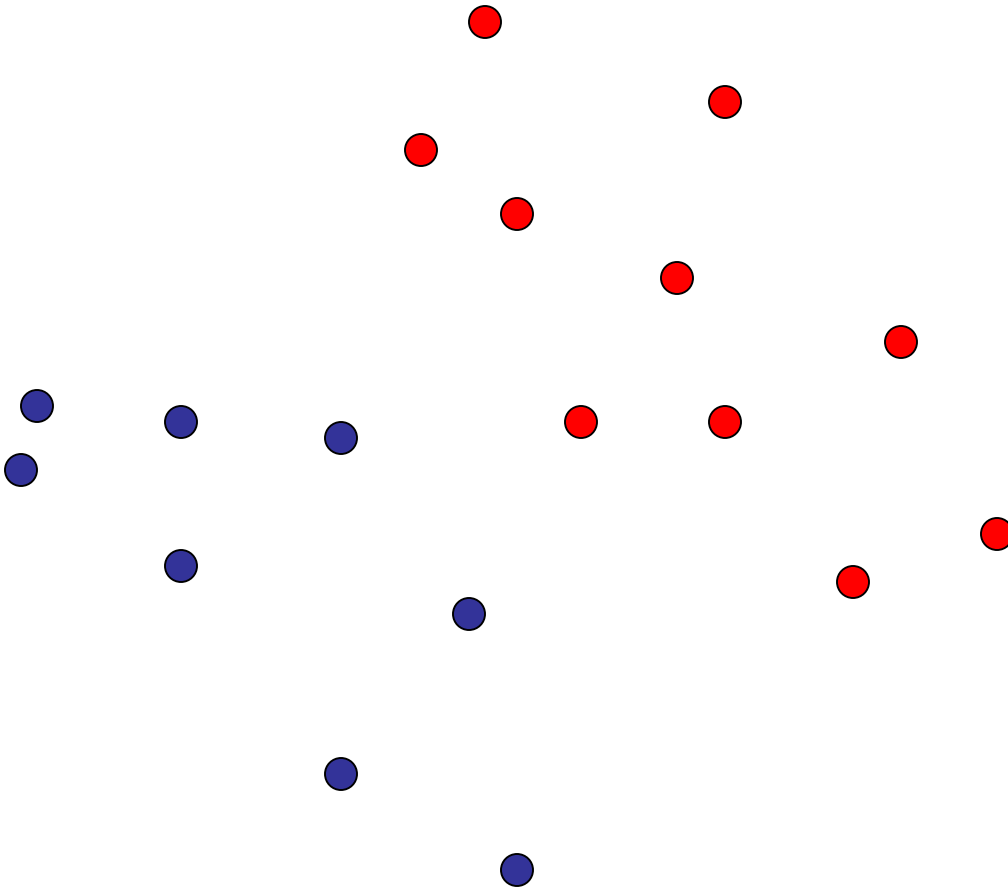
Model space



- For a linear classifier, the training data is used to learn w and then discarded
- Only w is needed for classifying new data

Linear classifiers

- We want to classify two classes of points
- Each point x_i can have two labels y_i : $\{-1, +1\}$

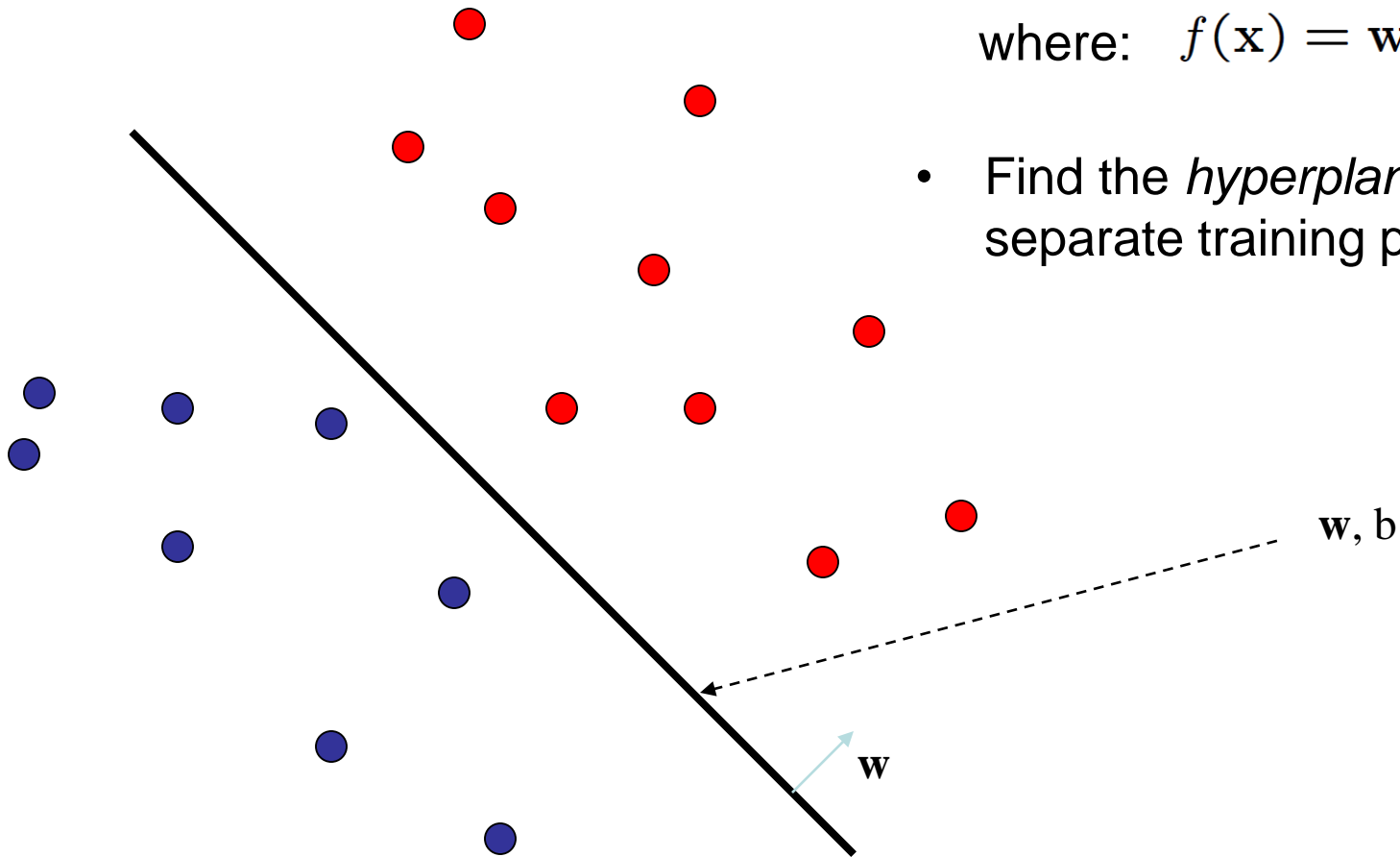


Linear classifiers

- GOAL: learn a classifier $f(\mathbf{x})$ such that:
$$f(\mathbf{x}_i) \begin{cases} \geq 0 & y_i = +1 \\ < 0 & y_i = -1 \end{cases}$$

where: $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$

- Find the *hyperplane* (\mathbf{w}, b) to separate training points

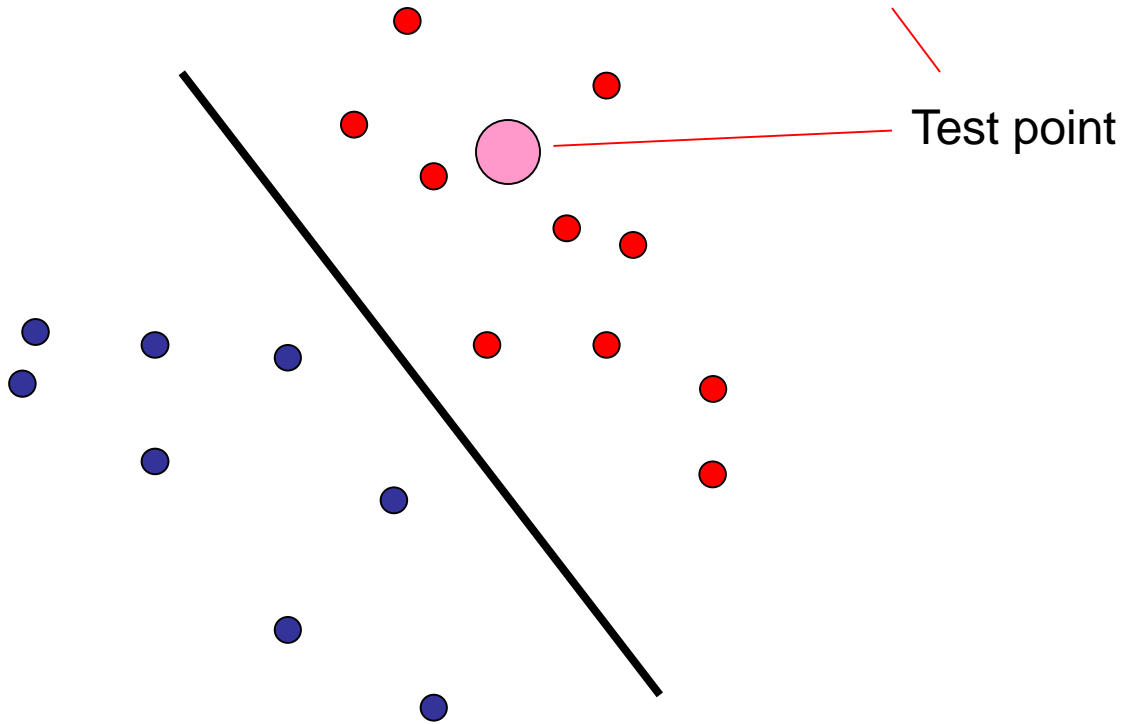


Linear classifiers

- Once w, b are learnt, we can do classification:

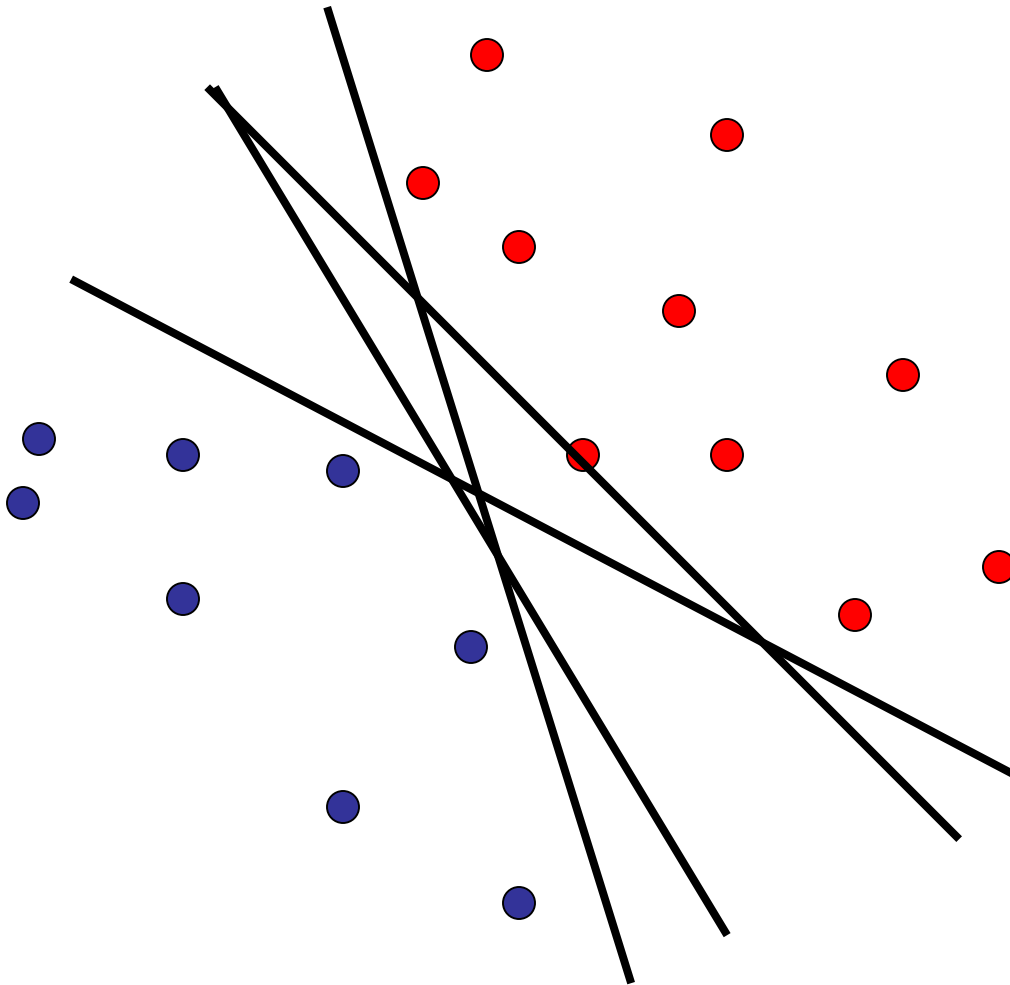
if $\mathbf{x} \cdot \mathbf{w} + b \geq 0 \rightarrow$ class 1

if $\mathbf{x} \cdot \mathbf{w} + b < 0 \rightarrow$ class 2



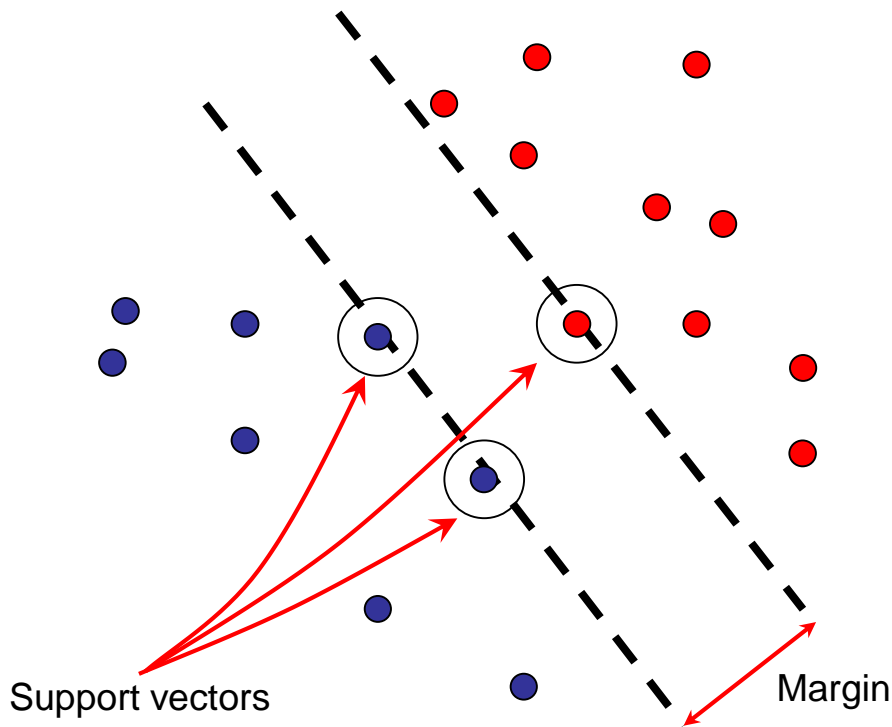
Linear classifiers

Which hyperplane is best?



Support vector machines

- Select two hyperplanes such:
 - They separate the training points
 - There are no points between them
 - Their distance is maximized



Support vectors: $\mathbf{x}_i \cdot \mathbf{w} + b = \pm 1$

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

Margin = $2 / \|\mathbf{w}\|$

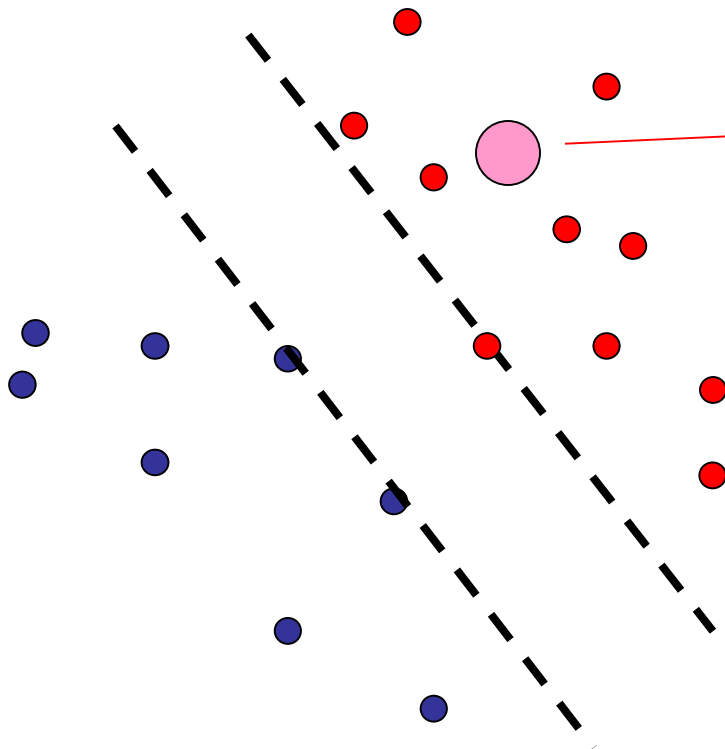
Solution:

$$\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$$
$$b = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i \cdot \mathbf{w} - y_i)$$

- The region bounded by them is called "the margin".
- **Maximum margin** solution: most stable under perturbations of the inputs

Support vector machines

- Classification: $f(\mathbf{x}) = \sum_i \alpha_i y_i \mathbf{x}_i \cdot \mathbf{x} + b$



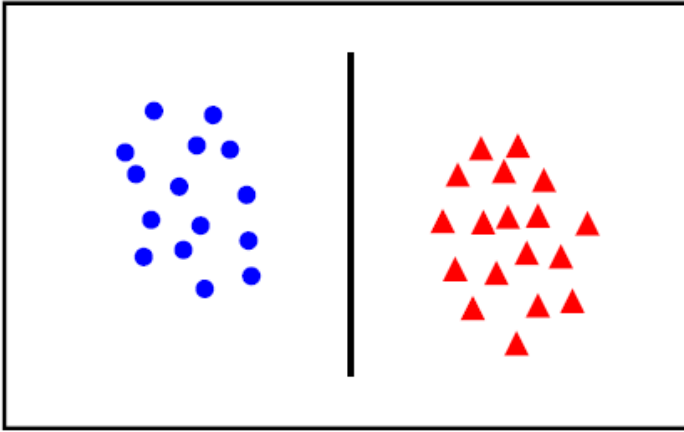
Test point

if $f(\mathbf{x}) \geq 0 \rightarrow \mathbf{x} \in \textit{class 1}$

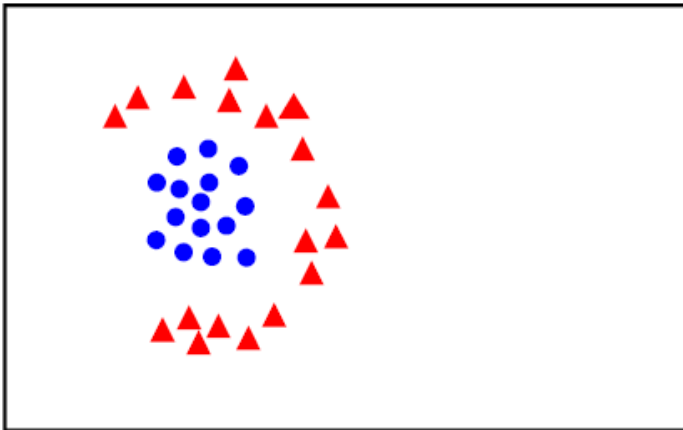
if $f(\mathbf{x}) < 0 \rightarrow \mathbf{x} \in \textit{class 2}$

Linear separability

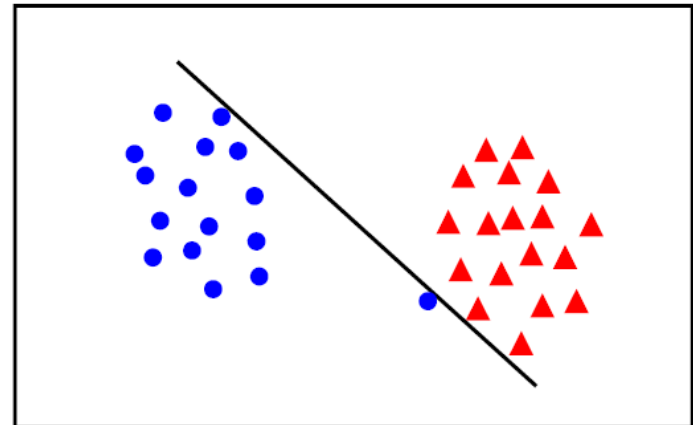
Linearly separable



Linearly **not** separable



Linearly separable
with small margins

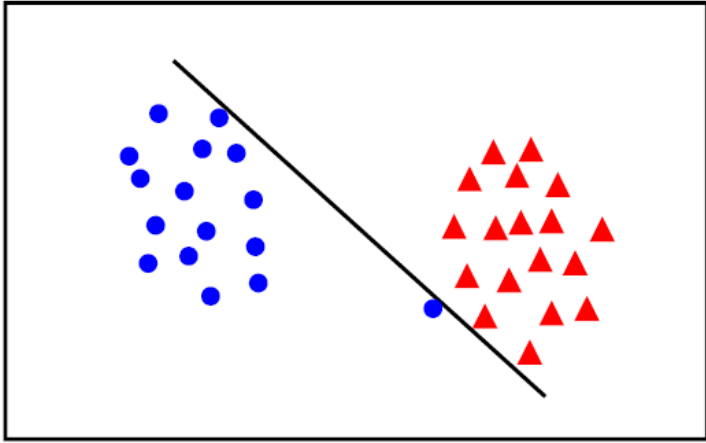


Linear separability

Two possible solutions:

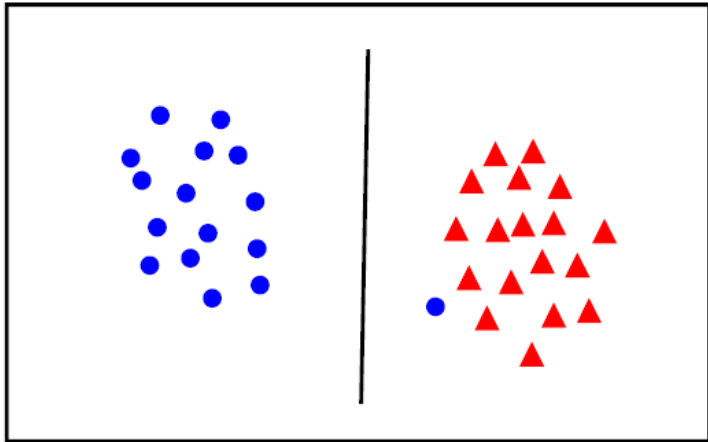
- Introduce soft variables (through slack variables)
- Non linear separation function (e.g., non linear SVM)

Soft margins



The points can be linearly separated but there is a very narrow margin

Or they cannot be separated at all



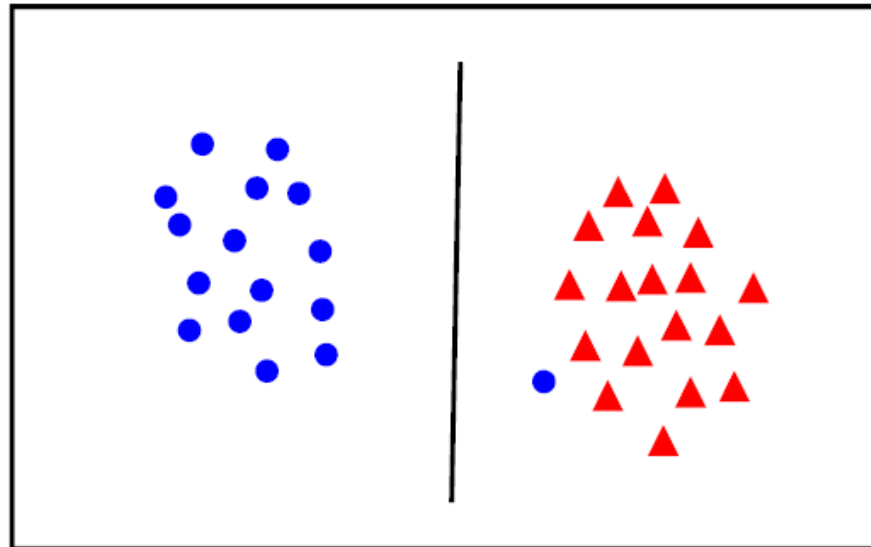
IDEA: still seek at large margin solution, even though one constraint is violated

In general there is a trade off between the margin and the number of mistakes on the training data

Soft margins

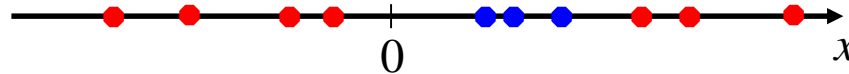
By Corinna Cortes and Vladimir N. Vapnik, 1995

- Use soft margin violations instead of the hard one: Find hyperplanes that split the examples as “cleanly” as possible, while still maximizing the distance to the nearest cleanly split examples.
- Introduce slack variables to measure the degree of misclassification of the data

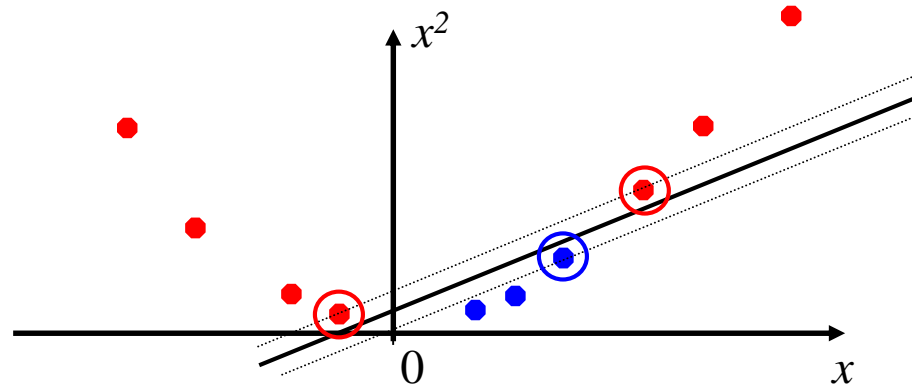


Nonlinear SVMs

- Given a non-linearly separable dataset:

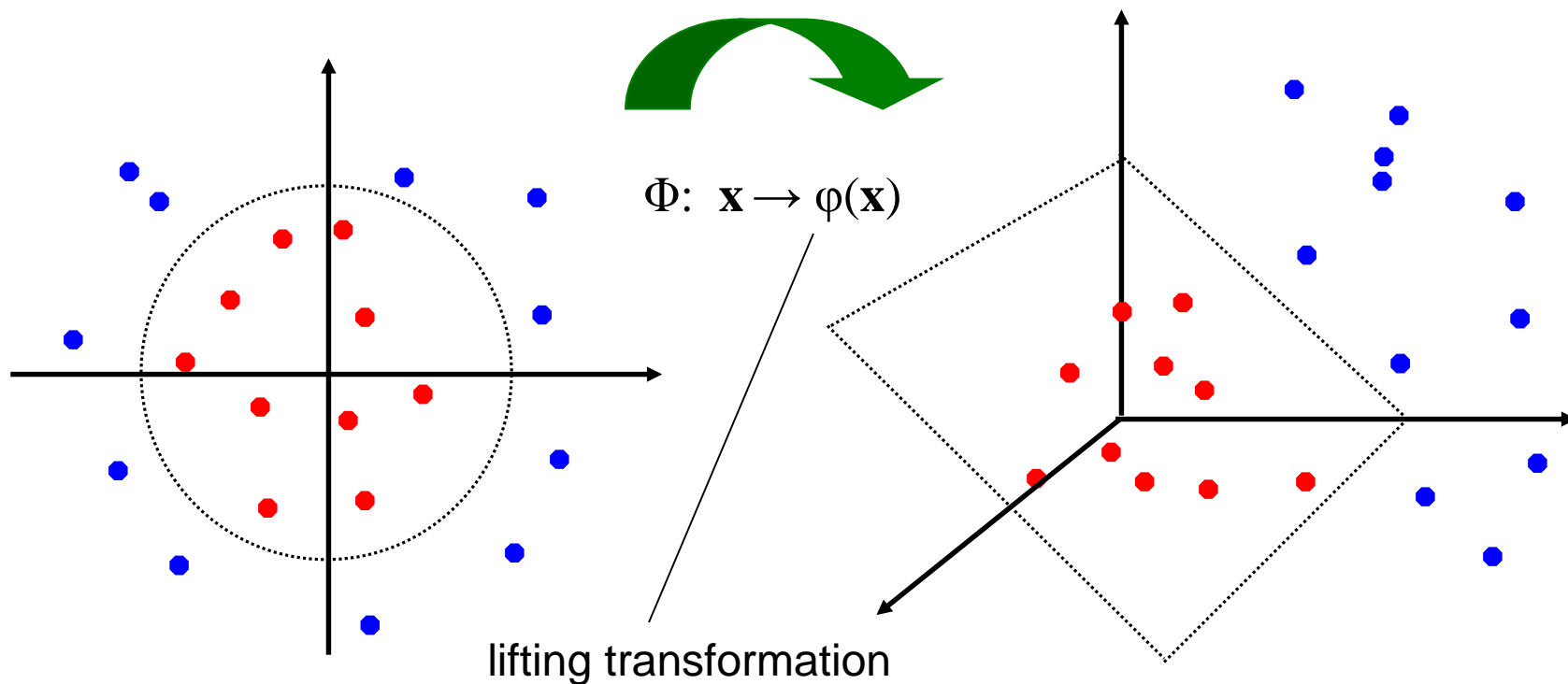


- Map it to a higher-dimensional space:



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:

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

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

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

NOTE:

- It is not required to compute $\boldsymbol{\varphi}(\mathbf{x})$ explicitly:
- The kernel must satisfy the “Mercer inequality”

Kernels for bags of features

- Histogram intersection kernel:

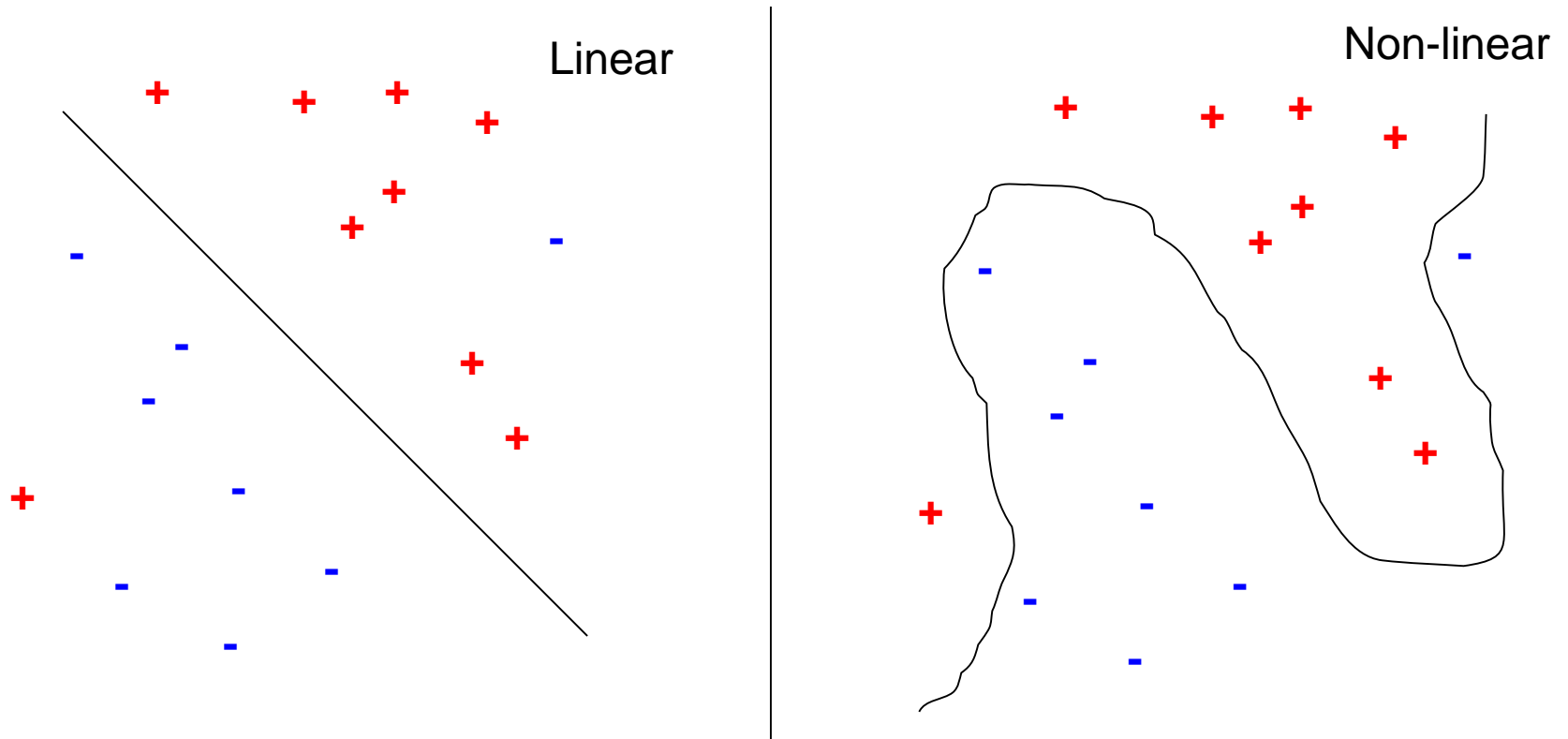
$$I(h_1, h_2) = \sum_{i=1}^N \min(h_1(i), h_2(i))$$

- Generalized Gaussian kernel:

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

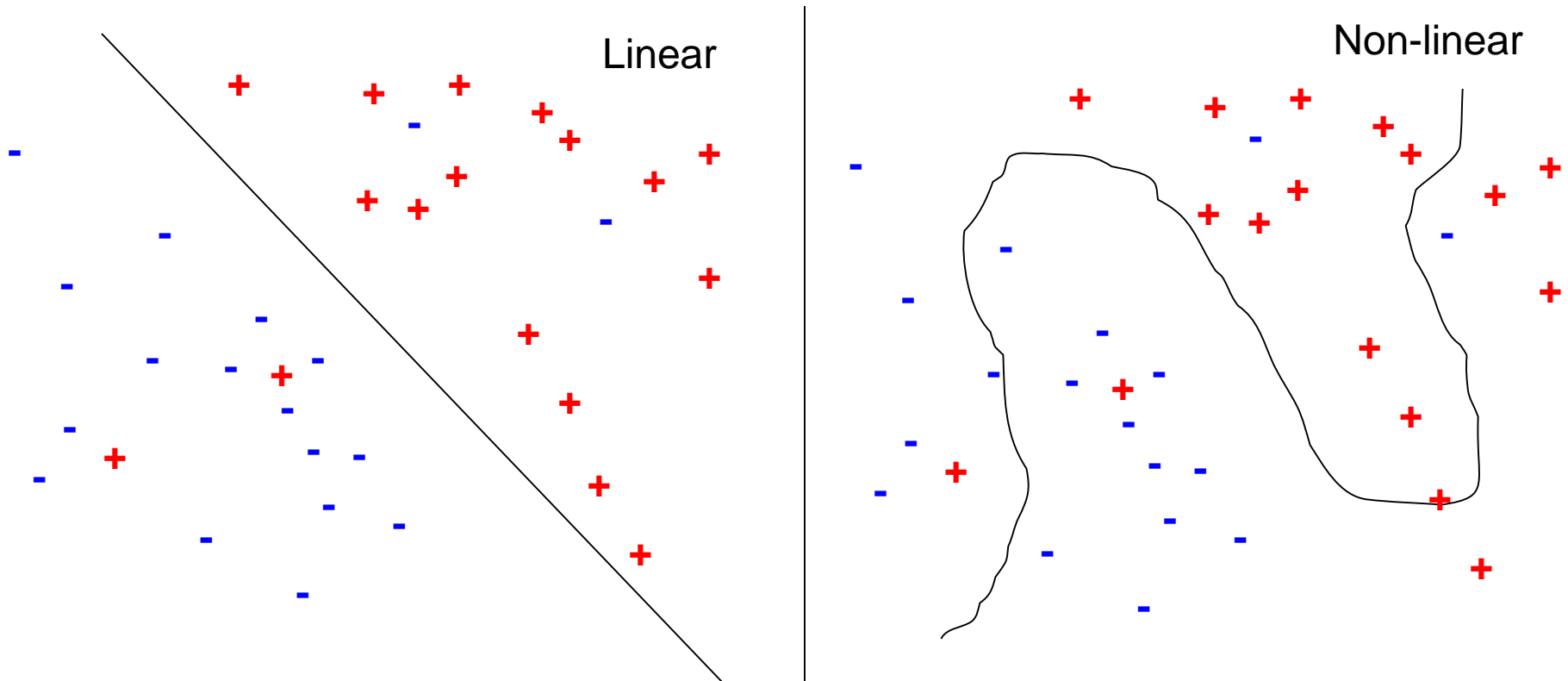
- D can be Euclidean distance, χ^2 distance etc...

Which classifier to use?



Which classifier to use?

Let's add more training data!



A more complex model can over fit the data if these are not enough!

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

SVMs: Pros and cons

- Pros

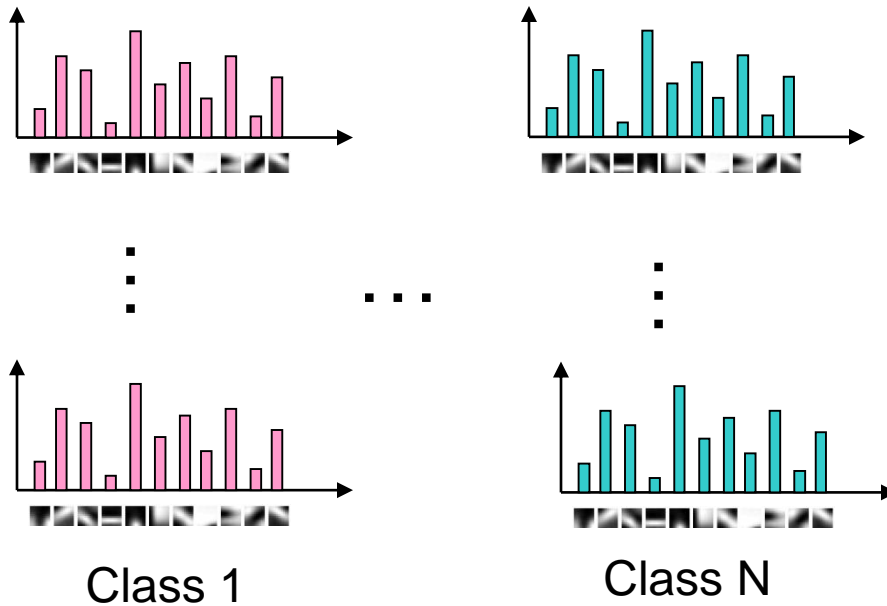
- Many publicly available SVM packages:
<http://www.kernel-machines.org/software>
- Kernel-based framework is very powerful, flexible
- SVMs work very well in practice, even with very small training sample sizes

- Cons

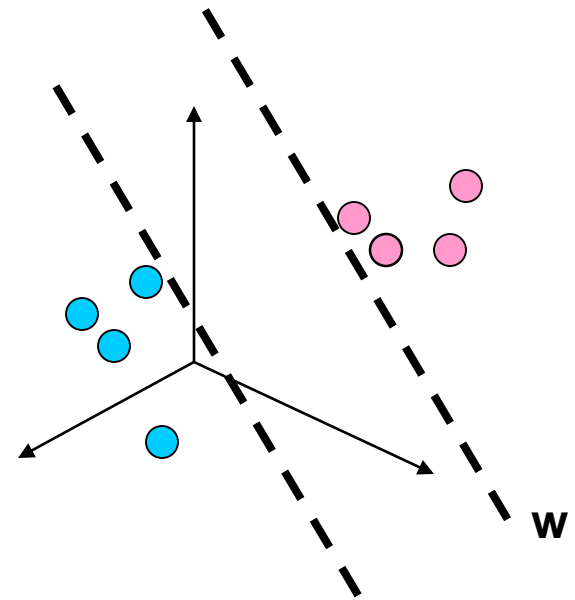
- No “direct” multi-class SVM, must combine two-class SVMs
- Computation, memory
 - During training time, must compute matrix of kernel values for every pair of examples
 - Learning can take a very long time for large-scale problems

Discriminative classifiers (linear classifier)

category models

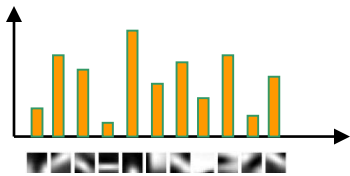


Model space



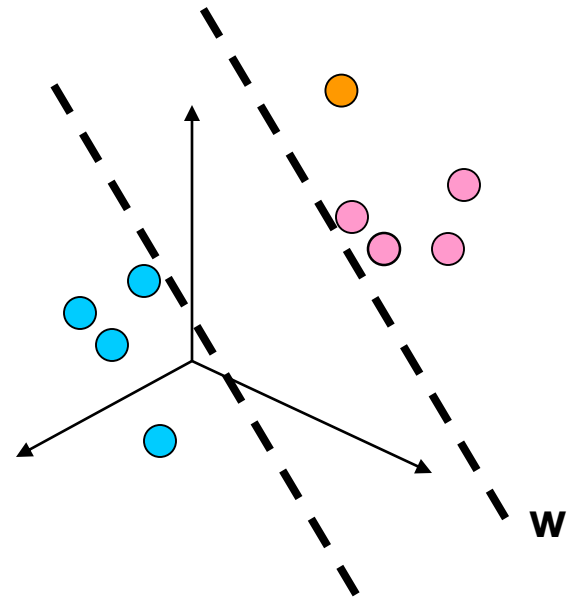
Discriminative classifiers (linear classifier)

Query image



Winning class: pink

Model space



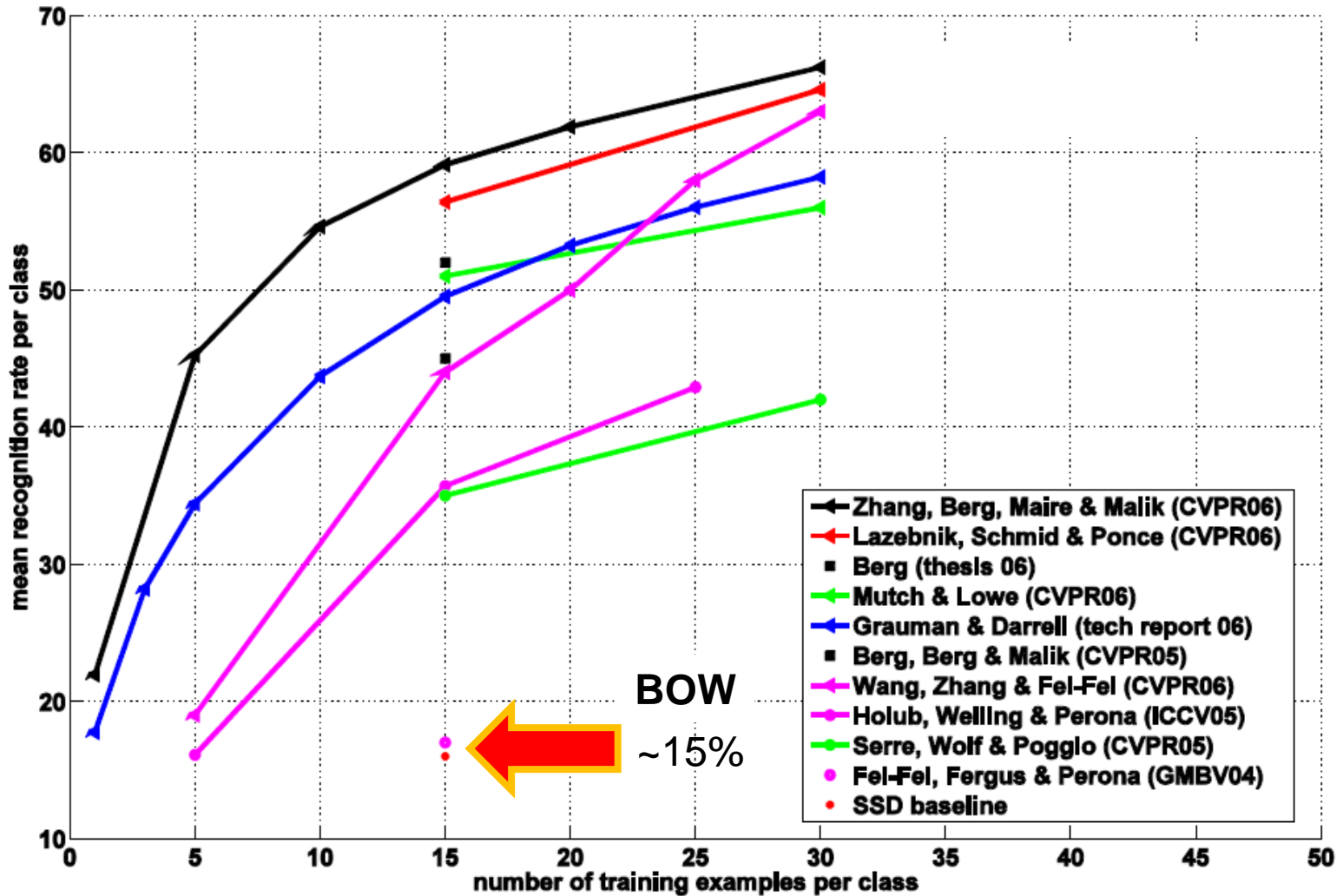
Caltech 101

Fei-Fei et al. (2004)

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

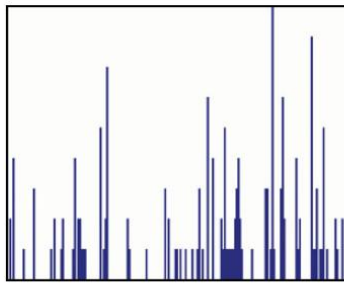
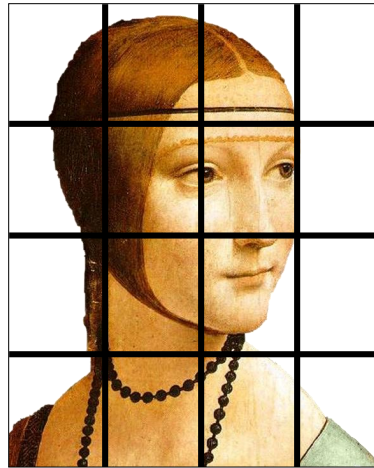


Caltech 101

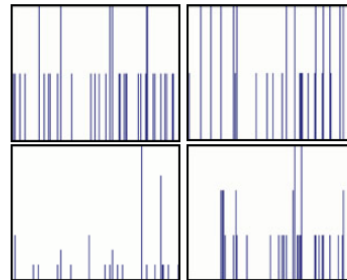


Spatial Pyramid Matching

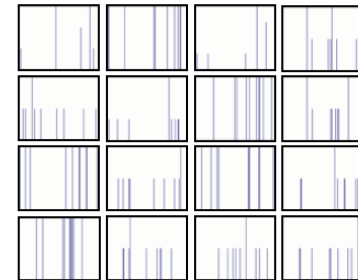
Beyond Bags of Features: Spatial Pyramid Matching for Recognizing Natural Scene Categories. S. Lazebnik, C. Schmid, and J. Ponce.. 2006



level 0



level 1



level 2

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

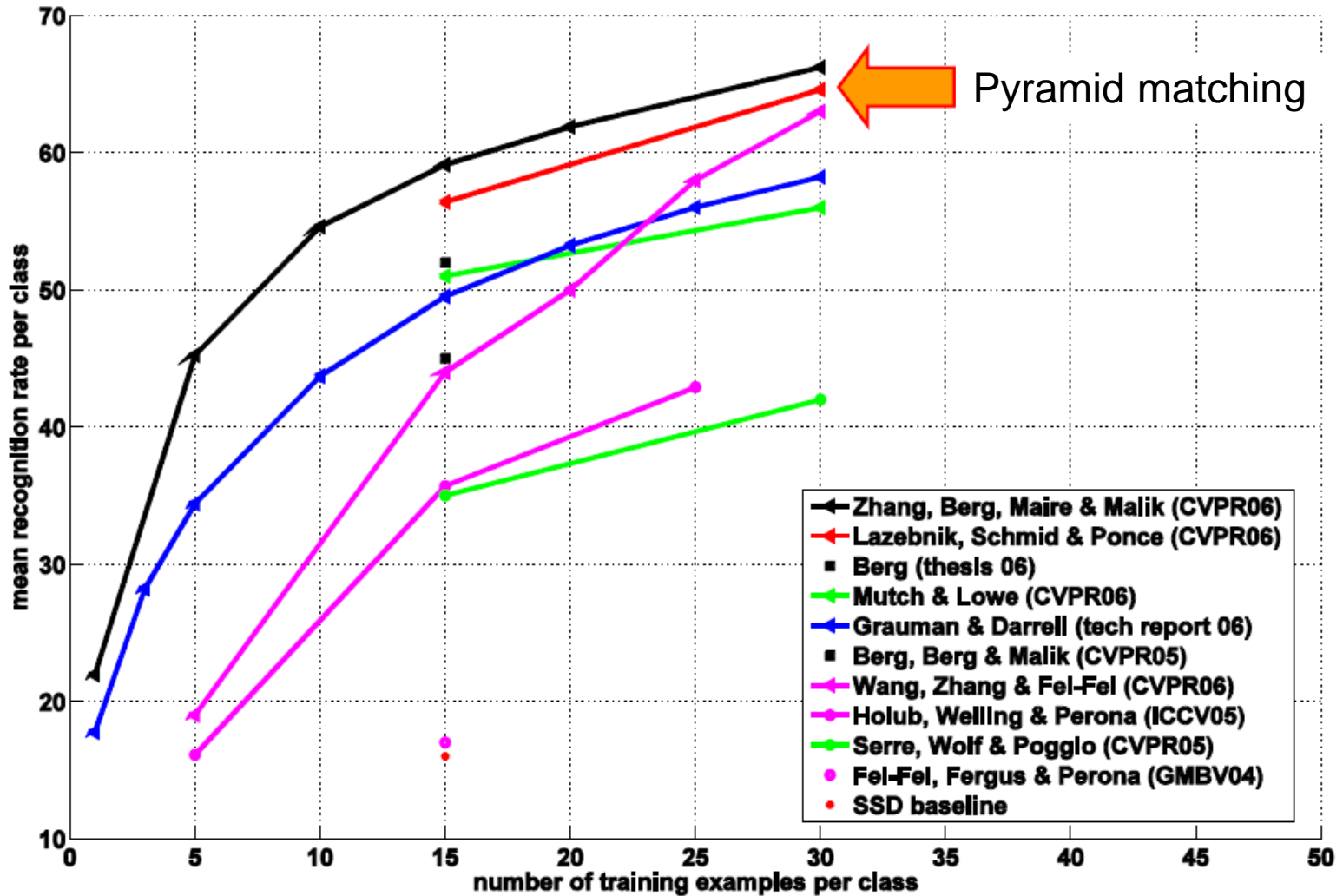
$$I(h_1, h_2) = \sum_{i=1}^N \min(h_1(i), h_2(i))$$

Caltech 101

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

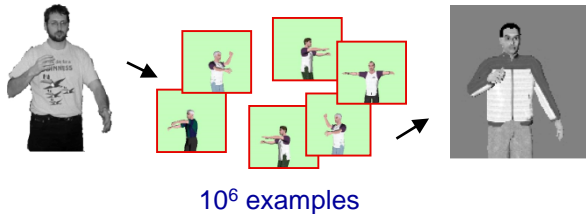
	Weak features (16)		Strong features (200)	
Level	Single-level	Pyramid	Single-level	Pyramid
0	15.5 \pm 0.9		41.2 \pm 1.2	
1	31.4 \pm 1.2	32.8 \pm 1.3	55.9 \pm 0.9	57.0 \pm 0.8
2	47.2 \pm 1.1	49.3 \pm 1.4	63.6 \pm 0.9	64.6 \pm 0.8
3	52.2 \pm 0.8	54.0 \pm 1.1	60.3 \pm 0.9	64.6 \pm 0.7

Caltech 101



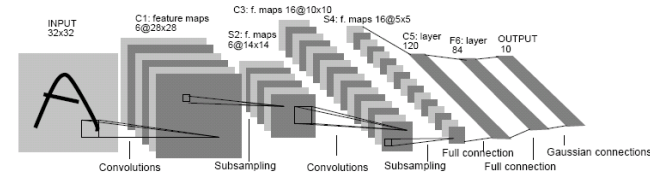
Discriminative models

Nearest neighbor



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

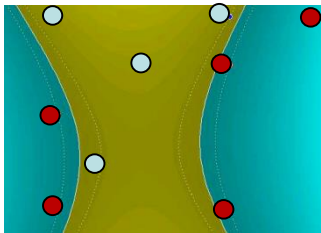
Neural networks



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

...

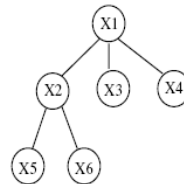
Support Vector Machines



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

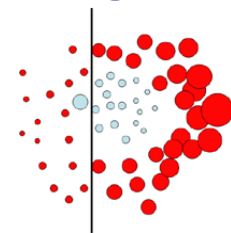
Latent SVM

Structural SVM



Felzenszwalb 00
Ramanan 03...

Boosting



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

Next Lecture

Bag of words models for object recognition and classification

- Generative methods