

Recognition and Classification in Images and Video

203.4780

http://www.cs.haifa.ac.il/~rita/visual_recog_course/Visual%20Recognition%20Spring%202014.htm

Introduction

The course is based on UT-Austin course: Special Topics in Computer Vision,
by Kristen Grauman.

Course overview

- Graduate course in computer vision.
- We will survey and discuss vision papers relating to object and activity recognition and scene understanding.
- The goal: understand classical and modern approaches to some important problems, analyzing their strengths and weaknesses, and identifying interesting open questions.

Requirements

- Paper review each week
- Participating in discussions
- Presenting one topic in a class (in pairs)
- Programming project (in pairs)
- Presentations are due one week before the slot your presentation is scheduled.
 - you will need to read the papers, create slides one week before the date you are signed up
 - **requires my approval for your presentation.**

Paper review

- Submit review for one of the assigned papers (with *) each week, except for the week of your presentation.
- **Reviews are due by 9 PM on the night before class (Sunday).**
- Email reviews to me, pasting the text directly into your mail (no attachments). Include [4738] in the subject header.

The review should address:

- Give a summary of the paper in your own words (very brief, 2-3 sentences)
- What is the main contribution of the paper?
- What are the primary strengths and weaknesses of the paper?
- How convincing are the experiments? If something specific is lacking, what should have been tested?
- Describe one specific way in which the work could be extended (bonus).
- Additional comments, including unclear points, connections you see between the papers.

Paper presentation

- Each team (of 2 students) will give a presentation in class covering 2 papers on a topic selected from the course syllabus list.
- The talk should be well-organized and polished, sticking to about 40 minutes (20 min. each student).
- Run through it beforehand and check the time (a good rule of thumb: a minute per slide).

Presentation should cover

- Clear statement of the problem
- Why the problem is interesting, important, difficult?
- Key technical ideas, how they work, main contributions, strengths and weaknesses
- Evaluation, summary of key experiments and data
- Open issues raised in the papers, likely extensions

Presentation guidelines

- Try to use applications to motivate the work when possible, and look for visual elements (images, videos) to put in the presentation.
- Check out the webpages linked on the class webpage, and also look at authors webpages for supplementary materials.
- It's ok to grab a few slides from conference talks etc. when available, but be sure to **clearly cite the source on each slide** that is not your own.

Project

- A project could be built around one of the topics of the syllabus
- It should be done with a partner.
- Experimental evaluation should be done on a benchmark data set (provided in the course page)
- You can use papers provided as an additional reading to choose your project.
- Initial project proposals will be due before the **middle of the term.**

Suggested Projects

- Object recognition using bag-of-words representation and discriminative learning.
- Face detection
- Pedestrian detection
- Face recognition
- Saliency in images
- Action recognition
- Image retrieval

Grading

- 30% participation, includes
 - attendance,
 - in-class discussions,
 - paper reviews.
- 40% presentations, includes
 - drafts submitted one week prior,
 - in-class presentation.
- 30% final project, includes
 - implementation,
 - presentation,
 - final report.

Syllabus

A. Recognizing specific objects

Global Features:

- Linear Subspaces
- Detection as a binary decision

Local Features:

- Local features, matching for object instances
- Visual Vocabularies and Bag of Words

Region-based Features:

- Mid-Level Representations

B. Beyond Single objects (using additional information)

- Saliency
- Attributes
- Context

C. Scalability problems

- Scaling with the large number of categories
- Large-scale search

D. Action recognition in video and images

Syllabus

A. Recognizing specific objects

Global Features:

- Linear Subspaces
- Detection as a binary decision

Local Features:

- Local features, matching for object instances
- Visual Vocabularies and Bag of Words

Region-based Features:

- ~~Mid-Level Representations~~

B. Beyond Single objects (using additional information)

- ~~Saliency~~
- Attributes
- Context

C. Scalability problems

- Scaling with the large number of categories
- Large-scale search

D. Action recognition in video and images

Two of three

Object Recognition

So what does object recognition involve?



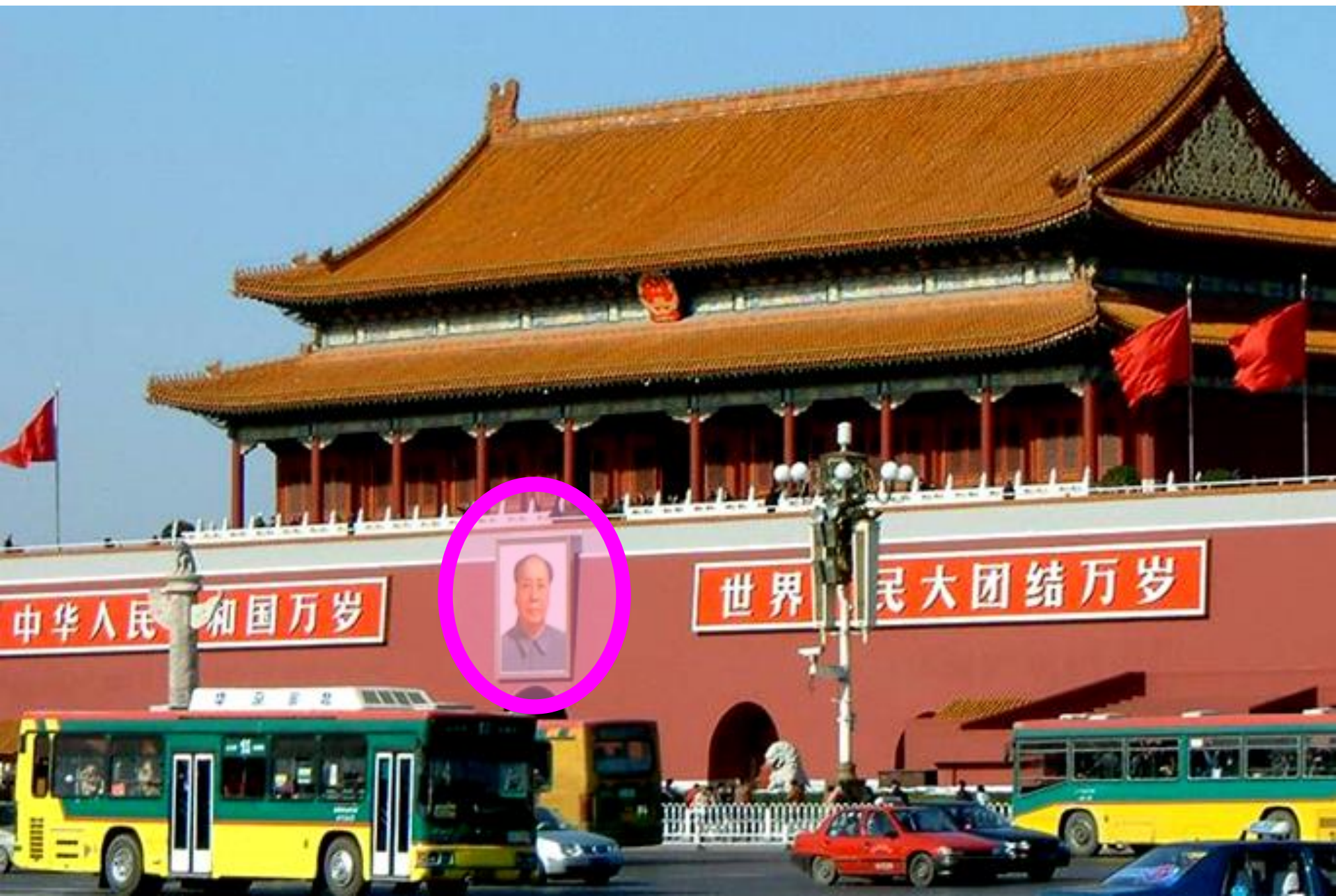
Verification: is that a bus?



Detection: locate the cars in the image



Identification: is that a picture of Mao?



Object categorization



sky

building

flag

face

banner

wall

street lamp

bus

bus

cars

Challenges 1: view point variation

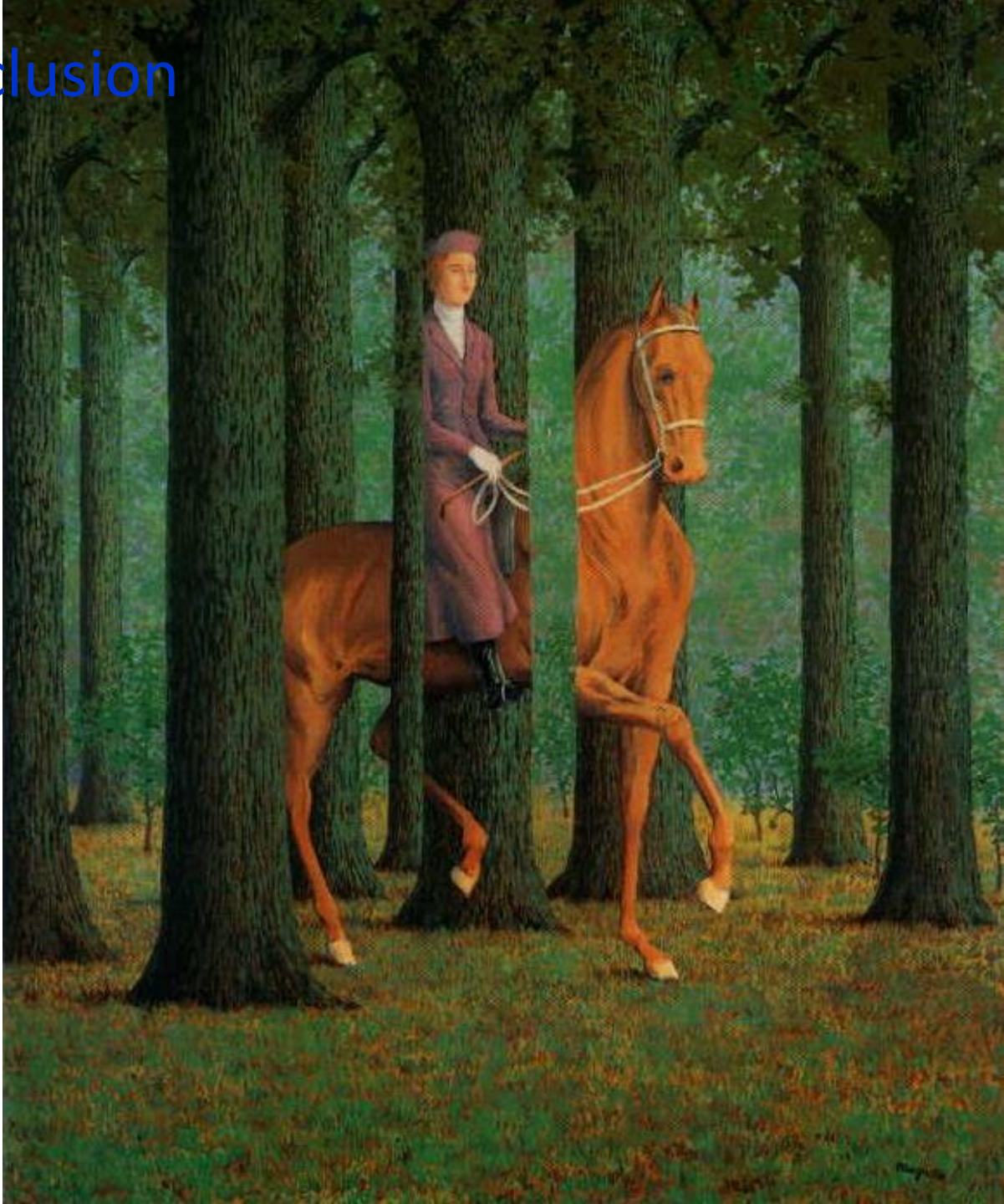


Michelangelo 1475-1564

Challenges 2: illumination



Challenges 3: occlusion

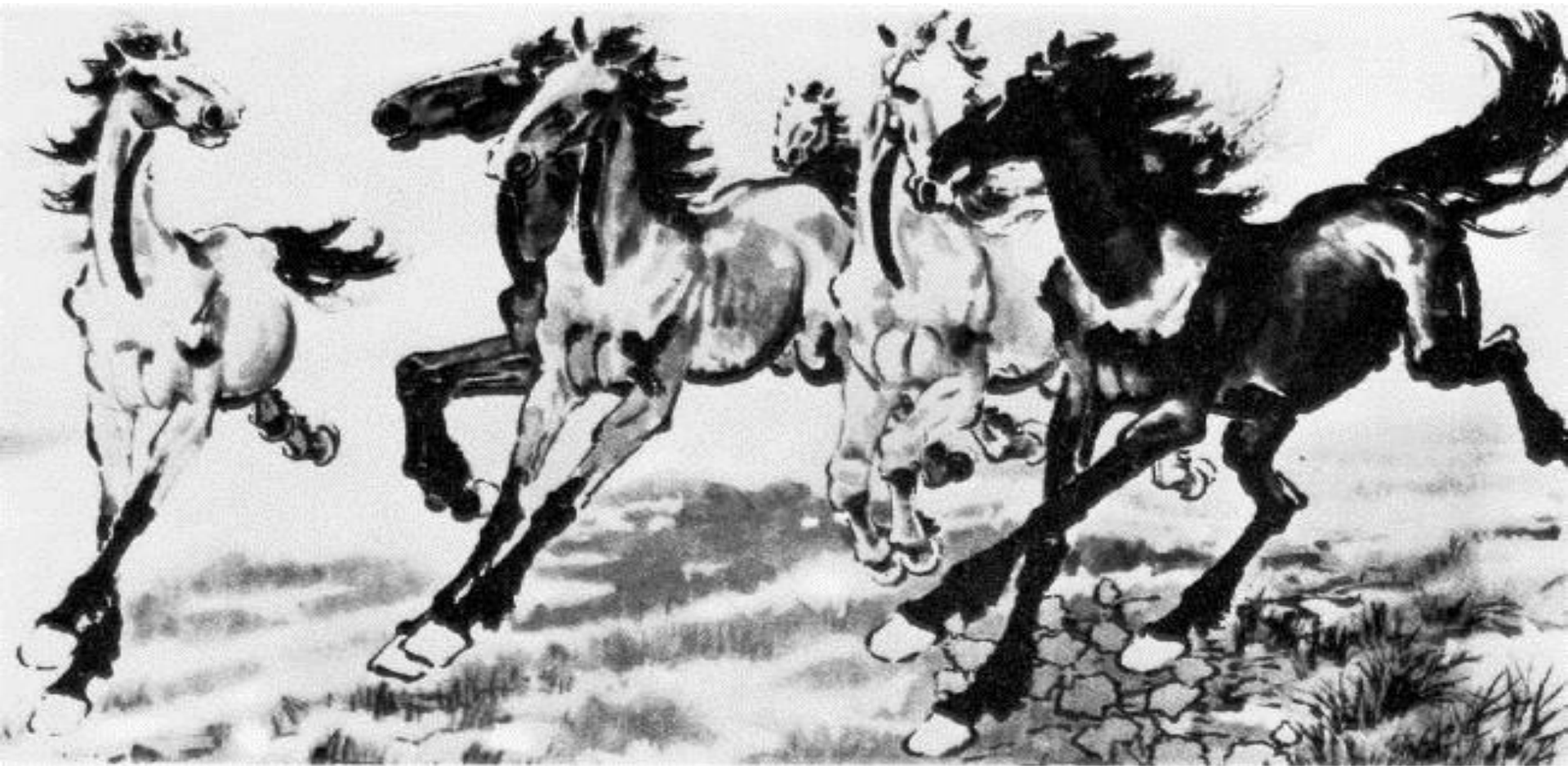


Magritte, 1957

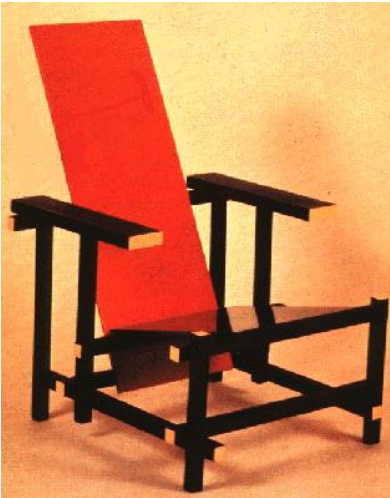
Challenges 4: scale



Challenges 5: deformation



Challenges 7: intra-class variation



Recognition Steps

Preprocessing

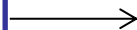
Train Images



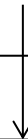
Features



Image
representation



Object Model



Recognition

Test Images



Features



Image
representation



Output:

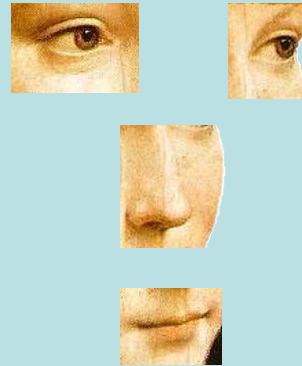
Yes/No
Class label
Position of the object

Features

Global

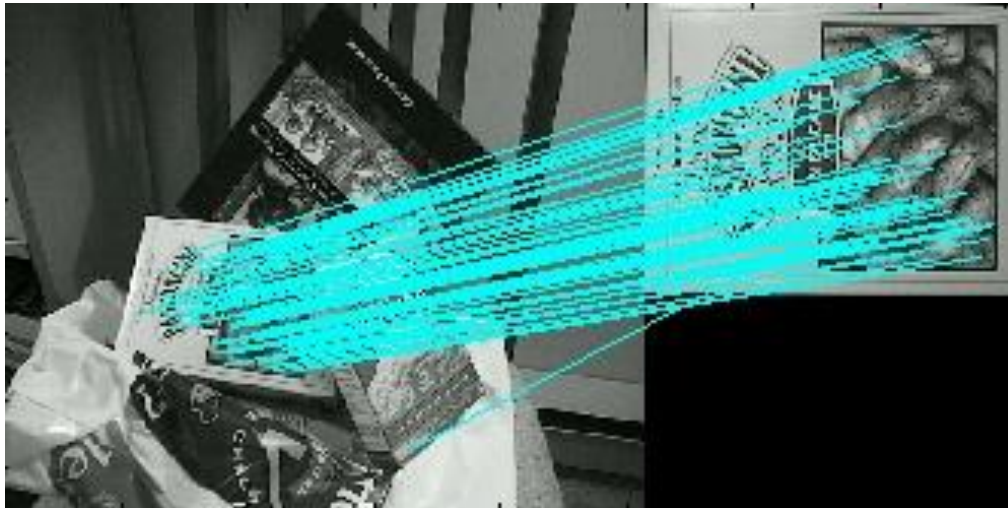


Part-based



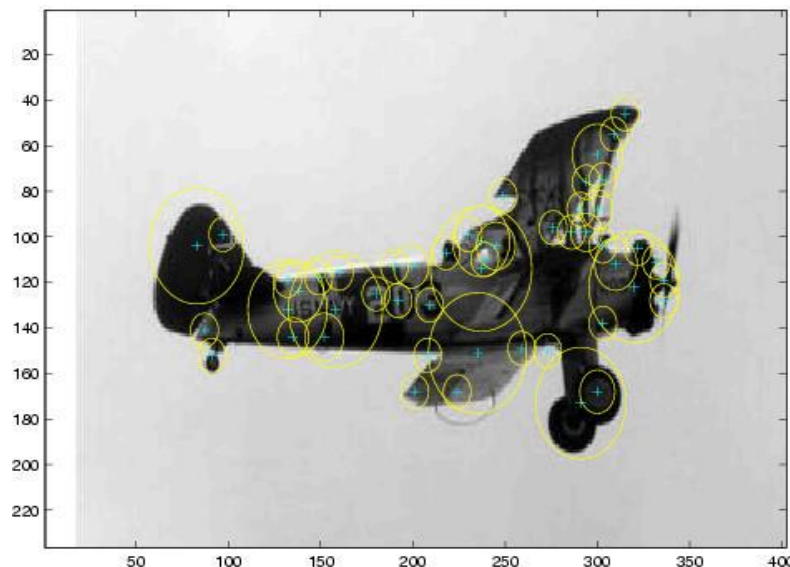
Recognition using local features

1. Find invariant local features
2. Match them with the model features
3. Vote on global geometric transformation



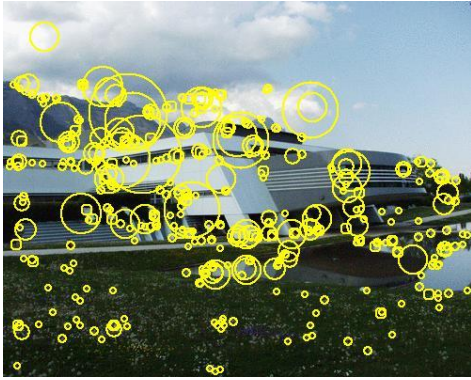
Interest Point Detectors

- Basic requirements:
 - Sparse
 - Informative
 - Repeatable
- Invariance
 - Rotation
 - Scale (Similarity)
 - Affine



Popular Detectors

Scale Invariant



Harris-Laplace



Difference of Gaussians



Laplace of Gaussians



Scale Saliency (Kadir-Braidy)

Affine Invariant



Harris-Laplace Affine



Difference of Gaussians
Affine



Laplace of Gaussians
Affine



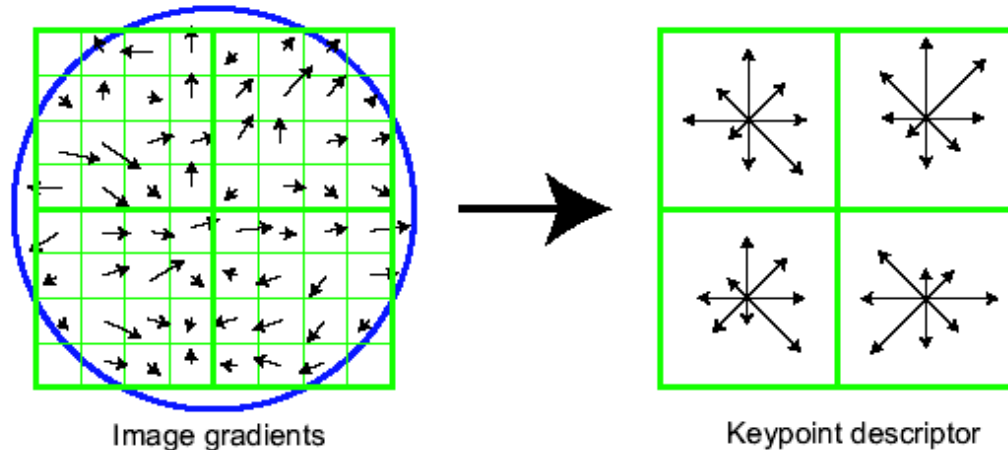
Affine Saliency (Kadir-Braidy)

Representation of appearance: Local Descriptors

- Invariance
 - Rotation
 - Scale
 - Affine
- Insensitive to small deformations
- Illumination invariance
 - Normalize out

SIFT – Scale Invariant Feature Transform

- Descriptor overview:
 - Determine **scale** (by maximizing DoG in scale and in space), **local orientation** (as the dominant gradient direction). Use this scale and orientation to make all further computations invariant to scale and rotation.
 - Compute **gradient orientation histograms** of several small windows (128 values for each point)
 - Normalize the descriptor to make it invariant to intensity change

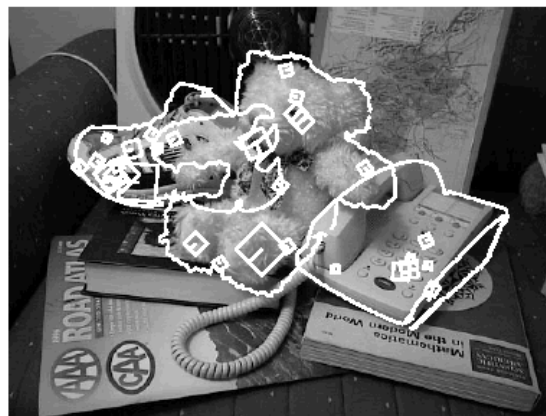
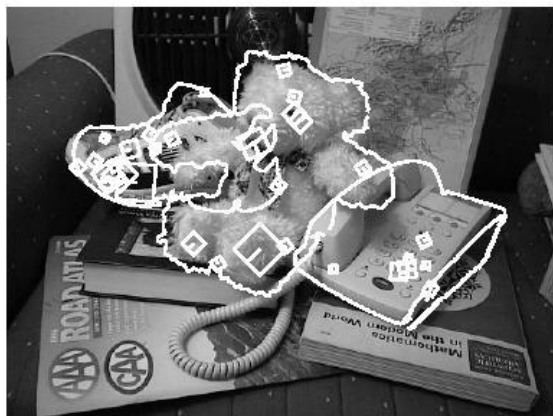


Recognizing Specific Objects

Learned models of local features, and got object outline from

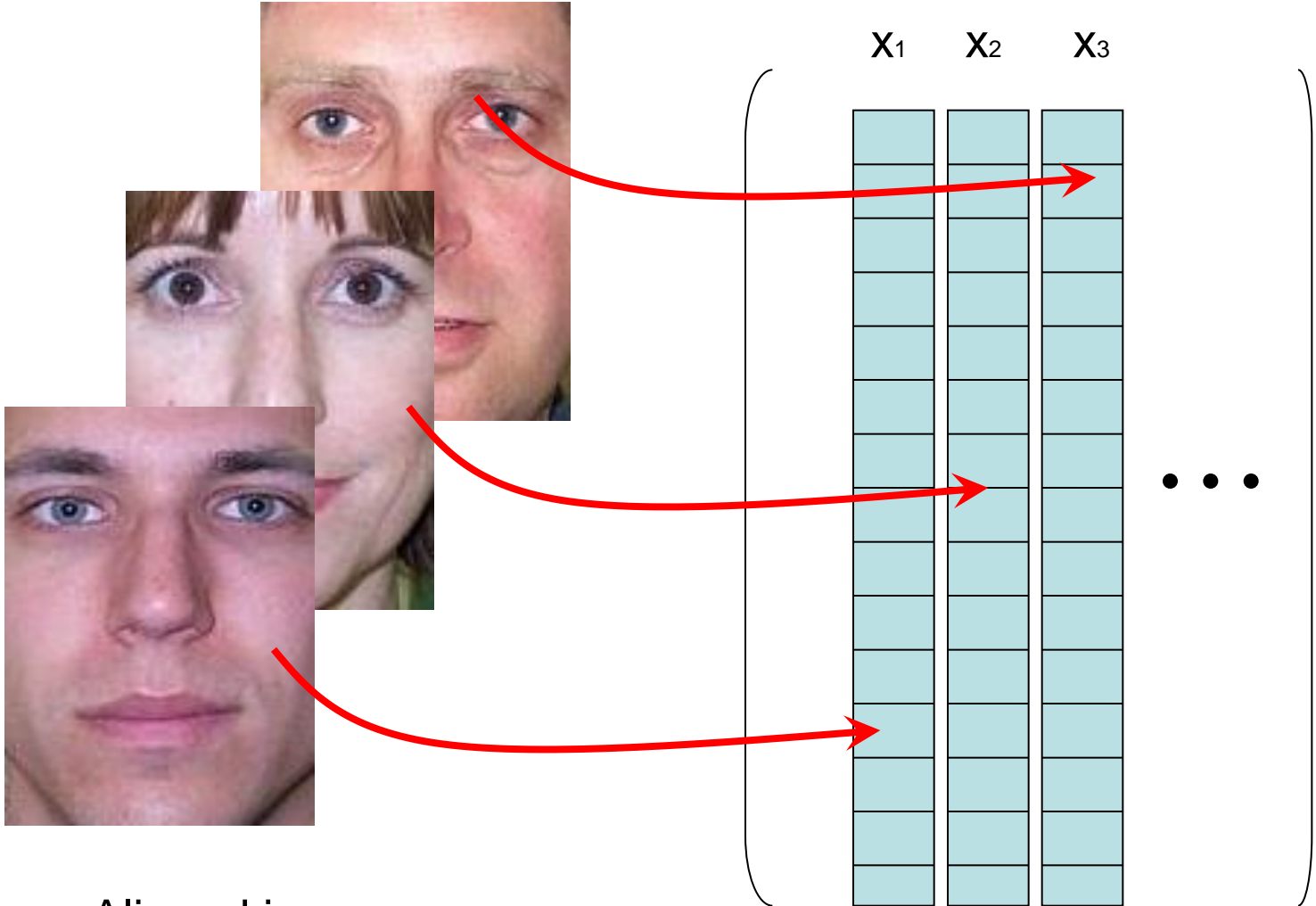


Objects may then be found under occlusion and 3D rotation



Recognition using global representations

...



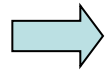
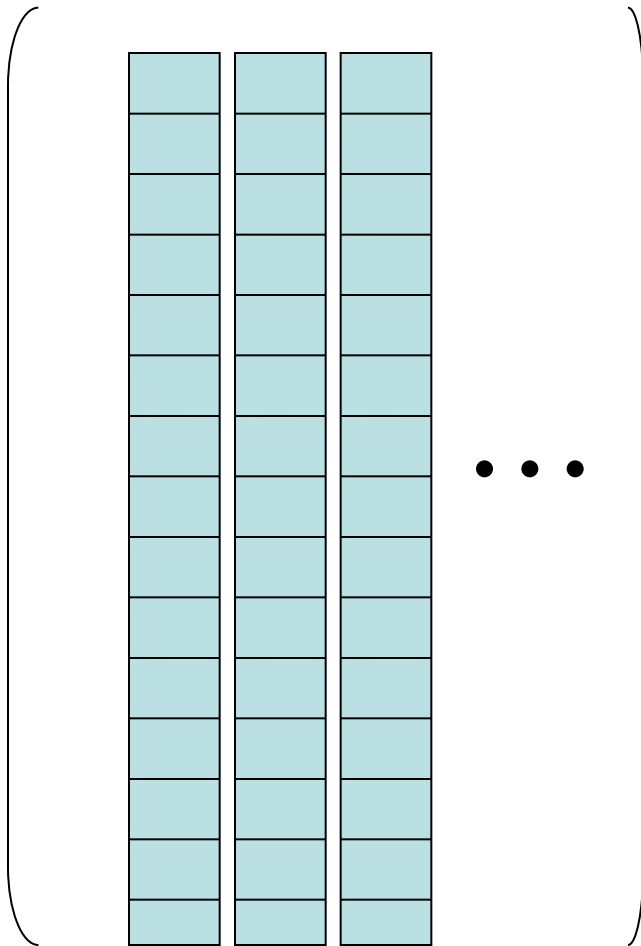
Aligned images

Vectors in high-dimensional space

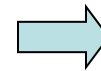
Global Approaches

Vectors in high-dimensional
space

X_1 X_2 X_3



Training
Involves some
dimensionality
reduction



Classifier

Dimensionality Reduction

- Inputs (**high dimensional**)
 $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ points in \mathbb{R}^d
- Outputs (**low dimensional**)
 $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n$ points in \mathbb{R}^k ($k \ll d$)
- Goal:
combine old features \mathbf{x} to create new features \mathbf{y}

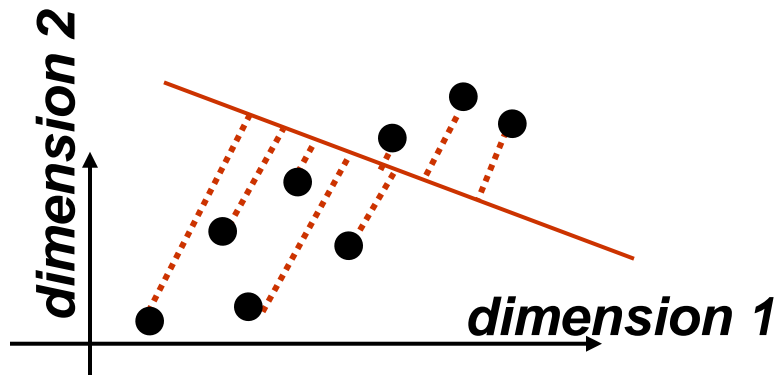
Linear Methods

$$\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_d \end{bmatrix} \Rightarrow \mathbf{W} \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_d \end{bmatrix} = \begin{bmatrix} \mathbf{w}_{11} & \cdots & \mathbf{w}_{1d} \\ \vdots & & \vdots \\ \mathbf{w}_{k1} & \cdots & \mathbf{w}_{kd} \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_d \end{bmatrix} = \begin{bmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_k \end{bmatrix} \quad \text{with } k < d$$

- Principle component analysis PCA
- Fisher Linear Discriminant (FLD)

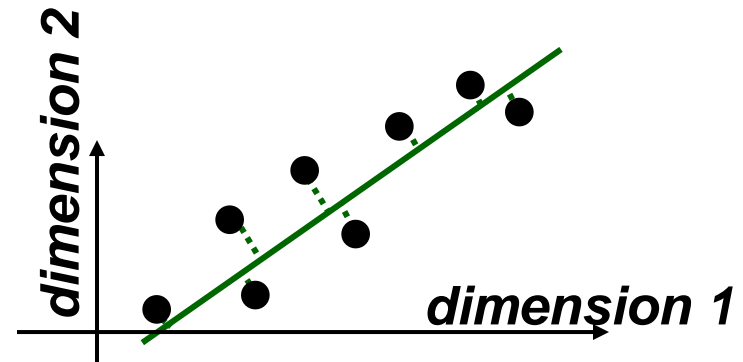
Principle Component Analysis

Good representation



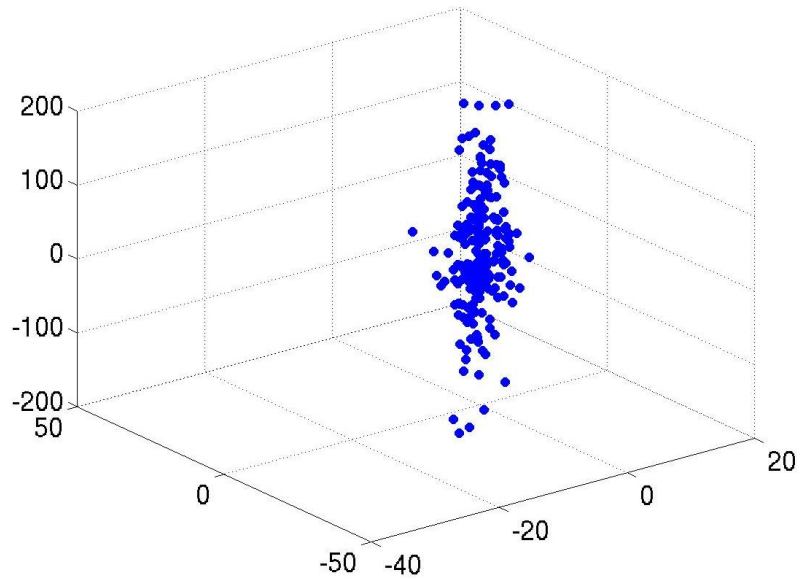
the projected data has a fairly large variance, and the points tend to be far from zero.

Poor representation

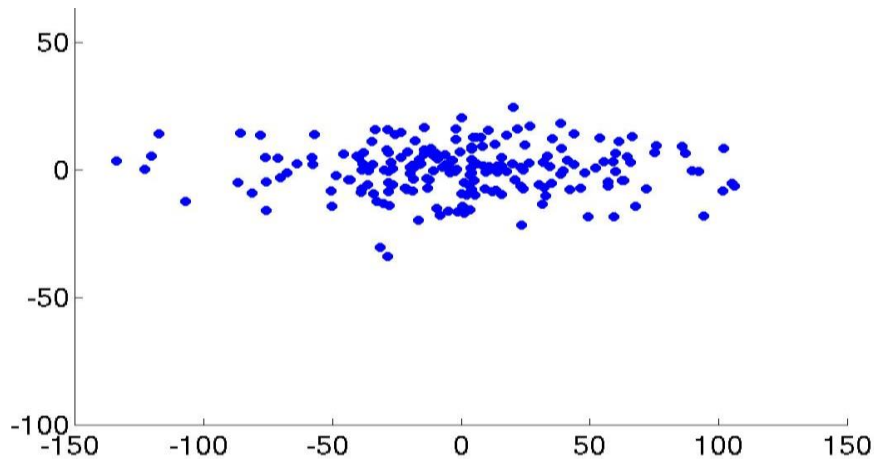


the projections have a significantly smaller variance, and are much closer to the origin.

PCA

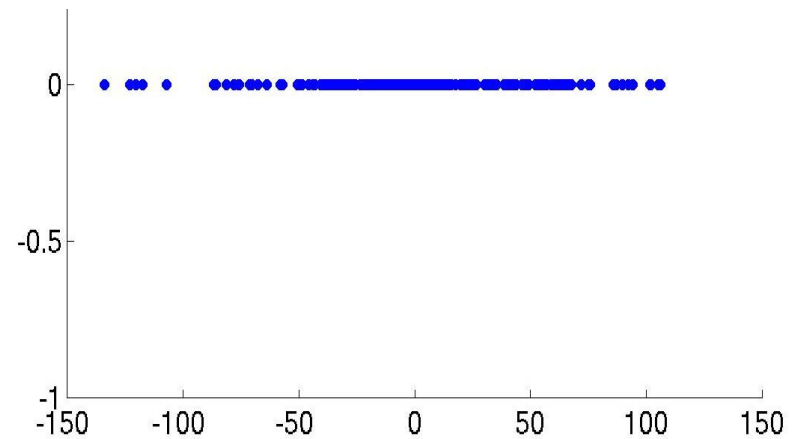


best 2D approximation



- Seek most accurate data representation in a lower dimensional space.
- The good direction/subspace to use for projection lies in the direction of largest variance.

best 1D approximation



Maximum Variance Subspace

- Assume inputs are centered: $\sum_{i=1}^n x_i = 0$
- Given a unit vector u and a point x , the length of the projection of x onto u is given by $x^T u$
- Maximize projected variance:

$$\begin{aligned}\text{var}(y) &= \frac{1}{n} \sum_i (x_i^T u)^2 = \frac{1}{n} \sum_i u^T x_i x_i^T u \\ &= u^T \left(\frac{1}{n} \sum_i x_i x_i^T \right) u\end{aligned}$$

1D Subspace

- Maximizing $u^T C u$ subject to $\|u\| = 1$
where $C = \frac{1}{n} \sum_i x_i x_i^T$ is the empirical covariance matrix of the data, gives the principle eigenvector of C .

d-dimensional subspace

- to project the data into a d-dimensional subspace ($k \ll d$), we should choose u_1, \dots, u_k to be the top k eigenvectors of C .
- u_1, \dots, u_k now form a new, orthogonal basis for the data.
- The low dimensional representation of x is given

by

$$y = \begin{bmatrix} u_1^T x \\ u_2^T x \\ \vdots \\ u_k^T x \end{bmatrix}$$

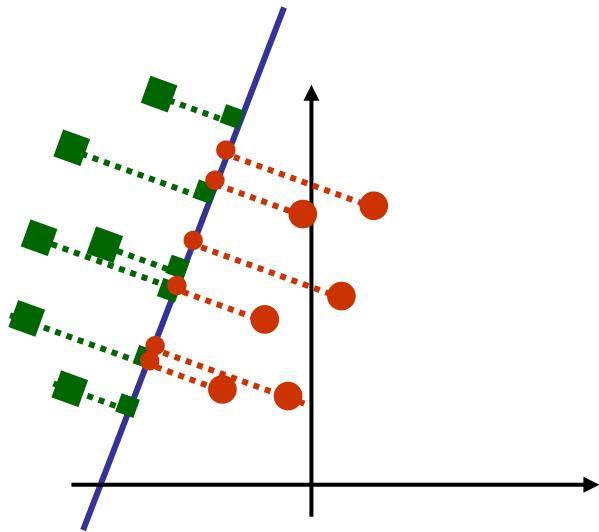
Interpreting PCA

- **Eigenvectors:**
principal axes of maximum variance subspace.
- **Eigenvalues:**
variance of projected inputs along principle axes.
- **Estimated dimensionality:**
number of significant (nonnegative) eigenvalues.

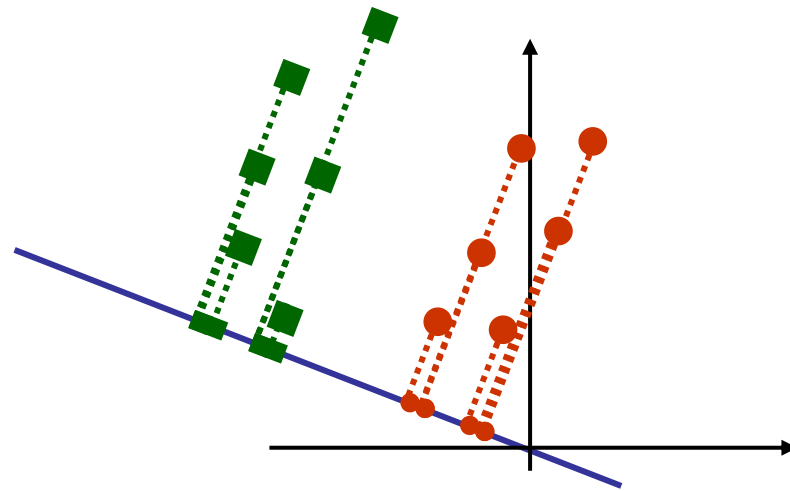
Fisher Linear Discriminant

- Main idea: find projection to a line s.t. samples from different classes are well separated.

Example in 2D



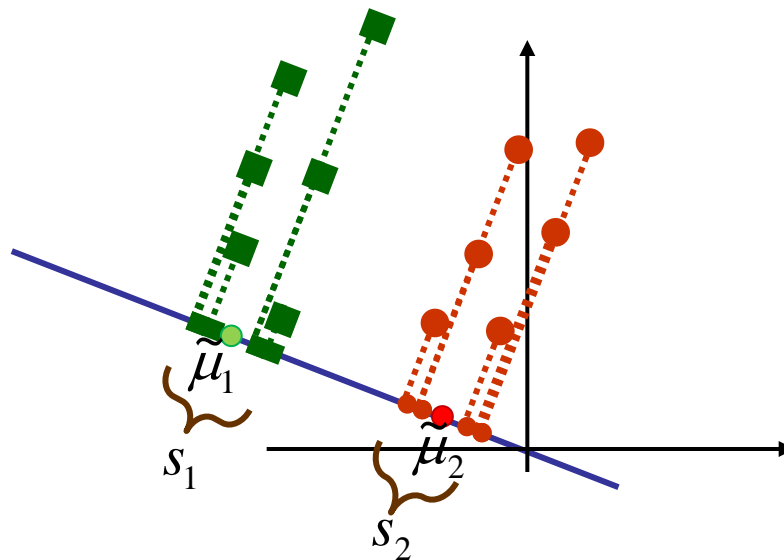
*bad line to project to,
classes are mixed up*



*good line to project to,
classes are well separated*

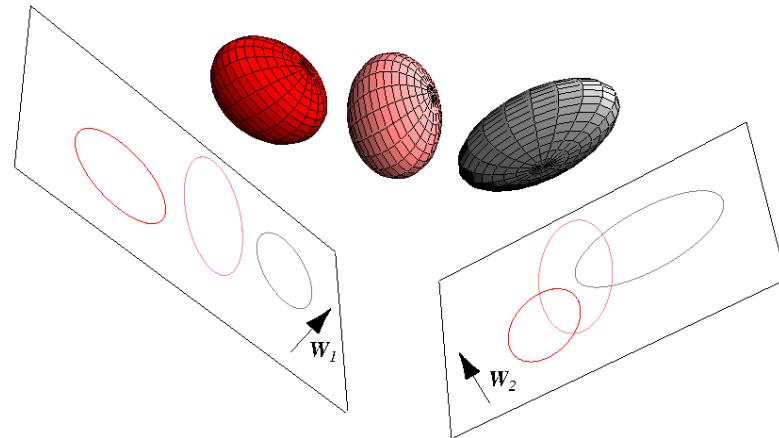
FLD

- Find projection that
 - Maximizes the distance between class means $|\tilde{\mu}_1 - \tilde{\mu}_2|$
 - Minimizes the scatter of the classes S_1, S_2



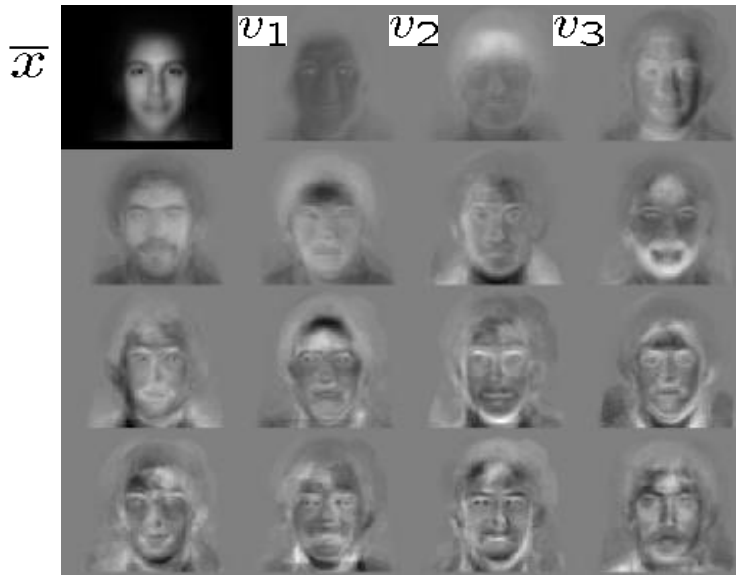
Multiple Discriminant Analysis

- Can generalize FLD to multiple classes
- In case of c classes, can reduce dimensionality to $c-1$ dimensions or less.



Linear Subspaces

Eigenfaces (PCA)



top left image is linear combination of the rest.

Fisherfaces (FLD)



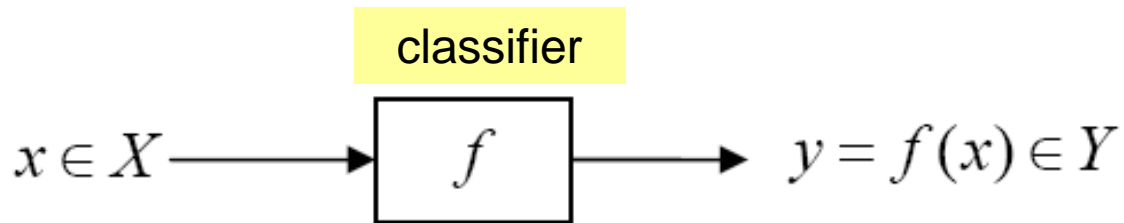
first four Fisherfaces from a set of 100 classes

Classification in Machine Learning

- A **classifier** is a function or an algorithm that maps every possible input (from a legal set of inputs) to a finite set of categories.
- X – **input space**, $x \in X$ **sample** from an input space.
- A typical input space is high-dimensional, for example $x = \{x_1, \dots, x_d\} \in R^d$, $d > 1$. We also call x a **feature vector**.
- Ω is a **finite set of categories** to which the input samples belong: $\Omega = \{1, 2, \dots, C\}$.
- $w_i \in \Omega$ are called **labels**.

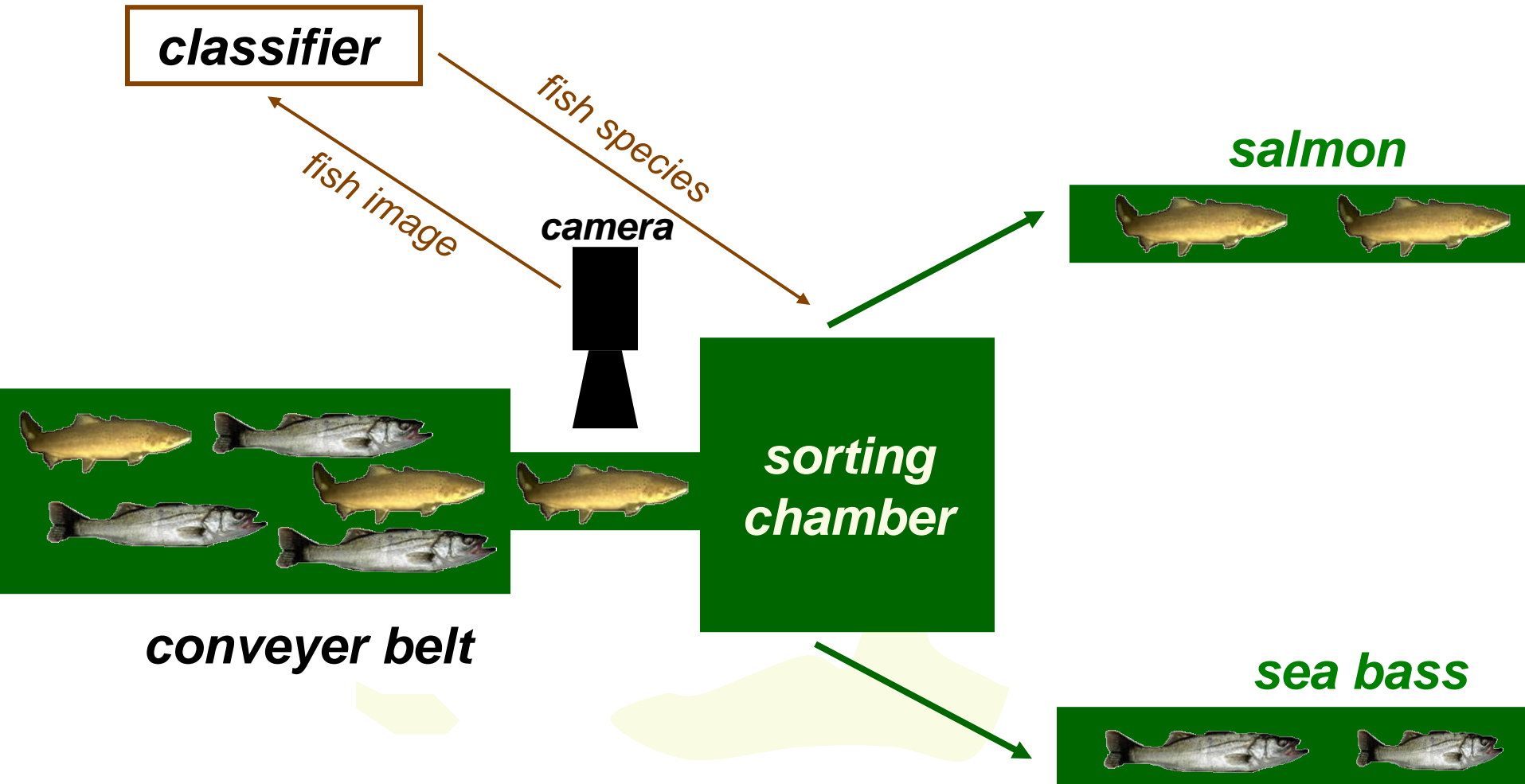
Definition of Classification

- Y is a finite **set of decisions** – the **output set** of the classifier.
- A classifier is a function $f : X \rightarrow Y$



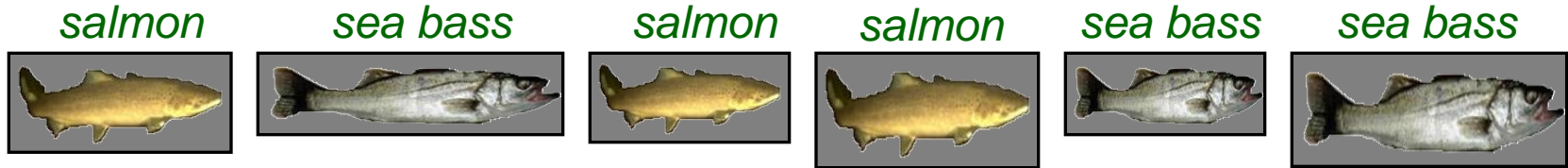
- Classification is also called **Pattern Recognition**.

Toy Application: fish sorting



How to design a PR system?

- **Collect data** and classify by hand



- **Preprocess** by segmenting fish from background



- **Extract** possibly discriminating **features**

- length, lightness, width, number of fins, etc.

- **Classifier design**

- **Choose model**

- **Train classifier** on part of collected data (**training** data)

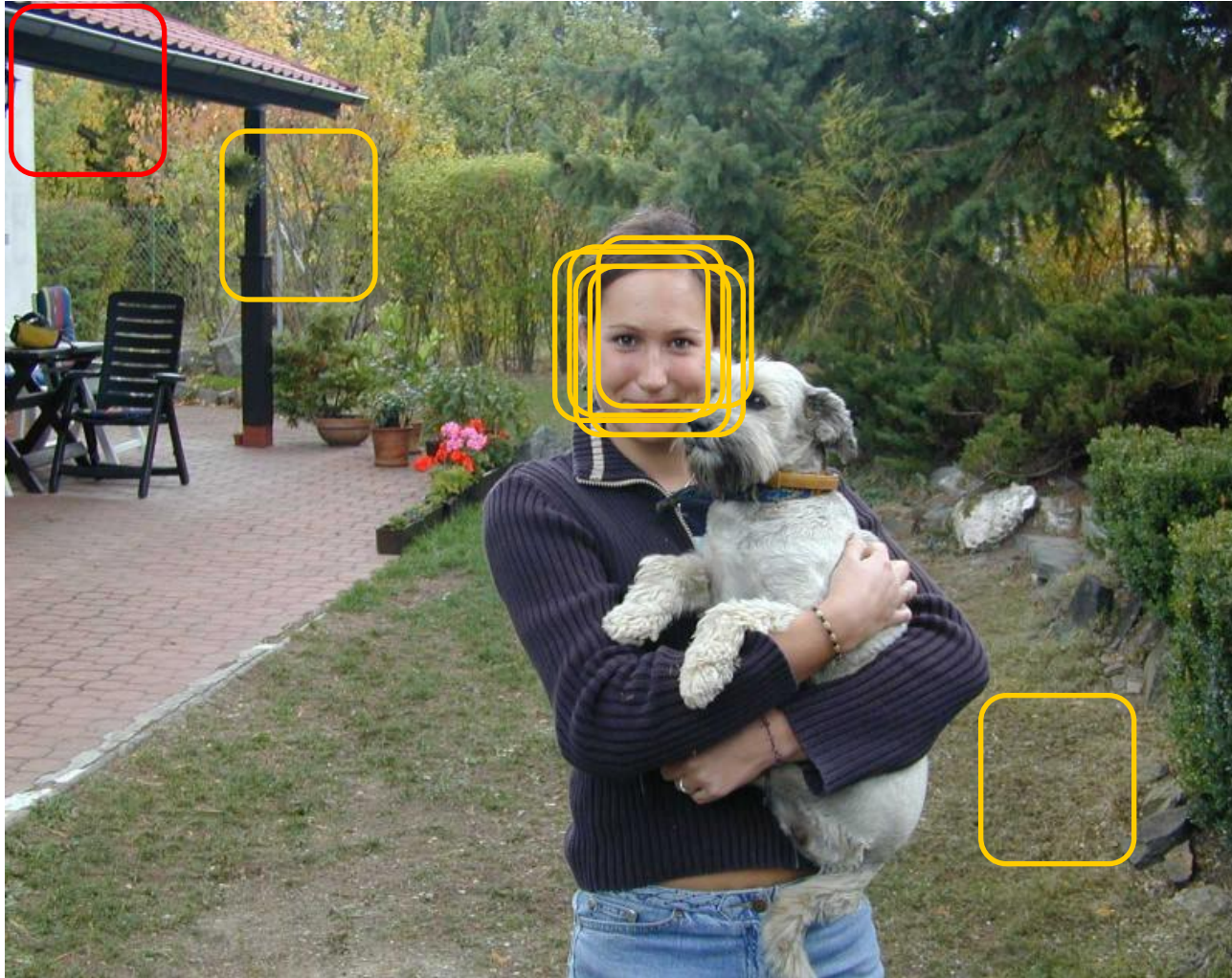
- **Test classifier** on the rest of collected data (**test** data)

i.e. the data not used for training

- Should classify **new** data (new fish images) well

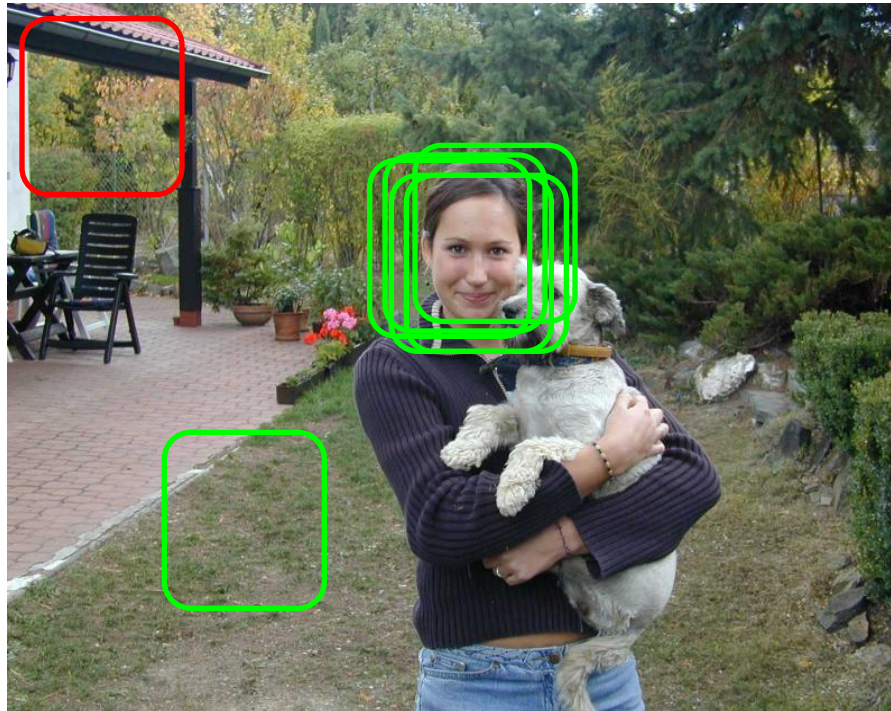
Detection

- Scale / position range to search over



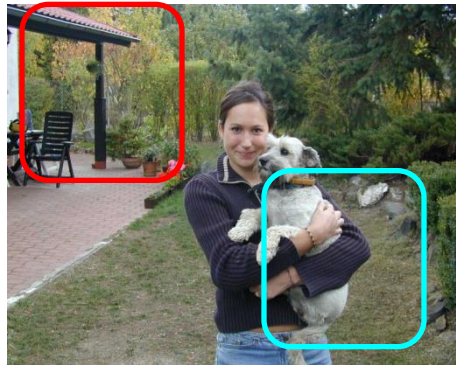
Detection

- Scale / position range to search over



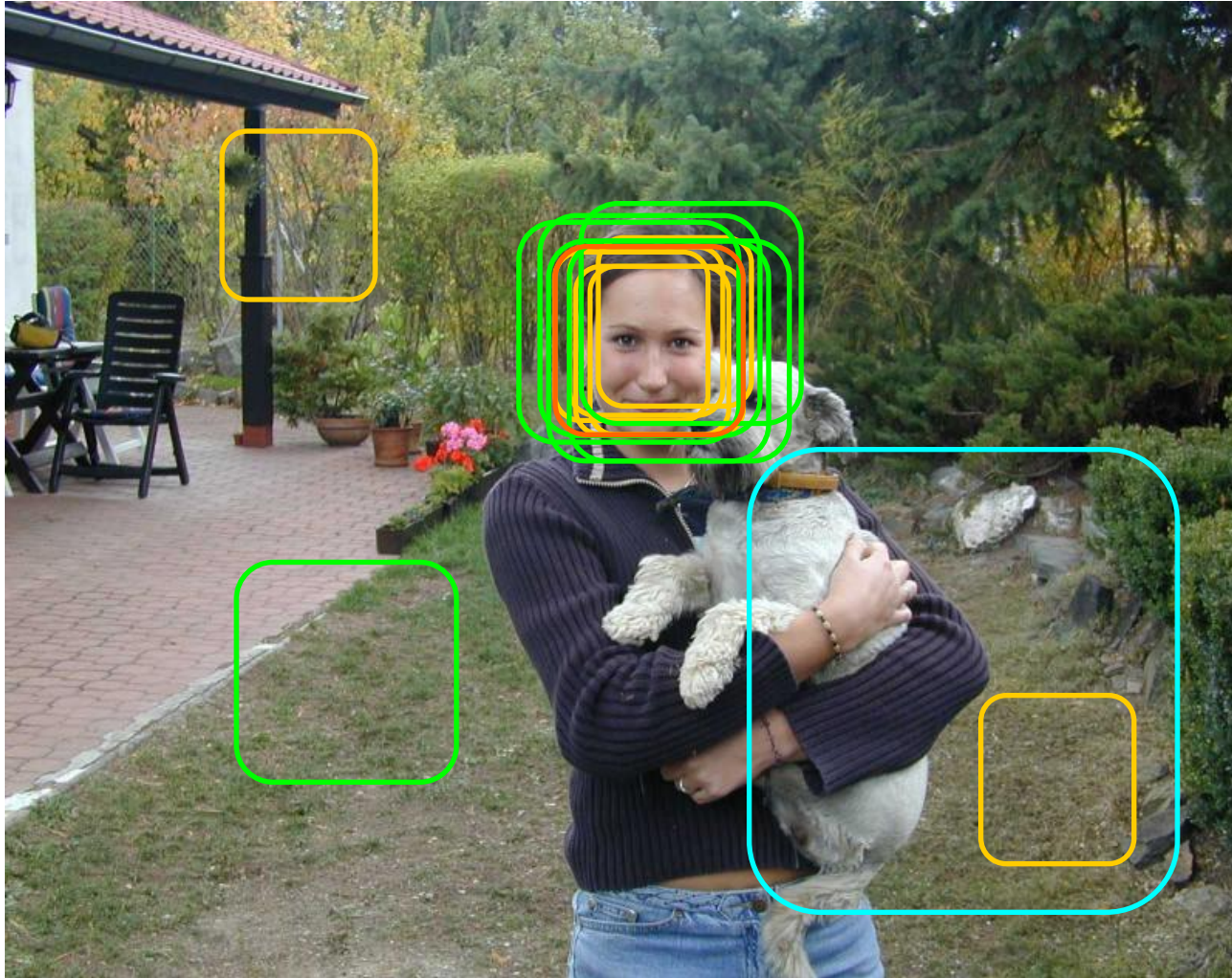
Detection

- Scale / position range to search over



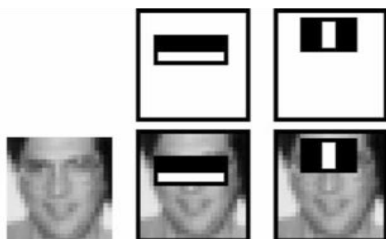
Detection

- Combine detection over space and scale.



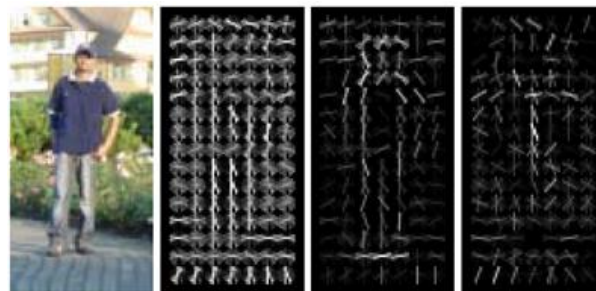
Detection as a binary decision

- Sliding window detection, detection as a binary decision problem.



Boosted Cascade
of Simple Features

Adaboost



Histograms of Oriented
Gradients for Human
Detection

SVM

Boosting, motivation

- It is usually hard to design an accurate classifier which generalizes well
- However it is usually easy to find many “rule of thumb” *weak* classifiers
 - A classifier is *weak* if it is only slightly better than random guessing
- Can we combine several weak classifiers to produce an accurate classifier?
 - Question people have been working on since 1980’s

AdaBoost

- Let's assume we have 2-class classification problem, with $x_i \in \mathbb{R}^n$, $y_i \in \{-1, 1\}$
- AdaBoost will produce a discriminant function:

$$g(x) = \sum_{t=1}^T \alpha_t h_t(x)$$

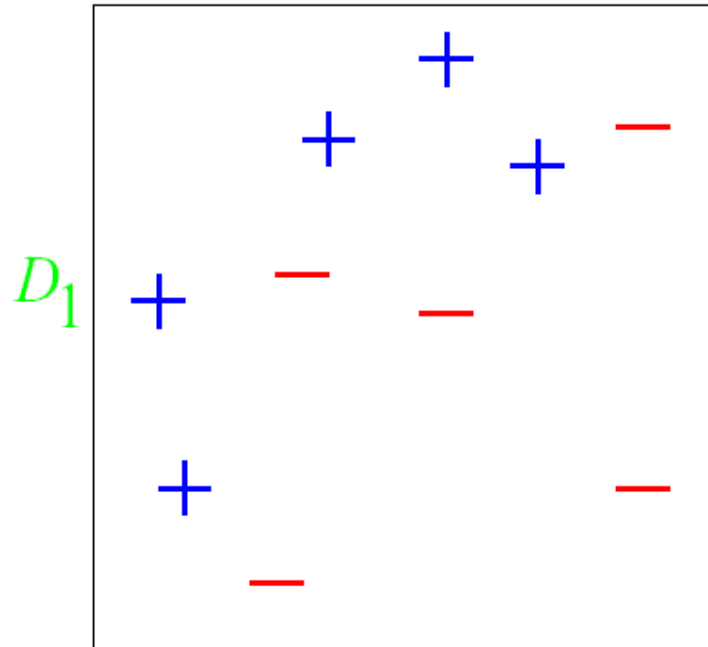
- where $h_t(x)$ is the “weak” classifier
- The final classifier is the sign of the discriminant function, that is $f_{\text{final}}(x) = \text{sign}[g(x)]$

AdaBoost

- $d(x)$ is the distribution of weights over the N training points $\sum d(x_i)=1$
- Initially assign uniform weights $d_0(x_i) = 1/N$ for all x_i
- At each iteration t :
 - Find best weak classifier $h_t(x)$ using weights $d_t(x)$
 - Compute the error rate ϵ_t as
$$\epsilon_t = \sum_{i=1 \dots N} d_t(x_i) \cdot I[y_i \neq h_t(x_i)]$$
 - assign weight α_t the classifier h_t 's in the final hypothesis
$$\alpha_t = \log ((1 - \epsilon_t)/\epsilon_t)$$
 - For each x_i , $d_{t+1}(x_i) = d_t(x_i) \cdot \exp[\alpha_t \cdot I(y_i \neq h_t(x_i))]$
 - Normalize $d_{t+1}(x_i)$ so that $\sum_{i=1} d_{t+1}(x_i) = 1$
- $f_{FINAL}(x) = \text{sign} [\sum \alpha_t h_t(x)]$

AdaBoost Example

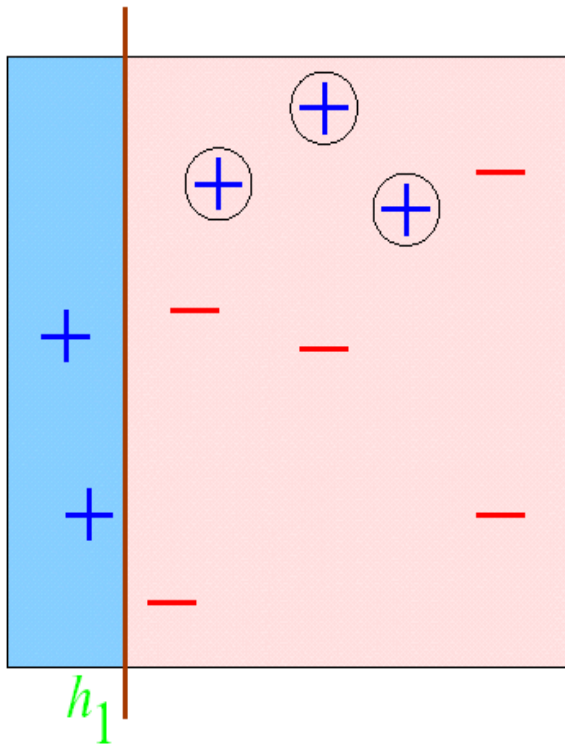
from "A Tutorial on Boosting" by Yoav Freund and Rob Schapire



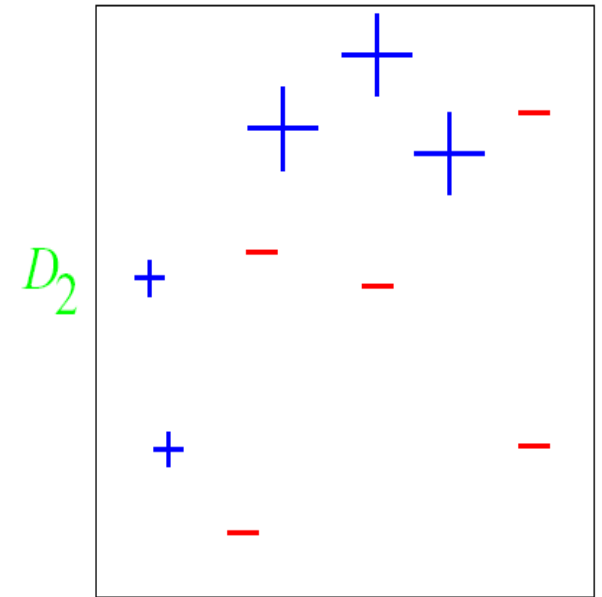
Original Training set : equal weights to all training samples

AdaBoost Example

ROUND 1

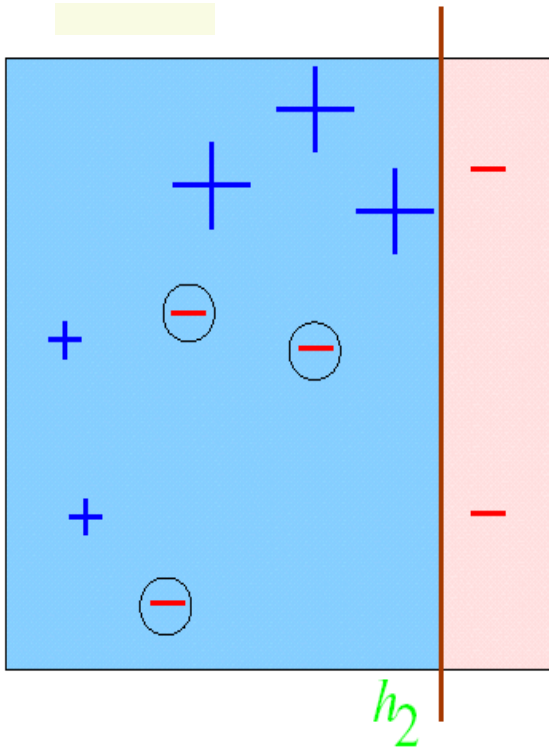


$$\begin{aligned}\epsilon_1 &= 0.30 \\ \alpha_1 &= 0.42\end{aligned}$$



AdaBoost Example

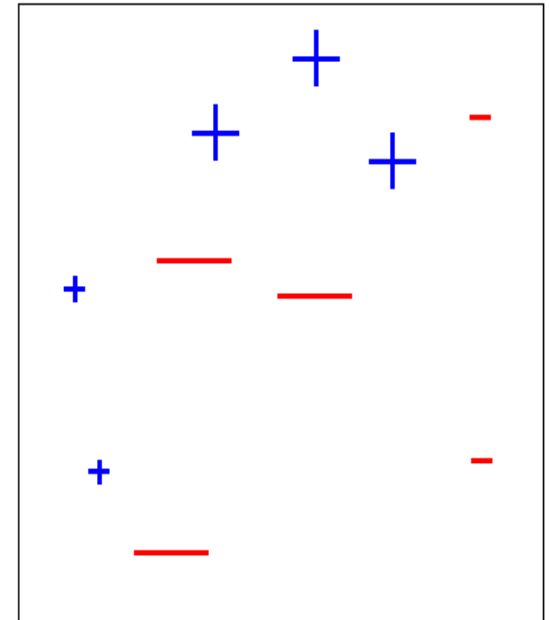
ROUND 2



$$\begin{aligned}\epsilon_2 &= 0.21 \\ \alpha_2 &= 0.65\end{aligned}$$

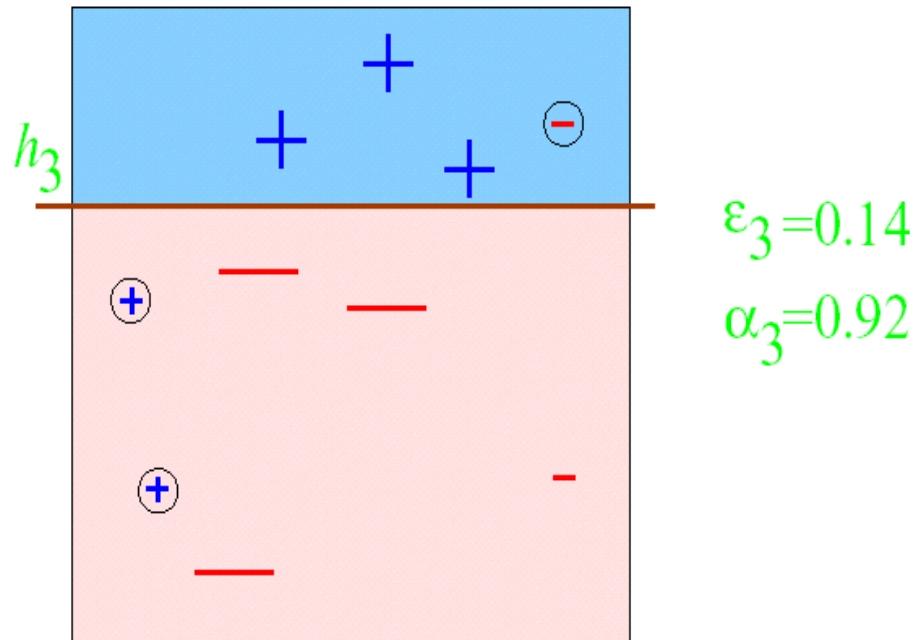


D_3



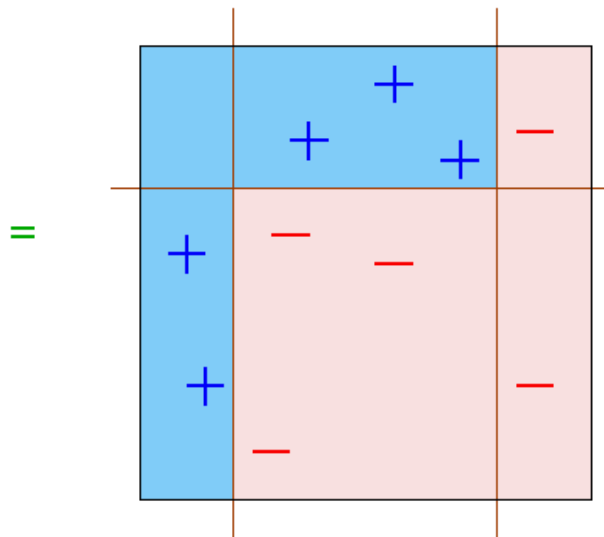
AdaBoost Example

ROUND 3



AdaBoost Example

$$H_{\text{final}} = \text{sign} \left(0.42 \begin{array}{|c|} \hline \text{blue} \\ \hline \end{array} + 0.65 \begin{array}{|c|} \hline \text{blue} \\ \hline \end{array} + 0.92 \begin{array}{|c|} \hline \text{blue} \\ \hline \end{array} \right)$$



SVM Problem Definition

Consider a training set of n iid samples

$$(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_n, \mathbf{y}_n)$$

where \mathbf{x}_i is a vector of length m and

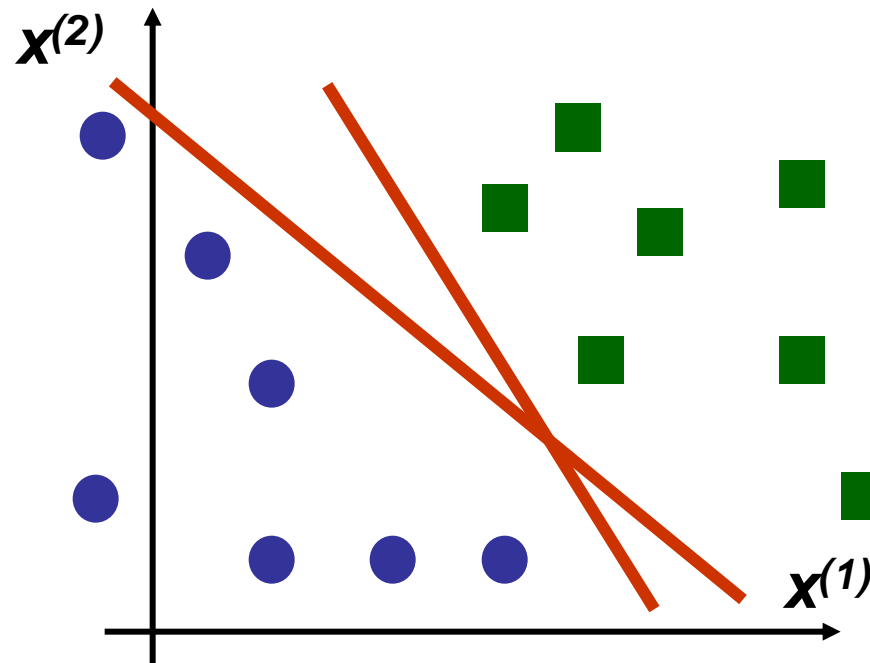
$\mathbf{y}_i \in \{+1, -1\}$ is the class label for data point \mathbf{x}_i .

Find a separating hyperplane $\mathbf{w} \cdot \mathbf{x} + b = 0$

corresponding to the decision function

$$f(\mathbf{x}) = \mathit{sign}(\mathbf{w} \cdot \mathbf{x} + b)$$

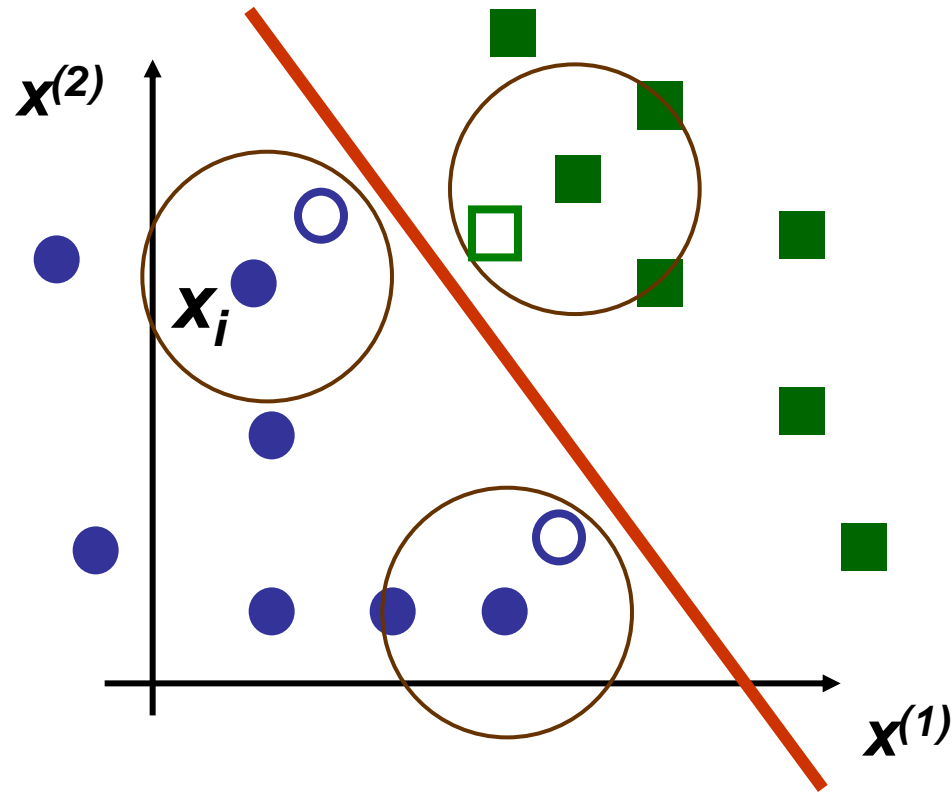
Separating Hyperplanes



- which separating hyperplane should we choose?

Separating Hyperplanes

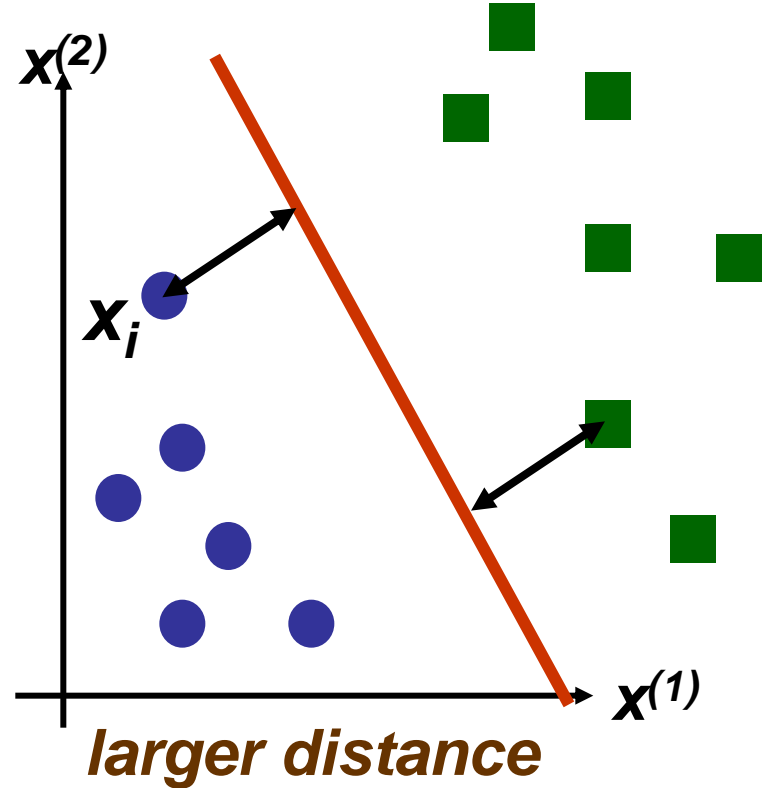
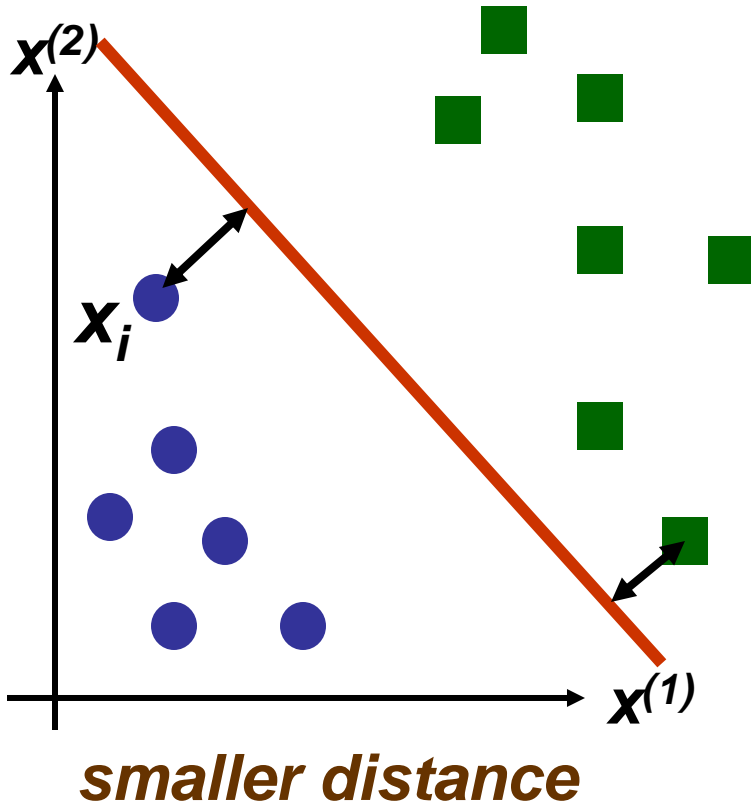
- Hyperplane as far as possible from any sample



- New samples close to the old samples will be classified correctly => Good generalization

SVM

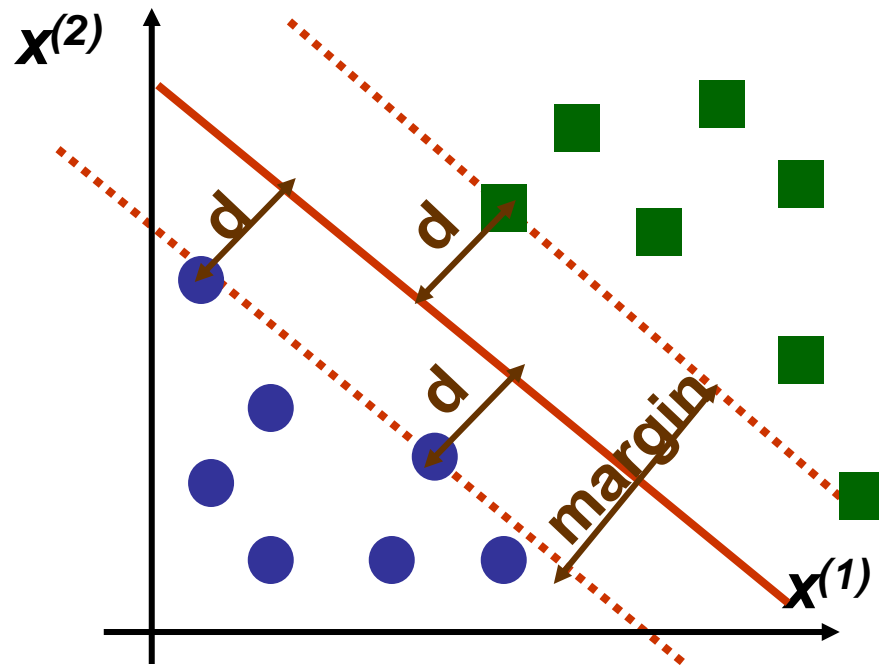
- Idea: maximize distance to the closest example



- For the optimal hyperplane
 - distance to the closest negative example = distance to the closest positive example

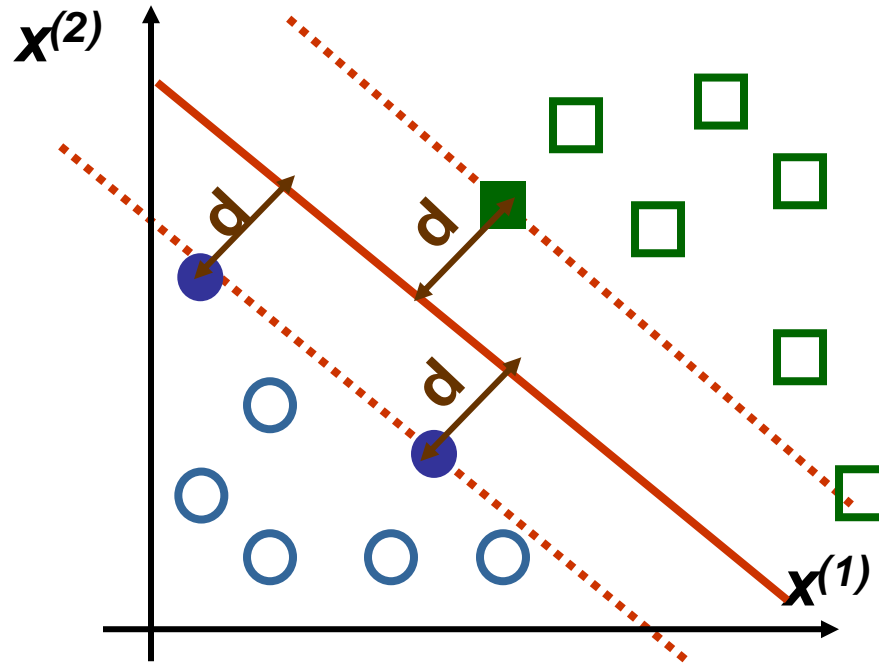
SVM: Linearly Separable Case

- SVM: maximize the *margin*



- Better generalization in theory and practice

SVM: Linearly Separable Case



- **Support vectors** are the samples closest to the separating hyperplane

SVM: Optimal Hyperplane

- Maximize margin $m = \frac{2}{\|\mathbf{w}\|}$

subject to constraints

$$\begin{cases} \mathbf{w}^t \mathbf{x}_i + \mathbf{b} \geq 1 & \mathbf{y}_i = 1 \\ \mathbf{w}^t \mathbf{x}_i + \mathbf{b} \leq -1 & \mathbf{y}_i = -1 \end{cases}$$

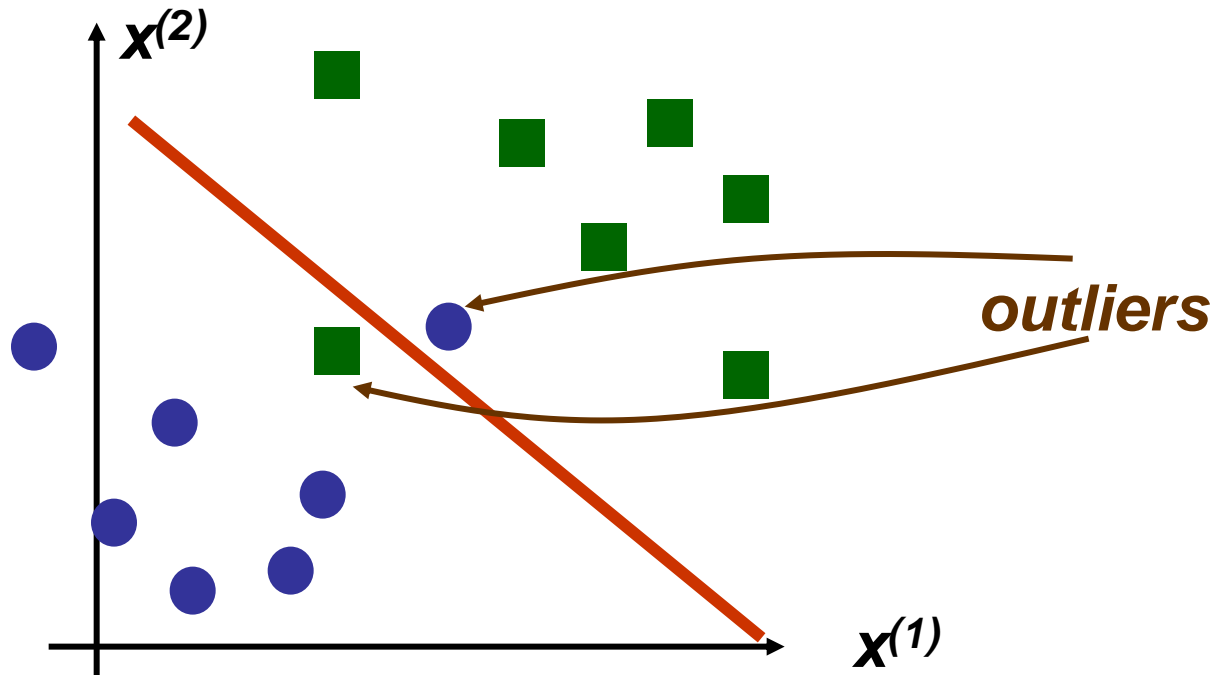
- Can convert our problem to

$$J(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{s.t.} \quad \mathbf{y}_i (\mathbf{w} \cdot \mathbf{x}_i + \mathbf{b}) \geq 1$$

- $J(\mathbf{w})$ is a quadratic function, thus there is a single global minimum

SVM: Non Separable Case

- Data is most likely to be not linearly separable, but linear classifier may still be appropriate



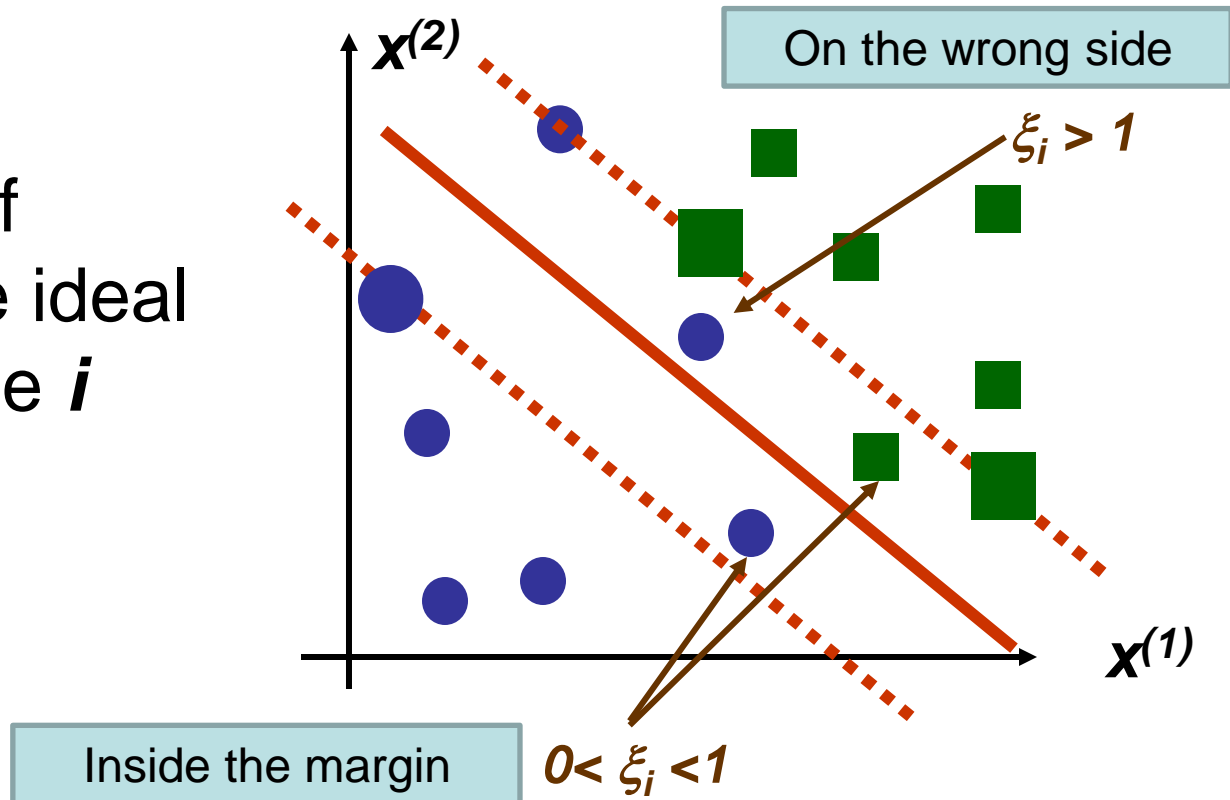
SVM with slacks

- Use nonnegative “slack” variables ξ_1, \dots, ξ_n (one for each sample)

- Change constraints from $y_i(\mathbf{w}^t \mathbf{x}_i + b) \geq 1 \quad \forall i$ to

$$y_i(\mathbf{w}^t \mathbf{x}_i + b) \geq 1 - \xi_i \quad \forall i$$

- ξ_i is a measure of deviation from the ideal position for sample i



SVM with slacks

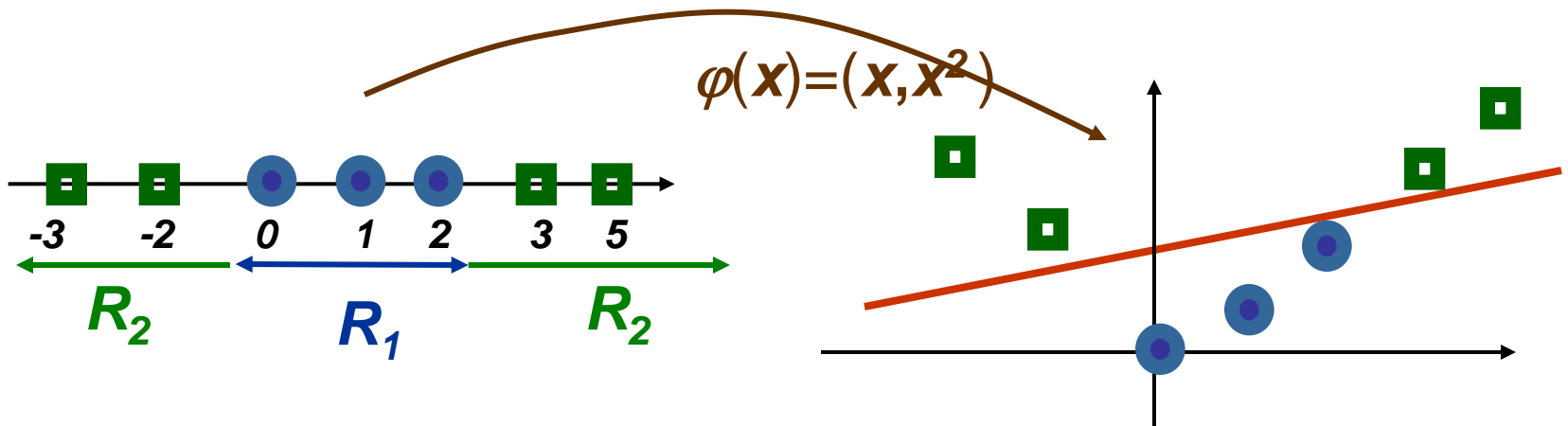
$$\text{minimize} \quad \mathcal{J}(\mathbf{w}, \xi_1, \dots, \xi_n) = \frac{1}{2} \|\mathbf{w}\|^2 + \mathbf{C} \sum_{i=1}^n \xi_i$$

constrained to $\mathbf{y}_i (\mathbf{w}^t \mathbf{x}_i + \mathbf{b}) \geq 1 - \xi_i$ and $\xi_i \geq 0 \quad \forall i$

- $\mathbf{C} > 0$ is a constant which measures relative weight of the first and second terms
 - if \mathbf{C} is small, we allow a lot of samples not in ideal position
 - if \mathbf{C} is large, we want to have very few samples not in ideal position

Non Linear Mapping

- Solve a non linear classification problem with a linear classifier
 1. Project data \mathbf{x} to high dimension using function $\varphi(\mathbf{x})$
 2. Find a linear discriminant function for transformed data $\varphi(\mathbf{x})$
 3. Final nonlinear discriminant function is $\mathbf{g}(\mathbf{x}) = \mathbf{w}^t \varphi(\mathbf{x}) + w_0$

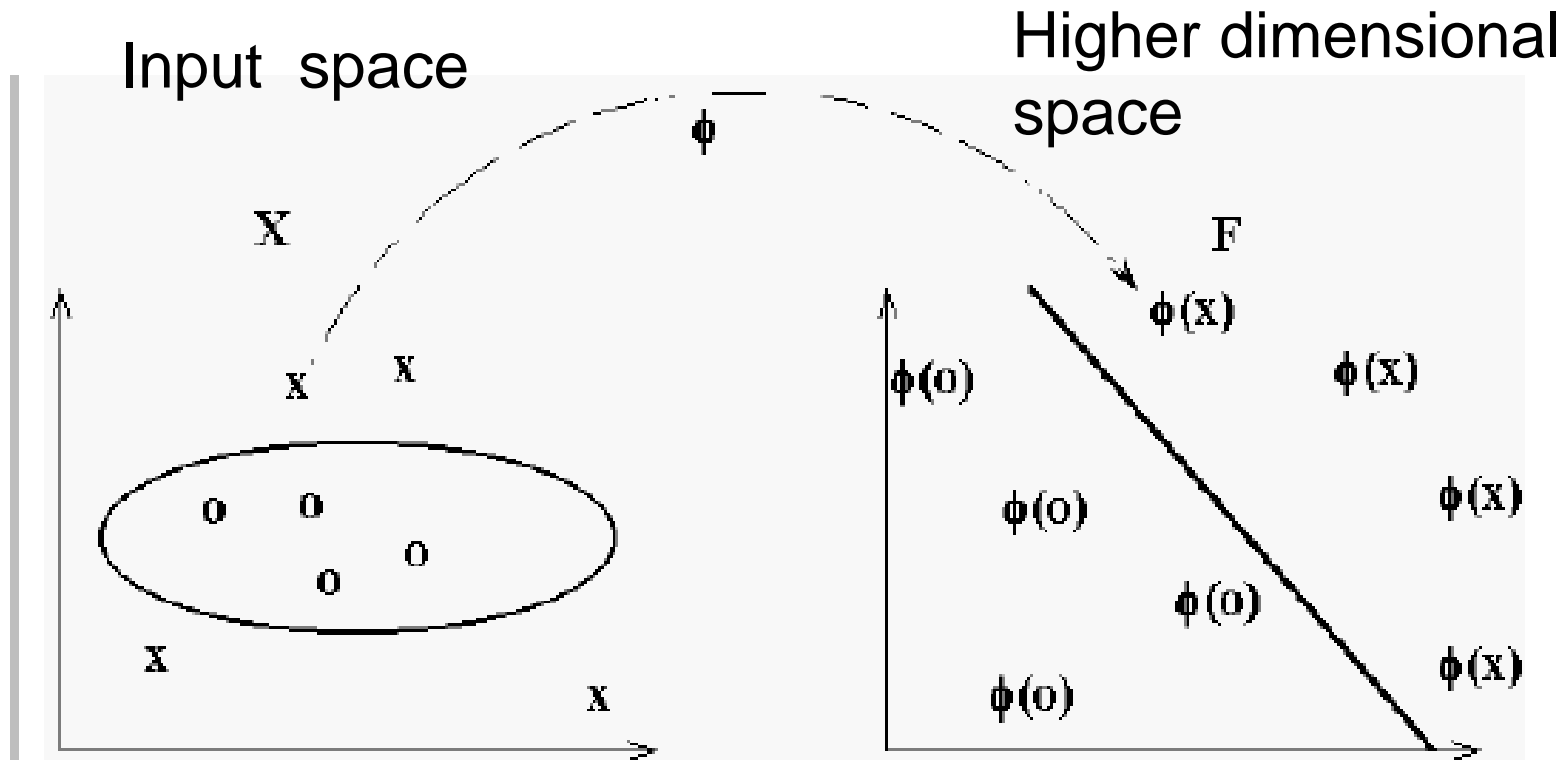


- In 2D, discriminant function is linear

$$\mathbf{g}\left(\begin{bmatrix} \mathbf{x}^{(1)} \\ \mathbf{x}^{(2)} \end{bmatrix}\right) = [\mathbf{w}_1 \quad \mathbf{w}_2] \begin{bmatrix} \mathbf{x}^{(1)} \\ \mathbf{x}^{(2)} \end{bmatrix} + w_0$$

- In 1D, discriminant function is not linear $\mathbf{g}(\mathbf{x}) = \mathbf{w}_1 \mathbf{x} + \mathbf{w}_2 \mathbf{x}^2 + w_0$

Kernel Trick



$$K(x_1, x_2) = \langle \phi(x_1), \phi(x_2) \rangle$$

$$x_1, x_2 \in R^n$$

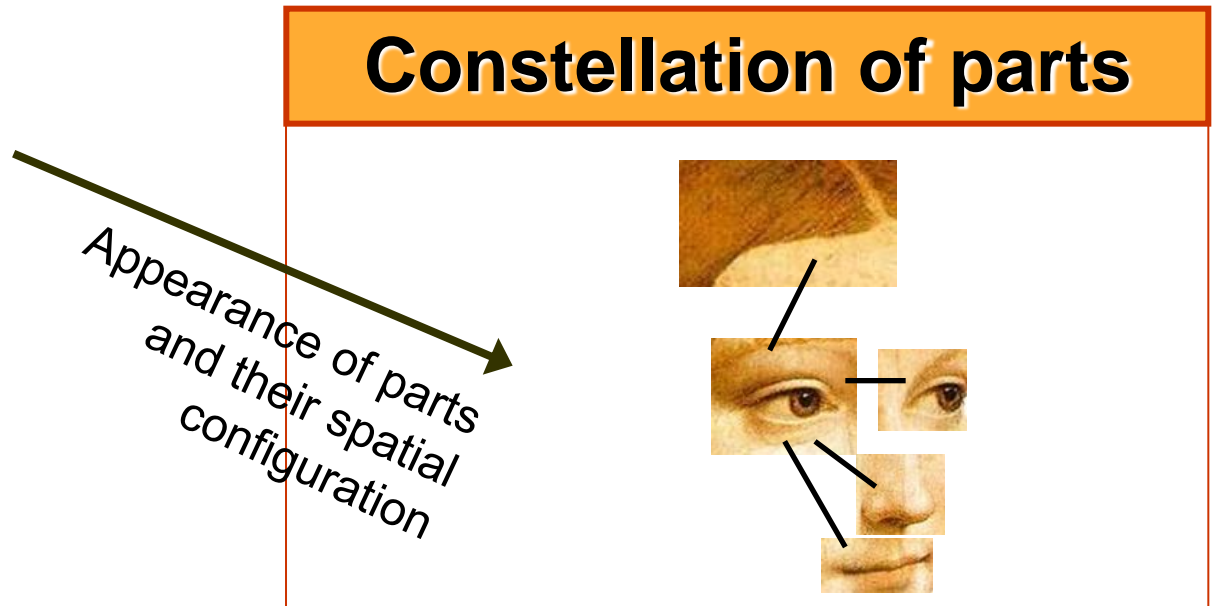
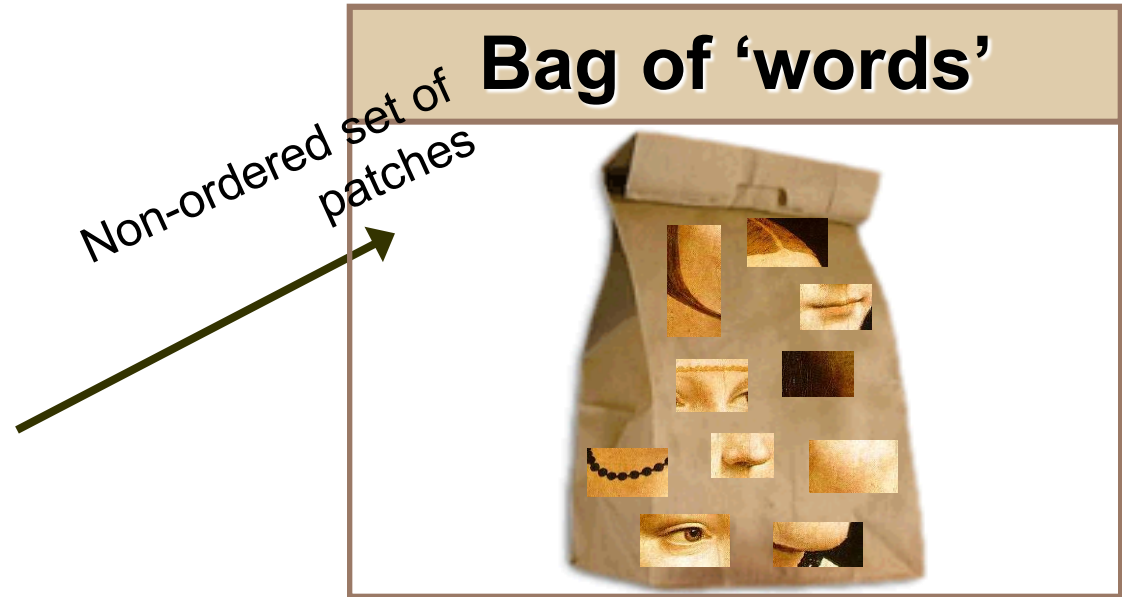
$$\phi(x_1), \phi(x_2) \in F$$

The Kernel Trick

- Choose kernel $K(\mathbf{x}_i, \mathbf{x}_j)$ corresponding to some function $\phi(\mathbf{x}_i)$ which takes sample \mathbf{x}_i to a higher dimensional space (don't need to know $\phi(\mathbf{x}_i)$)
- Replace dot products in the SVM formulation with kernel values.
- Need to compute the kernel matrix for the training data
- Need to compute $K(\mathbf{x}_i, \mathbf{x})$ for all SVs \mathbf{x}_i

Part-Based Approaches

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.

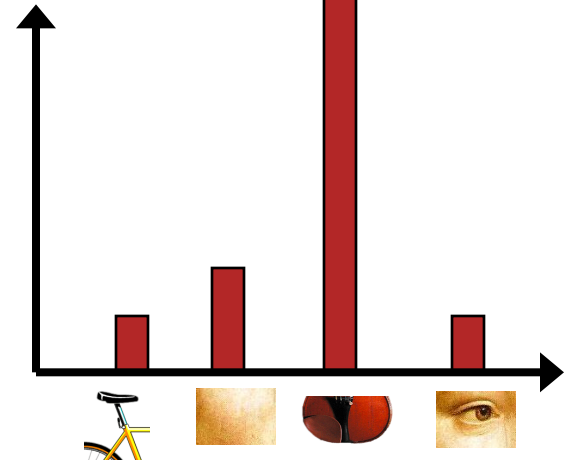
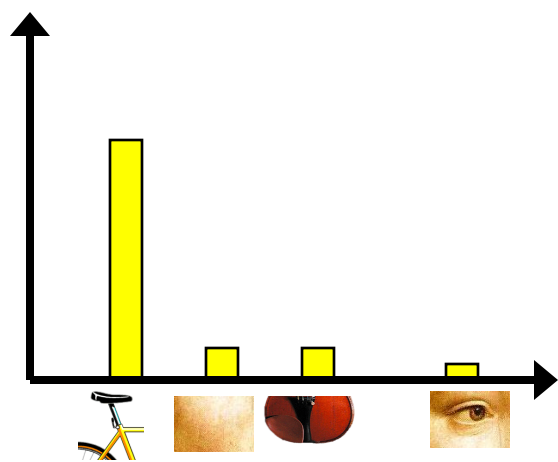
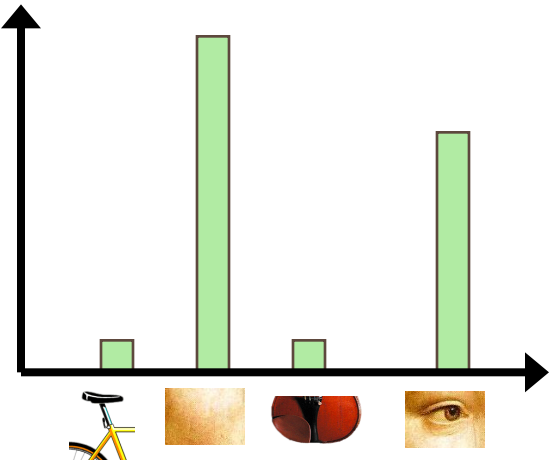
For a long time, the retinal image was considered as a visual centers in the brain. It is a moving image that is perceived and known. Hubel and Wiesel demonstrated that the *message of the image falling on the retina undergoes a step-wise analysis in a system of nerve cells stored in columns. In this 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, and a 10% increase in imports to \$600bn.

China's government also noted that the demand for the yuan against the US dollar is increasing. Beijing has made it clear that it will take its time and tread carefully in allowing the yuan to rise further in value.

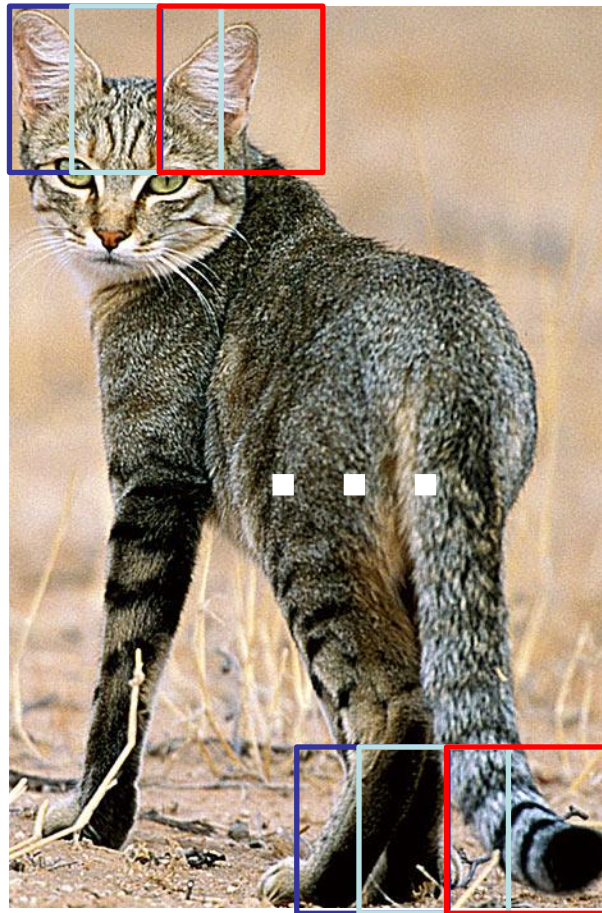
China, trade, surplus, commerce, exports, imports, US, yuan, bank, domestic, foreign, increase, trade, value



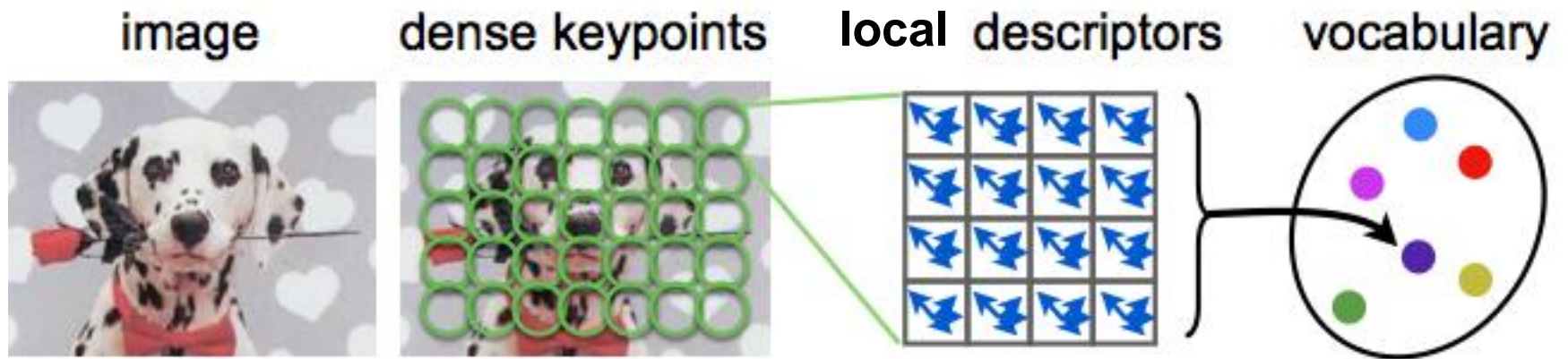
Regular dense grid over space and scale



Dense keypoints over space and scale

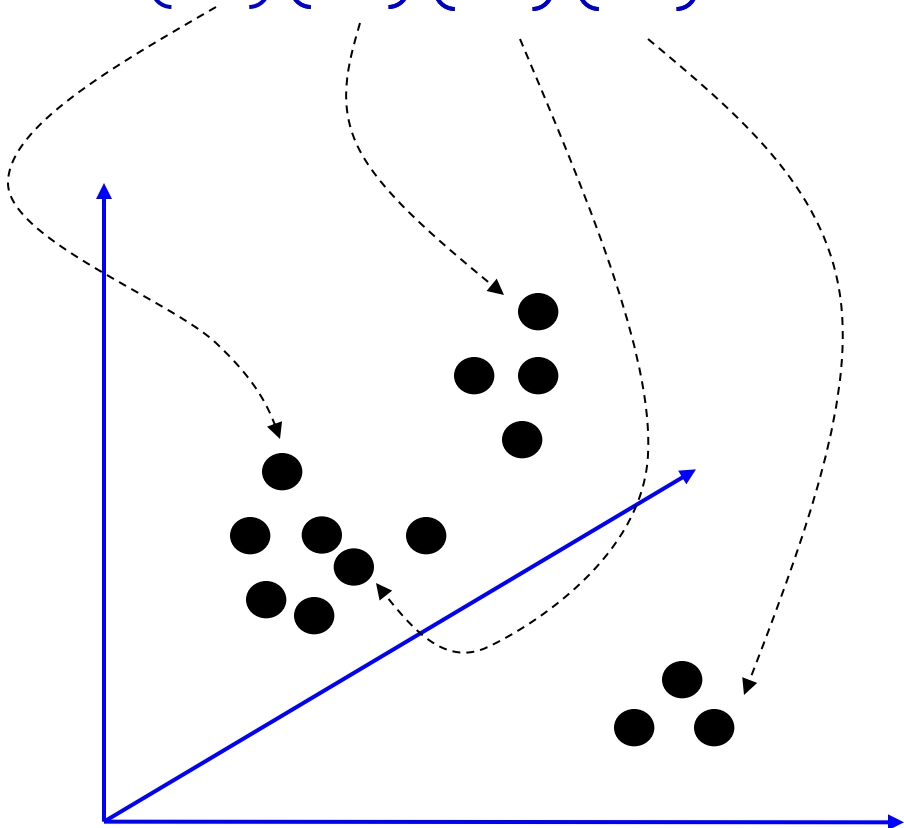
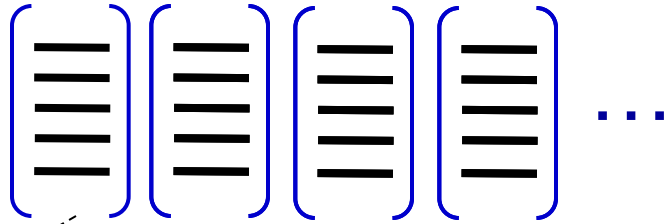


Representation Scheme



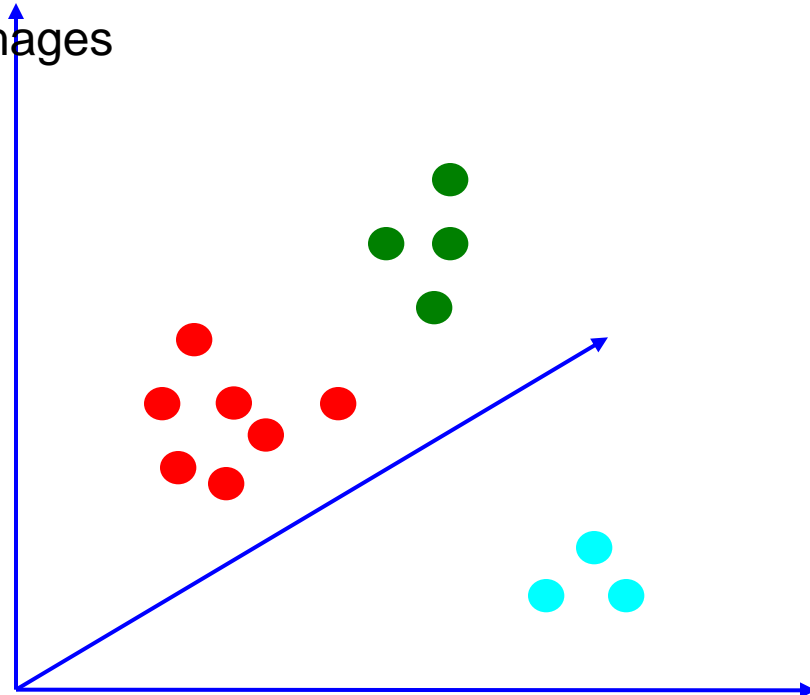
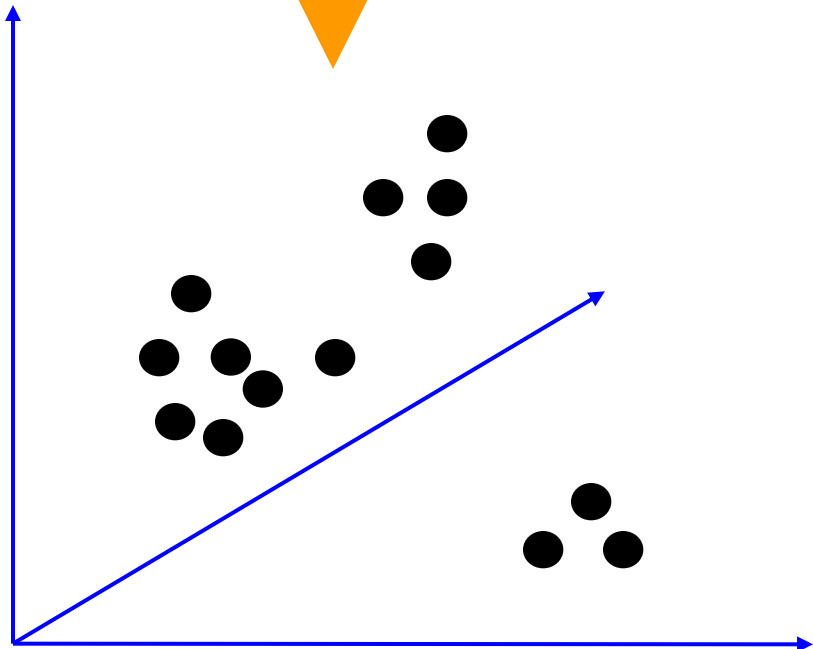
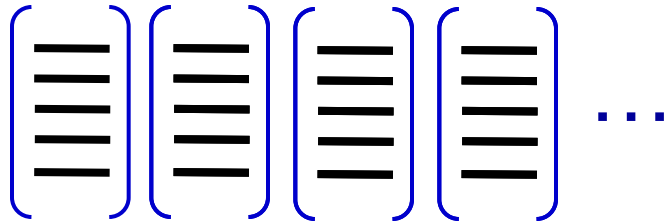
Vocabulary formation

Local descriptors from all (subset of) images



Vocabulary formation

Local descriptors from all (subset of) images



Clustering

Clustering

- **Input:**
 - Training samples $x_1, \dots, x_n \in R^n$
 - No labels are given.
- **Goal:** group input samples into classes of similar objects – cohesive “clusters.”
 - high intra-class similarity
 - low inter-class similarity
- **Algorithms:** many, most common k-means

k-means, definitions

- k – the number of clusters
a parameter of the algorithm
- μ_i cluster centroids
represent our current guesses for the positions
of the centers of the clusters
- Initialization: pick k random training
samples.
Other initialization methods are also possible.

k-means, algorithm

1. Initialize cluster centroids: $\mu_1, \dots, \mu_k \in R^n$
2. Repeat until convergence:

For every i , set

Assign point to the closest cluster centroid

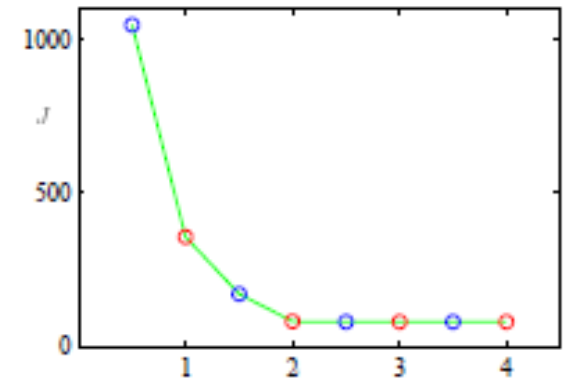
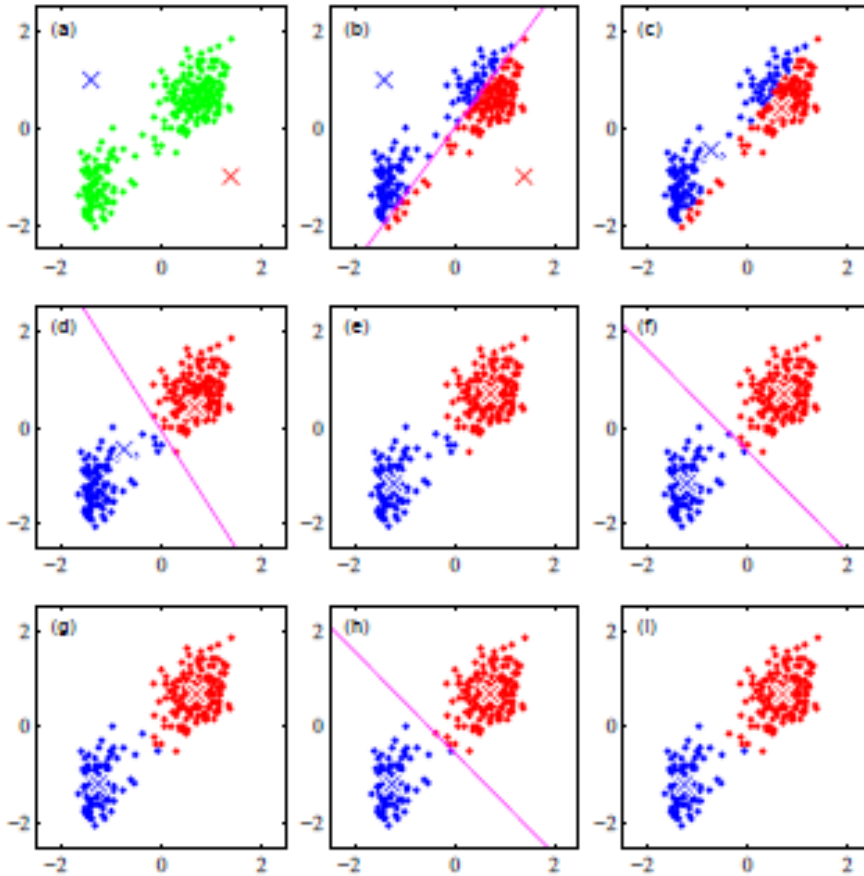
$$c_i = \arg \min_j \|x_i - \mu_j\|^2$$

For each j , set

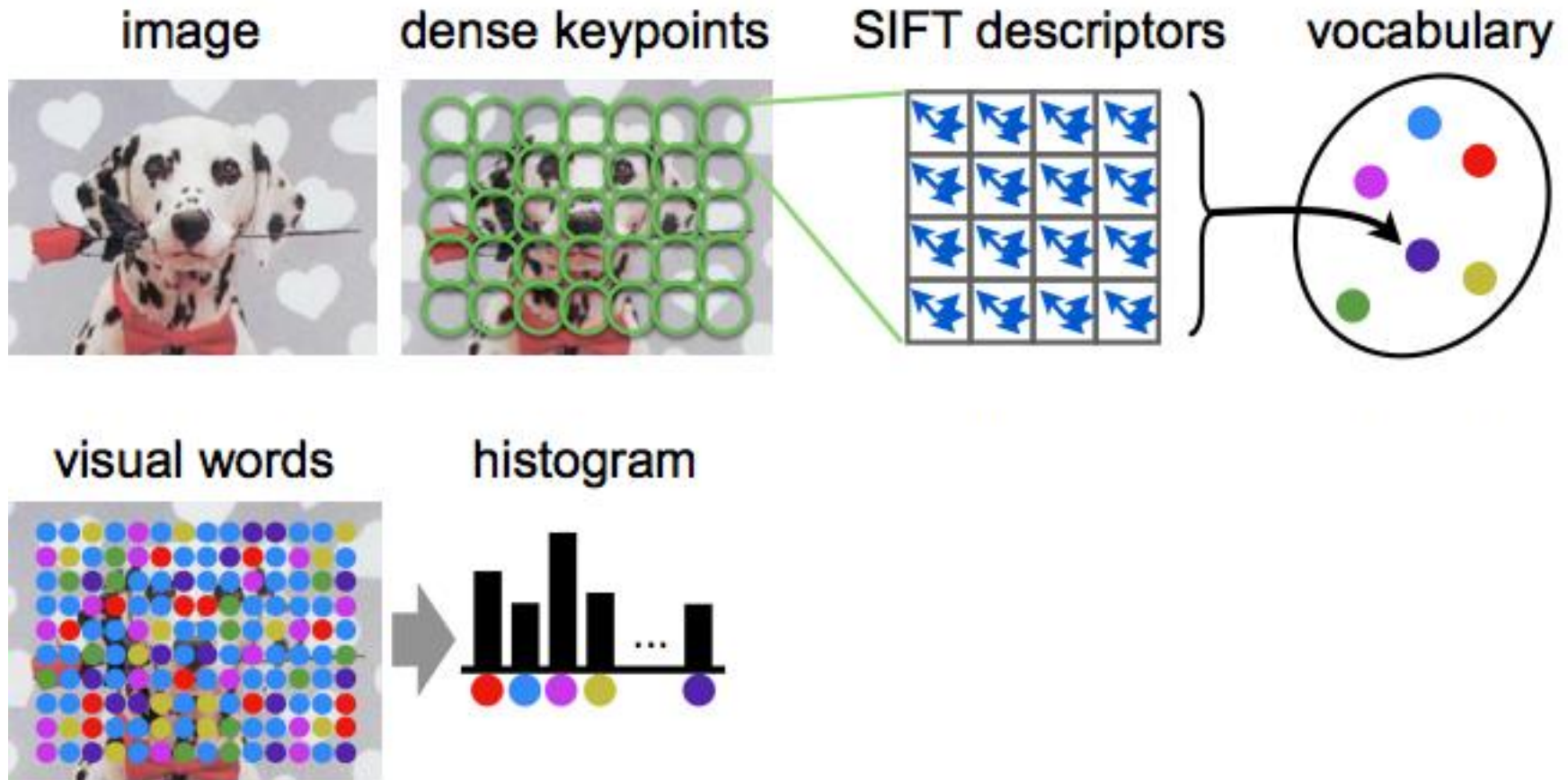
Update centroids to be the mean of the points assigned to it

$$\mu_j = \frac{\sum_{i=1}^n 1\{c_i = j\} x_i}{\sum_{i=1}^n 1\{c_i = j\}}$$

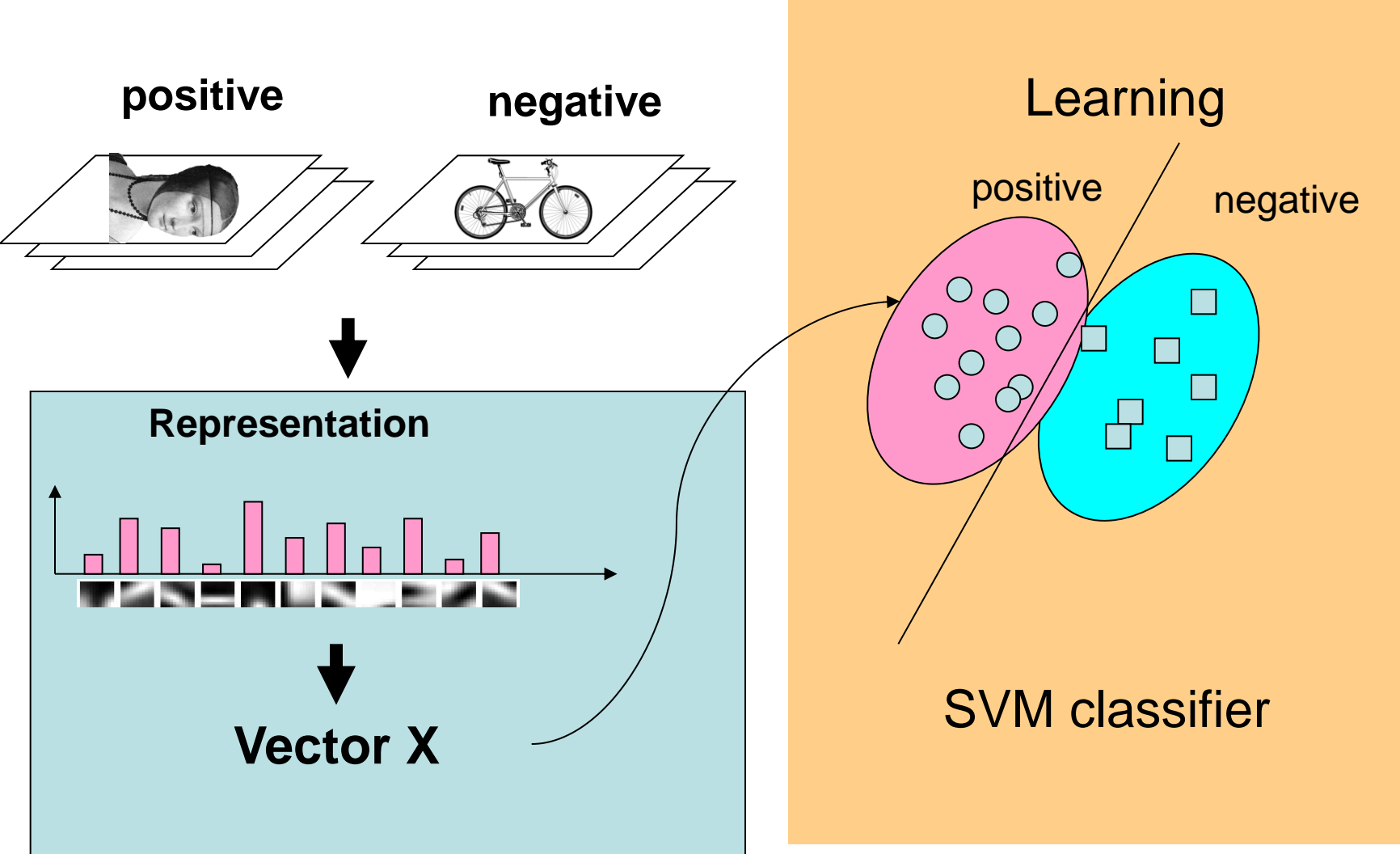
K-means, example



Representation Scheme



SVM classification

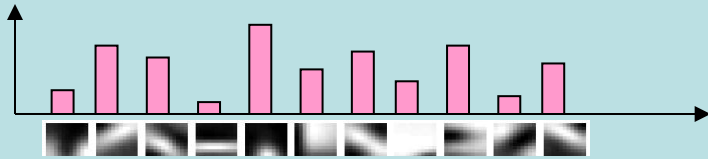


SVM classification

Test image



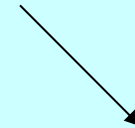
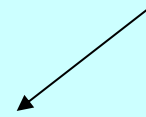
Representation



Vector X

Recognition

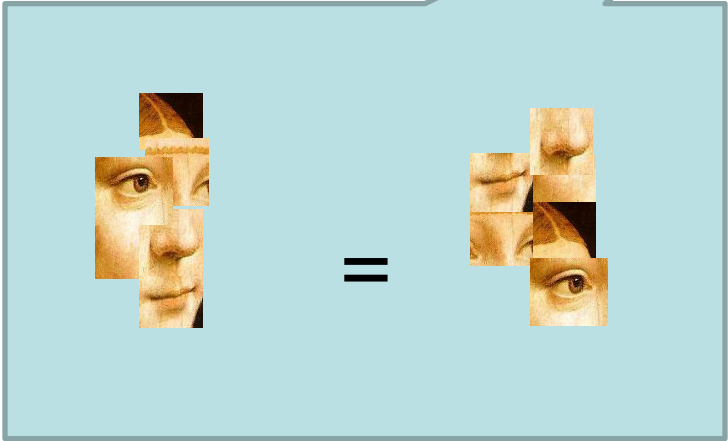
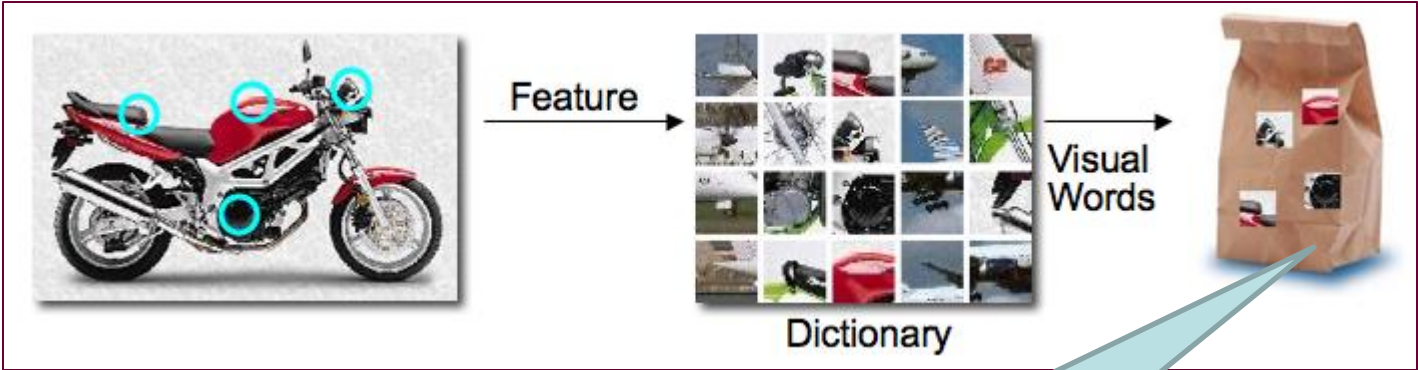
SVM(X)



**Contains
object**

Doesn't
contain
object

Bag of Features



Bag of Features



Pros: fairly flexible and computationally efficient

Cons: problems with large clutter



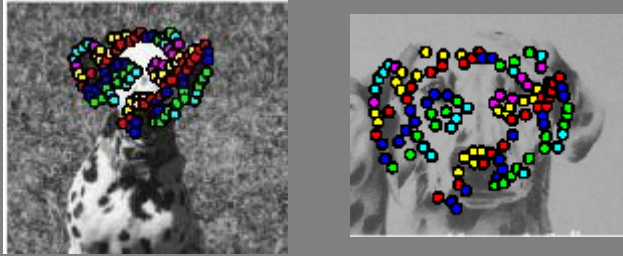
Different objects, but Similar representations;



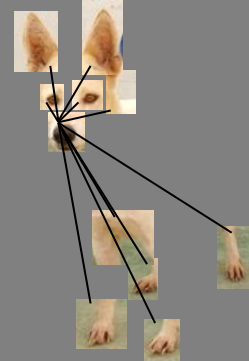
Similar objects, different representations;

Fusion of appearance and shape

Geometric
correspondence search



Constellation model

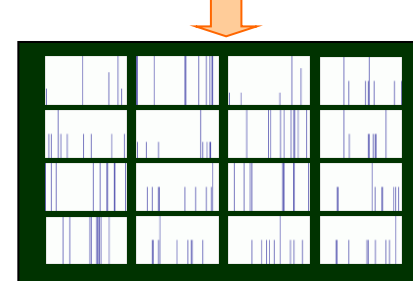
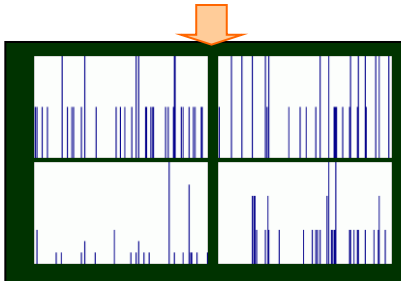
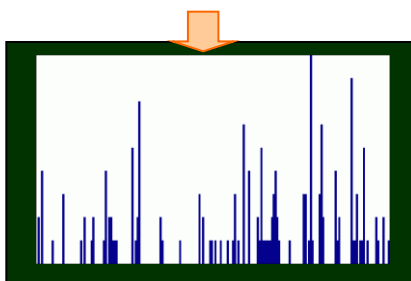
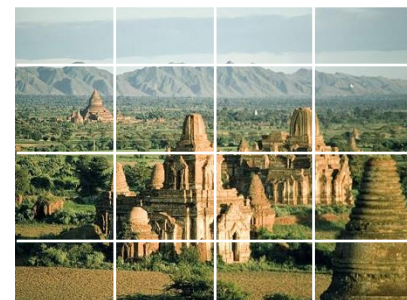
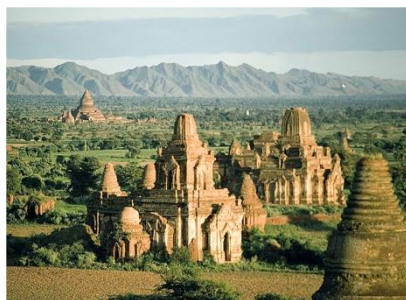


Pros: Structural information is important for recognition.

Cons: computationally expensive, restricted to small variation in shape within the category.

Beyond Bags of Features

- Computing bags of features on sub-windows of the whole image.



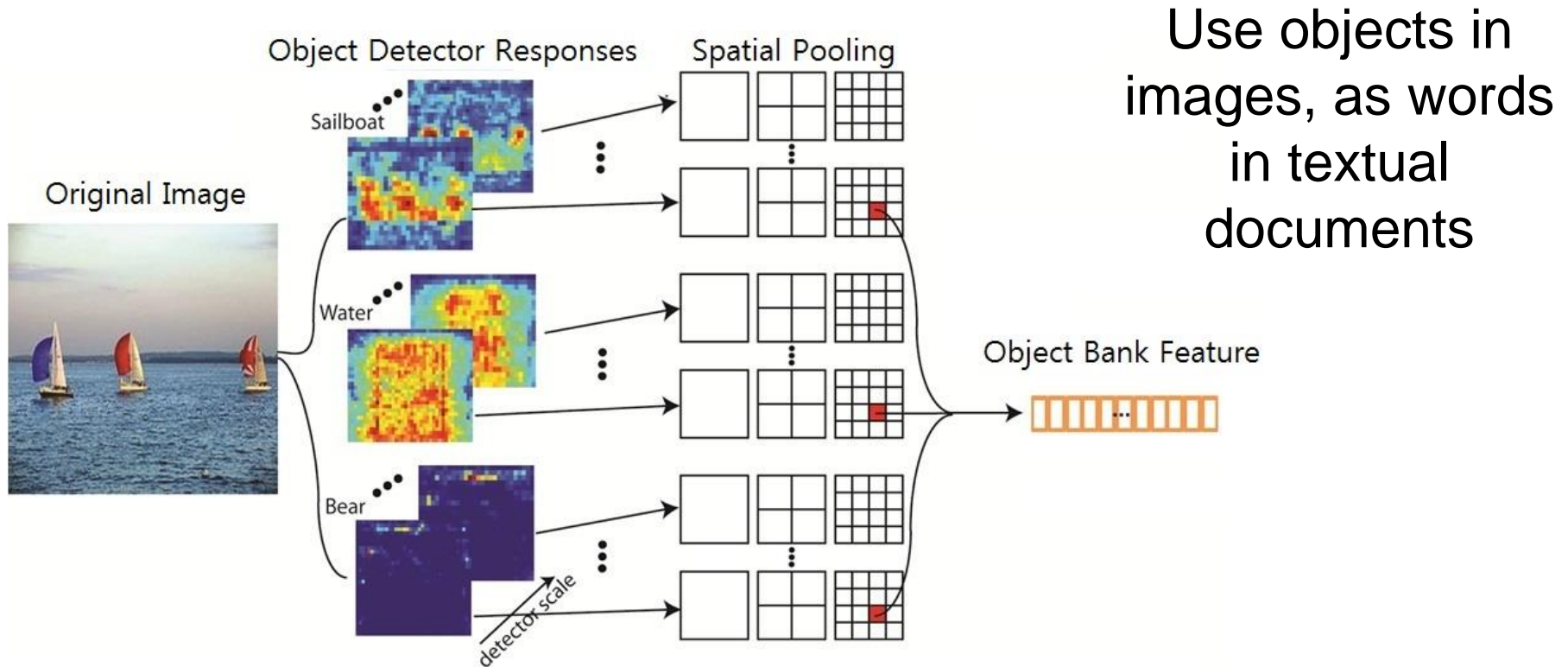
Context and Scenes

- What is behind the red rectangle?



Context is correlated
with the object

Semantically Meaningful Feature



[Li-Jia Li](#), [Hao Su](#), [Yongwhan Lim](#), [Robert Cosgriff](#), Daniel Goodwin, and [Li Fei-Fei](#) [Vision Lab, Stanford University](#)

Describing Objects with Attributes

- Shift the goal of recognition from naming to describing

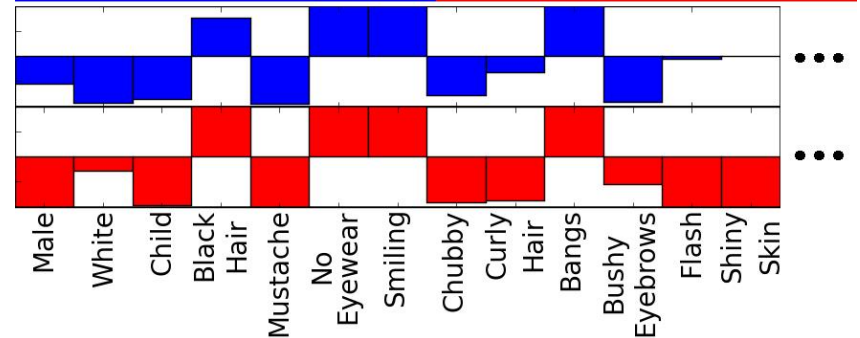
otter

black: yes
 white: no
 brown: yes
 stripes: no
 water: yes
 eats fish: yes

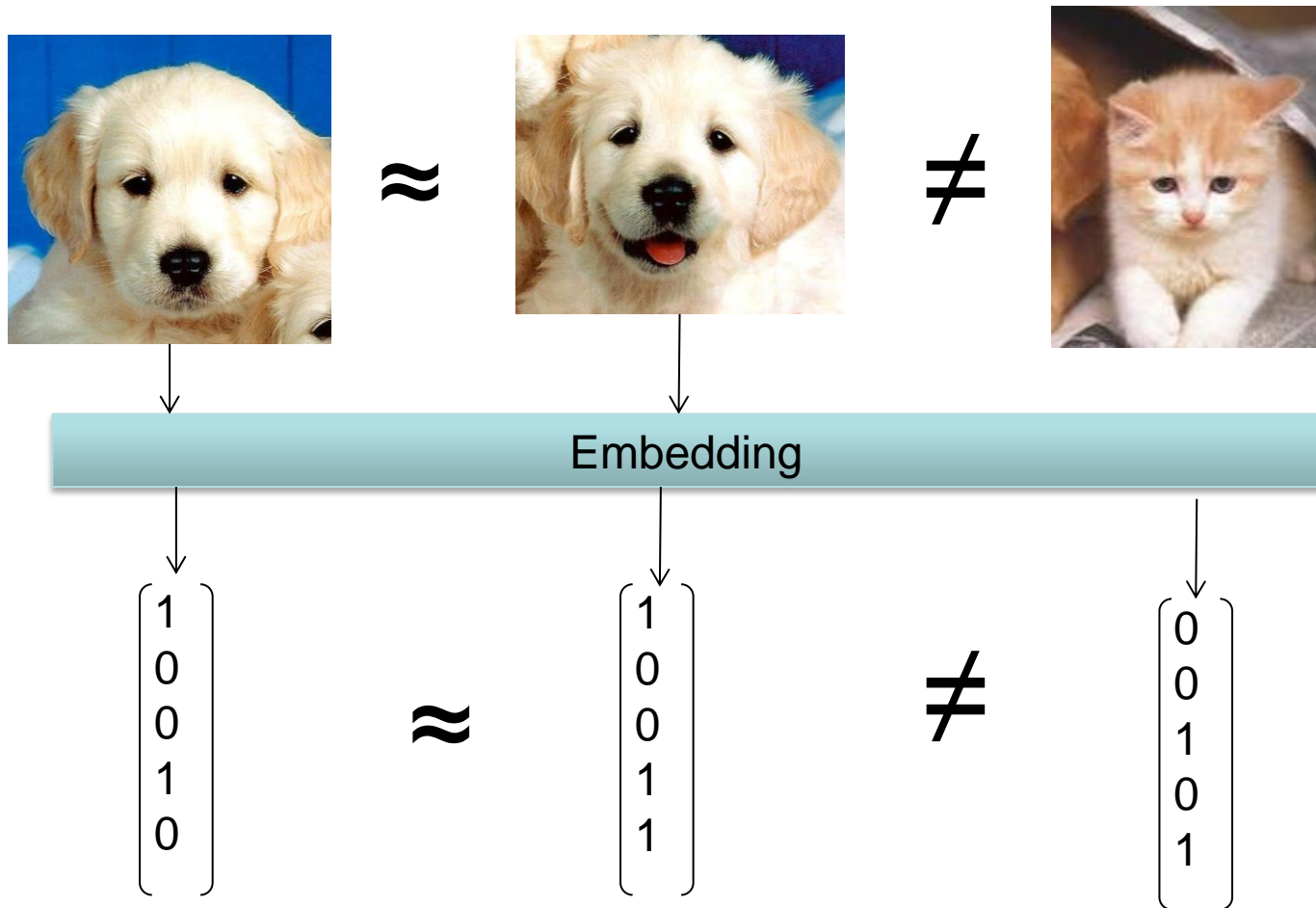


polar bear

black: no
 white: yes
 brown: no
 stripes: no
 water: yes
 eats fish: yes



Large-scale image/object search



Large-scale image/object search

General form of linear projection-based hashing function:

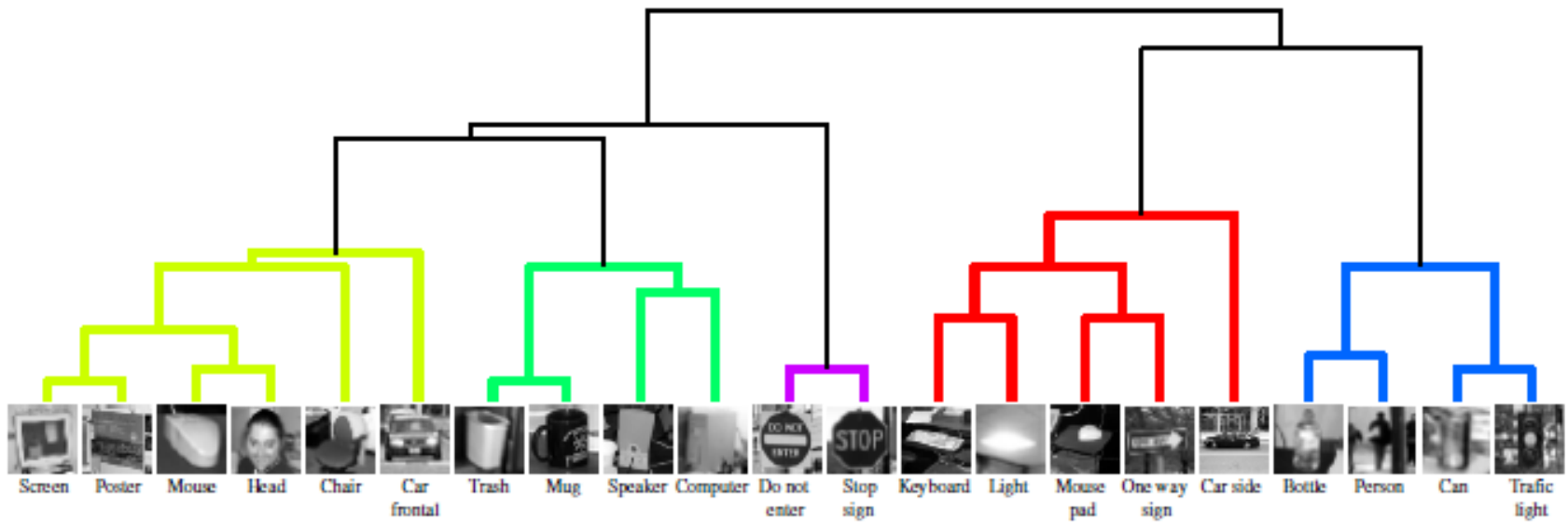
$$h_i(I) = \text{sign}\left(f\left(w_i^T I - b_i\right)\right)$$

Learn embedding $W = [w_1 \dots w_k]$, such that:
If I and J are similar (for example $\|I - J\|^2 < \varepsilon$), then

$$d_{\text{Hamming}}(I, J) < k$$

Dealing with many categories

Sharing Features between classes



Dealing with many categories

Model Transfer

Learn Model on



...



Adjust Model to



Activity recognition

