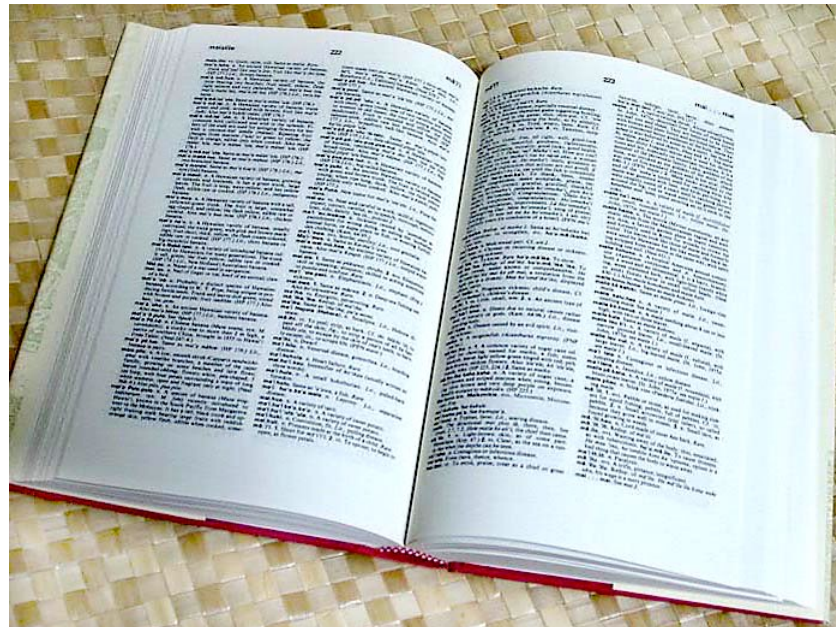# RECOGNITION

Thanks to Svetlana Lazebnik and Andrew Zisserman for the use of some slides

# How many categories?

# Variability makes recognition hard

Camera position
Illumination
Shape parameters

Within-class variations?

# Variations within the same class

# History

1960s – early 1990s: geometry

1990s: appearance
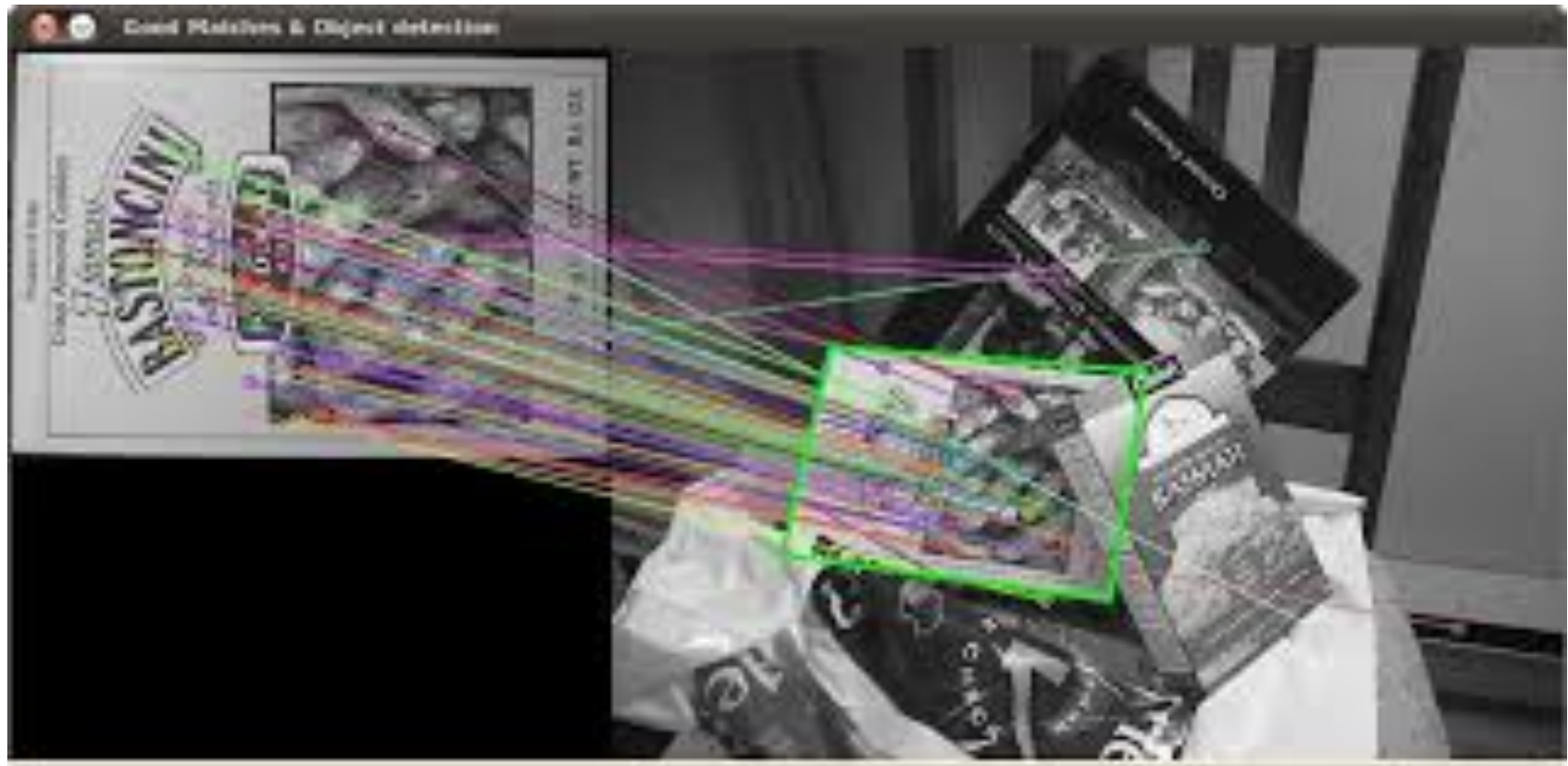
Mid-1990s: sliding window

Late 1990s: local features

Early 2000s: parts-and-shape models

Mid-2000s: bags of features

Present trends: data-driven methods, context

# 2D objects

# Eigenfaces (Turk & Pentland, 1991)



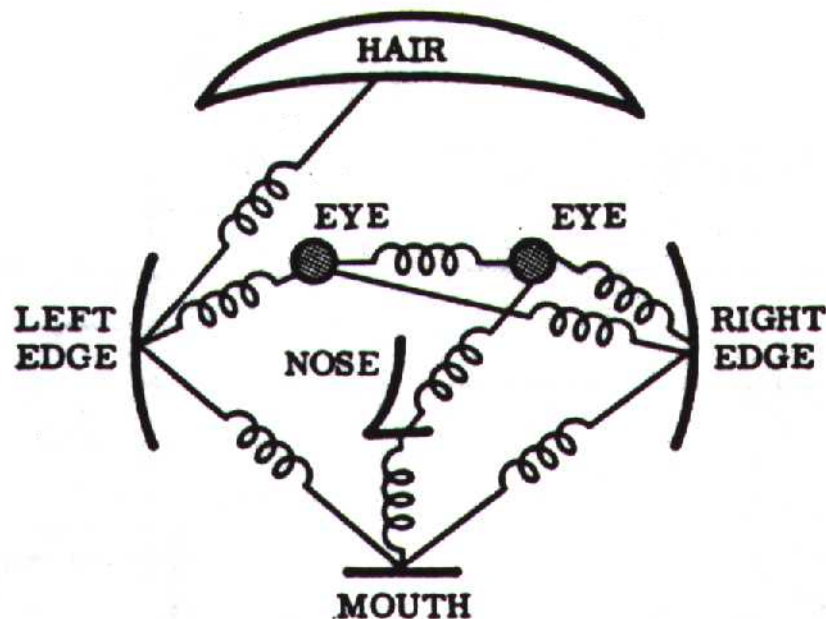| Experimental | Correct/Unknown Recognition Percentage | | |
|---|---|---|---|
| Condition | Lighting | Orientation | Scale |
| Forced classification | 96/0 | 85/0 | 64/0 |
| Forced 100% accuracy | 100/19 | 100/39 | 100/60 |
| Forced 20% unknown rate | 100/20 | 94/20 | 74/20 |

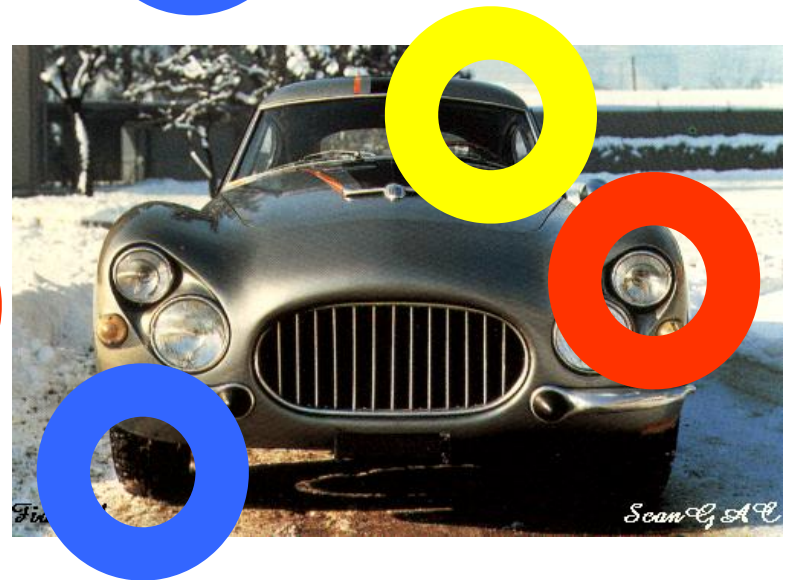# Local features



D. Lowe (1999, 2004)

# Parts-and-shape models

Model:

- Object as a set of parts
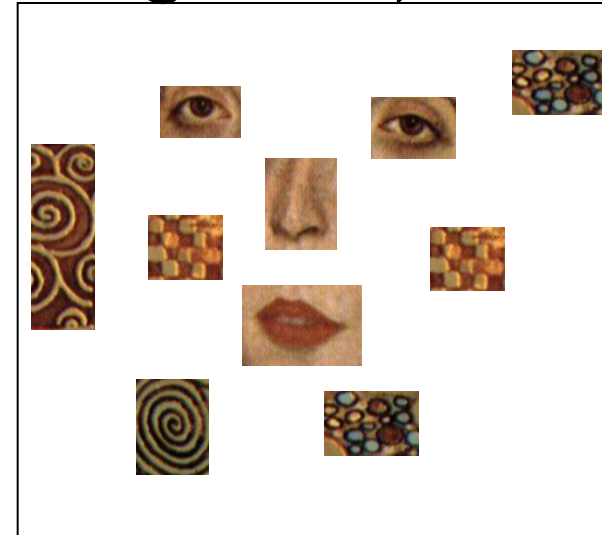- Relative locations between parts
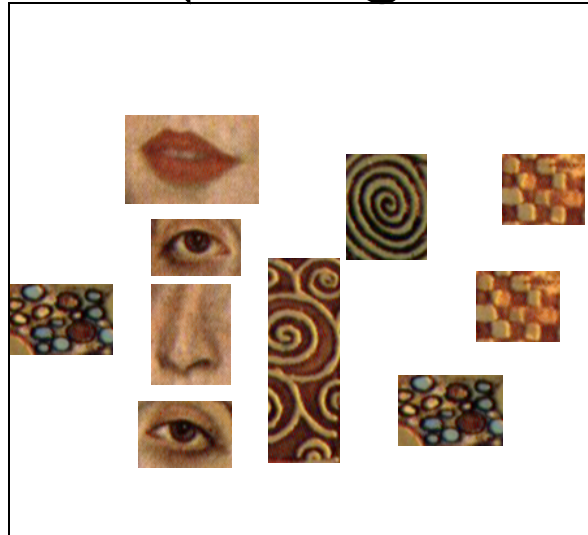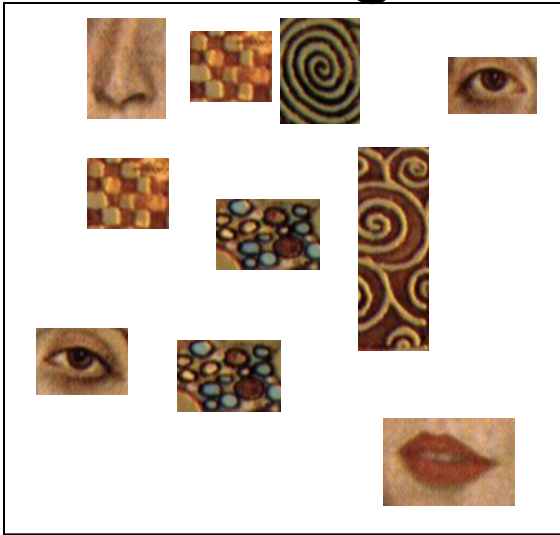- Appearance of part

# Constellation models



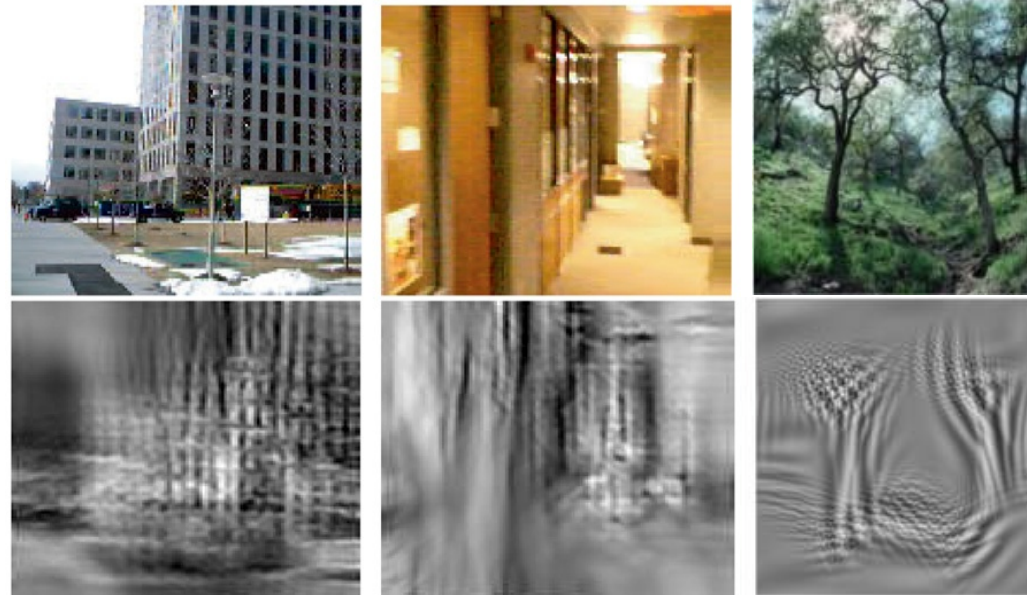Weber, Welling & Perona (2000), Fergus, Perona & Zisserman (2003)
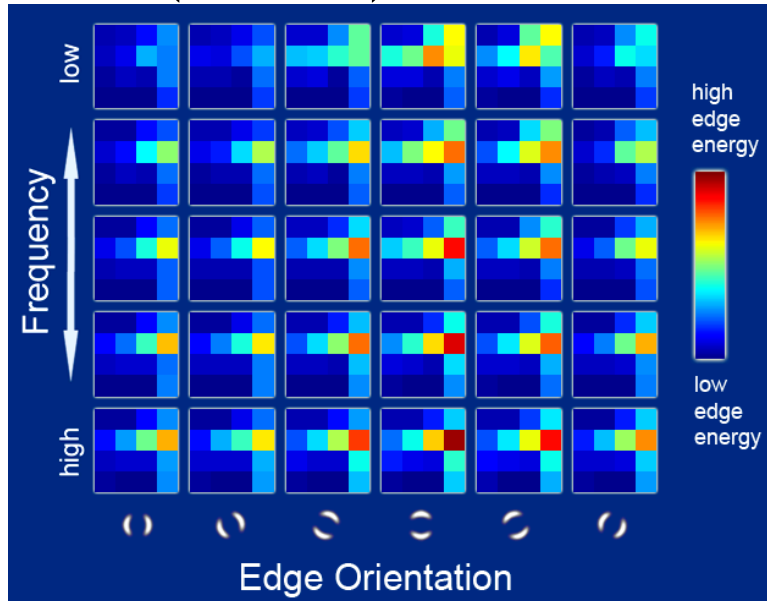
# Objects as texture

- All of these are treated as being the same – no segmentation (background, foreground)

# Global scene descriptors

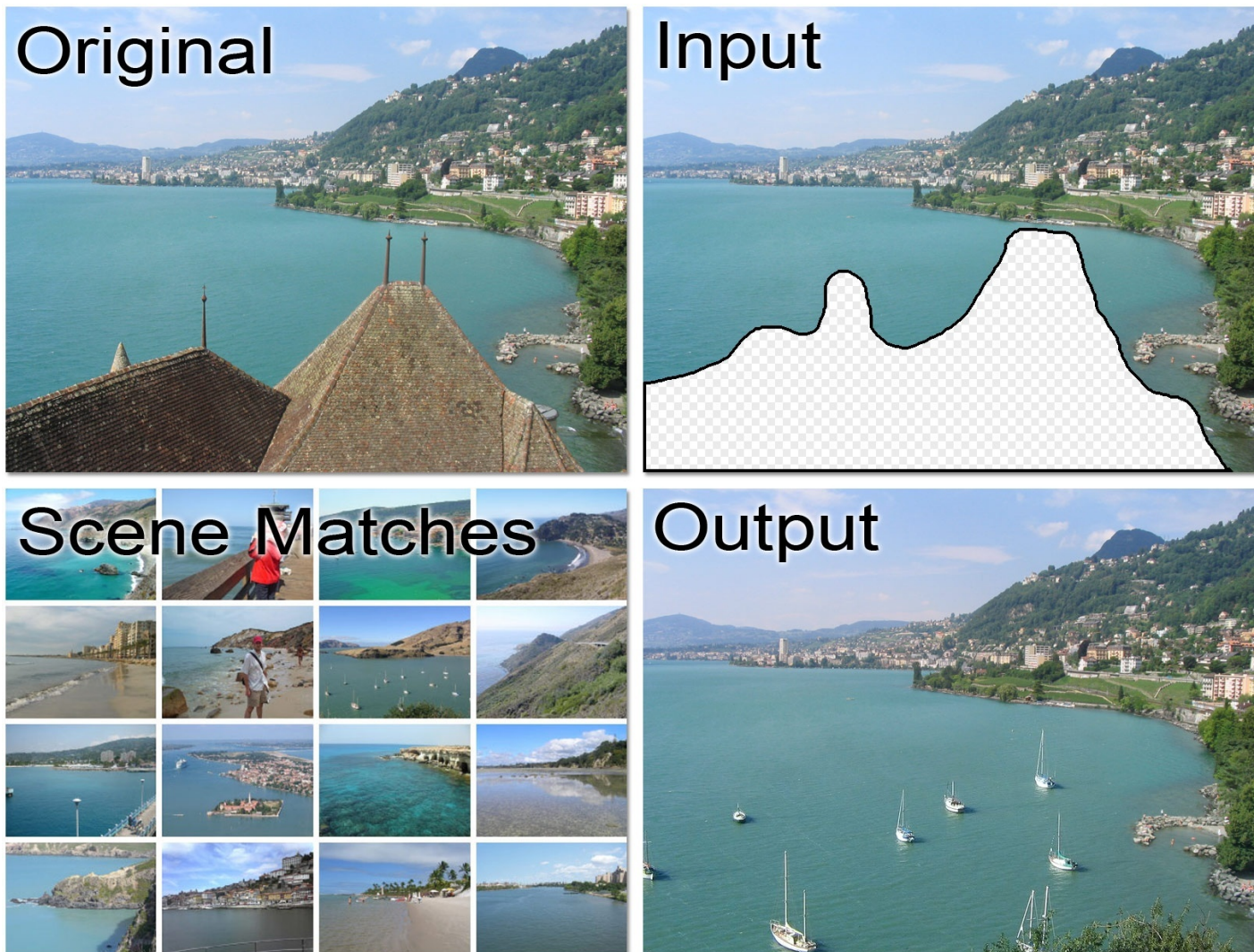- The "gist" of a scene: Oliva & Torralba (2001)



http://people.csail.mit.edu/torralba/code/spatialenvelope/

# Data-driven methods



Original

Input

Scene Matches

Output

J. Hays and A. Efros, Scene Completion using Millions of Photographs, SIGGRAPH 2007

# Data-driven methods



(a) Query Image

(b) Retrieval Set

(c) Superpixels

Building    Road

Car    Sky
(d) Per-class Likelihoods

Building    Sky

Car

Road

(e) Final Labeling

J. Tighe and S. Lazebnik, ECCV 2010

# Overview

- Basic recognition tasks

- A statistical learning approach

- Traditional or "shallow" recognition pipeline

  - Bags of features

  - Classifiers

- Currently best approaches: neural networks and "deep" recognition pipeline

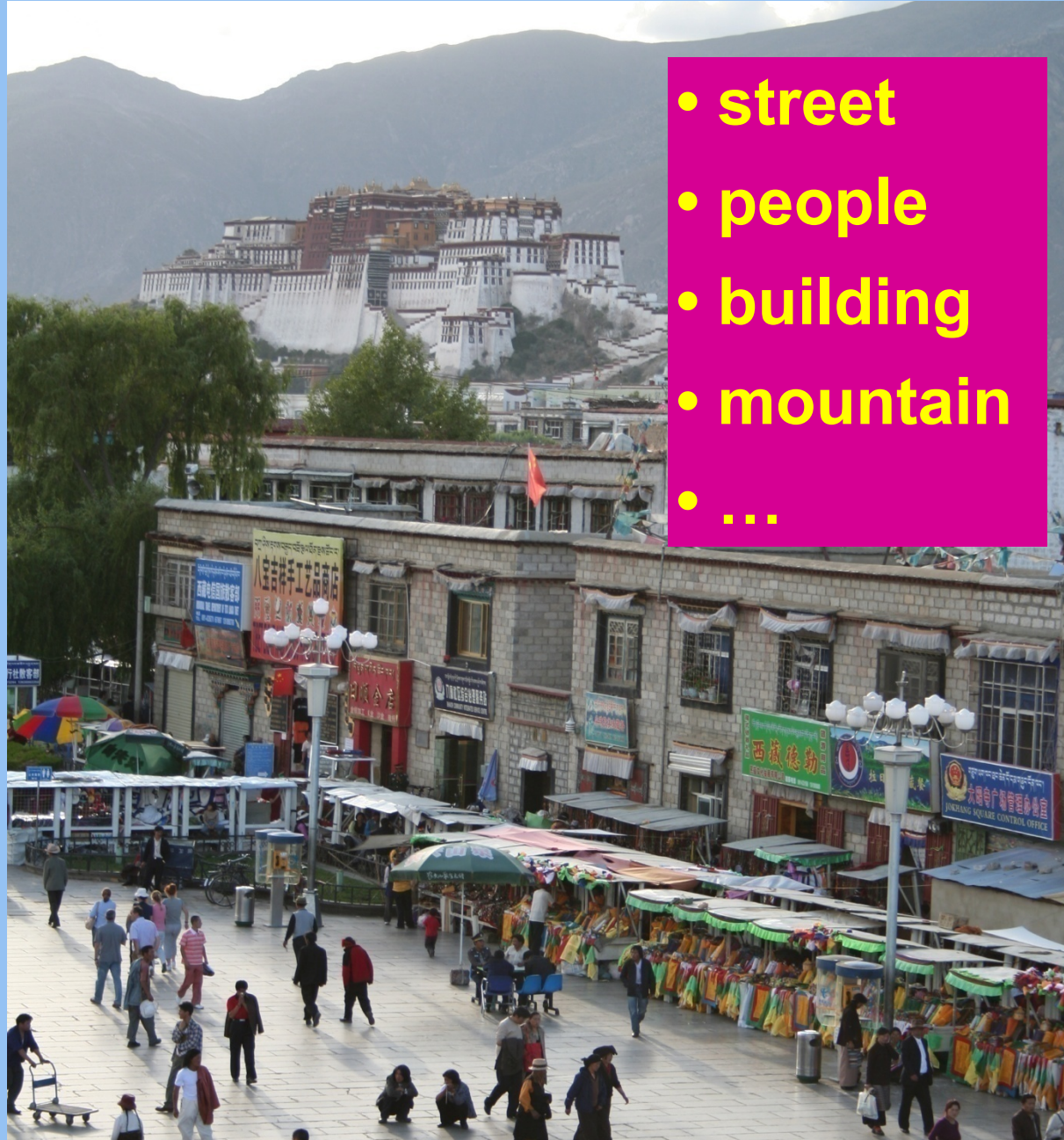# RECOGNITION AS AN APPLICATION OF MACHINE LEARNING

# Common recognition tasks

# Image classification



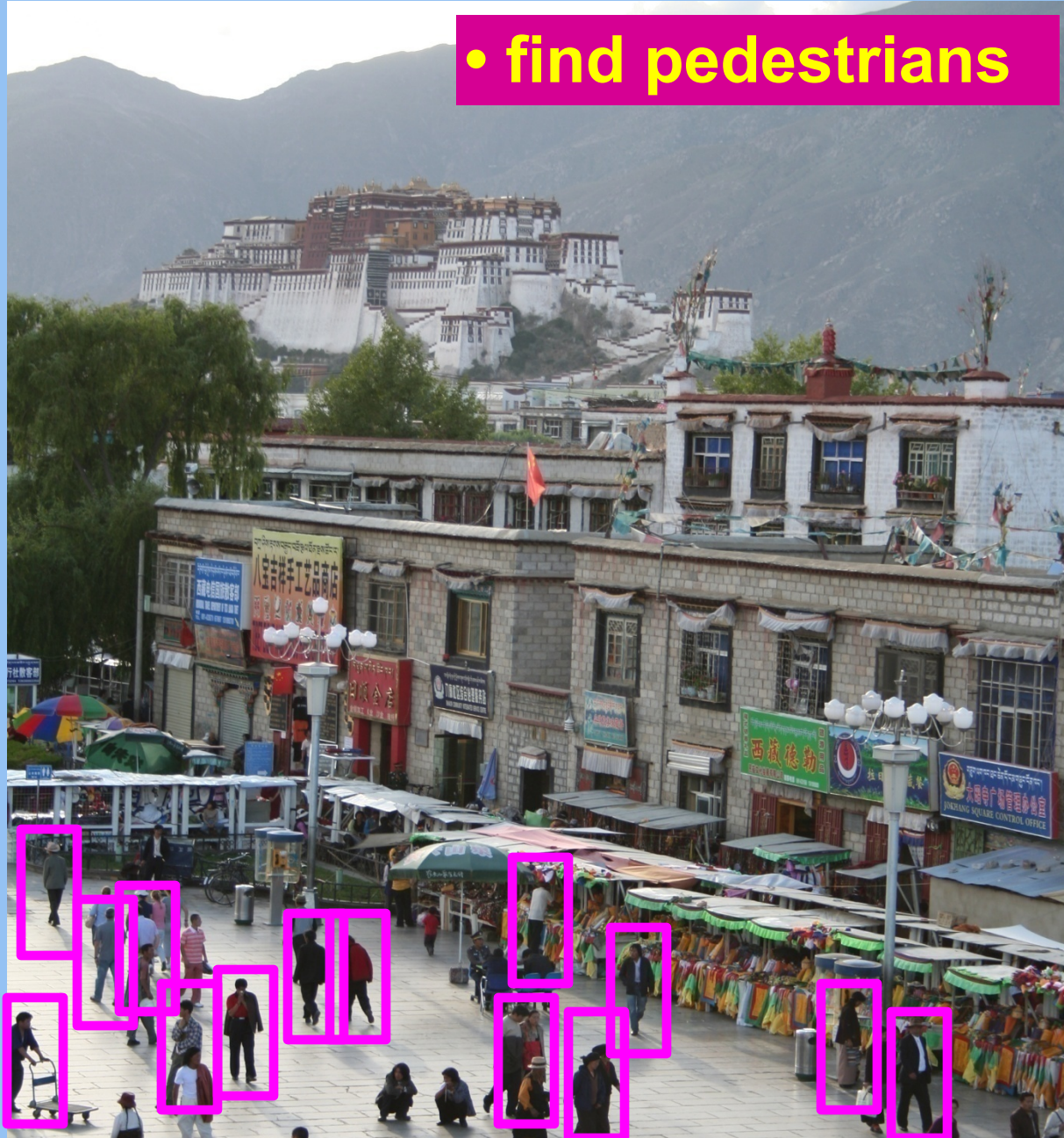- **outdoor/indoor**
- **city/forest/factory/etc.**

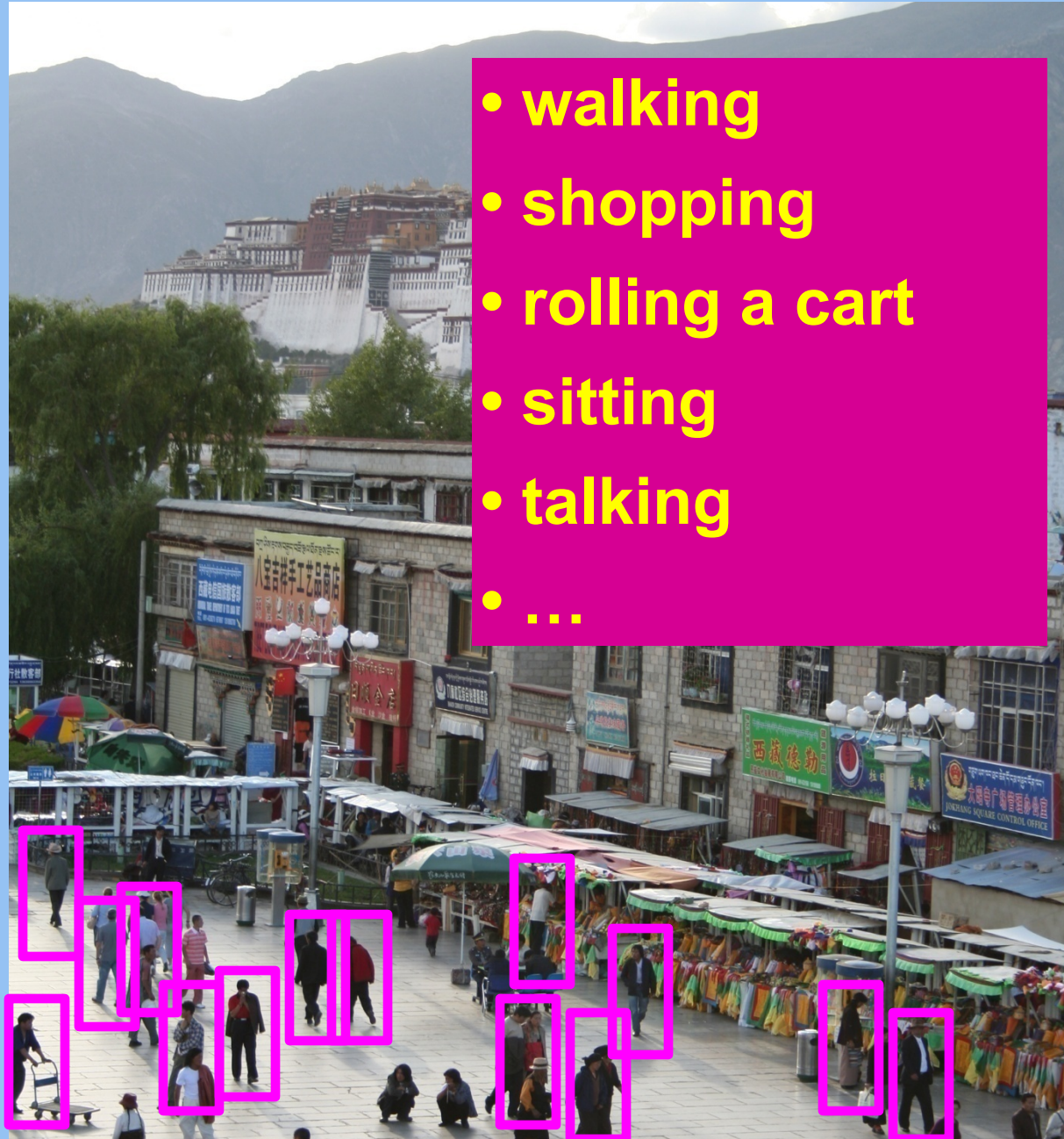# Image tagging



- **street**
- **people**
- **building**
- **mountain**
- **…**

# Object detection



- **find pedestrians**

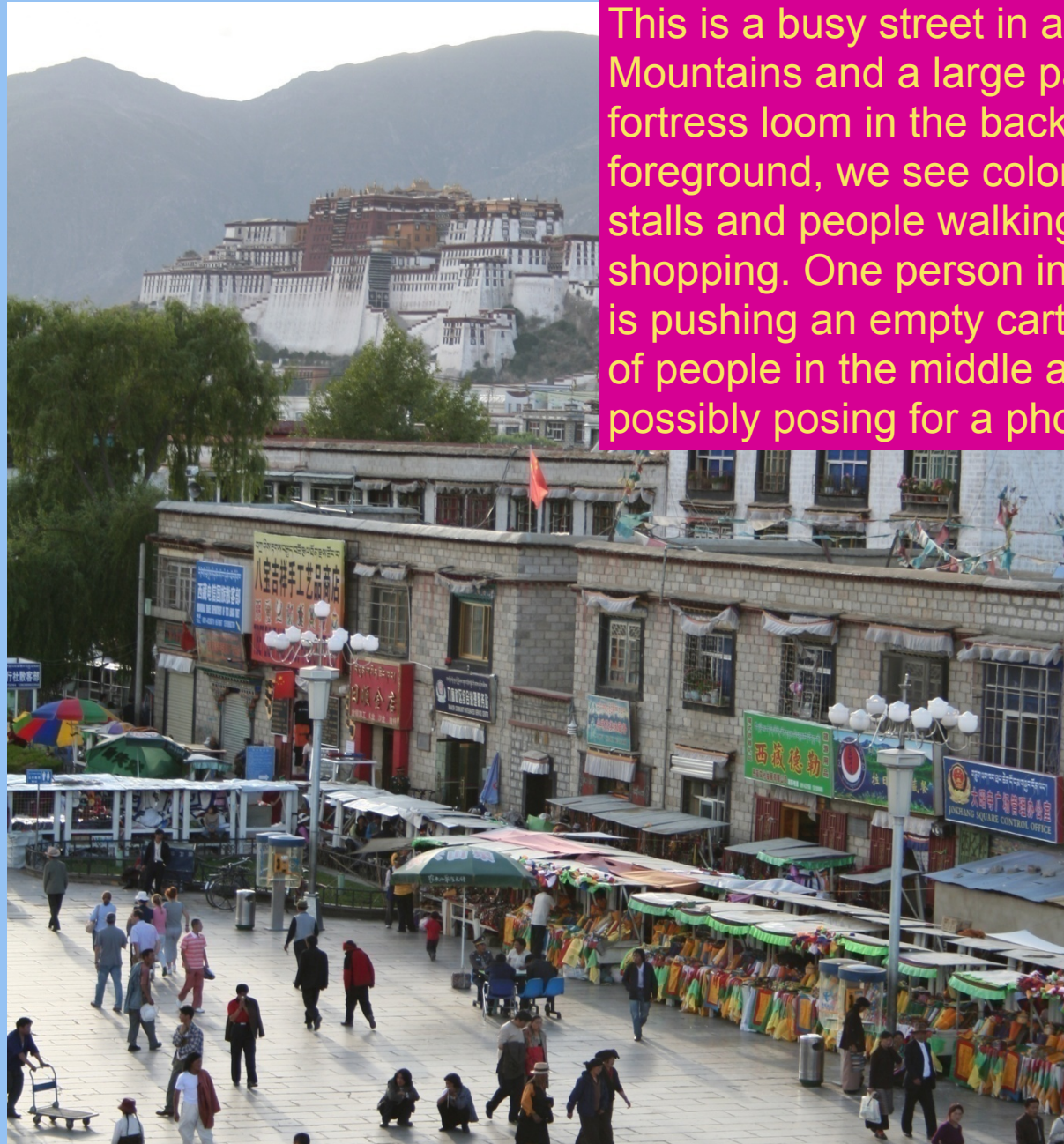# Activity recognition



- **walking**
- **shopping**
- **rolling a cart**
- **sitting**
- **talking**
- **…**

# Image parsing

# Image description



This is a busy street in an Asian city. Mountains and a large palace or fortress loom in the background. In the foreground, we see colorful souvenir stalls and people walking around and shopping. One person in the lower left is pushing an empty cart, and a couple of people in the middle are sitting, possibly posing for a photograph.

# Image classification

# The statistical learning framework

- Apply a prediction function to a feature representation of the image to get the desired output:

$$f(\text{ }) = \text{``apple''}$$

$$f(\text{ }) = \text{``tomato''}$$

$$f(\text{ }) = \text{``cow''}$$

# The statistical learning framework

$$y = f(\mathbf{x})$$

output     prediction function     Image feature

- **Training:** given a *training set* of labeled examples $\{(\mathbf{x}_1, y_1), \ldots, (\mathbf{x}_N, y_N)\}$, estimate the prediction function $f$ by minimizing the prediction error on the training set
- **Testing:** apply $f$ to a never before seen *test example* $\mathbf{x}$ and output the predicted value $y = f(\mathbf{x})$

# Steps

**Training**



**Testing**

Slide credit: D. Hoiem

# Traditional recognition pipeline

Image Pixels $\Rightarrow$ **Hand-designed feature extraction** $\Rightarrow$ **Trainable classifier** $\Rightarrow$ Object Class
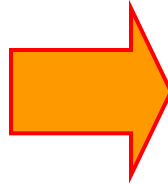
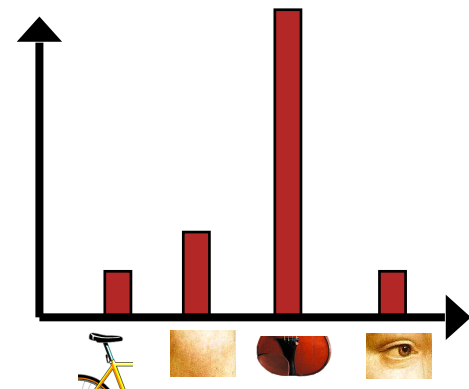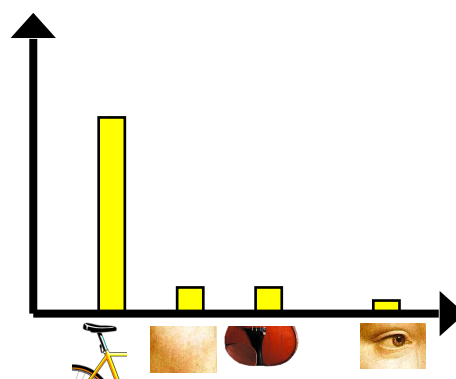- Features are not learned
- Trainable classifier is often generic (e.g. SVM)

# Bags of features
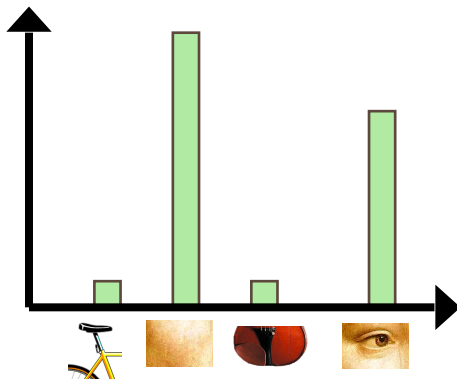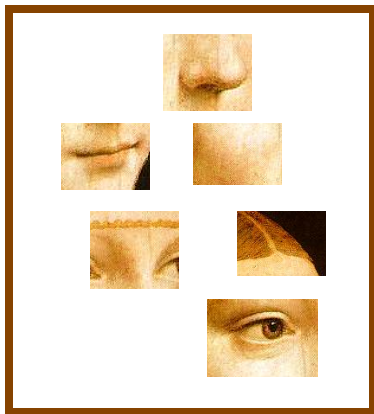
# Traditional features: Bags-of-features

1. Extract local features
2. Learn "visual vocabulary"
3. Quantize local features using visual vocabulary
4. Represent images by frequencies of "visual words"

# Texture recognition



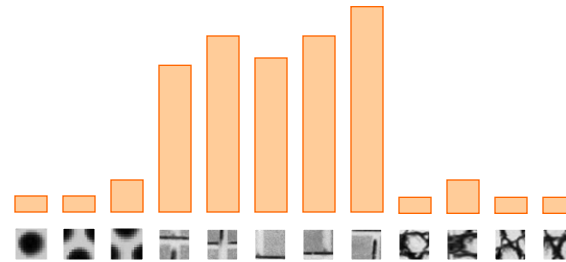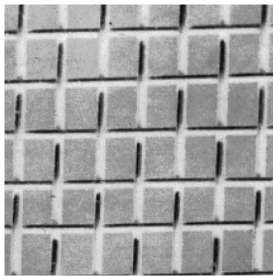Julesz 1981; Cula & Dana, 2001; Leung & Malik 2001; Mori, Belongie & Malik, 2001; Schmid 2001; Varma & Zisserman, 2002, 2003; Lazebnik, Schmid & Ponce, 2003

# 1. Local feature extraction

- Sample patches and extract descriptors

# Keypoints



patches surrounding keypoints

# 2. Learning the visual vocabulary

Extracted descriptors from the training set

# 2. Learning the visual vocabulary

# 2. Learning the visual vocabulary

Visual vocabulary
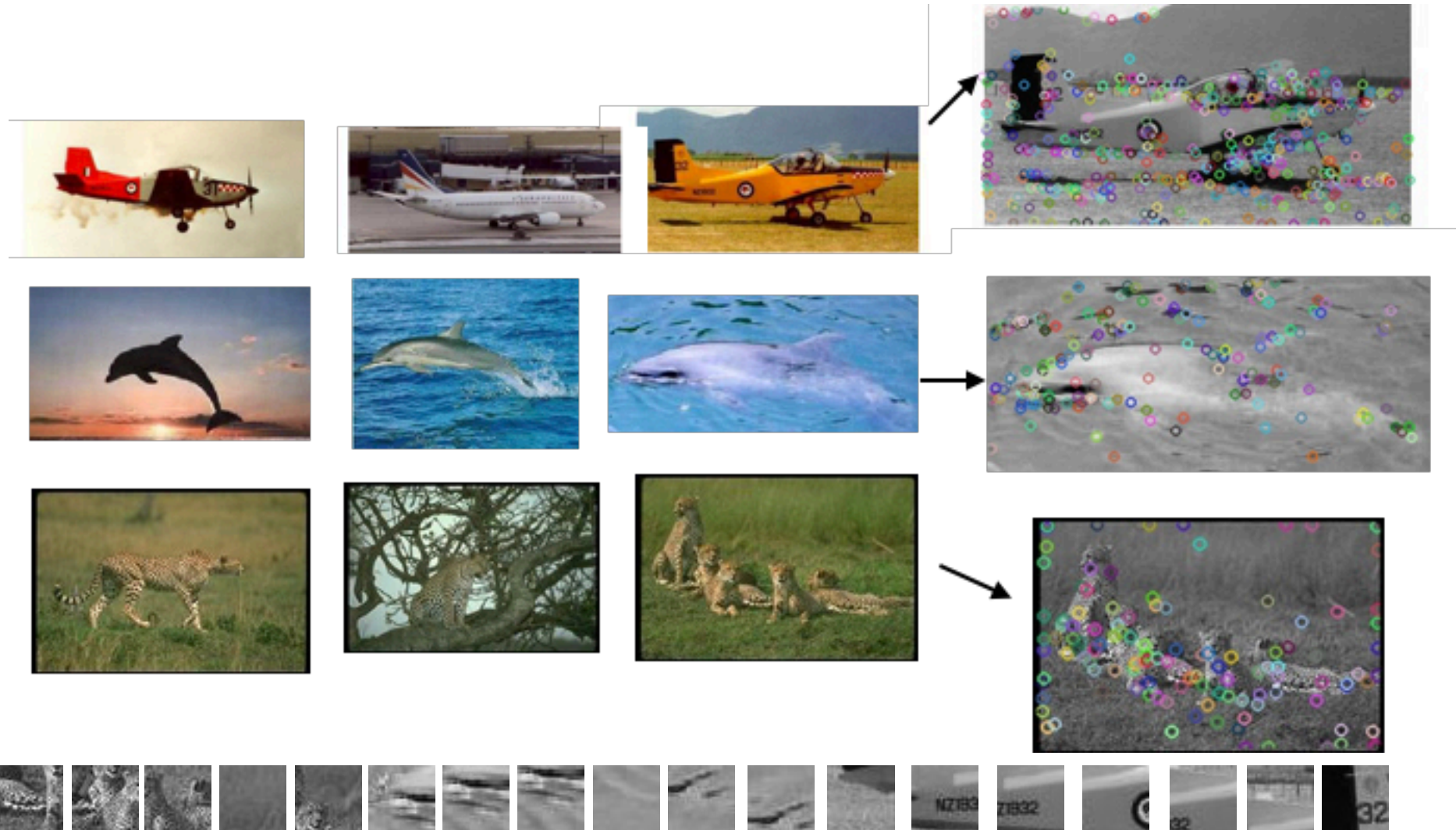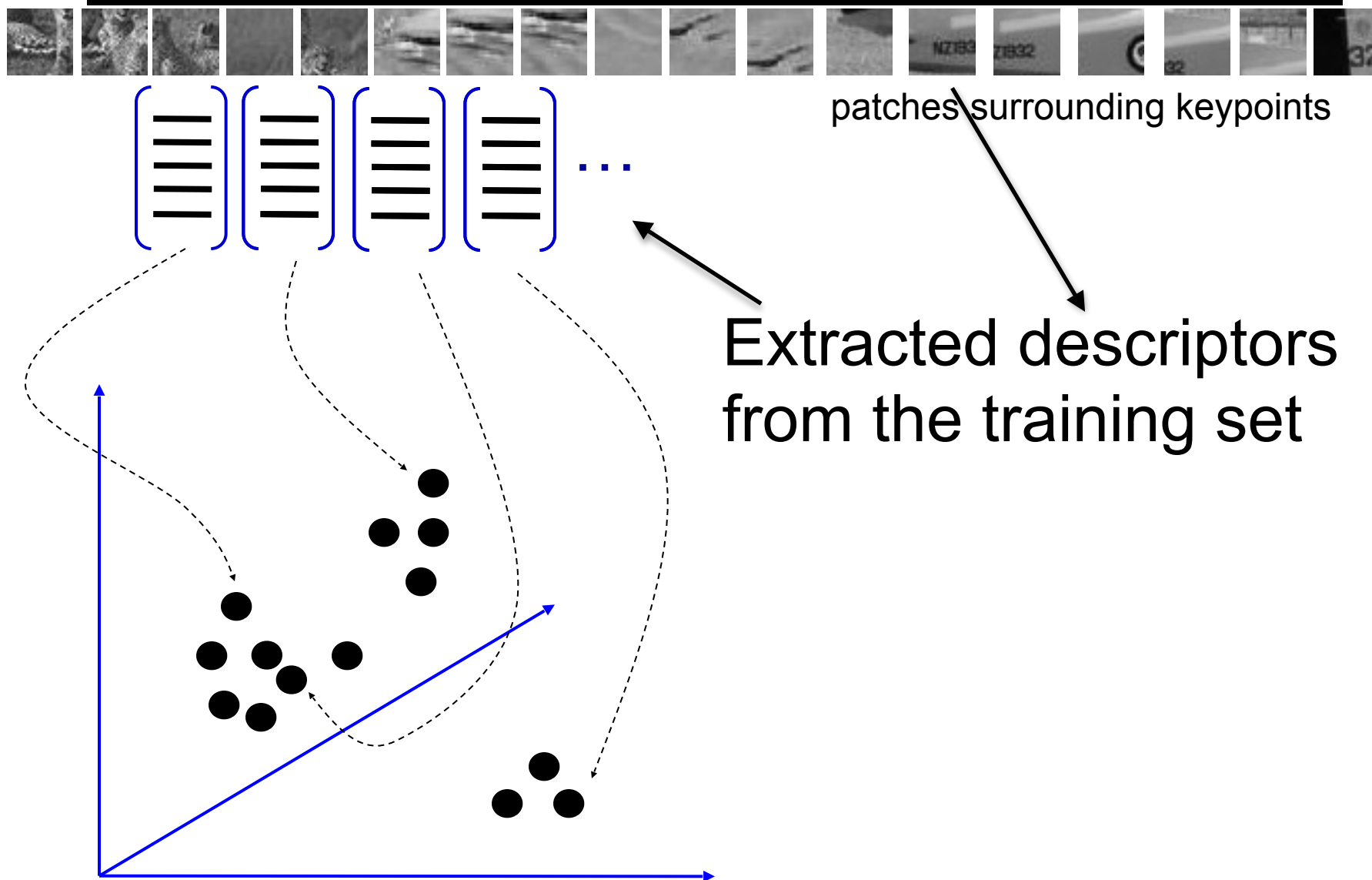
Clustering

# Review: K-means clustering

- Want to minimize sum of squared Euclidean distances between features $\mathbf{x}_i$ and their nearest cluster centers $\mathbf{m}_k$

$$D(X, M) = \sum_{\text{cluster } k} \sum_{\substack{\text{point } i \text{ in} \\ \text{cluster } k}} (\mathbf{x}_i - \mathbf{m}_k)^2$$

Algorithm:

- Randomly initialize K cluster centers

- Iterate until convergence:
  - Assign each feature to the nearest center
  - Recompute each cluster center as the mean of all features assigned to it

# Example visual vocabulary



**Appearance codebook**

Source: B. Leibe

# Bag-of-features steps

1. Extract local features
2. Learn "visual vocabulary"
3. **Quantize local features using visual vocabulary**
4. **Represent images by frequencies of "visual words"**

# Traditional recognition pipeline

Image Pixels → **Hand-designed feature extraction** → **Trainable classifier** → Object Class

# Classifiers: Nearest neighbor



Training examples from class 1

Test example

Training examples from class 2

f(**x**) = label of the training example nearest to **x**

All we need is a distance function for our inputs
No training required!

# K-nearest neighbor classifier

- For a new point, find the k closest points from training data
- Vote for class label with labels of the k points

# K-nearest neighbor classifier



Credit: Andrej Karpathy, http://cs231n.github.io/classification/

# K-nearest neighbor classifier



Left: Example images from the CIFAR-10 dataset. Right: first column shows a few test images and next to each we show the top 10 nearest neighbors in the training set according to pixel-wise difference.

Credit: Andrej Karpathy, http://cs231n.github.io/classification/

# Linear classifiers



Find a *linear function* to separate the classes:

$$f(\mathbf{x}) = \text{sgn}(\mathbf{w} \cdot \mathbf{x} + b)$$

# Visualizing linear classifiers



Source: Andrej Karpathy, http://cs231n.github.io/linear-classify/

# Nearest neighbor vs. linear classifiers

- **NN pros:**
  - Simple to implement
  - Decision boundaries not necessarily linear
  - Works for any number of classes
  - *Nonparametric* method
- **NN cons:**
  - Need good distance function
  - Slow at test time
- **Linear pros:**
  - Low-dimensional *parametric* representation
  - Very fast at test time
- **Linear cons:**
  - Works for two classes
  - How to train the linear function?
  - What if data is not linearly separable?

# Support vector machines

- When the data is linearly separable, there may be more than one separator (hyperplane)



Which separator is best?

# Support vector machines

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

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

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

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

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

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

Support vectors

Margin

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

# Finding the maximum margin hyperplane

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

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

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

*Quadratic optimization problem*:

$$\min_{\mathbf{w},b} \frac{1}{2}\|\mathbf{w}\|^2 \quad \text{subject to} \quad y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$$

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

# SVM parameter learning

- Separable data:
$$\min_{\mathbf{w},b} \frac{1}{2}\|\mathbf{w}\|^2 \quad \text{subject to} \quad y_i(\mathbf{w}\cdot\mathbf{x}_i + b) \geq 1$$

Maximize margin

Classify training data correctly

- Non-separable data:
$$\min_{\mathbf{w},b} \quad \frac{1}{2}\|\mathbf{w}\|^2 + C\sum_{i=1}^{n}\max\left(0, 1 - y_i(\mathbf{w}\cdot\mathbf{x}_i + b)\right)$$
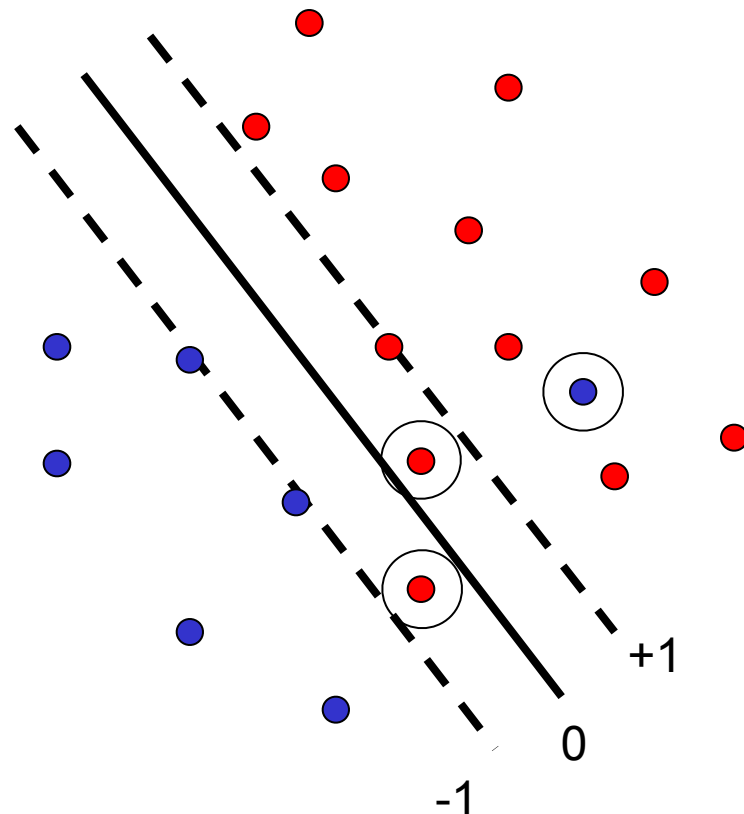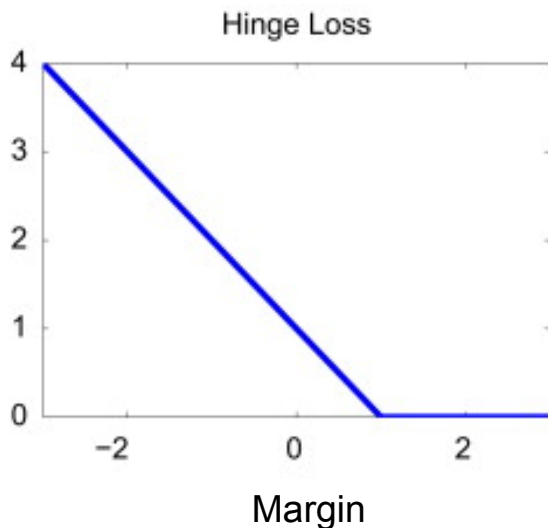
Maximize margin

Minimize classification mistakes

# SVM parameter learning

$$\min_{\mathbf{w},b} \quad \frac{1}{2}\|\mathbf{w}\|^2 + C\sum_{i=1}^{n}\max\left(0, 1 - y_i(\mathbf{w}\cdot\mathbf{x}_i + b)\right)$$



Hinge Loss

Margin

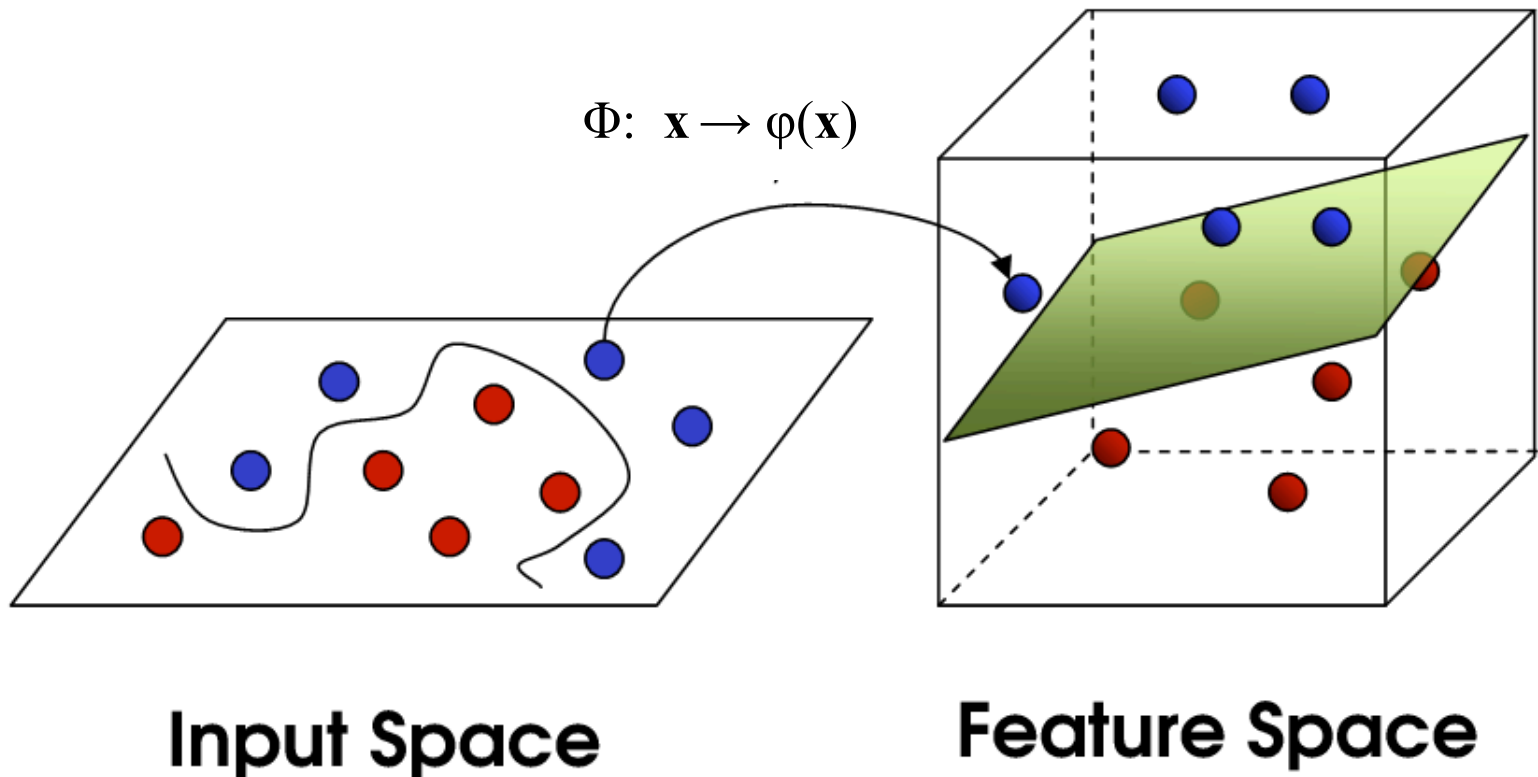+1

0

-1

Demo: http://cs.stanford.edu/people/karpathy/svmjs/demo

# Nonlinear SVMs

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

$$\Phi: \ \mathbf{x} \to \varphi(\mathbf{x})$$

**Input Space**　　　**Feature Space**

Image source

# Nonlinear SVMs

- Linearly separable dataset in 1D:

- Non-separable dataset in 1D:

- We can map the data to a *higher-dimensional space*:
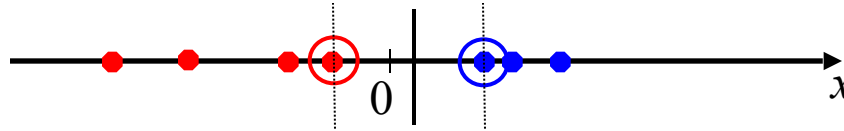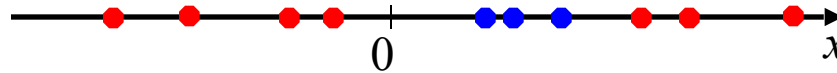
Slide credit: Andrew Moore

# The kernel trick

- General idea: the original input space can always be mapped to some higher-dimensional feature space where the training set is separable
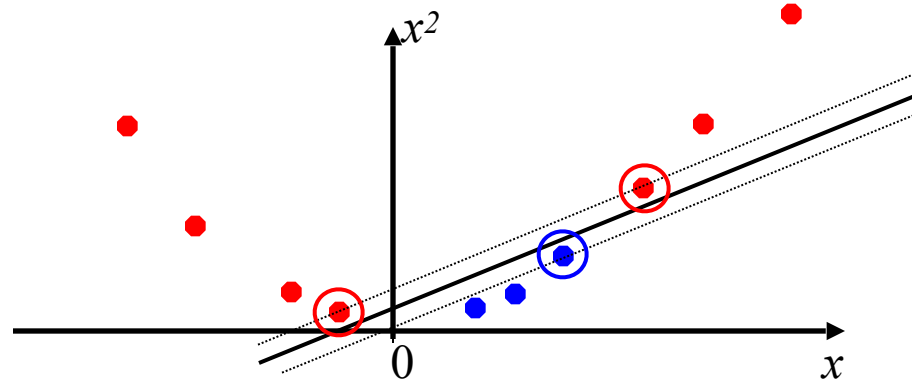
- **The kernel trick:** instead of explicitly computing the lifting transformation $\varphi(\mathbf{x})$, define a kernel function K such that

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

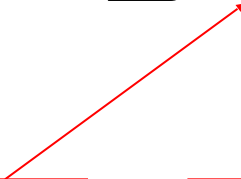(to be valid, the kernel function must satisfy *Mercer's condition*)

# The kernel trick

- Linear SVM decision function:

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

learned weight

Support vector

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

# The kernel trick

- Linear SVM decision function:

$$\mathbf{w} \cdot \mathbf{x} + b = \sum_i \alpha_i y_i \mathbf{x}_i \cdot \mathbf{x} + b$$
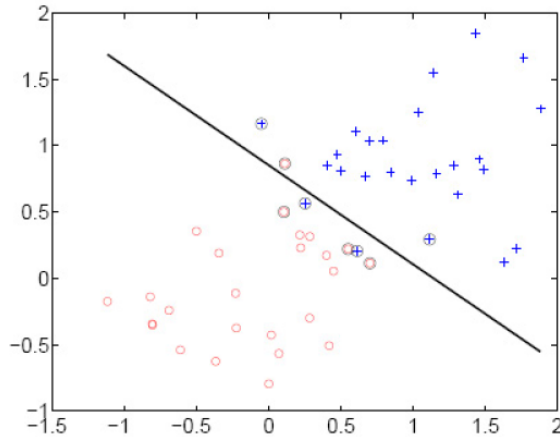
- Kernel SVM decision function:

$$\sum_i \alpha_i y_i \varphi(\mathbf{x}_i) \cdot \varphi(\mathbf{x}) + b = \sum_i \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + b$$
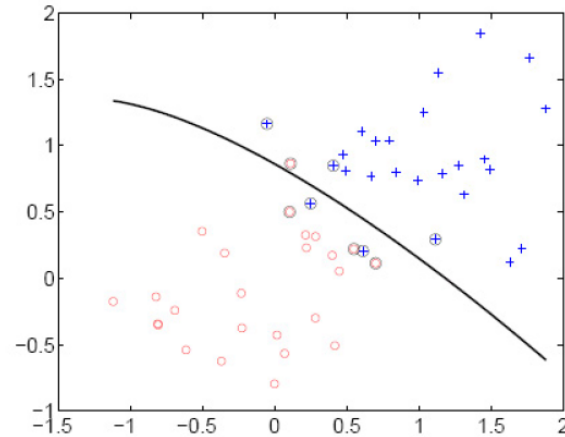
- This gives a nonlinear decision boundary in the original feature space

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

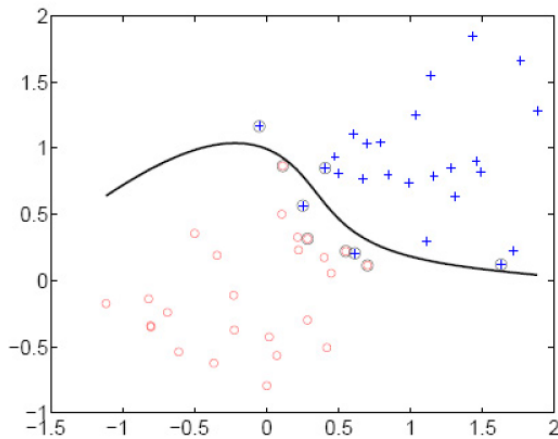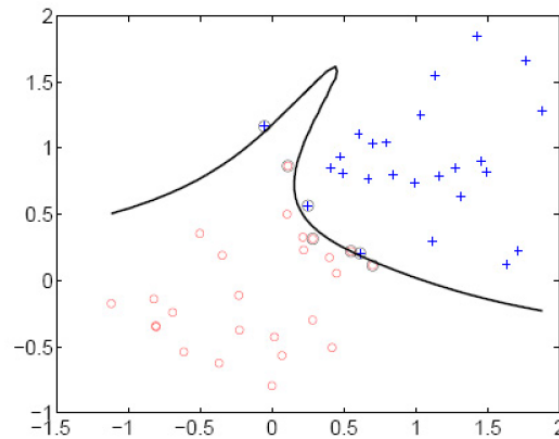# Polynomial kernel: $K(\mathbf{x}, \mathbf{y}) = (c + \mathbf{x} \cdot \mathbf{y})^d$



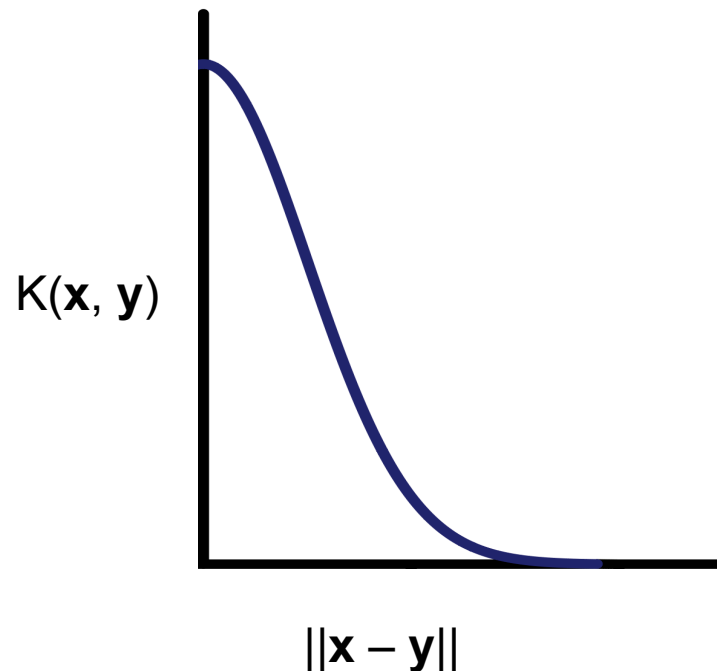linear

$2^{nd}$ order polynomial

$4^{th}$ order polynomial
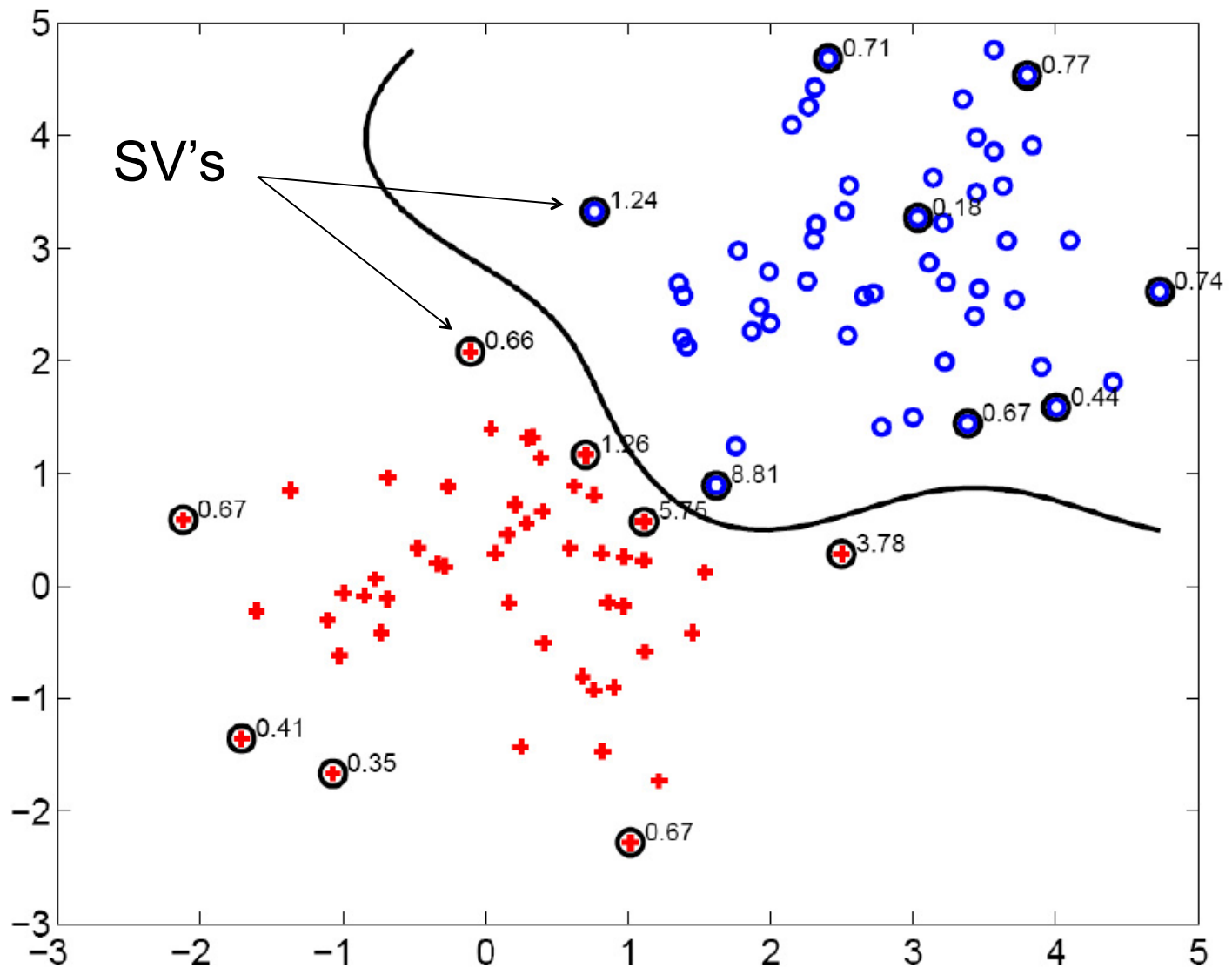
$8^{th}$ order polynomial

# Gaussian kernel

- Also known as the radial basis function (RBF) kernel:

$$K(\mathbf{x}, \mathbf{y}) = \exp\left( -\frac{1}{\sigma^2} \|\mathbf{x} - \mathbf{y}\|^2 \right)$$

K(**x**, **y**)

||**x** − **y**||

# Gaussian kernel

# Kernels for histograms

- Histogram intersection:

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

- Square root (Bhattacharyya kernel):

$$K(h_1, h_2) = \sum_{i=1}^{N} \sqrt{h_1(i) \, h_2(i)}$$

# SVMs: Pros and cons

- Pros
  - Kernel-based framework is very powerful, flexible
  - Training is convex optimization, globally optimal solution can be found
  - Amenable to theoretical analysis
  - SVMs work very well in practice, even with very small training sample sizes

- Cons
  - No "direct" multi-class SVM, must combine two-class SVMs (e.g., with one-vs-others)
  - Computation, memory (esp. for nonlinear SVMs)

# Generalization

- Generalization refers to the ability to correctly classify never before seen examples

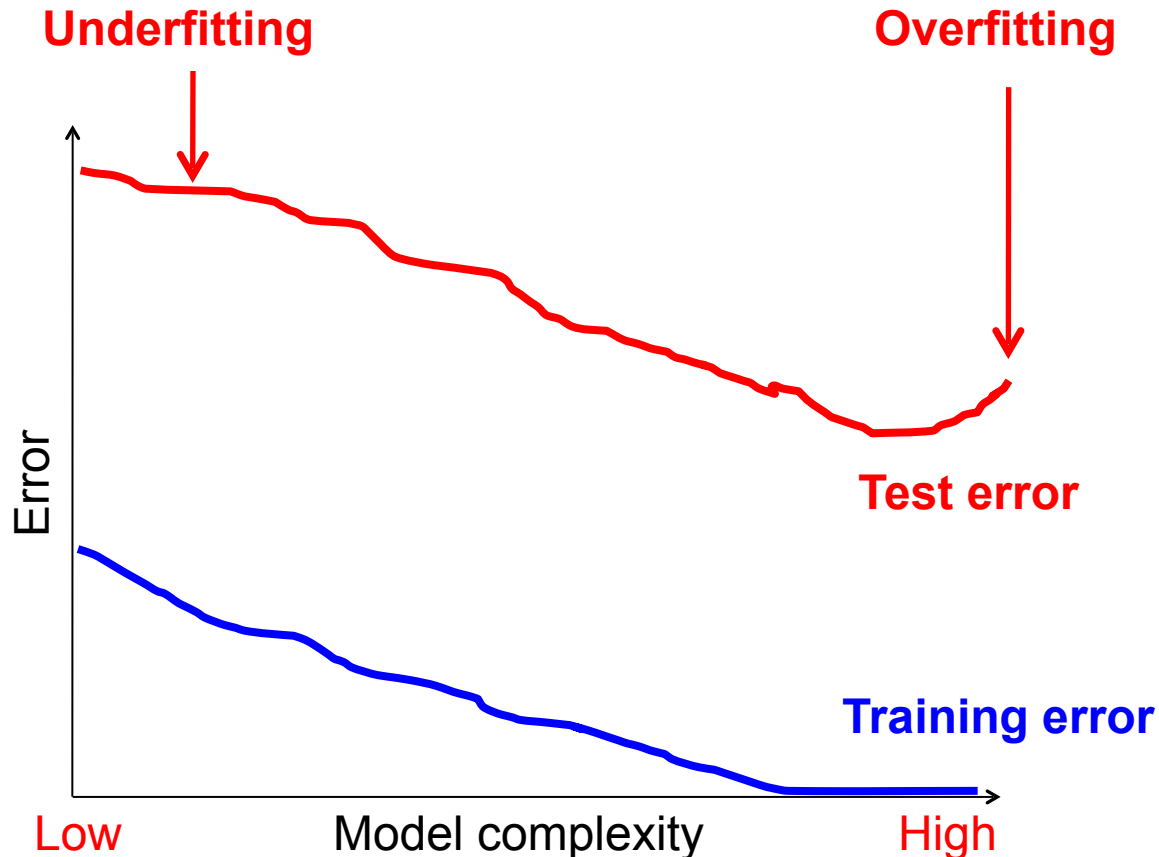- Can be controlled by turning "knobs" that affect the complexity of the model



Training set (labels known)

Test set (labels unknown)

# Diagnosing generalization ability

- **Training error:** how does the model perform on the data on which it was trained?

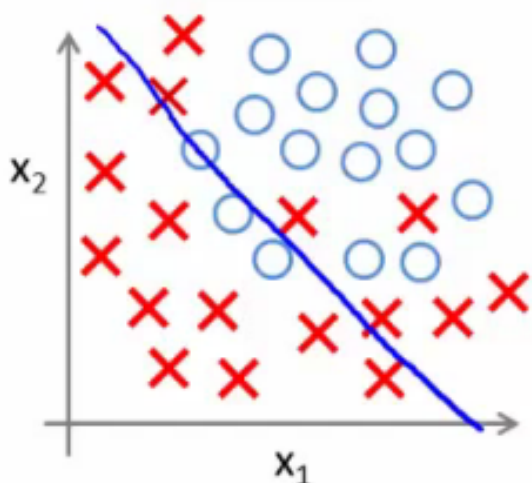- **Test error:** how does it perform on never before seen data?
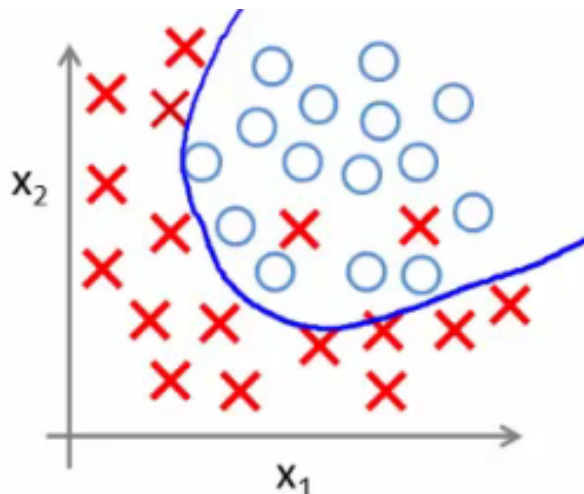


Source: D. Hoiem

# Underfitting and overfitting

- **Underfitting:** training and test error are both *high*
  - Model does an equally poor job on the training and the test set
  - Either the training procedure is ineffective or the model is too "simple" to represent the data
- **Overfitting:** Training error is *low* but test error is *high*
  - Model fits irrelevant characteristics (noise) in the training data
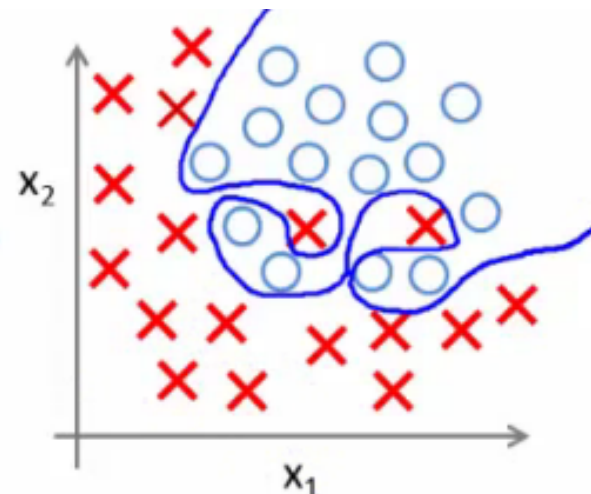  - Model is too complex or amount of training data is insufficient

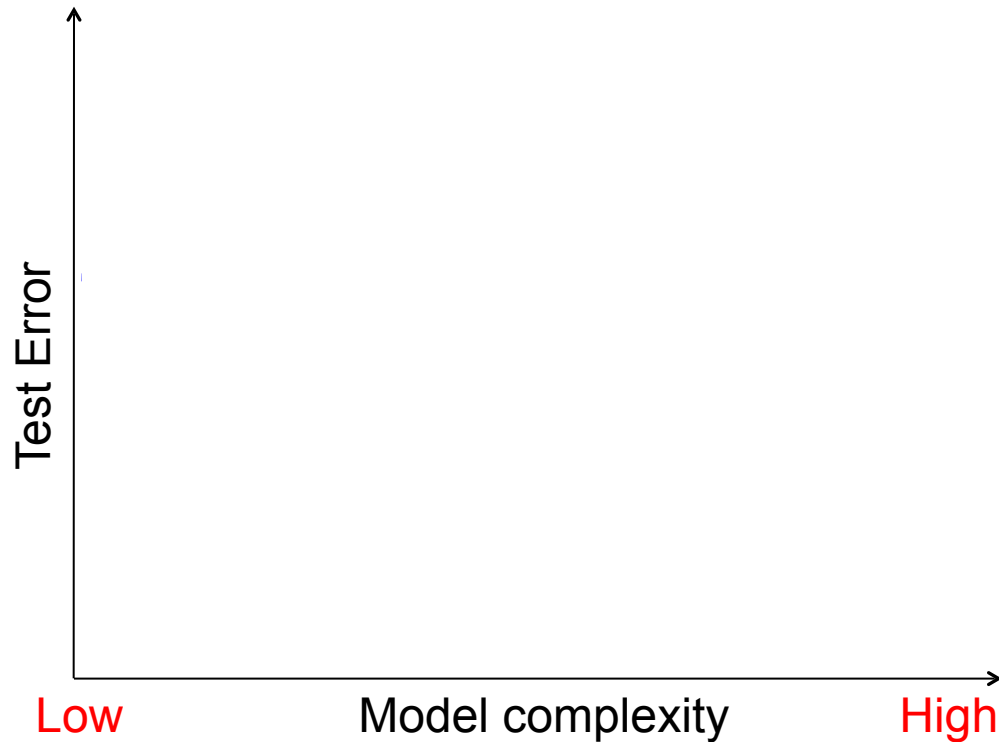Underfitting          Good generalization          Overfitting

# Effect of training set size
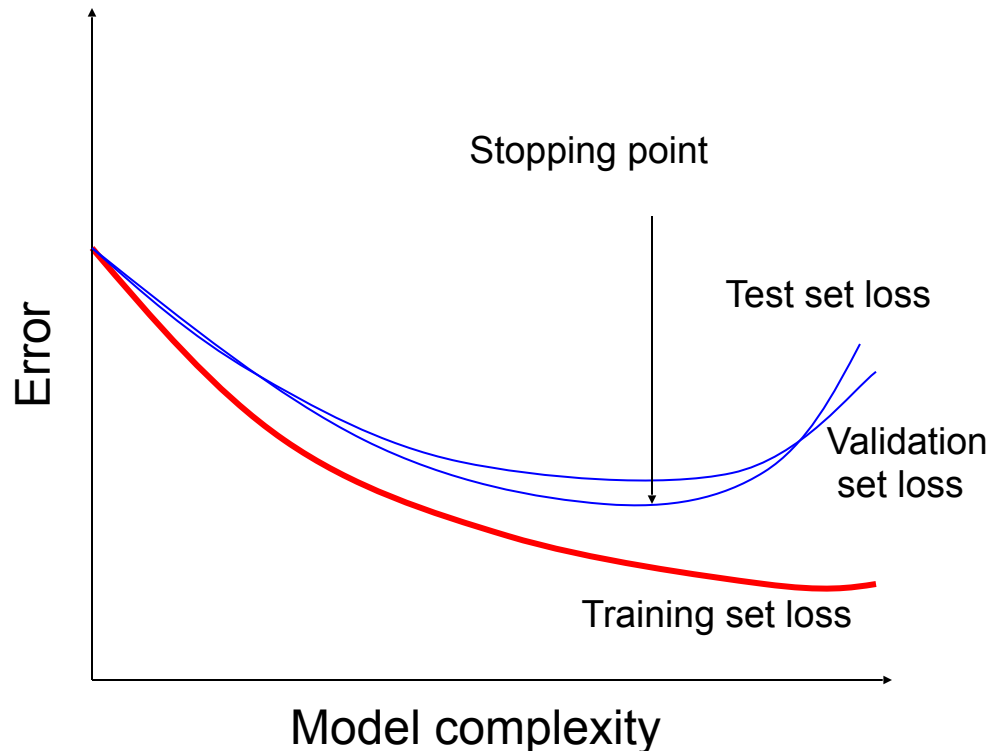


Test Error

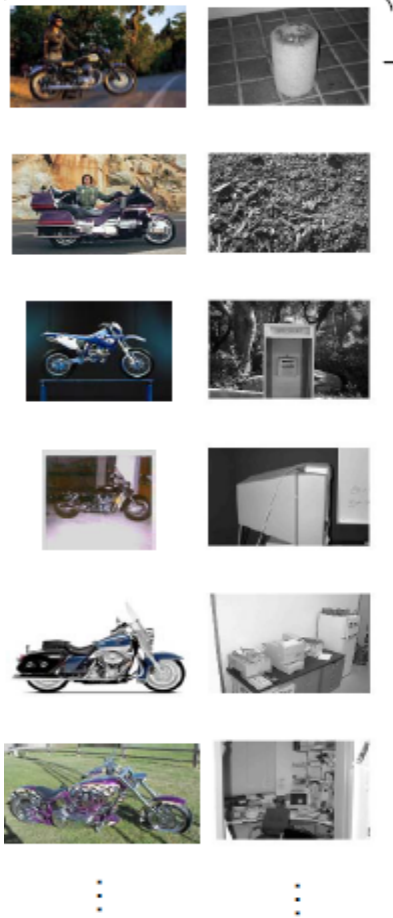Low    Model complexity    High

# Validation

- Split the data into **training**, **validation**, and **test** subsets
- Use training set to **optimize model parameters**
- Use validation test to **choose the best model**
- Use test set only to **evaluate performance**

# Summary

**The different steps**