

# Recognition, Detection, and Classification

(Part 1)

Computer Vision I

CSE 252A

Lecture 13

# Announcements

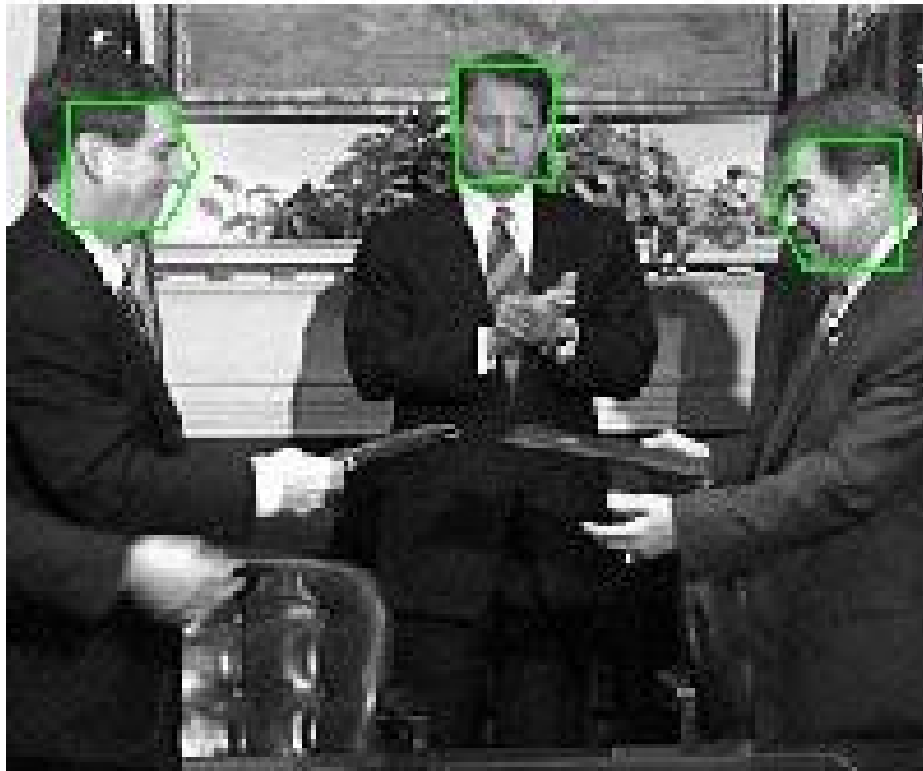
- Assignment 3 is due Nov 22, 11:59 PM
- Assignment 4 will be released Nov 22
  - Due Dec 6, 11:59 PM

# Recognition



Given a database of objects and an image determine what, if any of the objects are present in the image.

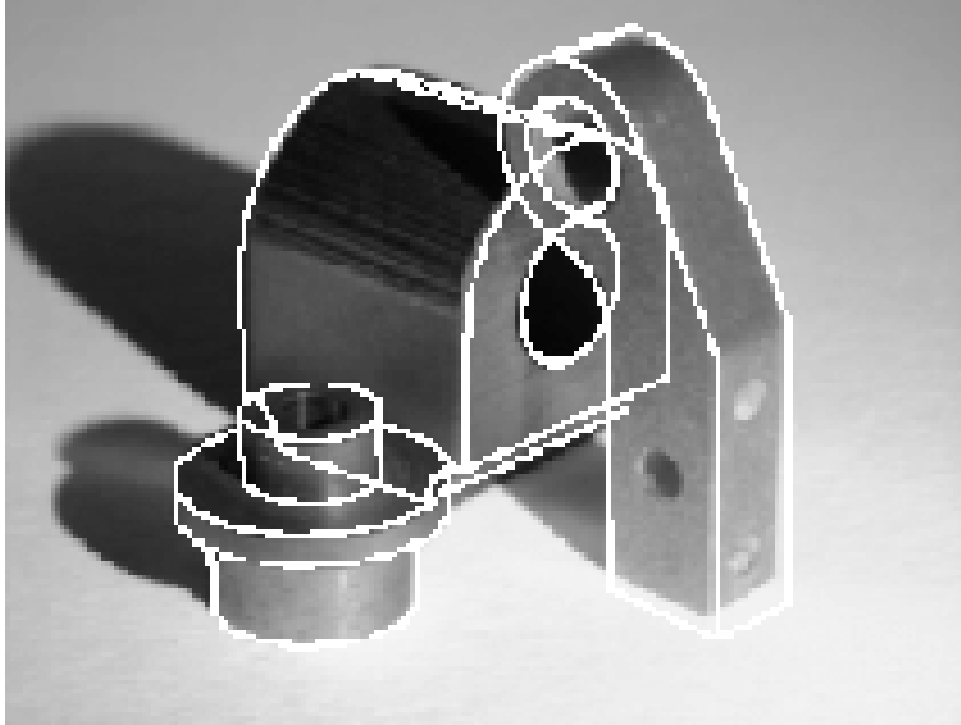
# Recognition



Given a database of objects and an image determine what, if any of the objects are present in the image.

... and where are they and how many

# Recognition

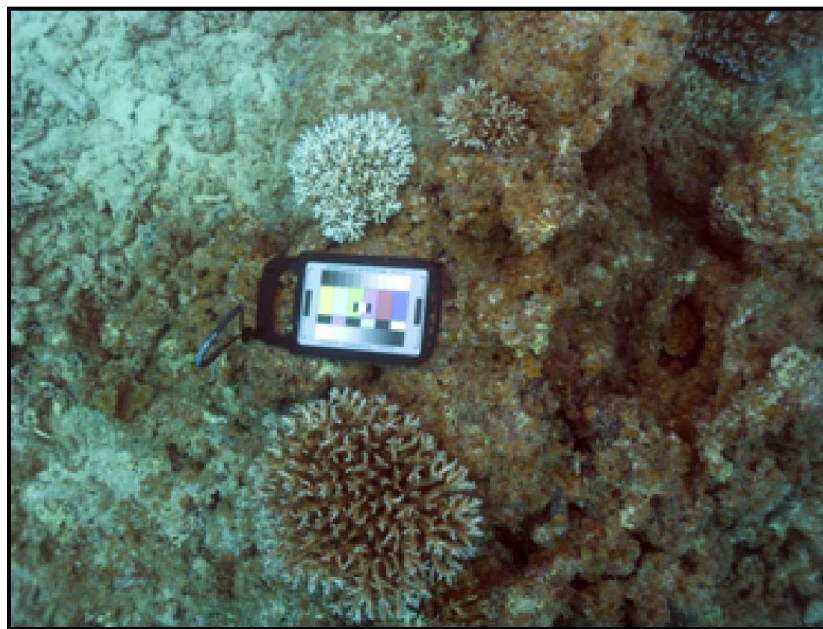


Given a database of objects and an image determine what, if any of the objects are present in the image.

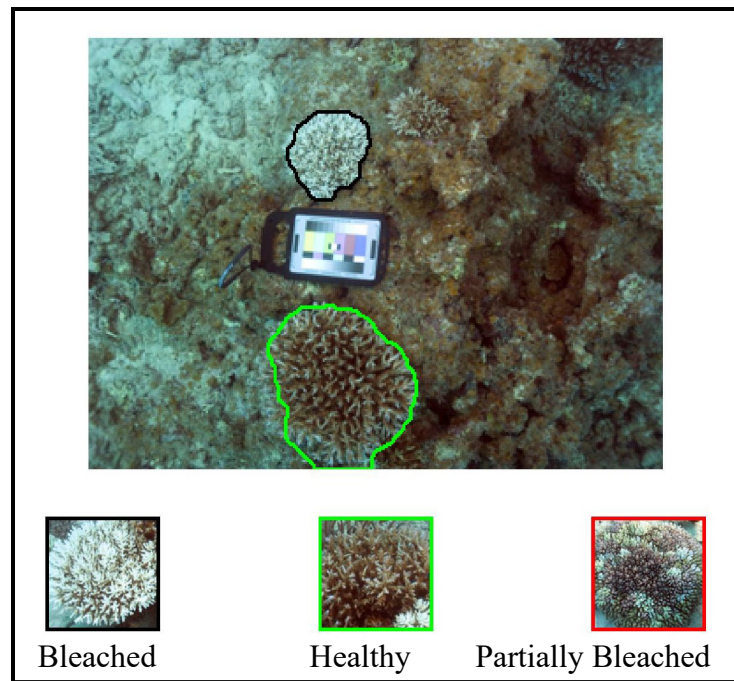
... and where are they and how many

... and what is 3D pose

**Example: where are the coral heads and which ones are healthy and which are bleached?**

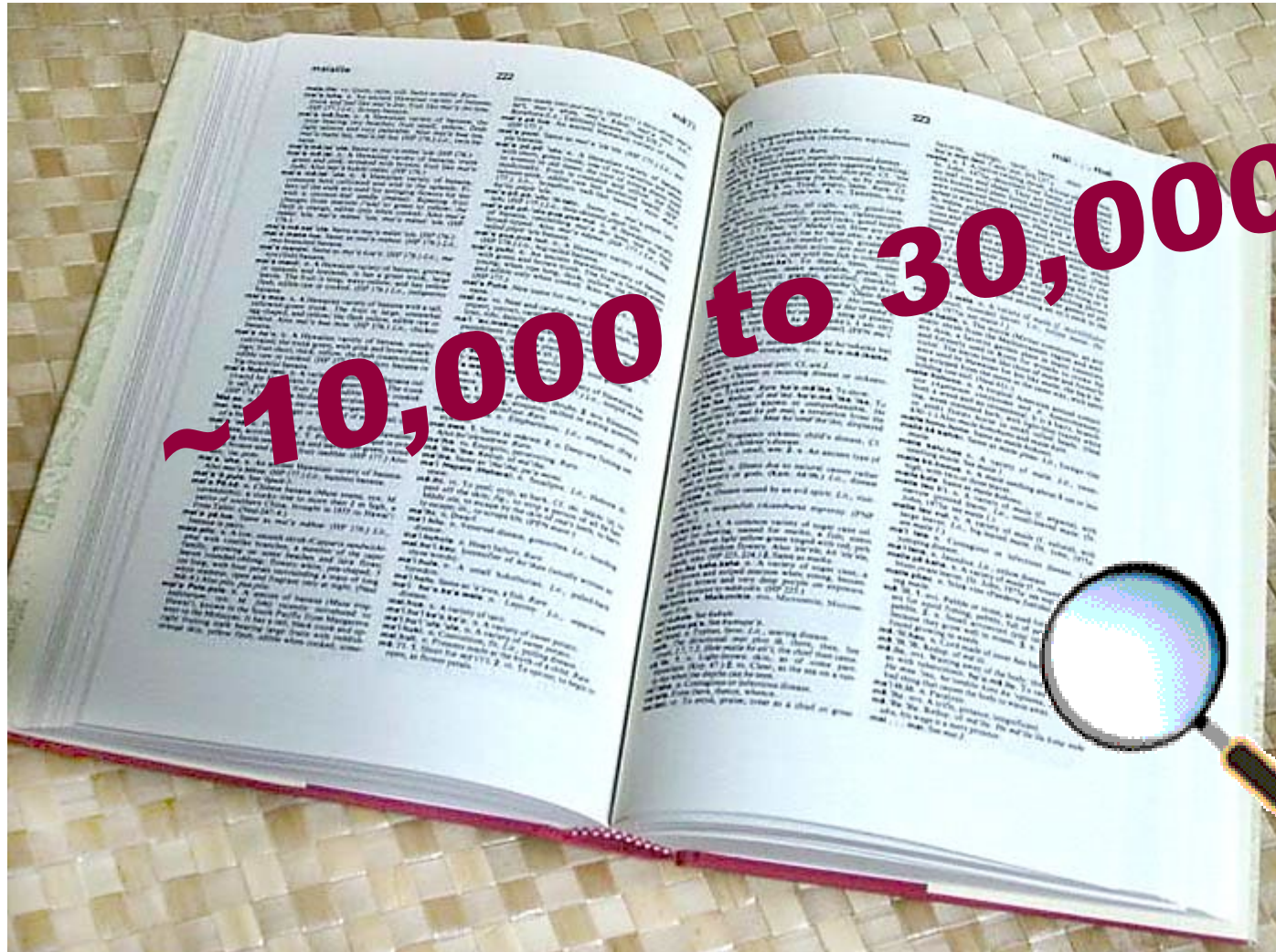


**Input Image**



**Segmented/labelled Image**

# How many visual object categories are there?



Biederman 1987

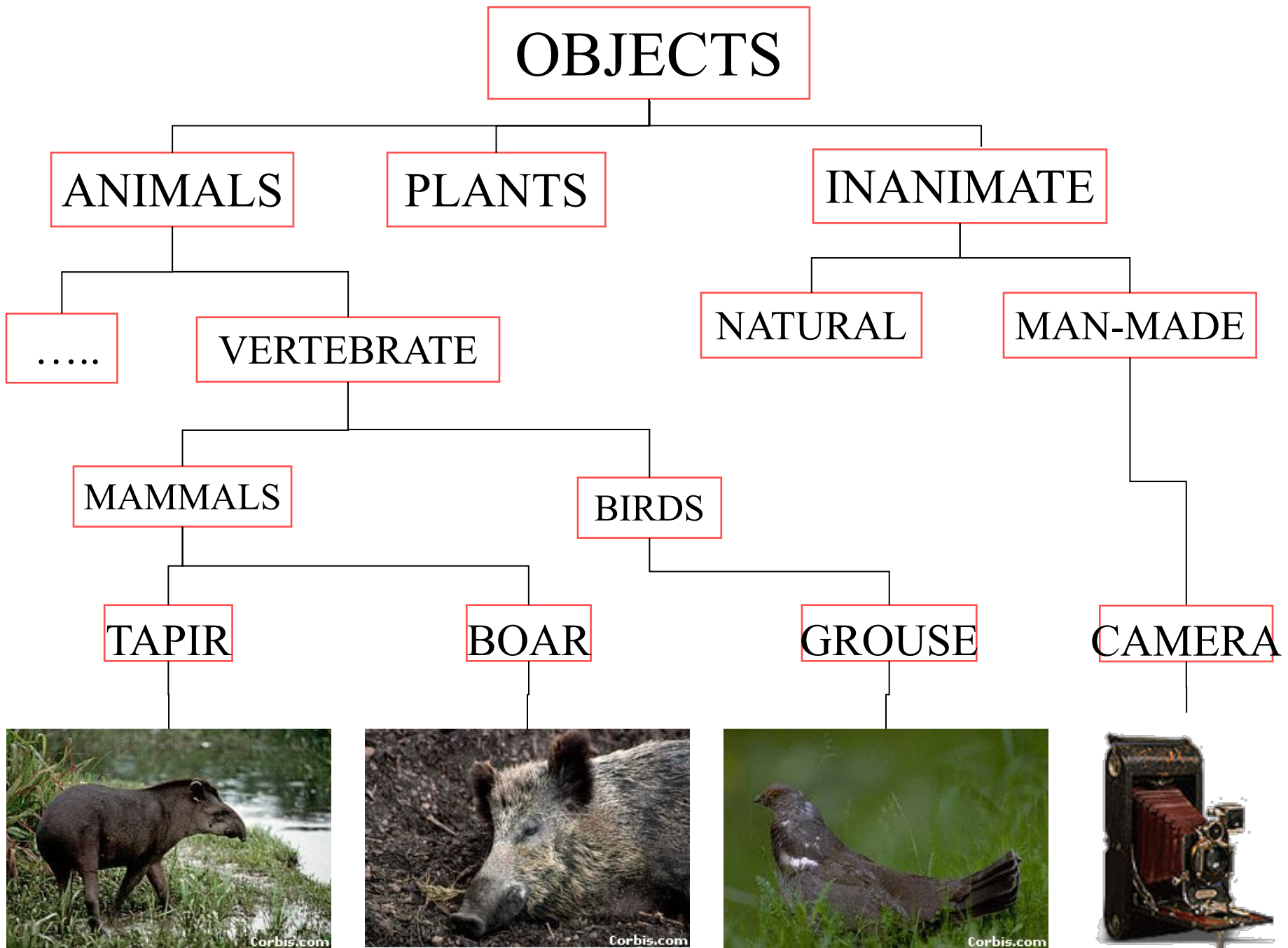
Computer Vision I



~10,000 to 30,000





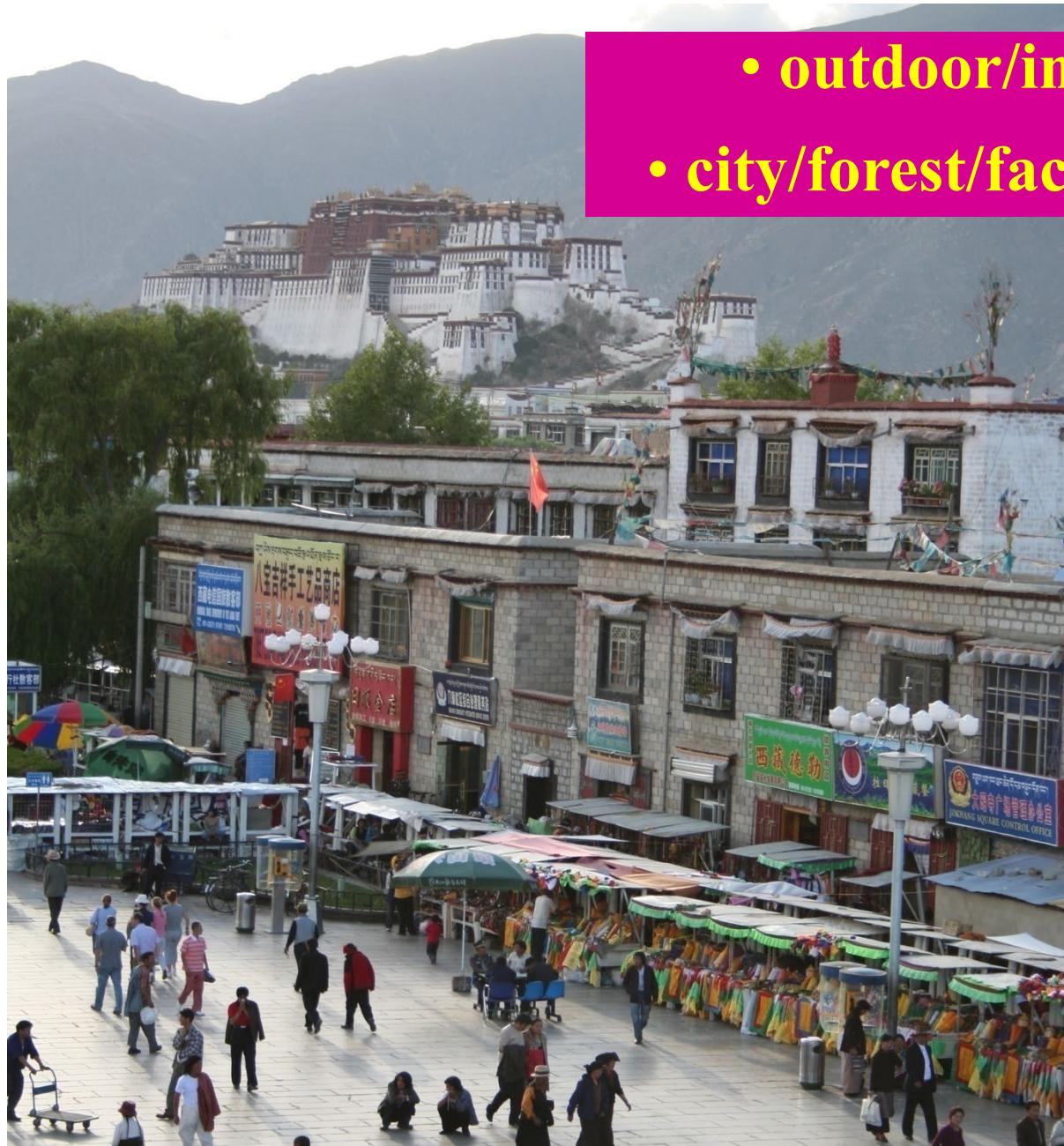


# Specific recognition tasks

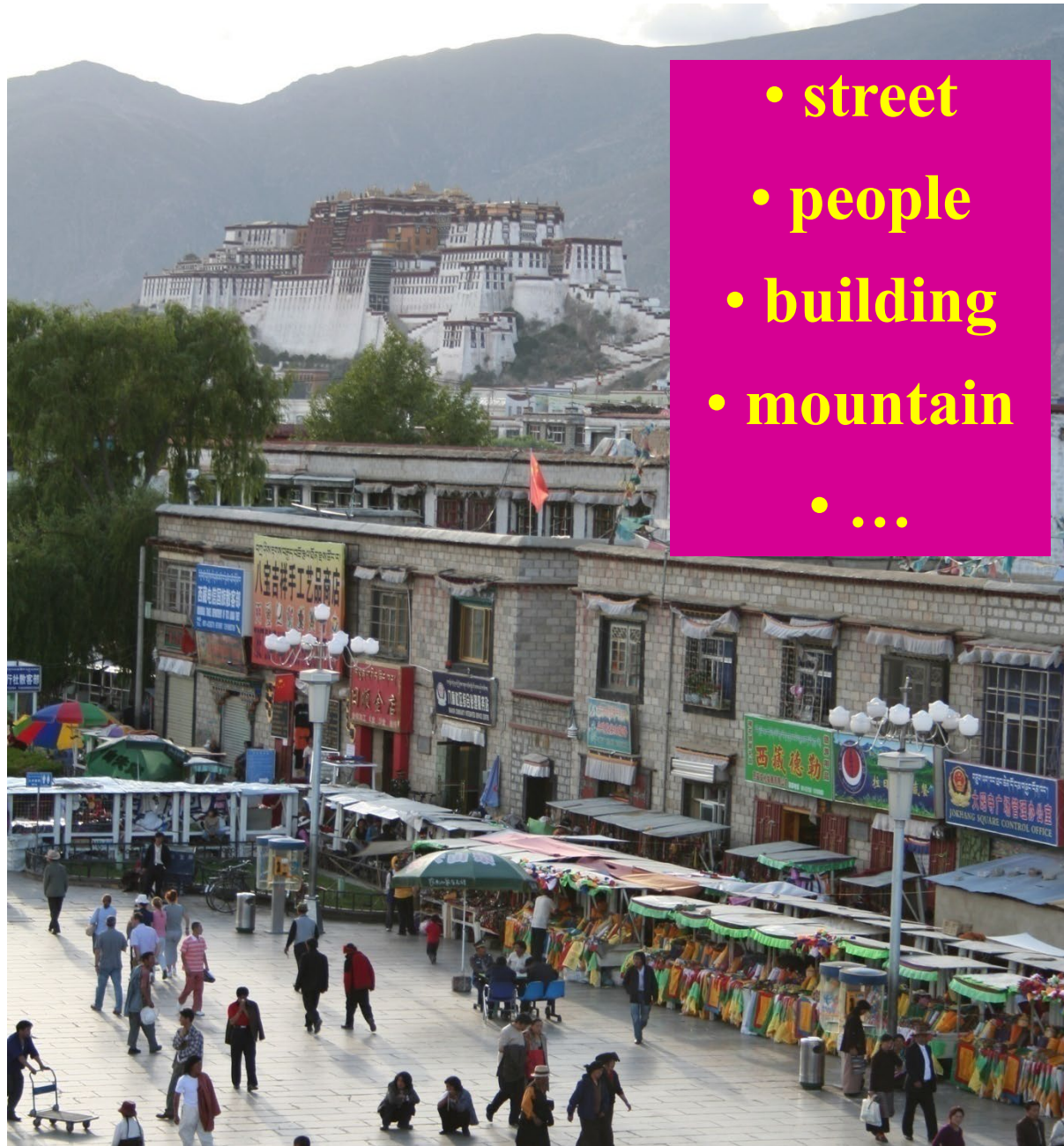


# Scene categorization

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



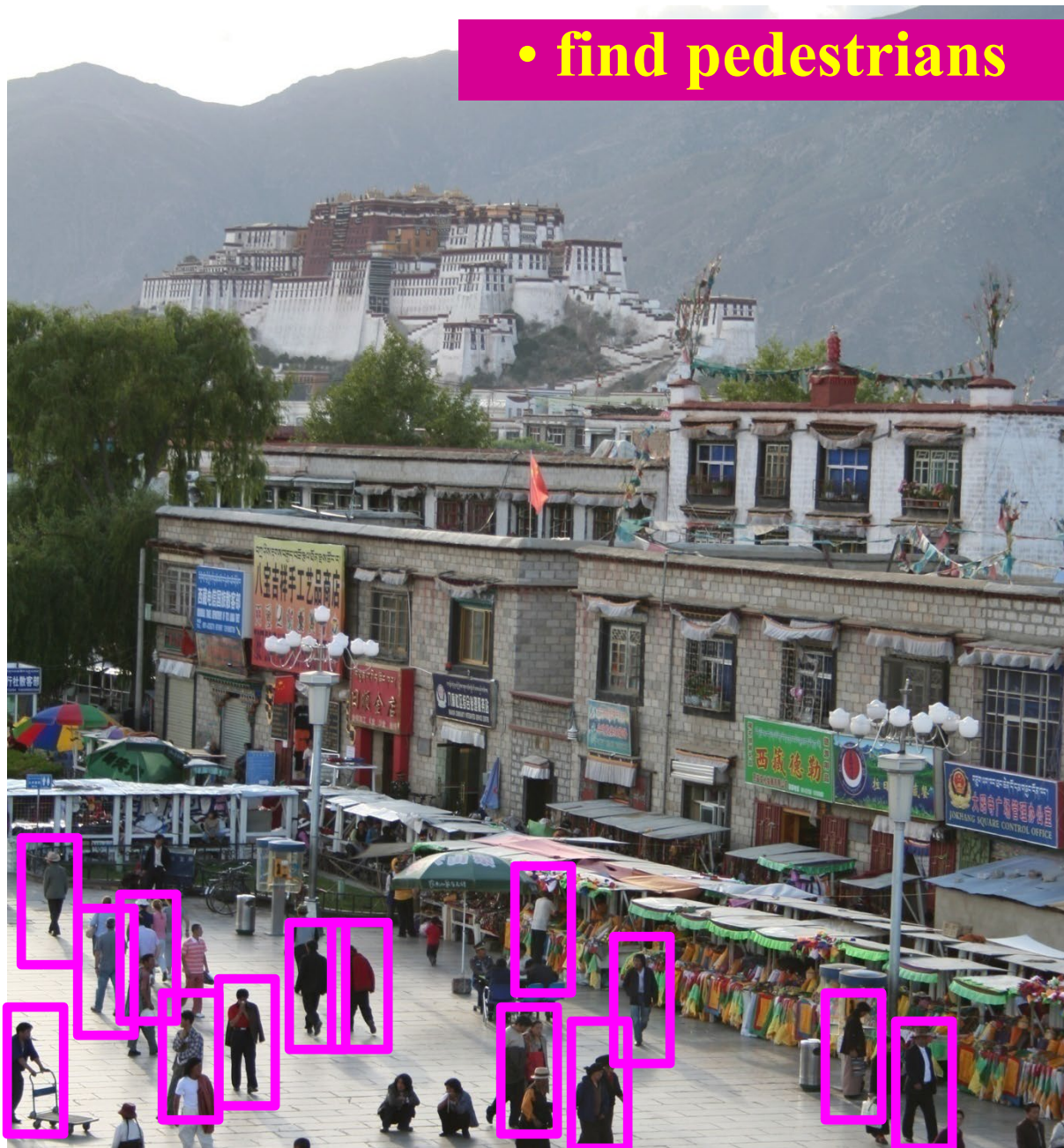
# Image annotation/tagging



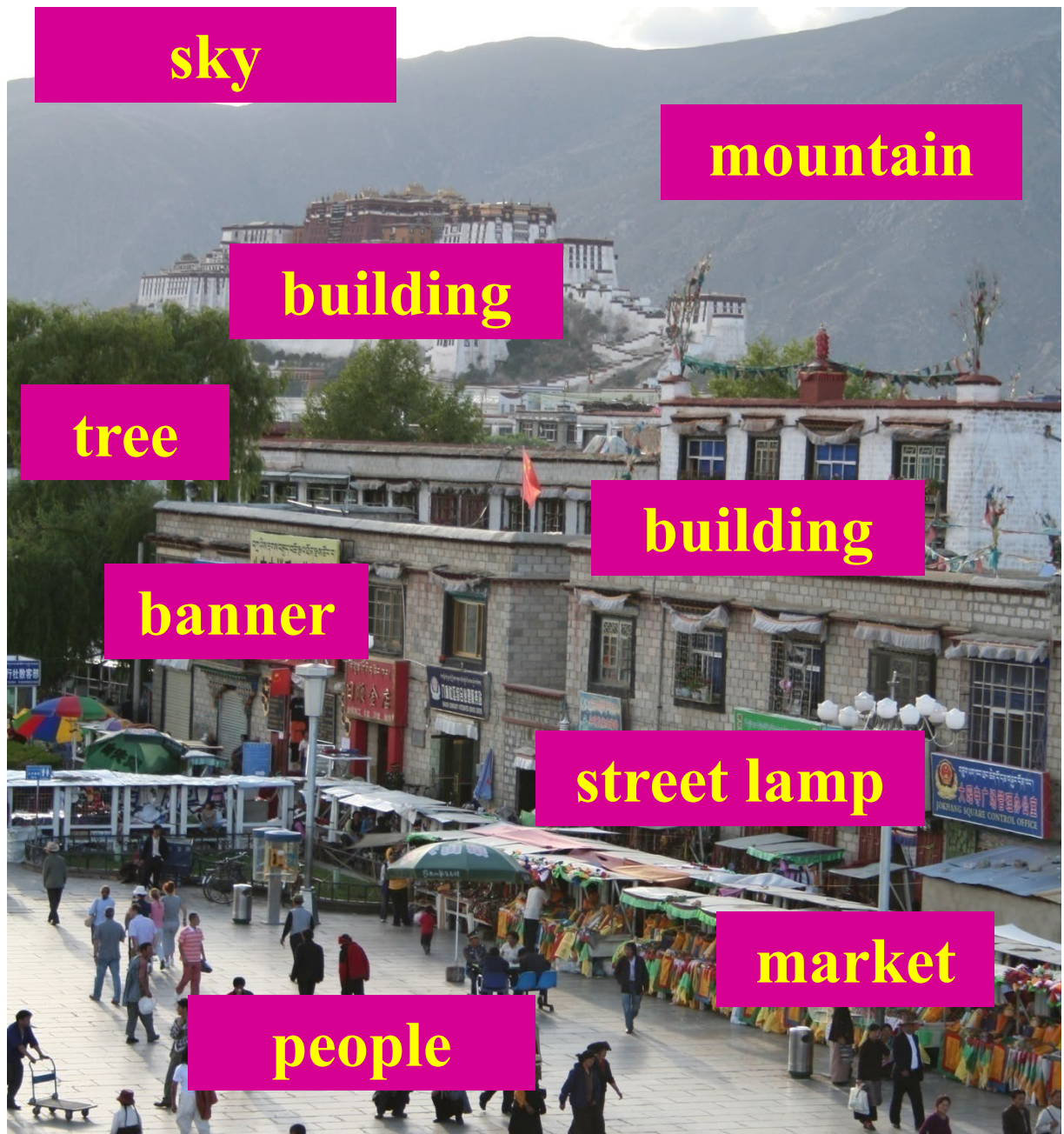
- street
- people
- building
- mountain
- ...

# Object detection

- find pedestrians



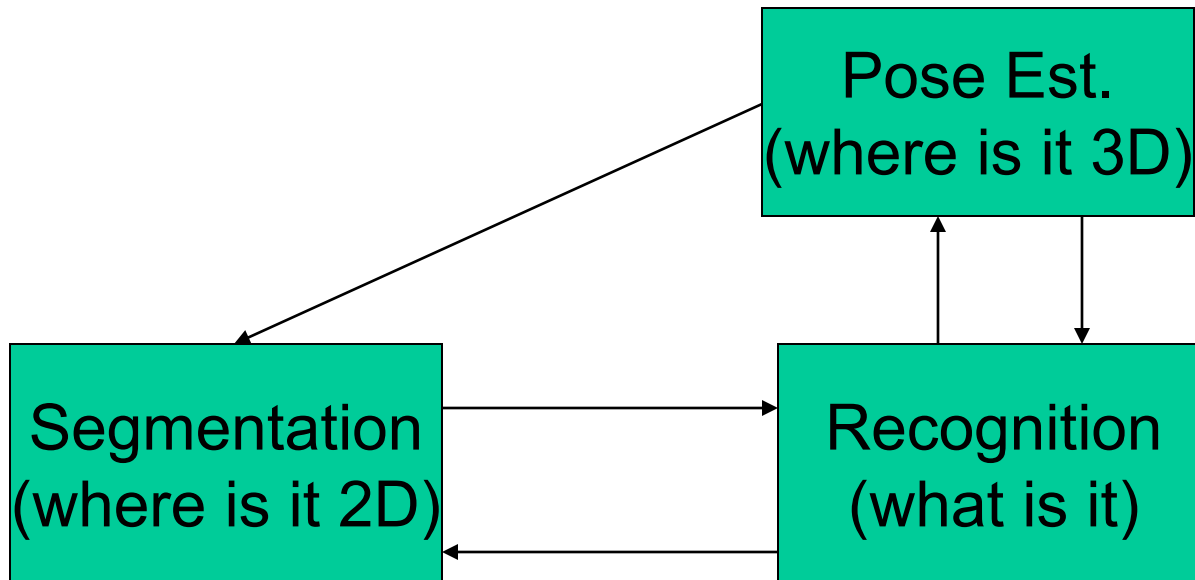
# Image parsing/Semantic Segmentation



# Object Recognition: The Problem

Given: A database  $D$  of “known” objects and an image  $I$ :

1. Determine which (if any) objects in  $D$  appear in  $I$
2. Determine the pose (rotation and translation) of the object



**WHAT AND WHERE!!!**

# Recognition Challenges

- Within-class variability

- Different objects within the class have different shapes or different material characteristics
- Deformable
- Articulated
- Compositional



- Pose variability:

- 2-D Image transformation (translation, rotation, scale)
- 3-D Pose Variability and perspective projection

- Lighting

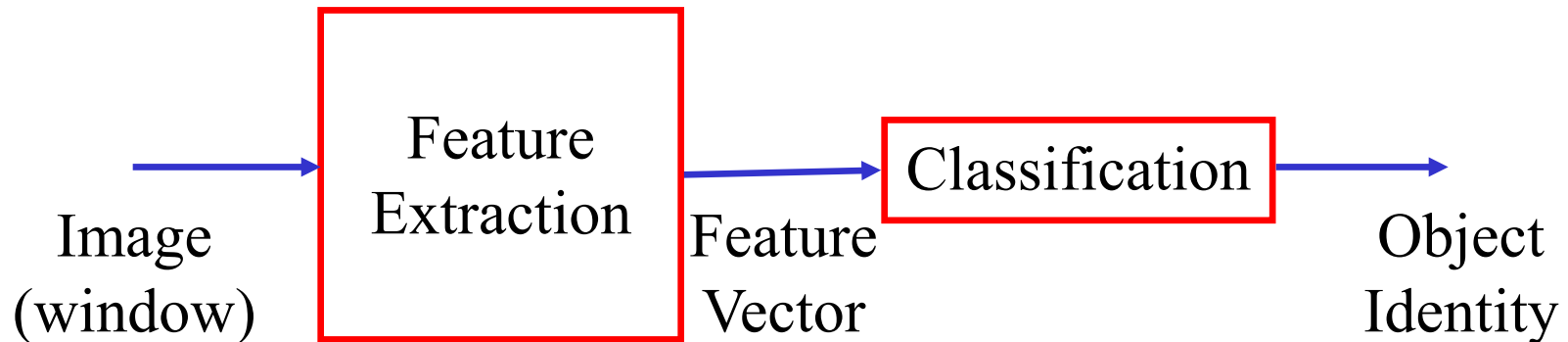
- Direction (multiple sources & type)
- Color
- Shadows

- Occlusion – partial occlusion, self occlusion

- Clutter in background -> false positives



# Sketch of a Pattern Recognition Architecture



Features might be:

- Histograms of colors
- SIFT descriptors
- Outputs of Convolutional Neural Network
- ...

Classifiers might be:

- Nearest Neighbor
- Bayesian classifier
- Logistic regression
- Support vector machine
- Softmax
- ...

# Sliding window approach to face detection

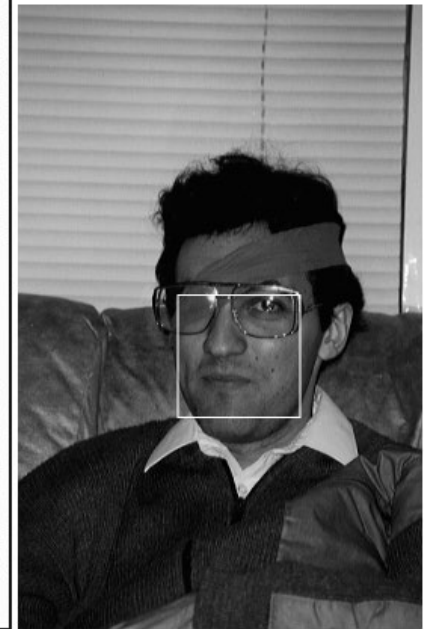
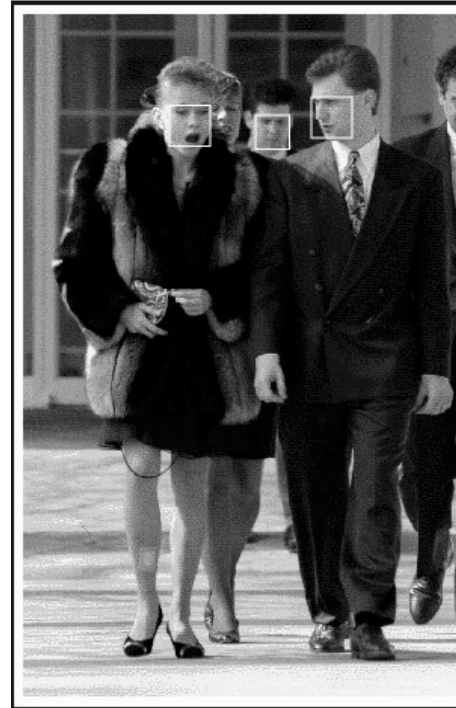
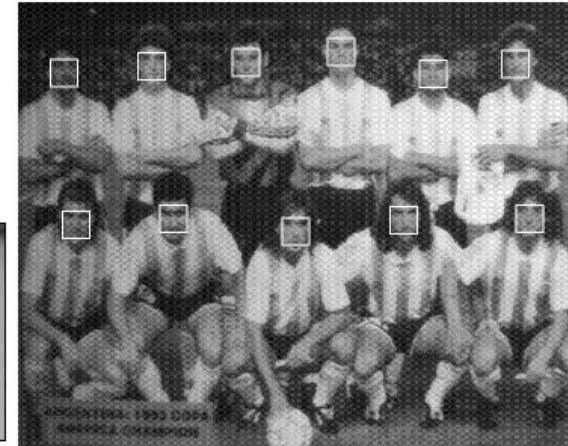
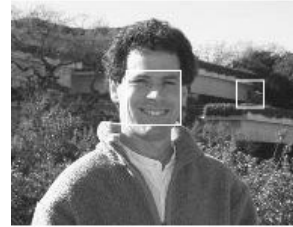


# Binary classification example: Face Detection

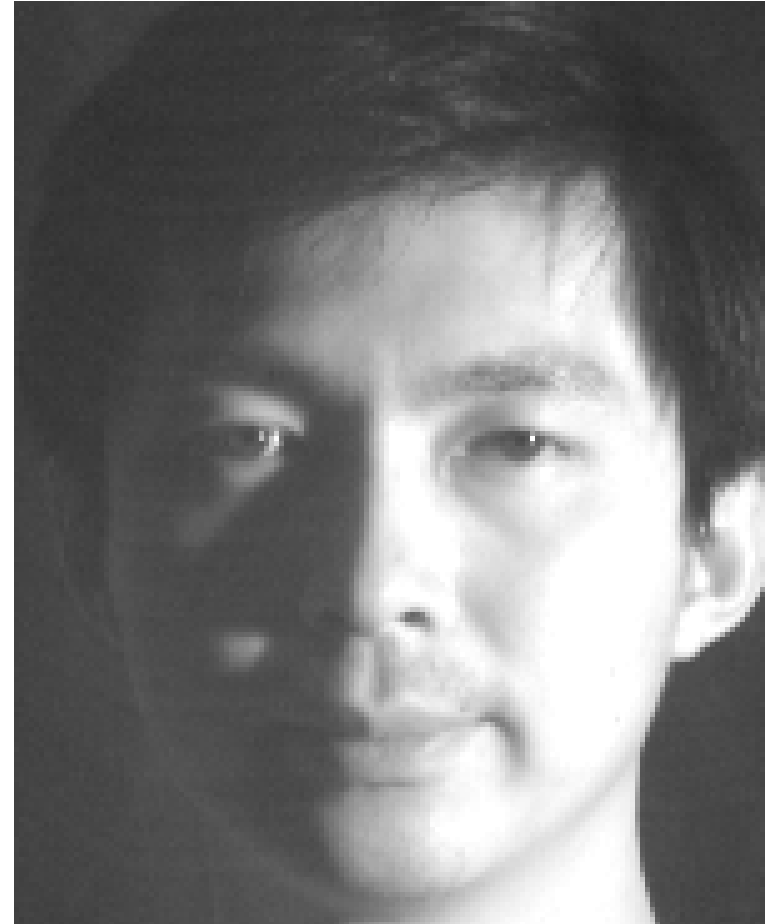
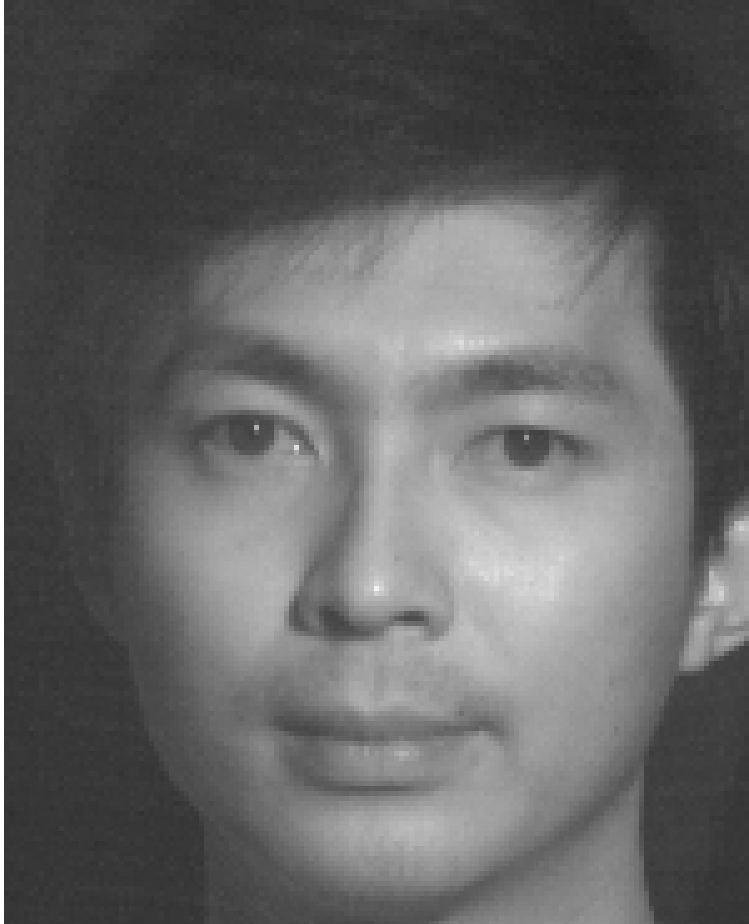
- Scan window over image.
- Classify window as either:
  - Face
  - Non-face



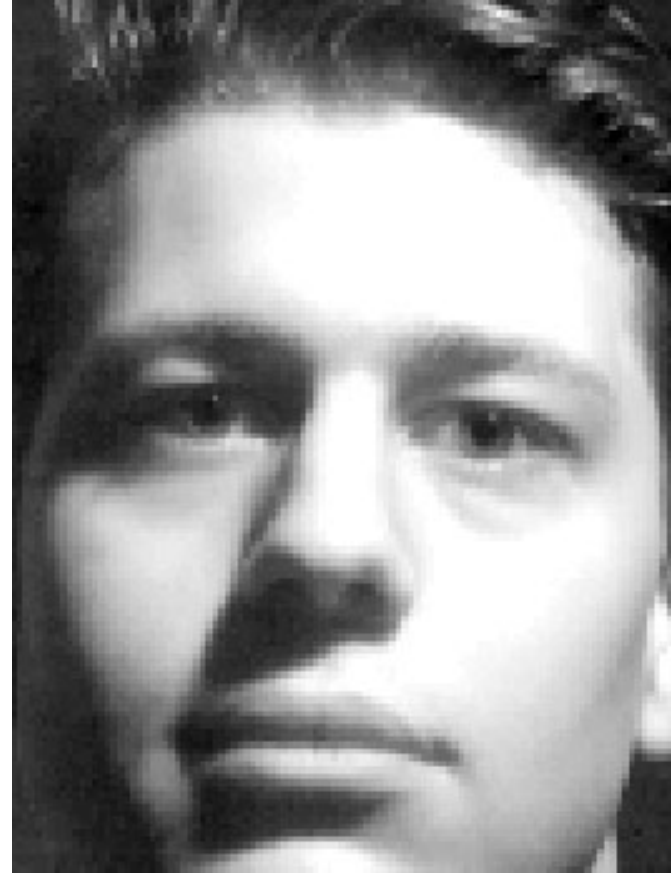
See for example, Viola-Jones face detector in OpenCV



Binary classification example:  
are these two faces from the same person?



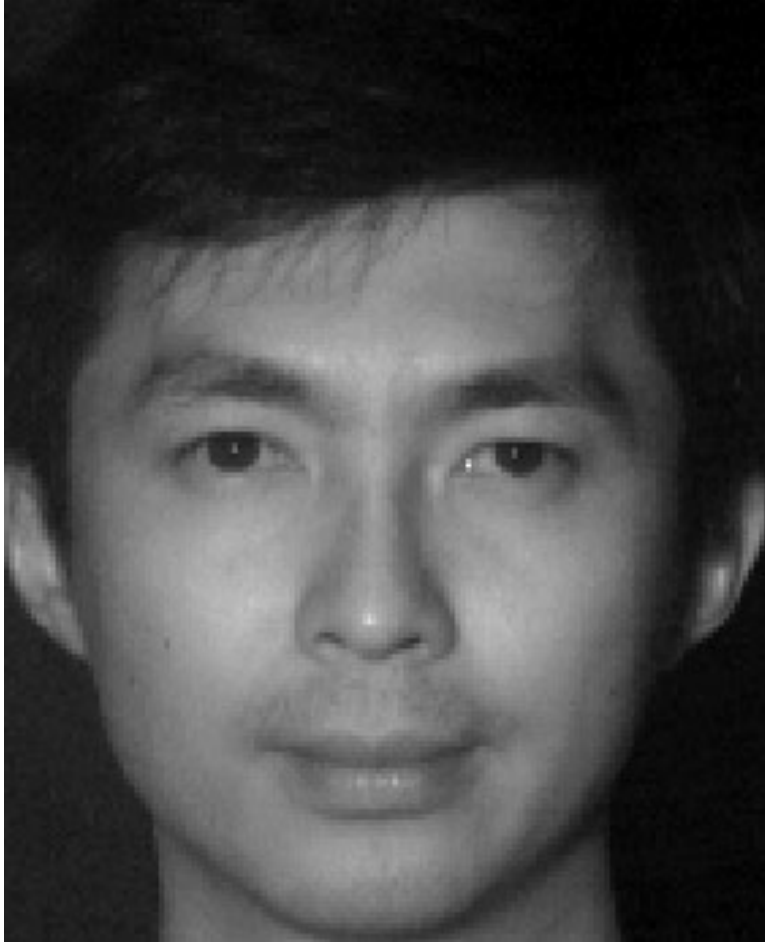
Binary classification example:  
are these two faces from the same person?



Binary classification example:  
are these two faces from the same person?



Binary classification example:  
are these two faces from the same person?



# Evaluating a Binary Classifier

Same face or different face?



		<u>Classifier Output</u>	
		Same	Different
<u>Truth</u>	Same	True Positive Rate	False Negative Rate
	Different	False Positive Rate	True Negative Rate

True Positive Rate (TPR) + False Negative Rate (FNR) = 1

False Positive Rate (FPR) + True Negative Rate (TNR) = 1

i.e., sum of rows is 1



# Evaluating a Binary Classifier

Same face or different face?



		<u>Classifier Output</u>	
		Same	Different
<u>Truth</u>	Same	0.98	0.02
	Different	0.05	0.95

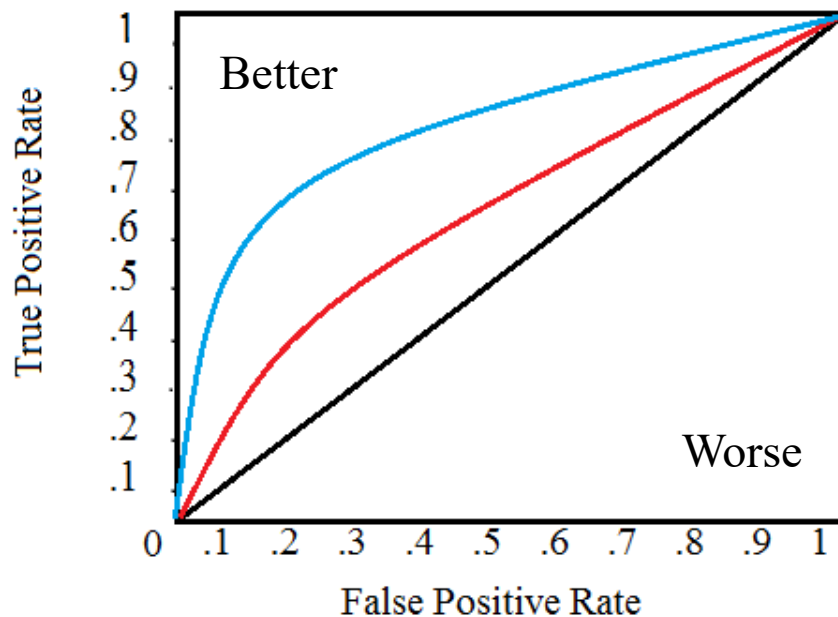
True Positive Rate (TPR) + False Negative Rate (FNR) = 1

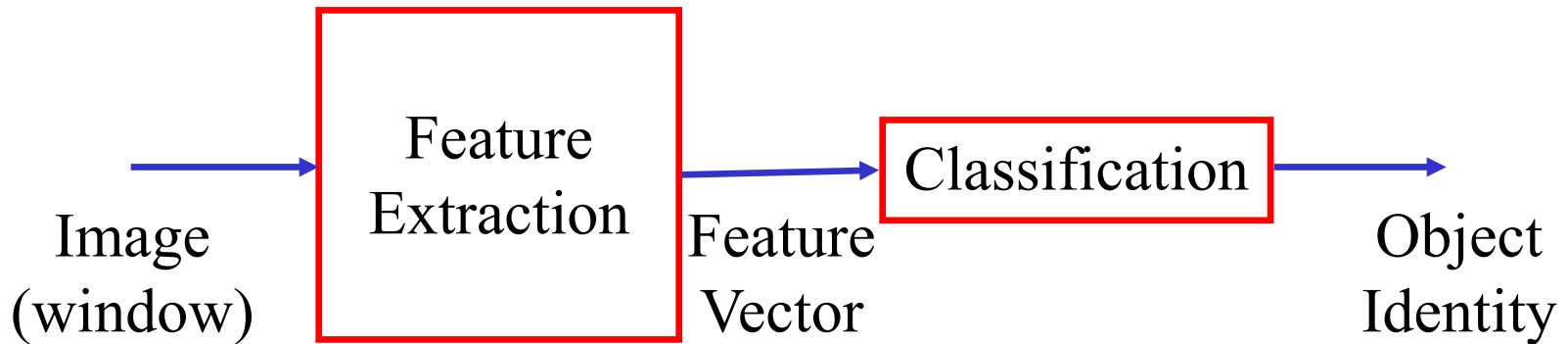
False Positive Rate (FPR) + True Negative Rate (TNR) = 1

i.e., sum of rows is 1

# ROC Curve: Evaluating a binary classifier

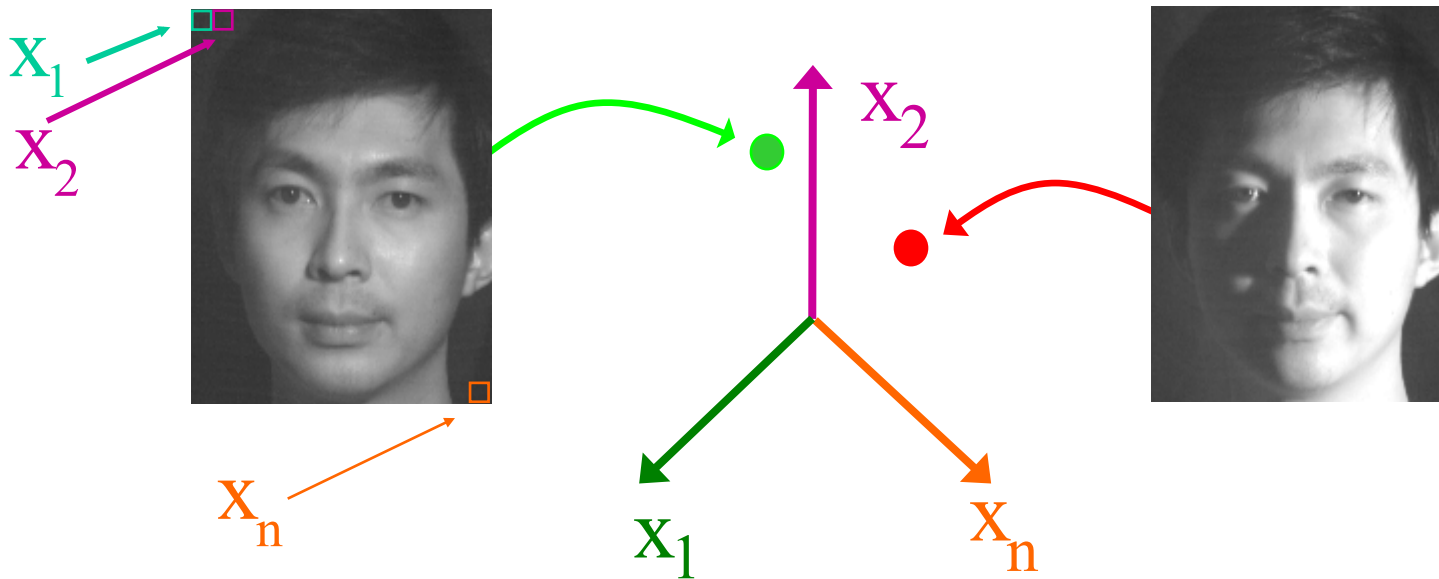
- For a face detector, there are two types of errors:
  - False Positives (e.g., non-face is detected as a face)
  - False Negatives (e.g., face is missed)
- Most classifiers have a **threshold** that can be varied, which leads to
- ROC Curve (Receiver Operator Characteristic) - Plot of tradeoff between False Positive Rate and True Positive Rate as a result of varying the threshold.





- So, what are the features?
- So, what is the classifier?

# Simplest feature: Image as a Feature Vector



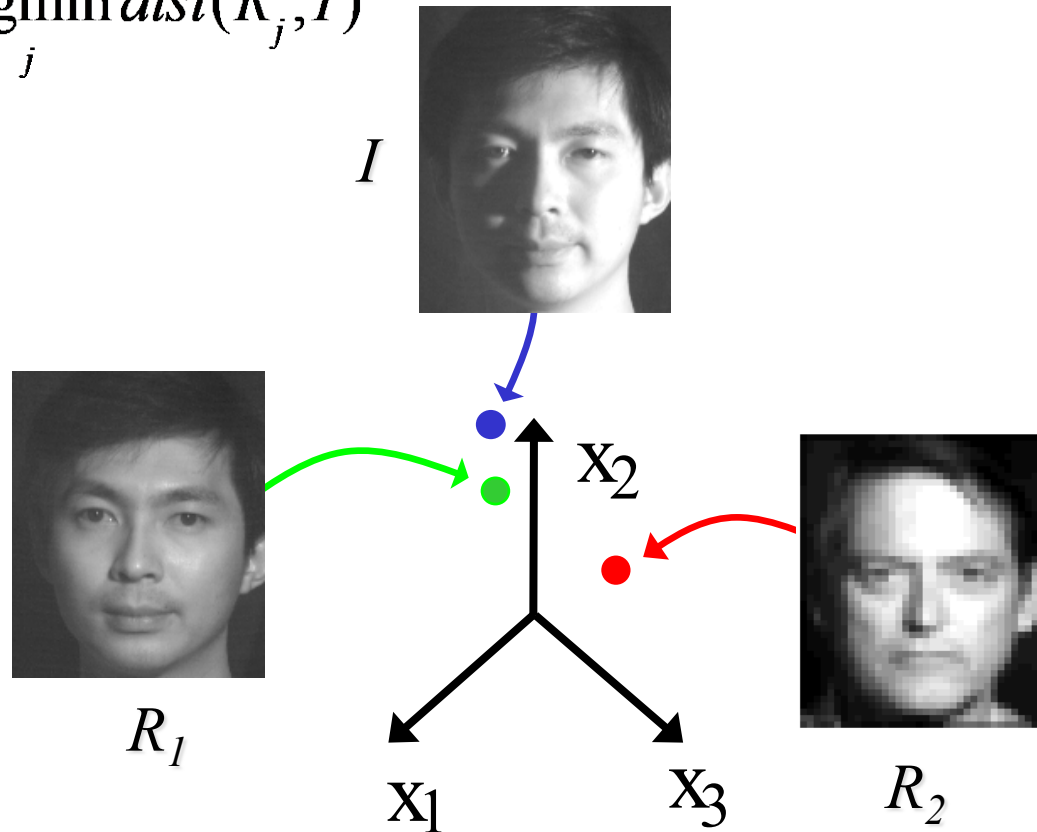
- Consider an  $n$ -pixel image to be a point in an  $n$ -dimensional space,  $\mathbf{x} \in \mathbf{R}^n$
- Each pixel value  $x_i$  is a coordinate of  $\mathbf{x}$

# Nearest Neighbor Classifier

n classes (different people)

$\{R_j : j=1, \dots, n\}$  : set of training images, one per person

$$\text{label} = \underset{j}{\operatorname{argmin}} \operatorname{dist}(R_j, I)$$

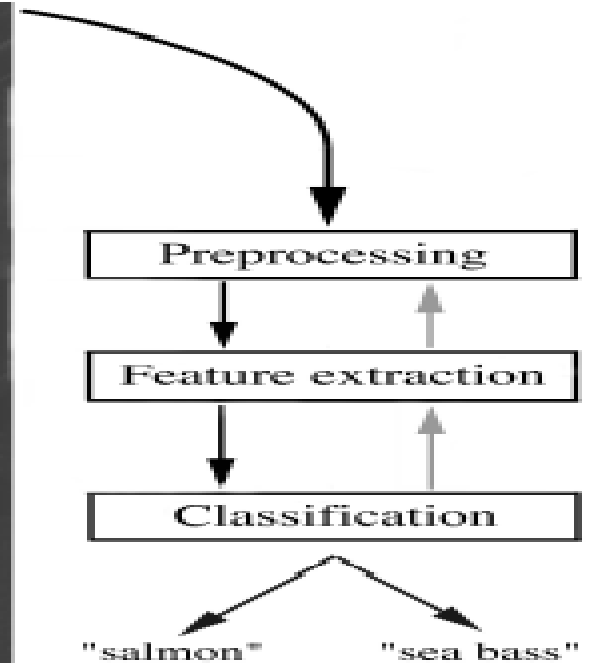


# Your summer internship

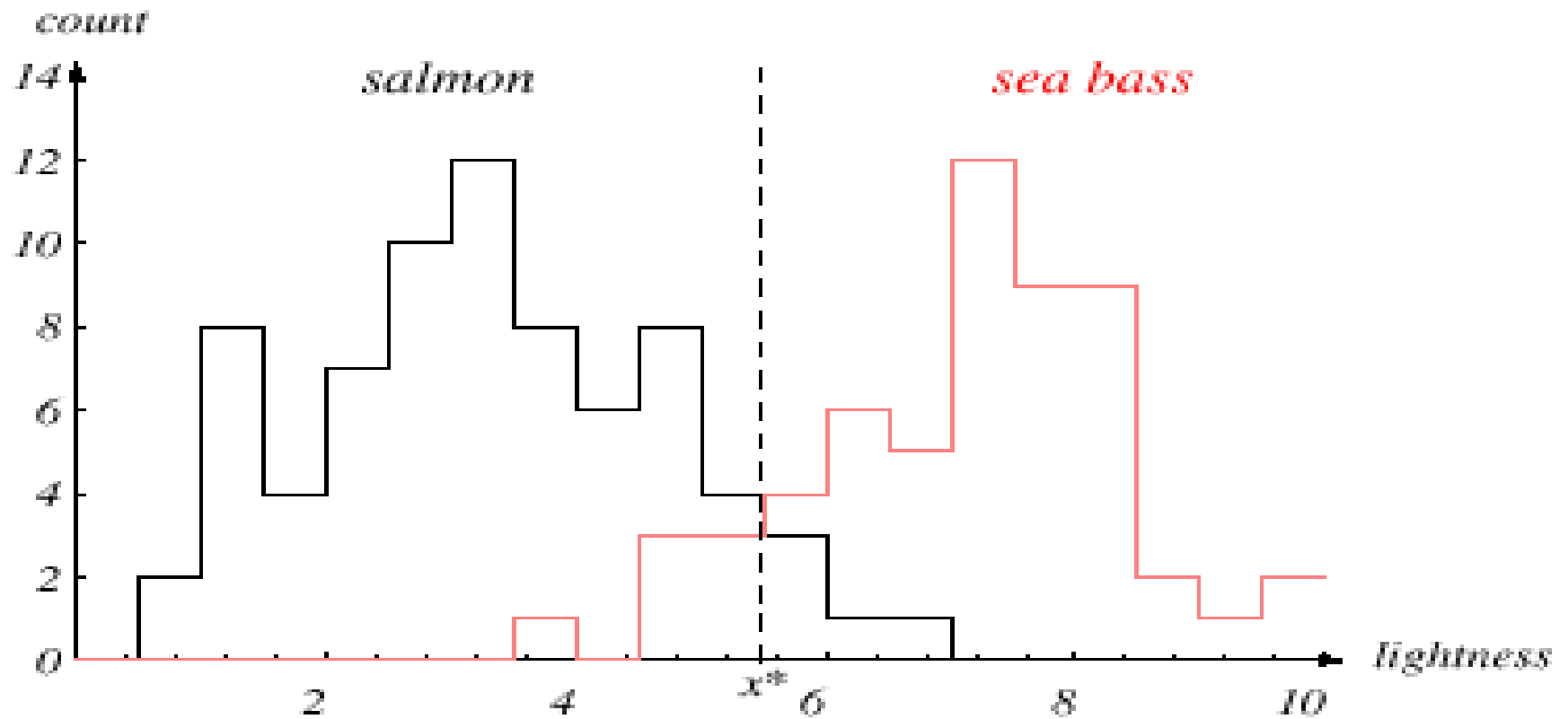


# Bayesian Classification Example

- Sorting incoming Fish on a conveyor according to species using optical sensing



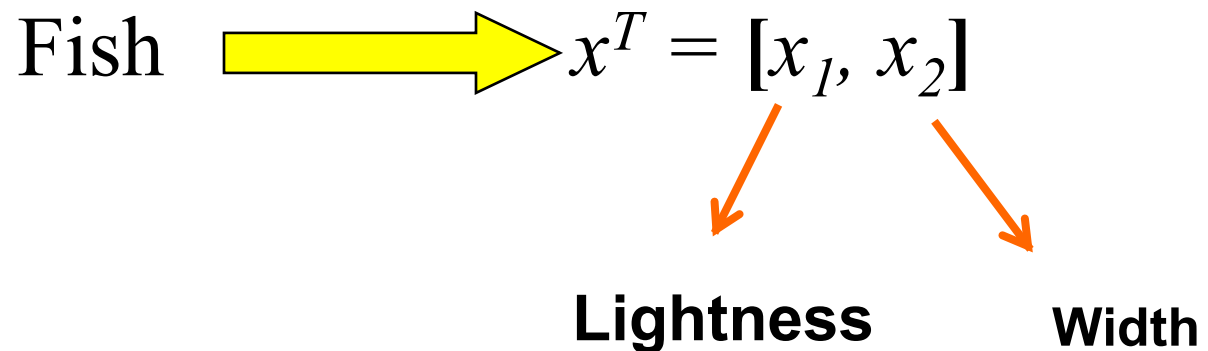
# Bayesian Classification Example: salmon or sea bass



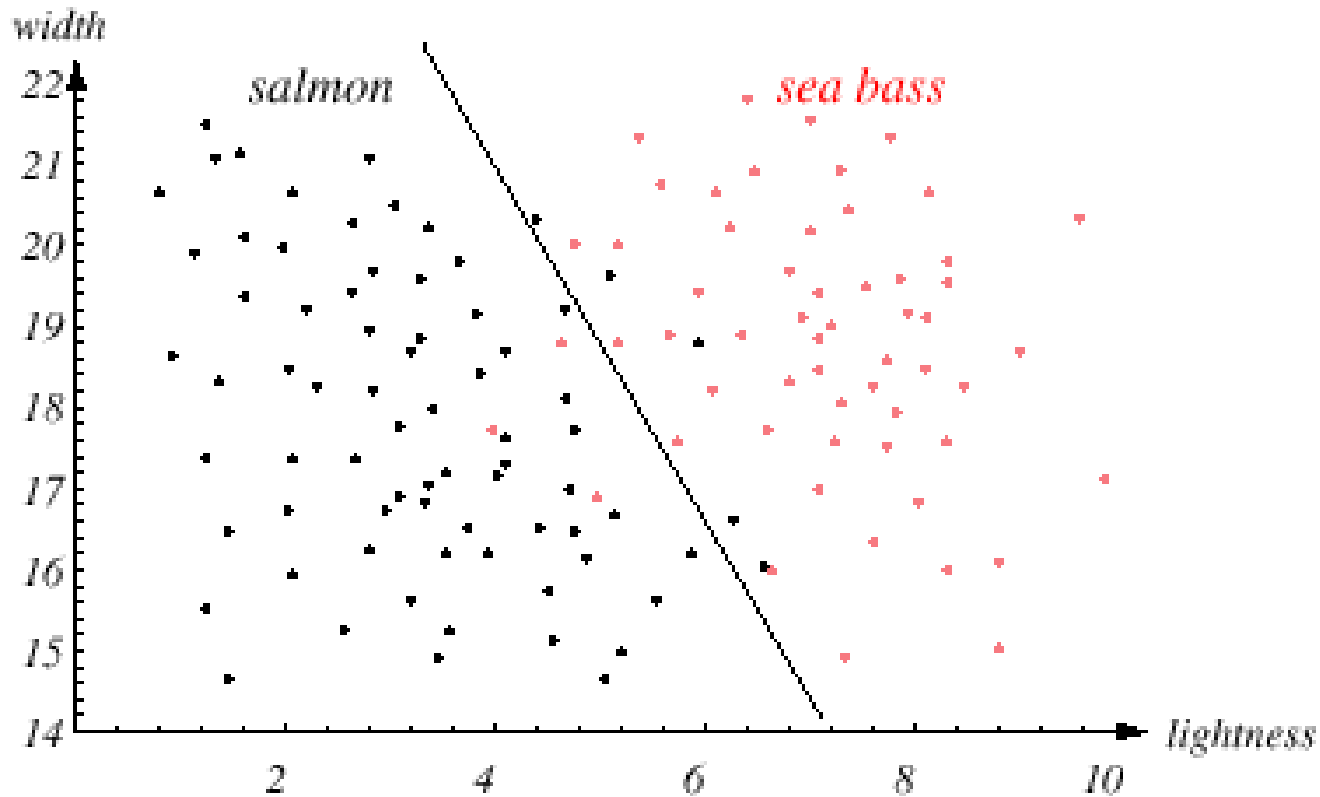


# Bayesian Classification Example: salmon or sea bass

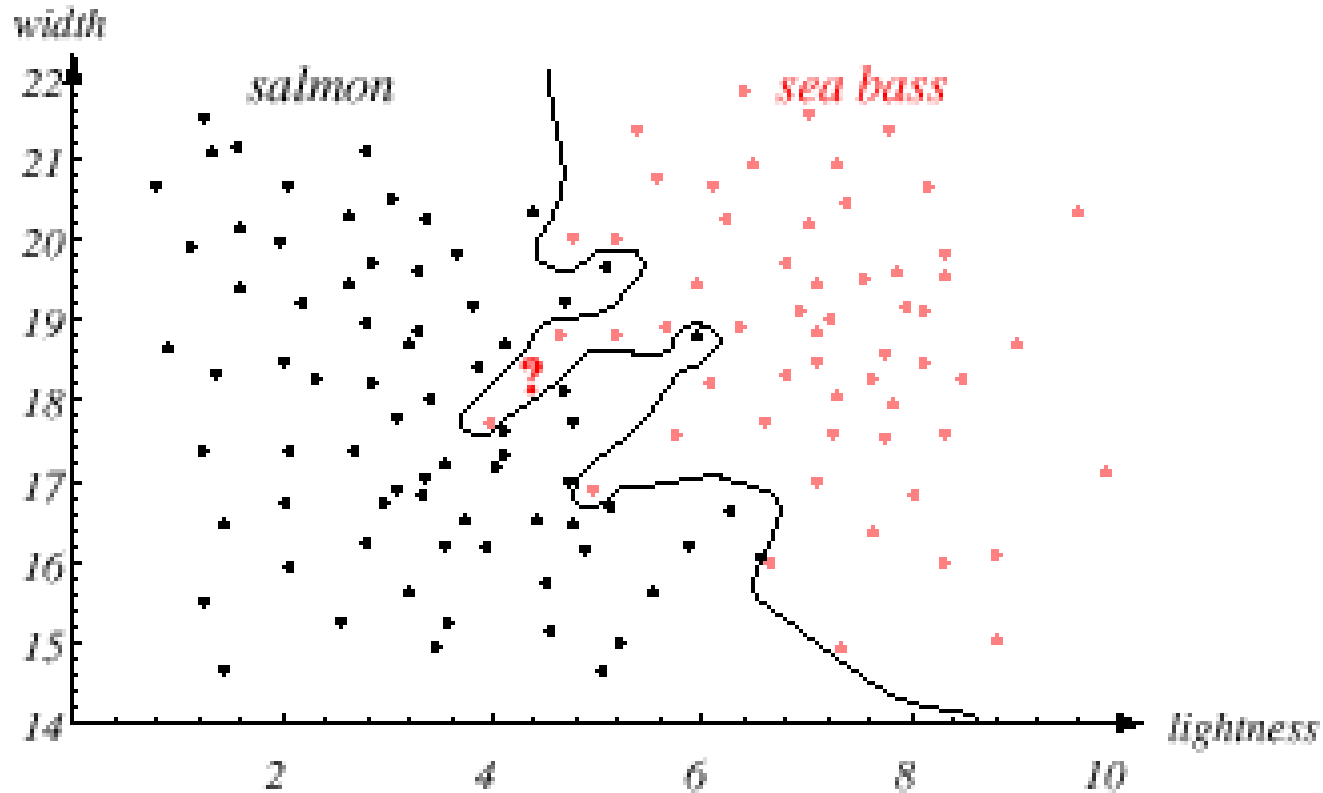
- Adopt the lightness and add the width of the fish



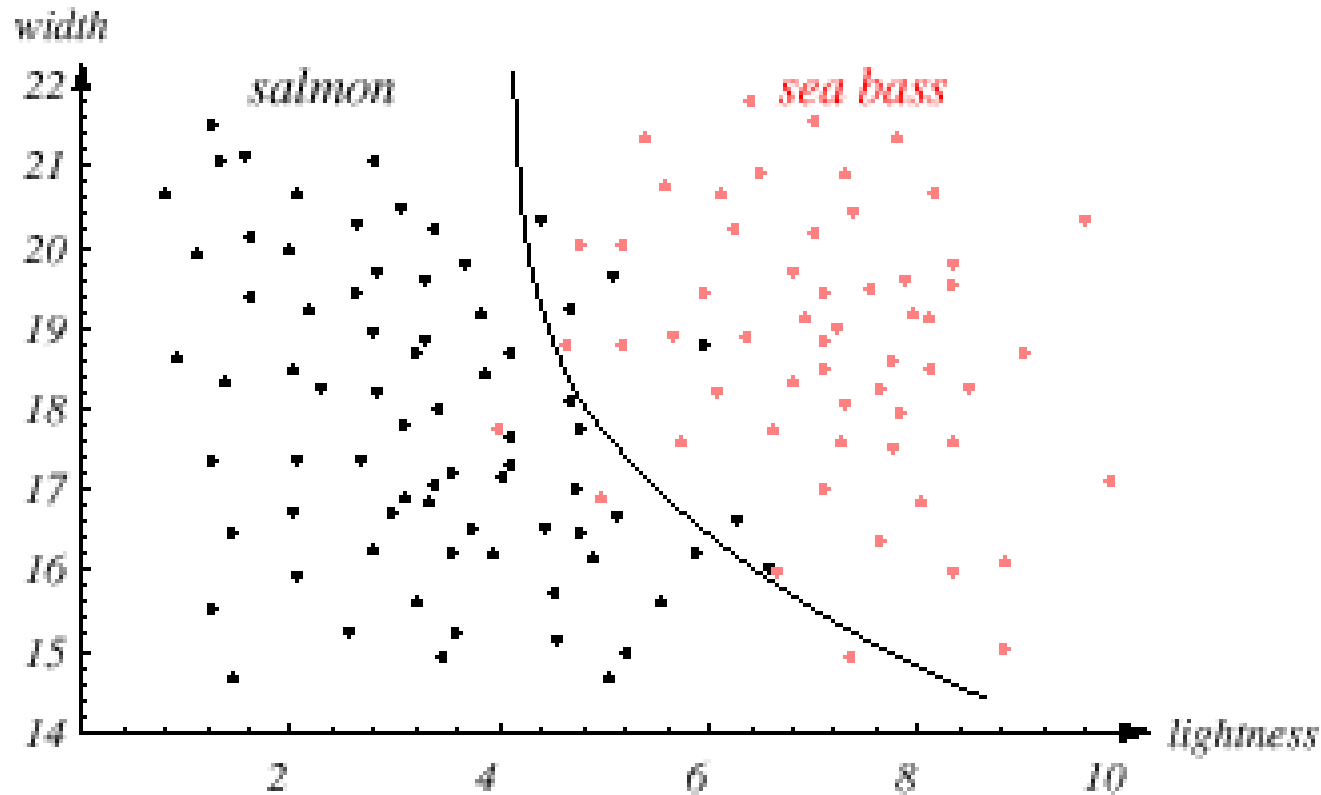
# Bayesian Classification Example: salmon or sea bass



# Bayesian Classification Example: salmon or sea bass



# Bayesian Classification Example: salmon or sea bass



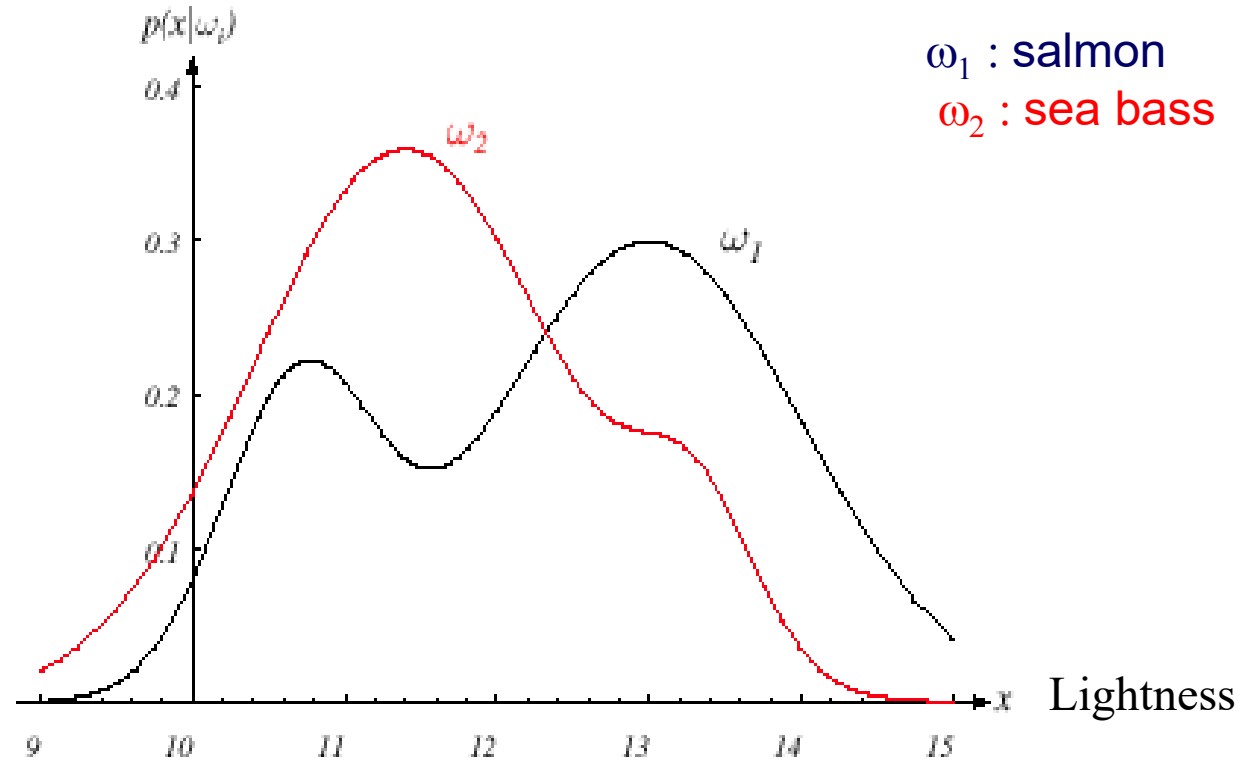
# Bayesian Classification

- The sea bass/salmon example
  - State of nature is a random variable,  $\omega_I$ 
    - $\omega_1$  – *the fish is a salmon*
    - $\omega_2$  – *the fish is a sea bass*
  - *Prior Probabilities*
    - $P(\omega_1), P(\omega_2)$
    - $P(\omega_i) > 0$
    - $P(\omega_1) + P(\omega_2) = 1$  (exclusivity and exhaustivity)
  - *Example prior: bass & salmon are equally likely*
    - $P(\omega_1) = P(\omega_2) = 1/2$  (uniform priors)

# Bayesian Classification

- Decision rule with only the prior information
  - Decide  $\omega_1$  if  $P(\omega_1) > P(\omega_2)$  otherwise decide  $\omega_2$
- Use of the class–conditional information
- *Let  $\mathbf{x}$  be a vector of features*
- $P(\mathbf{x} | \omega_1)$  and  $P(\mathbf{x} | \omega_2)$  describe the distribution in features (say lightness) between populations of sea-bass and salmon.

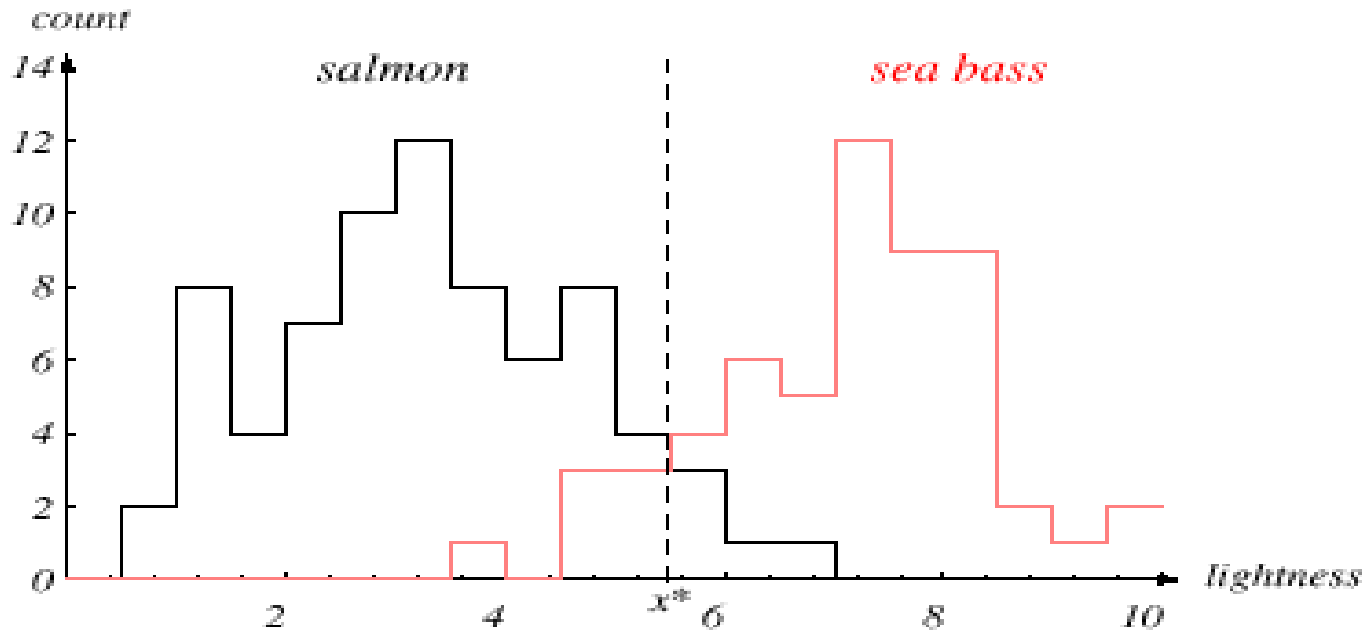
# Likelihood: $P(x | \omega_i)$



**FIGURE 2.1.** Hypothetical class-conditional probability density functions show the probability density of measuring a particular feature value  $x$  given the pattern is in category  $\omega_i$ . If  $x$  represents the lightness of a fish, the two curves might describe the difference in lightness of populations of two types of fish. Density functions are normalized, and thus the area under each curve is 1.0. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

Remember this histogram?

If the histogram is normalized by dividing by the number of samples, the result is sampled  $P(x | \omega_i)$





# Computing the Posterior

- Posterior  $P(\omega_j | \mathbf{x})$ , likelihood  $P(\mathbf{x} | \omega_j)$ , prior  $P(\mathbf{x})$

$$P(\omega_j | \mathbf{x}) = \frac{P(\mathbf{x} | \omega_j)P(\omega_j)}{P(\mathbf{x})} \quad \text{(BAYES RULE)}$$

where in case of two categories the probability density  $P(\mathbf{x})$  is

$$P(\mathbf{x}) = \sum_{j=1}^{j=2} P(\mathbf{x} | \omega_j)P(\omega_j)$$

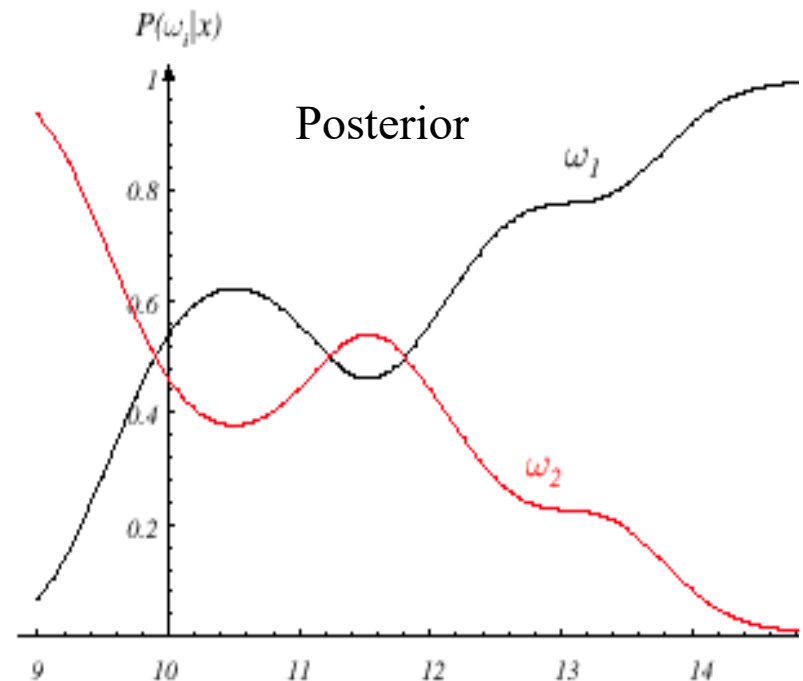
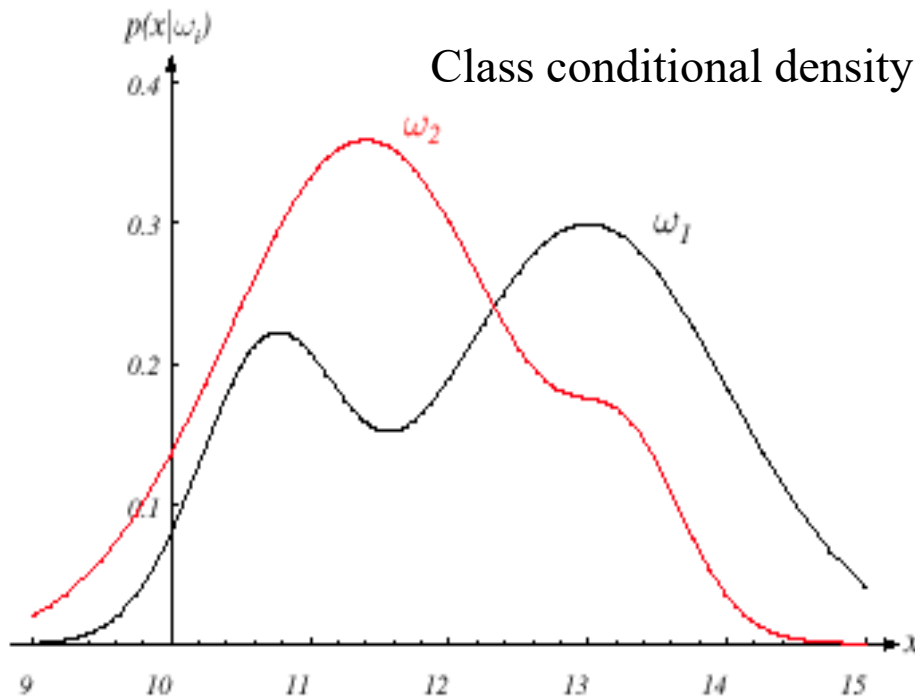
- Given a measurement  $\mathbf{x}$ , one can evaluate  $P(\omega_j | \mathbf{x})$  for all classes

# Maximum a posteriori classifier (MAP)

$$g_j(x) = P(\omega_j | x) = \frac{P(x | \omega_j)P(\omega_j)}{P(x)}$$

$P(x|\omega_j)$  : Class conditional density  
 $P(\omega_j)$  : Prior of class  $j$

Classification:  $\hat{j} = \arg \max_j g_j(x)$



# Bayesian Classification

- Think of a classifier as set of scalar functions  $g_i(\mathbf{x})$  – one for each class  $i$  – that assigns a score to the feature vector  $\mathbf{x}$ , and then chooses the class  $i$  with the highest score
- The Bayesian classifier assigns a score based on the *a posteriori* probabilities:  $g_i(\mathbf{x}) = P(\omega_i | \mathbf{x})$

# Decision Regions and Decision Boundary

**Decision region**  $R_i$  of class  $\omega_i$

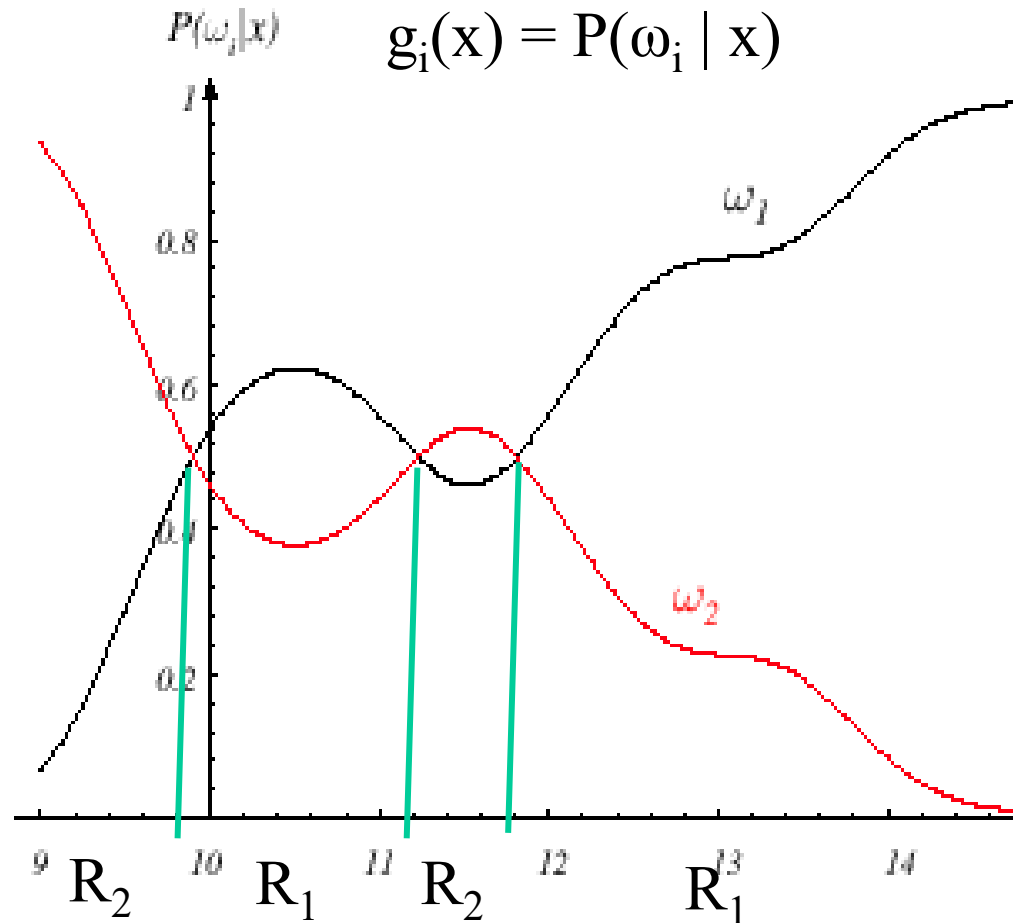
$R_i$  is the region of the feature space that will be classified as class  $\omega_i$

$$R_i = \{\mathbf{x} : g_i(\mathbf{x}) > g_j(\mathbf{x}), \forall j, i \neq j\}$$

**Decision boundary** between decision regions between  $R_i$  and  $R_j$  for classes  $\omega_i$  and  $\omega_j$

$$H_{ij} = \{\mathbf{x} : g_i(\mathbf{x}) = g_j(\mathbf{x}), g_i(\mathbf{x}) > g_k(\mathbf{x}), \forall k, i \neq k\}$$

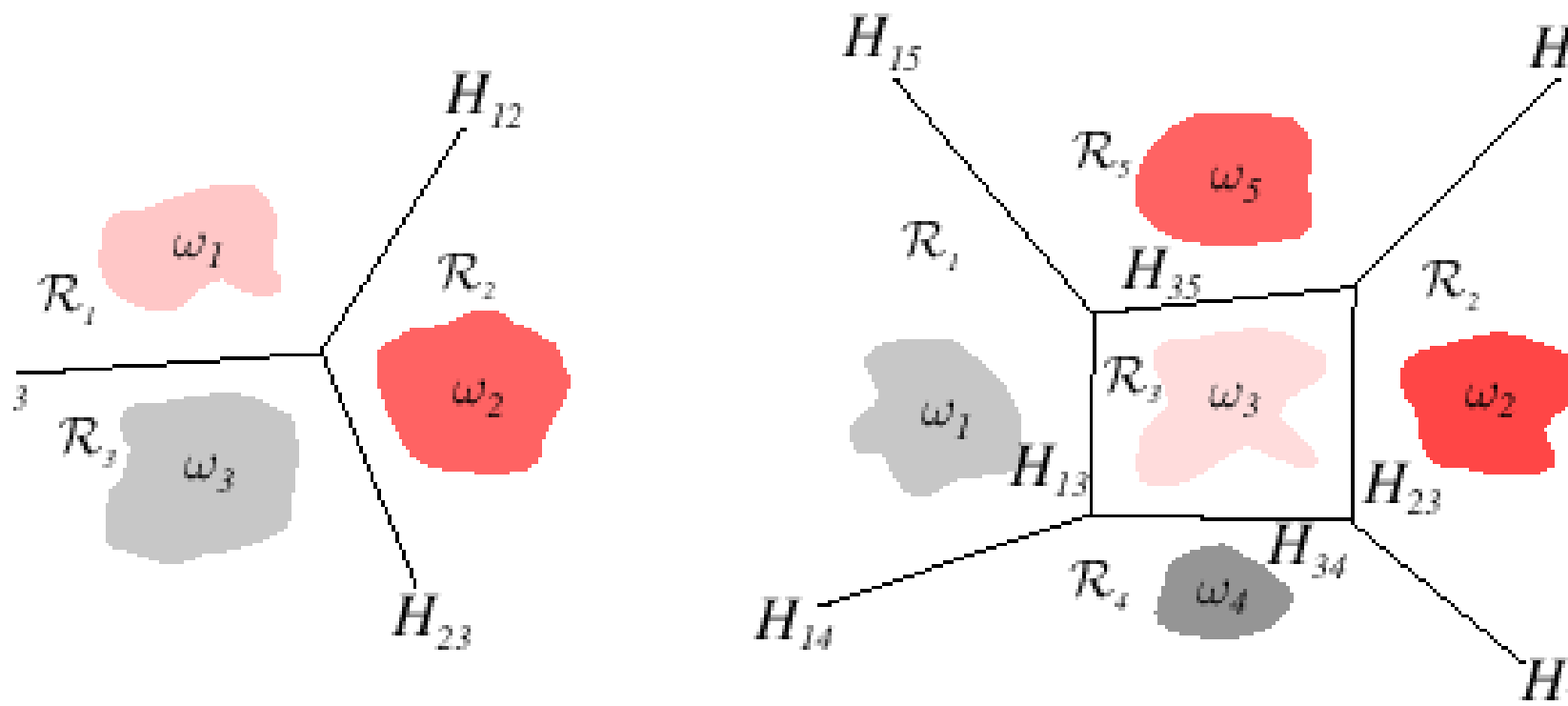
# Decision Regions and Boundaries



$$R_i = \{\mathbf{x} : g_i(\mathbf{x}) > g_j(\mathbf{x}), \forall j, i \neq j\}$$

$$H_{ij} = \{\mathbf{x} : g_i(\mathbf{x}) = g_j(\mathbf{x}), g_i(\mathbf{x}) > g_k(\mathbf{x}), \forall k, i \neq k\}$$

# Multiclass – 2D Feature Space



**IE 5.4.** Decision boundaries produced by a linear machine for a three-class and a five-class problem. From: Richard O. Duda, Peter E. Hart, and David G. Stork. *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

# Evaluating Multi-class classifiers

- Confusion matrix
  - True positive (TP): diagonal
  - False negative (FN): along row, excluding diagonal
  - False positive (FP): along column, excluding diagonal

	airplane	automobile	bird	cat	deer	dog	frog	horse	ship	truck	Actual	Recall
airplane	880	6	40	3	3	2	3	9	8	6	960	91.67%
automobile	0	982	0	2	1	11	2	0	1	20	1019	96.37%
bird	4	1	954	6	9	3	19	1	3	3	1003	95.11%
cat	2	1	9	935	6	29	18	4	1	9	1014	92.21%
deer	2	0	7	17	920	5	9	28	2	1	991	92.84%
dog	0	0	8	86	8	895	15	10	1	3	1026	87.23%
frog	0	2	3	7	3	0	993	0	0	1	1009	98.41%
horse	0	1	4	3	6	5	3	956	0	3	981	97.45%
ship	17	7	0	2	1	0	4	0	940	16	987	95.24%
truck	9	12	1	1	0	37	1	2	5	942	1010	93.27%
Predicted	914	1012	1026	1062	957	987	1067	1010	961	1004	<b>10000</b>	93.98%
Precision	96.28%	97.04%	92.98%	88.04%	96.13%	90.68%	93.06%	94.65%	97.81%	93.82%	94.05%	93.9...

$$\text{Recall} = \frac{tp}{tp + fn}$$

Total population

Average recall

$$\frac{\text{Total true positives}}{\text{Total population}}$$

$$\text{Precision} = \frac{tp}{tp + fp}$$

Average precision

- The Bayes decision rule is optimal in that it minimizes the probability of error.
- The challenge – To implement it you need:
  - $P(\omega_j)$  – this is pretty easy to get
  - $P(\mathbf{x} | \omega_j)$  – this is really hard to measure and approximate, especially when  $\mathbf{x}$  is higher than a few dimensions.
- What can we do?
  - Parametric form for  $P(\mathbf{x} | \omega_j)$
  - K-th nearest neighbor classifier
  - Learn decision rules  $g_i(\mathbf{x})$  directly, rather than trying to model and estimate  $P(\mathbf{x} | \omega_j)$  and using Bayes.
  - Reduce dimension of  $\mathbf{x}$



# Plug-in classifiers

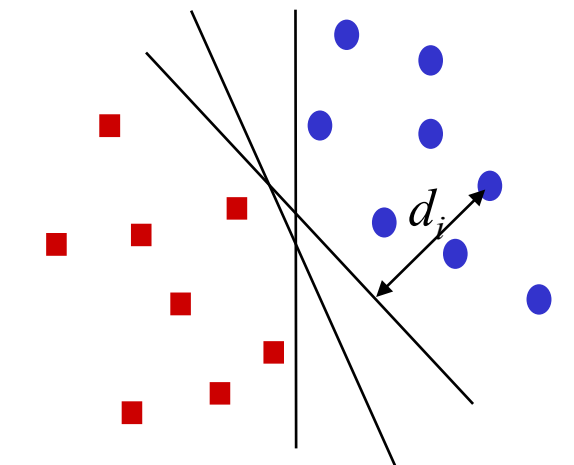
- Assume that class conditional distributions  $P(\mathbf{x}|\omega_i)$  have some parametric form with parameters  $W$
- Given training samples  $\mathbf{x}$  for class  $\omega_i$ , estimate parameters  $W$
- Common parametric form:
  - Uniform over region
  - Normal distribution with shared covariance matrix, different means
  - Normal distribution but with different covariance matrices

# Support Vector Machines

- Bayes classifiers and generative approaches in general try to model of the posterior,  $p(\omega_i|\mathbf{x})$
- Instead, try to obtain the decision boundary directly
  - Potentially easier, because we need to encode only the geometry of the boundary, not any irrelevant wiggles in the posterior.
  - Not all points affect the decision boundary

# Support Vector Machines

- Set  $S$  of points  $\mathbf{x}_i \in \mathbb{R}^n$ , each  $\mathbf{x}_i$  belongs to one of two classes  $y_i \in \{-1, 1\}$
- The goal is to find a hyperplane that divides  $S$  in these two classes



Separating hyperplanes

$$\mathbf{w}^T \mathbf{x} + b = 0$$

$S$  is separable if  $\exists \mathbf{w} \in \mathbb{R}^n, b \in \mathbb{R}$

$$y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$

$$d_i = \frac{\mathbf{w}^T \mathbf{x}_i + b}{\|\mathbf{w}\|} \quad \|\mathbf{w}\| = w$$

Closest point

$$y_i d_i = \frac{1}{w}$$

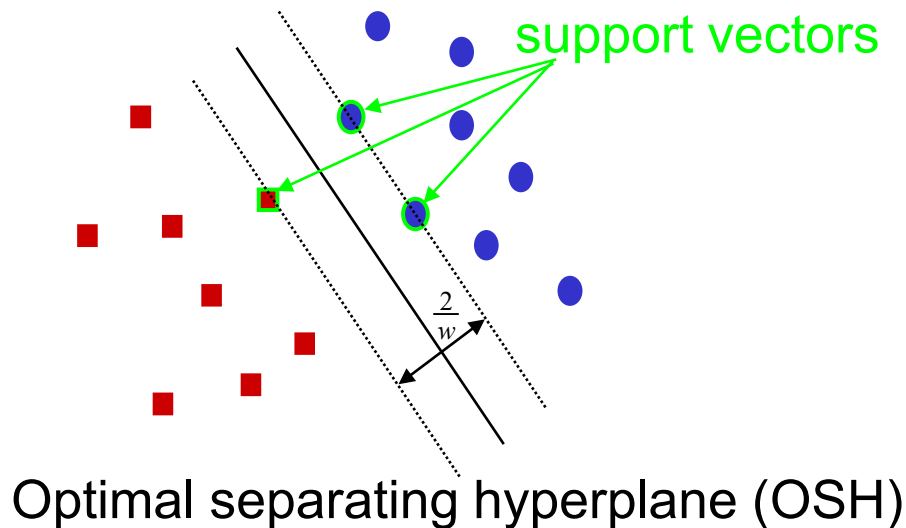
# Support Vector Machines

- Optimal separating hyperplane maximizes  $\frac{1}{w}$

Problem 1:

$$\text{Minimize } \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

$$\text{Subject to } y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1, \quad i = 1, 2, \dots, N$$



# Solve using Lagrange multipliers

- Lagrangian

$$L(\mathbf{w}, b, \boldsymbol{\alpha}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{i=1}^N \alpha_i \{y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1\}$$

– at solution

$$\frac{\partial L}{\partial b} = \sum_{i=1}^N y_i \alpha_i = 0$$

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i = 0$$

$$L(\boldsymbol{\alpha}) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j$$

– therefore

$$\boldsymbol{\alpha} \geq 0$$

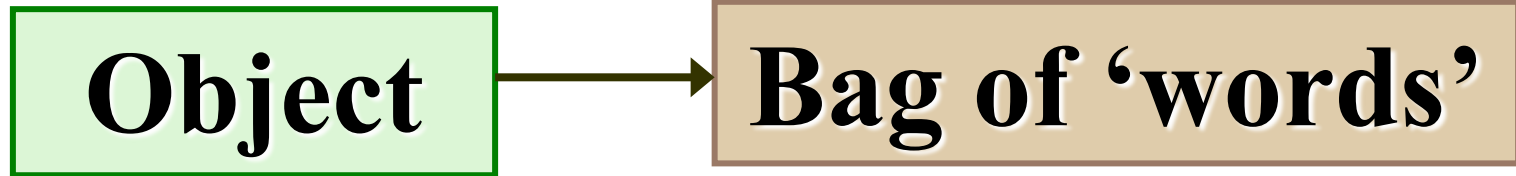
# Decision function

- Once  $w$  and  $b$  have been computed the classification decision for input  $x$  is given by

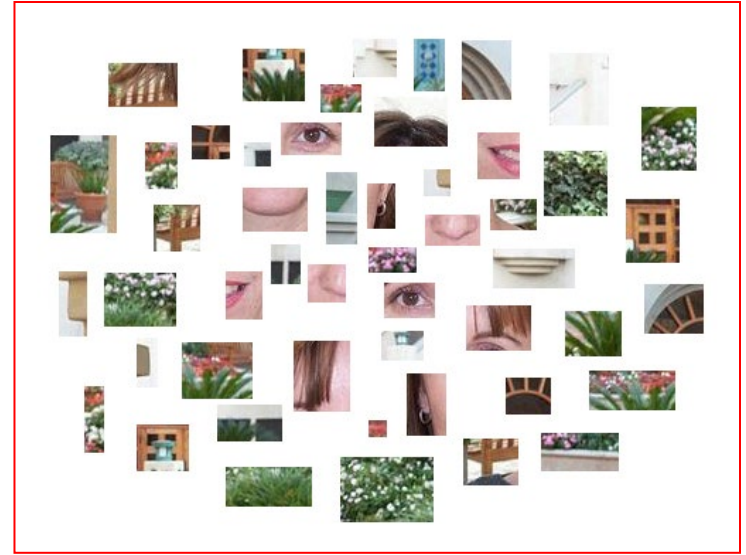
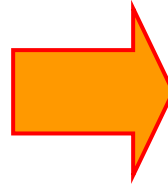
$$f(x) = \text{sign}(w^T x + b)$$

- Note the globally optimal solution can always be obtained (convex problem)

# Bag-of-features models



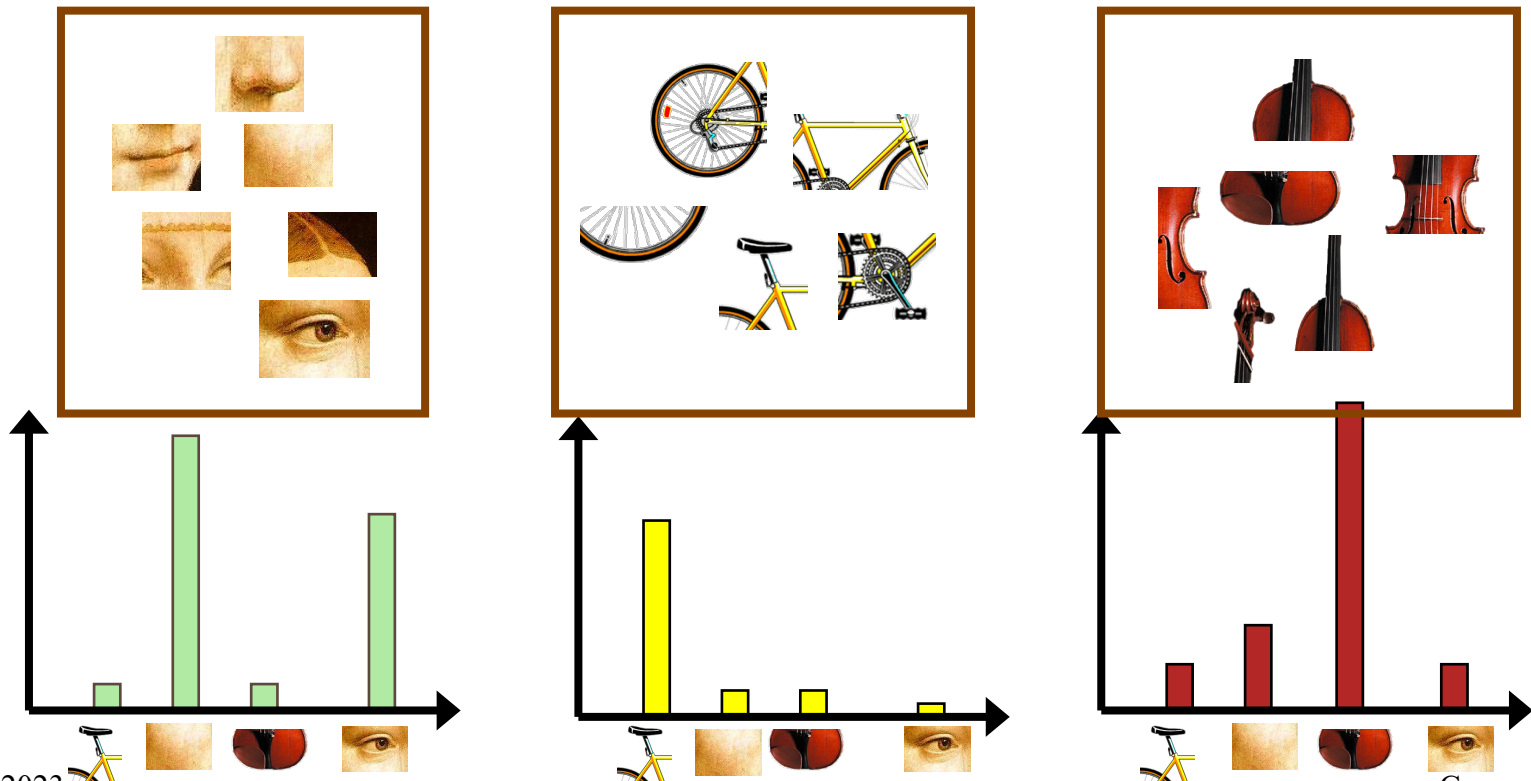
# Bag-of-features models





# Bag-of-features steps

1. Extract features
2. Learn “visual vocabulary”
3. Quantize features using visual vocabulary
4. Represent images by frequencies (histogram) of “visual words”
5. Recognition using histograms as input to classifier



# Feature extraction

- Regular grid or interest regions



# Pattern Classification Summary

- Supervised vs. Unsupervised: Do we have labels?
- Supervised
  - Nearest Neighbor
  - Bayesian
    - Plug in classifier
    - Distribution-based
    - Projection Methods (Fisher's, LDA)
  - Neural Network
  - Support Vector Machine
  - Kernel methods
- Unsupervised
  - Clustering
  - Reinforcement learning

# Next Lectures

- Recognition, detection, and classification