

## Recognition

---



The "Margaret Thatcher Illusion", by Peter Thompson

### Readings

- C. Bishop, "Neural Networks for Pattern Recognition", Oxford University Press, 1998, Chapter 1.
- [Szeliski, Chapter 14.2.1 \(eigenfaces\)](#)

## Recognition

---



The "Margaret Thatcher Illusion", by Peter Thompson

## What do we mean by "recognition"?

Next 15 slides adapted from Li, Fergus, & Torralba's excellent [short course](#) on category and object recognition



## Verification: is that a lamp?



Detection: are there people?



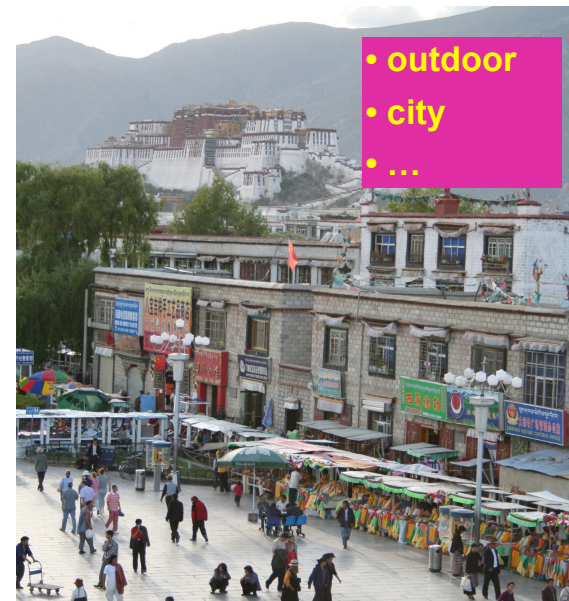
Identification: is that Potala Palace?



Object categorization



Scene and context categorization

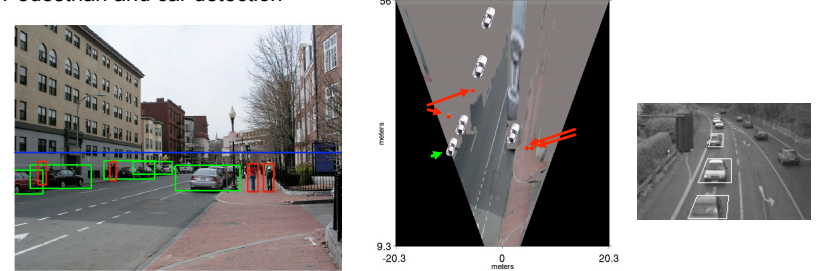


## Applications: Computational photography



## Applications: Assisted driving

### Pedestrian and car detection



### Lane detection



- Collision warning systems with adaptive cruise control,
- Lane departure warning systems,
- Rear object detection systems,

## Applications: image search



**Places**  
[London](#)  
[New York](#)  
[Egypt](#)  
[Forbidden City](#)

**Celebrities**  
[Michael Jordan](#)  
[Angelina Jolie](#)  
[Halle Berry](#)  
[Seth Rogan](#)  
[Rihanna](#)

**Art**  
[impressionism](#)  
[Keith Haring](#)  
[cubism](#)  
[Salvador Dali](#)  
[pointillism](#)

**Shopping**  
[evening gown](#)  
[necklace](#)  
[shoes](#)

### Refine your image search with visual similarity

Similar Images allows you to search for images using pictures rather than words. Click the "Similar images" link under an image to find other images that look like it. Try a search of your own or click on an example below.

#### paris



[Similar images](#)



[Similar images](#)



[Similar images](#)



[Similar images](#)

#### temple



[Similar images](#)



[Similar images](#)



[Similar images](#)



[Similar images](#)

## Challenges: viewpoint variation



Michelangelo 1475-1564

## Challenges: illumination variation

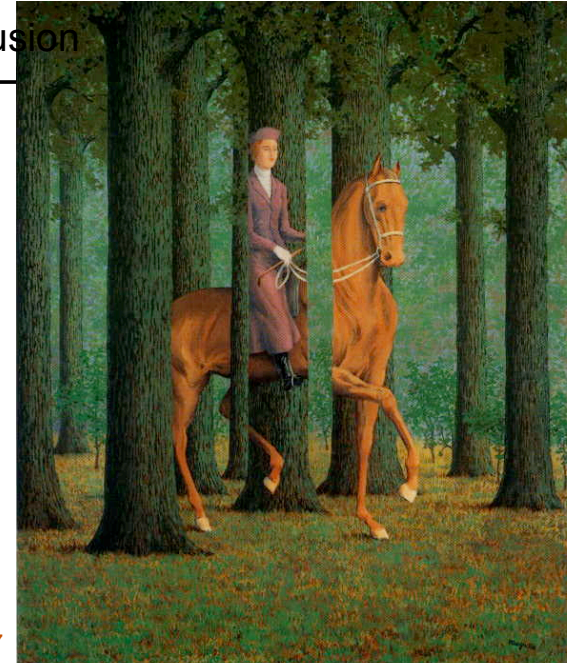
---



slide credit: S. Ullman

## Challenges: occlusion

---



Magritte, 1957

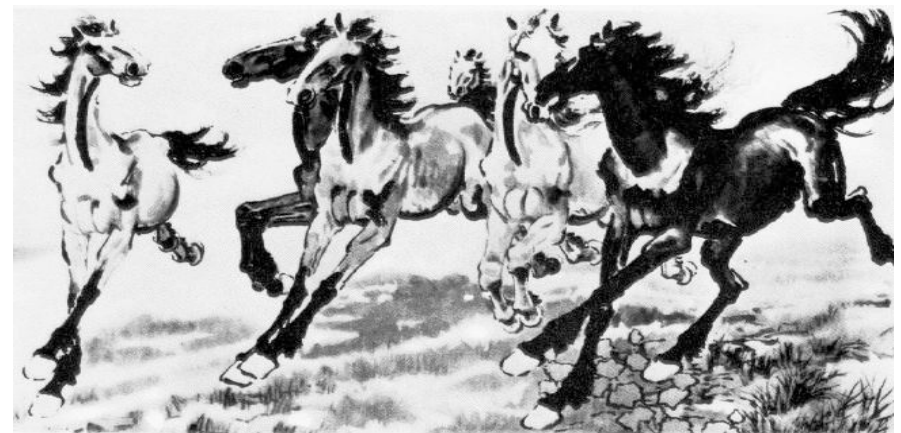
## Challenges: scale

---



## Challenges: deformation

---

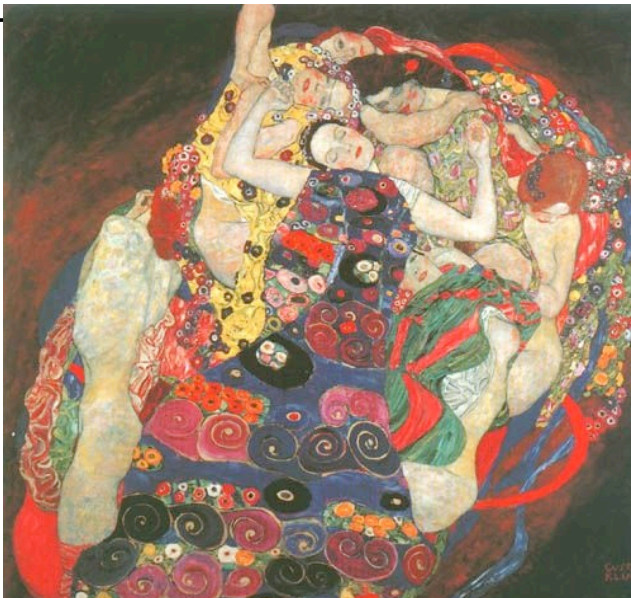


Xu, Beihong 1943

## Challenges: background clutter

---

Klimt, 1913



## Challenges: intra-class variation

---



## Face detection

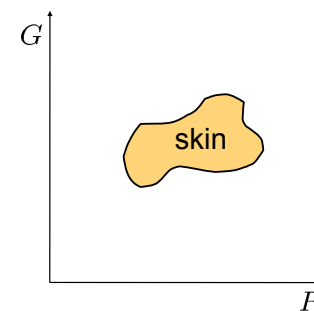
---



How to tell if a face is present?

## One simple method: skin detection

---



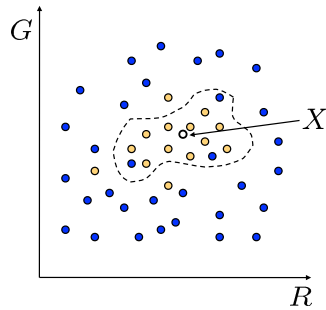
Skin pixels have a distinctive range of colors

- Corresponds to region(s) in RGB color space
  - for visualization, only R and G components are shown above

Skin classifier

- A pixel  $X = (R, G, B)$  is skin if it is in the skin region
- But how to find this region?

## Skin detection



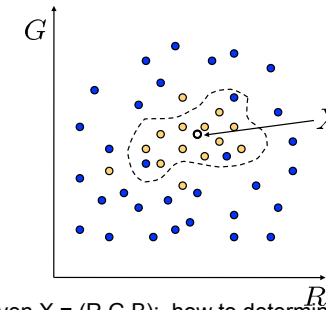
### Learn the skin region from examples

- Manually label pixels in one or more “training images” as skin or not skin
- Plot the training data in RGB space
  - skin pixels shown in orange, non-skin pixels shown in blue
  - some skin pixels may be outside the region, non-skin pixels inside. Why?

### Skin classifier

- Given  $X = (R,G,B)$ : how to determine if it is skin or not?

## Skin classification techniques



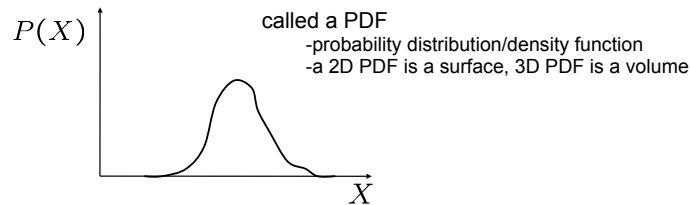
Skin classifier: Given  $X = (R,G,B)$ : how to determine if it is skin or not?

- Nearest neighbor
  - find labeled pixel closest to  $X$
- Find plane/curve that separates the two classes
  - popular approach: Support Vector Machines (SVM)
- Data modeling
  - fit a model (curve, surface, or volume) to each class
  - probabilistic version: fit a probability density/distribution model to each class

## Probability

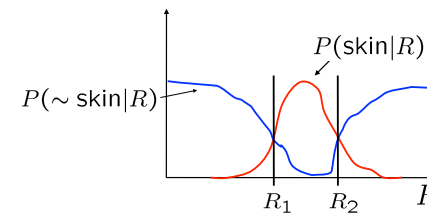
### Basic probability

- $X$  is a random variable
- $P(X)$  is the probability that  $X$  achieves a certain value



- $0 \leq P(X) \leq 1$
- $\int_{-\infty}^{\infty} P(X)dX = 1$  or  $\sum P(X) = 1$   
 continuous  $X$  or discrete  $X$
- Conditional probability:  $P(X | Y)$ 
  - probability of  $X$  given that we already know  $Y$

## Probabilistic skin classification



Now we can model uncertainty

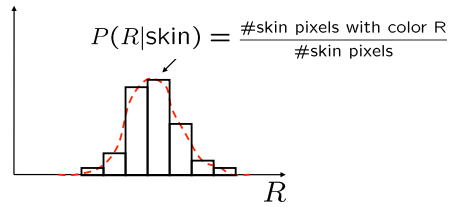
- Each pixel has a probability of being skin or not skin
  - $P(\sim \text{skin}|R) = 1 - P(\text{skin}|R)$

### Skin classifier

- Given  $X = (R,G,B)$ : how to determine if it is skin or not?
- Choose interpretation of highest probability
  - set  $X$  to be a skin pixel if and only if  $R_1 < X \leq R_2$

Where do we get  $P(\text{skin}|R)$  and  $P(\sim \text{skin}|R)$  ?

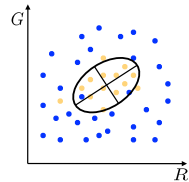
# Learning conditional PDF's



We can calculate  $P(R | \text{skin})$  from a set of training images

- It is simply a histogram over the pixels in the training images
  - each bin  $R_i$  contains the proportion of skin pixels with color  $R_i$

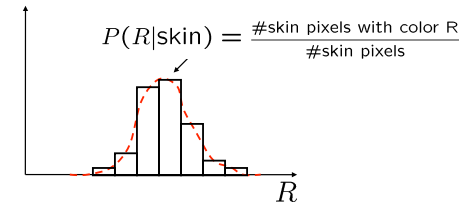
This doesn't work as well in higher-dimensional spaces. Why not?



Approach: fit parametric PDF functions

- common choice is rotated Gaussian
  - center  $c = \bar{X}$
  - covariance  $\sum_X (X - \bar{X})(X - \bar{X})^T$ 
    - orientation, size defined by eigenvecs, eigenvals

# Learning conditional PDF's



We can calculate  $P(R | \text{skin})$  from a set of training images

- It is simply a histogram over the pixels in the training images
  - each bin  $R_i$  contains the proportion of skin pixels with color  $R_i$

But this isn't quite what we want

- Why not? How to determine if a pixel is skin?
- We want  $P(\text{skin} | R)$  not  $P(R | \text{skin})$
- How can we get it?

# Bayes rule

$$P(X|Y) = \frac{P(Y|X)P(X)}{P(Y)}$$

In terms of our problem:

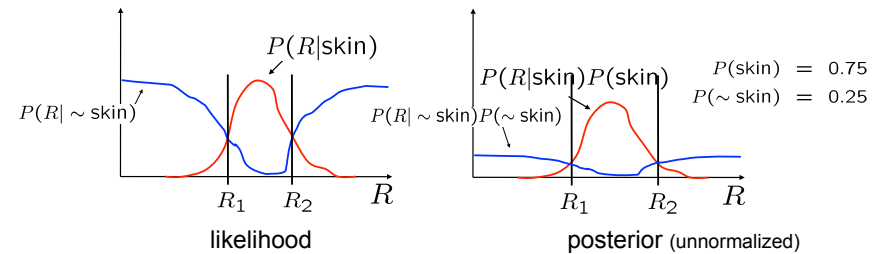
$$P(\text{skin}|R) = \frac{\overset{\substack{\text{what we measure} \\ \text{(likelihood)}}}{P(R|\text{skin})} \overset{\substack{\text{domain knowledge} \\ \text{(prior)}}}{P(\text{skin})}}{\underset{\substack{\text{normalization term} \\ P(R) = P(R|\text{skin})P(\text{skin}) + P(R|\sim\text{skin})P(\sim\text{skin})}}{P(R)}}$$

what we want (posterior)

What could we use for the prior  $P(\text{skin})$ ?

- Could use domain knowledge
  - $P(\text{skin})$  may be larger if we know the image contains a person
  - for a portrait,  $P(\text{skin})$  may be higher for pixels in the center
- Could learn the prior from the training set. How?
  - $P(\text{skin})$  may be proportion of skin pixels in training set

# Bayesian estimation

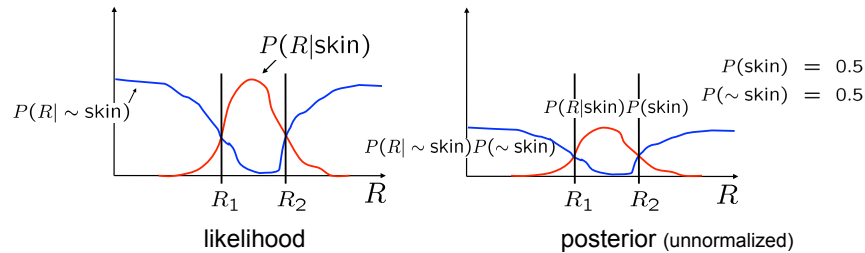


Bayesian estimation

= minimize probability of misclassification

- Goal is to choose the label (skin or  $\sim\text{skin}$ ) that maximizes the posterior
  - this is called **Maximum A Posteriori (MAP) estimation**

## Bayesian estimation



Bayesian estimation = minimize probability of misclassification

- Goal is to choose the label (skin or ~skin) that maximizes the posterior
  - this is called **Maximum A Posteriori (MAP) estimation**
- Suppose the prior is uniform:  $P(\text{skin}) = P(\sim\text{skin}) = 0.5$ 
  - in this case  $P(\text{skin}|R) = cP(R|\text{skin})$   $P(\sim\text{skin}|R) = cP(R|\sim\text{skin})$
  - maximizing the posterior is equivalent to maximizing the likelihood
    - »  $P(\text{skin}|R) > P(\sim\text{skin}|R)$  if and only if  $P(R|\text{skin}) > P(R|\sim\text{skin})$
  - this is called **Maximum Likelihood (ML) estimation**

## Skin detection results

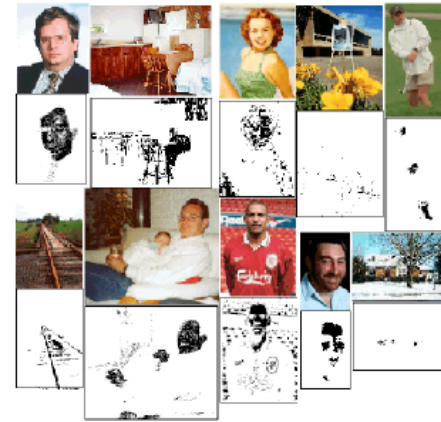
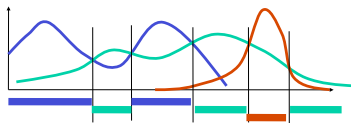


Figure 25.3. The figure shows a variety of images together with the output of the skin detector of Jones and Rehg applied to the image. Pixels marked black are skin pixels, and white are background. Notice that this process is relatively effective, and could certainly be used to focus attention on, say, faces and hands. Figure from "Statistical color models with application to skin detection," M.J. Jones and J. Rehg, Proc. Computer Vision and Pattern Recognition, 1999 © 1999, IEEE

## General classification

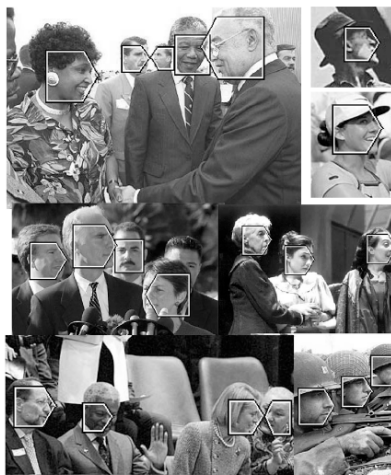
This same procedure applies in more general circumstances

- More than two classes
- More than one dimension



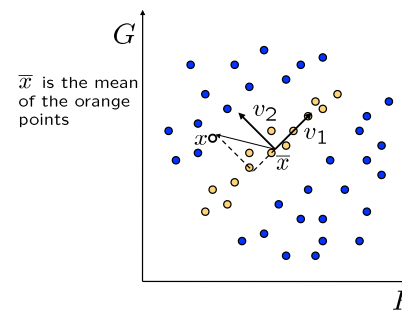
Example: face detection

- Here,  $X$  is an image region
  - dimension = # pixels
  - each face can be thought of as a point in a high dimensional space



H. Schneiderman and T.Kanade

## Linear subspaces



convert  $x$  into  $v_1, v_2$  coordinates

$$x \rightarrow ((x - \bar{x}) \cdot v_1, (x - \bar{x}) \cdot v_2)$$

What does the  $v_2$  coordinate measure?

- distance to line
- use it for classification—near 0 for orange pts

What does the  $v_1$  coordinate measure?

- position along line
- use it to specify which orange point it is

Classification can be expensive

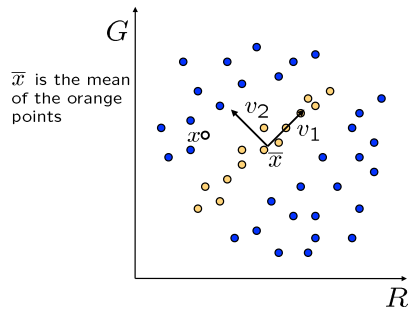
- Big search prob (e.g., nearest neighbors) or store large PDF's

Suppose the data points are arranged as above

- Idea—fit a line, classifier measures distance to line



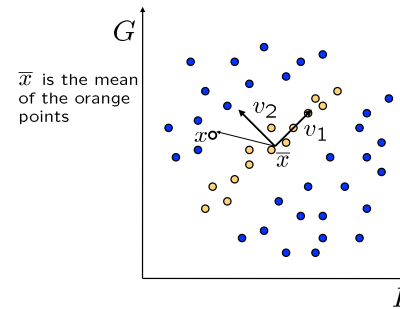
## Dimensionality reduction



### Dimensionality reduction

- We can represent the orange points with *only* their  $v_1$  coordinates
  - since  $v_2$  coordinates are all essentially 0
- This makes it much cheaper to store and compare points
- A bigger deal for higher dimensional problems

## Linear subspaces



Consider the variation along direction  $\mathbf{v}$  among all of the orange points:

$$var(\mathbf{v}) = \sum_{\text{orange point } \mathbf{x}} \|(\mathbf{x} - \bar{\mathbf{x}})^T \cdot \mathbf{v}\|^2$$

What unit vector  $\mathbf{v}$  minimizes  $var$ ?

$$\mathbf{v}_2 = \min_{\mathbf{v}} \{var(\mathbf{v})\}$$

What unit vector  $\mathbf{v}$  maximizes  $var$ ?

$$\mathbf{v}_1 = \max_{\mathbf{v}} \{var(\mathbf{v})\}$$

$$\begin{aligned} var(\mathbf{v}) &= \sum_{\mathbf{x}} \|(\mathbf{x} - \bar{\mathbf{x}})^T \cdot \mathbf{v}\|^2 \\ &= \sum_{\mathbf{x}} \mathbf{v}^T (\mathbf{x} - \bar{\mathbf{x}})(\mathbf{x} - \bar{\mathbf{x}})^T \mathbf{v} \\ &= \mathbf{v}^T \left[ \sum_{\mathbf{x}} (\mathbf{x} - \bar{\mathbf{x}})(\mathbf{x} - \bar{\mathbf{x}})^T \right] \mathbf{v} \\ &= \mathbf{v}^T \mathbf{A} \mathbf{v} \text{ where } \mathbf{A} = \sum_{\mathbf{x}} (\mathbf{x} - \bar{\mathbf{x}})(\mathbf{x} - \bar{\mathbf{x}})^T \end{aligned}$$

Solution:  $\mathbf{v}_1$  is eigenvector of  $\mathbf{A}$  with *largest* eigenvalue  
 $\mathbf{v}_2$  is eigenvector of  $\mathbf{A}$  with *smallest* eigenvalue

## Principal component analysis (PCA)

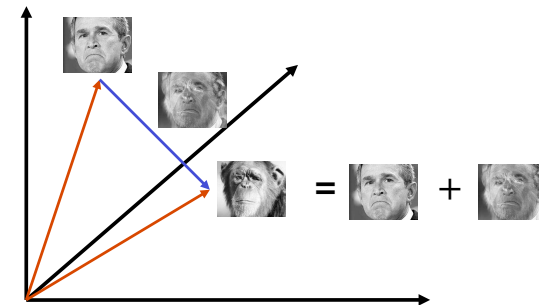
Suppose each data point is N-dimensional

- Same procedure applies:

$$\begin{aligned} var(\mathbf{v}) &= \sum_{\mathbf{x}} \|(\mathbf{x} - \bar{\mathbf{x}})^T \cdot \mathbf{v}\|^2 \\ &= \mathbf{v}^T \mathbf{A} \mathbf{v} \text{ where } \mathbf{A} = \sum_{\mathbf{x}} (\mathbf{x} - \bar{\mathbf{x}})(\mathbf{x} - \bar{\mathbf{x}})^T \end{aligned}$$

- The eigenvectors of  $\mathbf{A}$  define a new coordinate system
  - eigenvector with largest eigenvalue captures the most variation among training vectors  $\mathbf{x}$
  - eigenvector with smallest eigenvalue has least variation
- We can compress the data by only using the top few eigenvectors
  - corresponds to choosing a "linear subspace"
    - represent points on a line, plane, or "hyper-plane"
  - these eigenvectors are known as the **principal components**

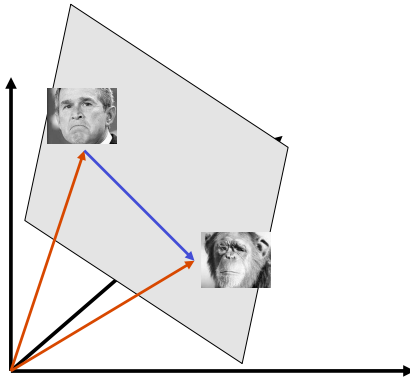
## The space of faces



An image is a point in a high dimensional space

- An  $N \times M$  image is a point in  $\mathbb{R}^{NM}$
- We can define vectors in this space as we did in the 2D case

## Dimensionality reduction



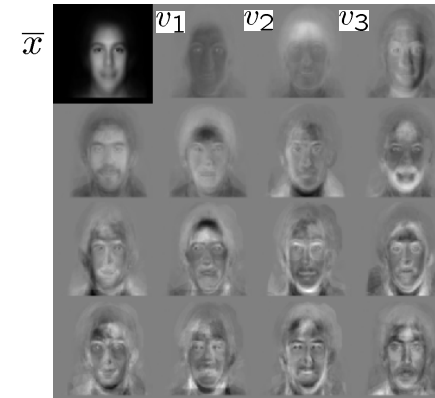
The set of faces is a “subspace” of the set of images

- Suppose it is  $K$  dimensional
- We can find the best subspace using PCA
- This is like fitting a “hyper-plane” to the set of faces
  - spanned by vectors  $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_K$
  - any face  $\mathbf{x} \approx \bar{\mathbf{x}} + a_1\mathbf{v}_1 + a_2\mathbf{v}_2 + \dots + a_k\mathbf{v}_k$

## Eigenfaces

PCA extracts the eigenvectors of  $\mathbf{A}$

- Gives a set of vectors  $\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3, \dots$
- Each one of these vectors is a direction in face space
  - what do these look like?



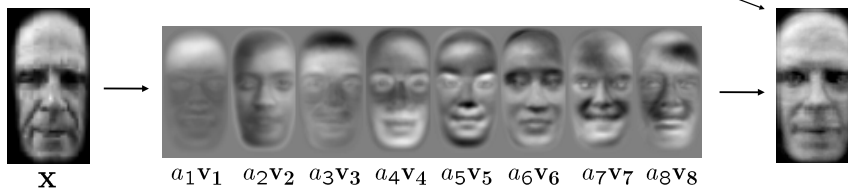
## Projecting onto the eigenfaces

The eigenfaces  $\mathbf{v}_1, \dots, \mathbf{v}_K$  span the space of faces

- A face is converted to eigenface coordinates by

$$\mathbf{x} \rightarrow \left( \underbrace{(\mathbf{x} - \bar{\mathbf{x}}) \cdot \mathbf{v}_1}_{a_1}, \underbrace{(\mathbf{x} - \bar{\mathbf{x}}) \cdot \mathbf{v}_2}_{a_2}, \dots, \underbrace{(\mathbf{x} - \bar{\mathbf{x}}) \cdot \mathbf{v}_K}_{a_K} \right)$$

$$\mathbf{x} \approx \bar{\mathbf{x}} + a_1\mathbf{v}_1 + a_2\mathbf{v}_2 + \dots + a_K\mathbf{v}_K$$



## Recognition with eigenfaces

Algorithm

1. Process the image database (set of images with labels)
  - Run PCA—compute eigenfaces
  - Calculate the  $K$  coefficients for each image
2. Given a new image (to be recognized)  $\mathbf{x}$ , calculate  $K$  coefficients

$$\mathbf{x} \rightarrow (a_1, a_2, \dots, a_K)$$

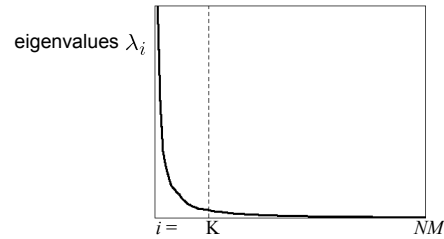
3. Detect if  $\mathbf{x}$  is a face

$$\|\mathbf{x} - (\bar{\mathbf{x}} + a_1\mathbf{v}_1 + a_2\mathbf{v}_2 + \dots + a_K\mathbf{v}_K)\| < \text{threshold}$$

4. If it is a face, who is it?
  - Find closest labeled face in database
    - nearest-neighbor in  $K$ -dimensional space

## Choosing the dimension K

---



How many eigenfaces to use?

Look at the decay of the eigenvalues

- the eigenvalue tells you the amount of variance “in the direction” of that eigenface
- ignore eigenfaces with low variance

## Object recognition

---

This is just the tip of the iceberg

- Better features:
  - edges (e.g., SIFT)
  - motion
  - depth/3D info
  - ...
- Better classifiers:
  - e.g., support vector machines (SVN)
- Speed (e.g., real-time face detection)
- Scale
  - e.g., Internet image search

Recognition is a very active research area right now