

# Machine Learning in Robotics

## Lecture 8: Dimensionality Reduction

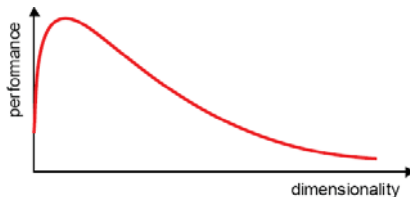
**Prof. Dongheui Lee**

*Institute of Automatic Control Engineering  
Technische Universität München*

dhlee@tum.de

# The curse of dimensionality

- The problems associated with multivariate data analysis as the dimensionality increases
- Toy example : 3-class pattern recognition or an approach which divide the sample space into equally spaced bins
  - 1D  $\rightarrow$  3 bins  $\rightarrow 3^2$  examples
  - 2D  $\rightarrow 3^2$  bins  $\rightarrow 3^3$  examples
  - 3D  $\rightarrow 3^3$  bins  $\rightarrow 3^4$  examples
- For a given sample size, there is a maximum number of features above which the performance of our classifier will degrade rather than improve



# The curse of dimensionality

## Implication of the curse of dimensionality

- Exponential growth in the number of examples, required to maintain a given sampling density
- Exponential growth in the complexity of the target function (a density estimate) with increasing dimensionality
- Humans have an extraordinary capacity to discern patterns in 1~3D. But these capability degrades drastically for higher dimensions

## Why dimensionality reduction?

- Learning a target function from data where some features are irrelevant → reduce variance. Improve accuracy
- Wish to visualize high dimensional data
- Intrinsic dimensionality of data is smaller than the number of features used to describe it

# Dimensionality Reduction

## Two approaches

- Feature extraction: creating a subset of new features by combinations of the existing features  $m > k$

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} \rightarrow \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_k \end{bmatrix} = f \left( \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} \right)$$

- Feature selection: choosing a subset of all the features (the ones more informative)

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} \rightarrow \begin{bmatrix} x_{i_1} \\ x_{i_2} \\ \vdots \\ x_{i_k} \end{bmatrix}$$

# Feature Extraction Problem

- Given a feature space, to find a mapping  $y = f(x)$  such that the reduced feature vector  $y$  preserves most of the information or structure of original feature  $x$
- The optimal mapping function  $f(x)$  will result in no increase in the minimum probability of error
- In general, the mapping function can be a nonlinear function. But, often feature extraction is limited to linear transforms.  $y = W^T x$

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} \rightarrow \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_k \end{bmatrix} = \begin{bmatrix} u_{11} & u_{12} & \cdots & u_{1m} \\ u_{21} & u_{22} & \cdots & u_{2m} \\ & & \ddots & \\ u_{k1} & u_{k2} & \cdots & u_{km} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}$$

# Two Approaches for Feature Extraction

## Unsupervised approaches

- Principal Components Analysis
- Singular Value Decomposition
- Independent Components Analysis

## Supervised approaches

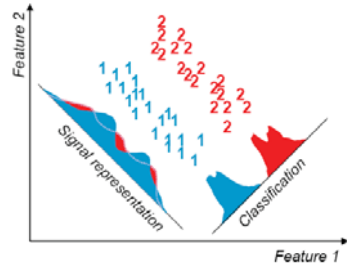
- Fisher Linear Discriminant Analysis (LDA)
- Hidden Layers of Neural Networks

# Feature Extraction

The selection of feature extraction mapping  $y = f(x)$  is guided by an objective function to maximize (or minimize)

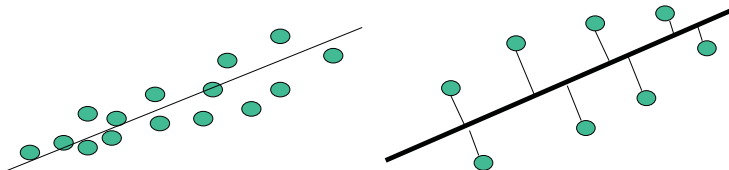
Depending on the criteria measures by objective function. Two categories

- Signal representation
  - Goal : to represent samples accurately
  - Principal Components Analysis (PCA)
- Classification
  - Goal : to <sup>增加</sup>enhance the <sup>有辨识力的</sup>class-discriminatory information
  - Linear Discriminant Analysis (LDA)



# Principal Components Analysis (PCA)

- Given data points in  $m$ -dimensional space, project into lower dimensional space while preserving as much information as possible  $\Rightarrow$  Choose a line that fits the data so the points are spread out well along the line
- Seeks a projection that best represent the data in a least-square sense  $\Rightarrow$  minimize sum of squares of distances to the line
- We wish to explain/summarize the underlying variance-covariance structure of a large set of variables through a few linear combinations of these variables





# Space Representation

- High-dimensional space representation  $\mathbf{x}$
- Low-dimensional space representation  $\mathbf{y}$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_k \end{bmatrix}$$

- PCA allows us to compute a linear transformation that maps data from a high dimensional space to a lower dimensional sub-space

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} \rightarrow \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_k \end{bmatrix} = \begin{bmatrix} u_{11} & u_{12} & \cdots & u_{1m} \\ u_{21} & u_{22} & \cdots & u_{2m} \\ & & \ddots & \\ u_{k1} & u_{k2} & \cdots & u_{km} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}$$

## Maximum variance formulation

- Goal: to project the data onto a space having dimensionality  $k < m$  while maximizing the variance of the projected data
- Begin with  $k = 1$
- $m$ -dimensional unit vector  $\mathbf{u}_1 = [u_{11} \quad u_{12} \quad \cdots \quad u_{1m}]^T$ ,  $\mathbf{u}_1^T \mathbf{u}_1 = 1$
- Each data point  $\mathbf{x}^{(i)}$  is then projected onto a scalar value  $y^{(i)} = \mathbf{u}_1^T \mathbf{x}^{(i)}$
- Mean, covariance of the original data  
 $\boldsymbol{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)}$ ,  $\mathbf{S} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}^{(i)} - \boldsymbol{\mu})(\mathbf{x}^{(i)} - \boldsymbol{\mu})^T$
- mean of the projected data is  $\mathbf{u}_1^T \boldsymbol{\mu}$  mean of y
- variance of the projected data  $\frac{1}{n} \sum_{i=1}^n \{\mathbf{u}_1^T \mathbf{x}^{(i)} - \mathbf{u}_1^T \boldsymbol{\mu}\}^2 = \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$  Variance of y
- Maximize the projected variance  $\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1$  wrt  $\mathbf{u}_1$  Lagrange Multiplier

$$\frac{\partial}{\partial \mathbf{u}_1} \mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 + \lambda_1 (1 - \mathbf{u}_1^T \mathbf{u}_1) = 0, \quad \mathbf{S} \mathbf{u}_1 = \lambda_1 \mathbf{u}_1$$

1- $\mathbf{u}^T \mathbf{u}$ : normalize  
避免  $\mathbf{u} = k\mathbf{u}$  多解

- Variance will be a maximum value when  $\mathbf{u}_1$  is equal to the eigenvector having the largest eigenvalue  $\lambda_1$
- When  $k$  increases, choose  $\mathbf{u}_1 \dots \mathbf{u}_k$  as corresponding eigenvectors to  $k$  largest eigenvalues  $\lambda_1, \lambda_2 \dots \lambda_k$



# PCA in General

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} \rightarrow \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_k \end{bmatrix} = \begin{bmatrix} u_{11} & u_{12} & \cdots & u_{1m} \\ u_{21} & u_{22} & \cdots & u_{2m} \\ & & \vdots & \\ u_{k1} & u_{k2} & \cdots & u_{km} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}$$

- $\{\mathbf{u}_j\}$  are orthogonal basis vectors.
- $\mathbf{u}_1 = [u_{11} \ u_{12} \ \cdots \ u_{1m}]^T$  is 1st eigenvector of covariance matrix, and known as the first principal component. This axis explains as much as possible of original variance in data set
- $\mathbf{u}_2 = [u_{21} \ u_{22} \ \cdots \ u_{2m}]^T$  is 2nd eigenvector of covariance matrix, and the 2nd principal component
- $\mathbf{u}_k = [u_{k1} \ u_{k2} \ \cdots \ u_{km}]^T$  is  $k$ th eigenvector of covariance matrix, and the  $k$ th principal component

# PCA algorithm

1. Compute the mean and covariance matrix

$$\mu = \frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)}, \quad S = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}^{(i)} - \mu)(\mathbf{x}^{(i)} - \mu)^T$$

x 为列向量  
S 为方阵

2. Find eigenvectors and eigenvalues of  $S$
3. Choose the  $k$  largest eigenvalues  $\lambda_1, \lambda_2 \dots \lambda_k$
4. Choose the corresponding eigenvectors  $\mathbf{u}_1 \dots \mathbf{u}_k$  and make transformation matrix  $\mathbf{W} = [\mathbf{u}_1 \mathbf{u}_2 \dots]$
5. Project original data  $\mathbf{y} = \mathbf{W}^T \mathbf{x}$

注意是S的特征向量构造的变换阵！不是X矩阵！

# Principal Components Analysis

- Principal Components Analysis is the oldest technique in multivariate analysis
- PCA is also known as the Karhunen-Loève transform
- PCA was first introduced by Pearson in 1901, and generalized by Loève in 1963
- The main limitation of PCA is that it does not consider class separability since it does not take into account the class label of the feature vector
  - PCA simply performs a coordinate rotation that aligns the transformed axes with the directions of maximum variance
  - There is no guarantee that the directions of maximum variance will contain good features for discrimination

匹配

最大方差对应的方向未必包含适合分类的特征

没有考虑到是否可分。（有些特征向量可能属于一类，但是其中对应较小特征根的却在变换中被舍弃。）

# PCA - simple example

Compute the principal components for the following two-dimensional dataset

$X = (x_1, x_2) = (1, 2), (3, 3), (3, 5), (5, 4), (5, 6), (6, 5), (8, 7), (9, 8)$

Solution

- The (biased) covariance estimate of the data is

$$S = \begin{bmatrix} 6.25 & 4.25 \\ 4.25 & 3.5 \end{bmatrix}$$

- The eigenvalues are the zeros of the characteristic equation

$$Su = \lambda u \Rightarrow |S - \lambda I| = 0 \Rightarrow \lambda_1 = 9.34, \lambda_2 = 0.41$$

- The eigenvectors are the solutions of the system

$$\begin{bmatrix} 6.25 & 4.25 \\ 4.25 & 3.5 \end{bmatrix} \begin{bmatrix} u_{11} \\ u_{12} \end{bmatrix} = \begin{bmatrix} \lambda_1 u_{11} \\ \lambda_1 u_{12} \end{bmatrix} \Rightarrow \begin{bmatrix} u_{11} \\ u_{12} \end{bmatrix} = \begin{bmatrix} 0.81 \\ 0.59 \end{bmatrix}$$

$$\begin{bmatrix} 6.25 & 4.25 \\ 4.25 & 3.5 \end{bmatrix} \begin{bmatrix} u_{21} \\ u_{22} \end{bmatrix} = \begin{bmatrix} \lambda_2 u_{21} \\ \lambda_2 u_{22} \end{bmatrix} \Rightarrow \begin{bmatrix} u_{21} \\ u_{22} \end{bmatrix} = \begin{bmatrix} -0.59 \\ 0.81 \end{bmatrix}$$

简化计算的方法：另  $u_1=1$ ，求出  $u_2$  然后再归一化。

# Linear Discriminant Analysis

- Multiple classes and PCA
  - Suppose there are  $C$  classes in the training data
  - PCA is based on the sample covariance which characterizes the scatter of the entire data set, irrespective of class-membership PCA缺点
  - The projection axes chosen by PCA might not provide good discrimination power
- What is the goal of LDA?
  - Perform dimensionality reduction while preserving as much of the class discriminatory information as possible
  - Seeks to find directions along which the classes are best separated
  - Takes into consideration the scatter within-classes but also the scatter between-classes 分散
  - More capable of distinguishing image variation due to identity from variation due to other sources such as illumination and expression  
不但考虑到class内的分散性，也考虑到class之间的分散性。适合图片分类。

# Linear Discriminant Analysis

In order to find a good projection vector, we need to define a measure of separation between the projected classes

- Projection  $y = W^T x$
- The mean of original and projected dataset

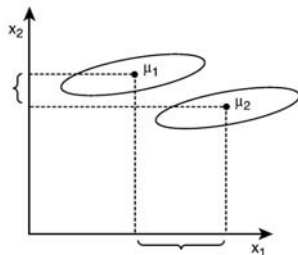
Original  $\mu_j = \frac{1}{n_j} \sum_{x \in X_j} x$  ,  $\tilde{\mu}_j = \frac{1}{n_j} \sum_{y \in Y_j} y = \frac{1}{n_j} \sum_{x \in X_j} w^T x = w^T \mu_j$  Projected

- The case of two classes, the distance between the projected means

$$|\tilde{\mu}_1 - \tilde{\mu}_2| = |w^T(\mu_1 - \mu_2)|$$

只考虑不同class之间的差别不够分类

⇒ Not so good measure, since it does not take into account the standard deviation within the class.





# Fisher Linear Discriminant Analysis

- Proposed by Ronald Aylmer Fisher
- To maximize a function that represents the difference between means normalized by a measure of the within-class scatter

$$J(\mathbf{w}) = \frac{|\tilde{\mu}_1 - \tilde{\mu}_2|^2}{S_1 + S_2}$$

- Define the scatter matrix for each class
  - Original space  $\mu_i = \frac{1}{n_i} \sum_{\mathbf{x} \in X_i} \mathbf{x}$  ,  $S_i = \frac{1}{n_i} \sum_{\mathbf{x} \in X_i} (\mathbf{x} - \mu_i)(\mathbf{x} - \mu_i)^T$
  - projected space  
 $\tilde{\mu}_i = \frac{1}{n_i} \sum_{\mathbf{x} \in X_i} \mathbf{w}^T \mathbf{x}$  ,  $\tilde{S}_i = \frac{1}{n_i} \sum_{\mathbf{x} \in X_i} \mathbf{w}^T (\mathbf{x} - \mu_i)(\mathbf{x} - \mu_i)^T \mathbf{w} = \mathbf{w}^T S_i \mathbf{w}$
- Within-class scatter matrix  $S_W = S_1 + S_2$  ,  $\tilde{S}_1 + \tilde{S}_2 = \mathbf{w}^T S_W \mathbf{w}$
- Between-class scatter matrix  
 $S_B = (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T$  ,  $(\tilde{\mu}_1 - \tilde{\mu}_2)^2 = \mathbf{w}^T S_B \mathbf{w}$
- The criterion function  $J(\mathbf{w}) = \frac{\mathbf{w}^T S_B \mathbf{w}}{\mathbf{w}^T S_W \mathbf{w}}$  SB尽量大, SW尽量小。最大化J。

Sy\_in : Sy\_between

# Linear Discriminant Analysis

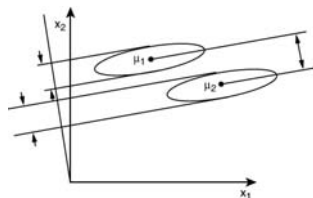
Maximize  $J(\mathbf{w})$

$$\frac{\partial}{\partial \mathbf{w}} [J(\mathbf{w})] = \frac{\partial}{\partial \mathbf{w}} \left[ \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} \right] = 0$$

$\vdots$

$$\mathbf{S}_W^{-1} \mathbf{S}_B \mathbf{w} = J \mathbf{w}$$

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \left[ \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} \right] = \mathbf{S}_W^{-1} (\mu_1 - \mu_2) \quad \text{参考 Bishop.189}$$



Properties of LDA

- LDA is a parametric method since it assumes unimodal Gaussian likelihoods  
LDA假设统一高斯分布likelihood!
- If the distributions are significantly non-Gaussian, the LDA projections will not be able to preserve any complex structure of the data that may be needed for classification

# PCA application: Face recognition

Eigenfaces for face detection/recognition

- Template matching problem
- Difficult to perform recognition in a high-dimensional space
- Improvements can be achieved by mapping data into a lower dimensionality space

Data is projected into a lower dimensional space using PCA.

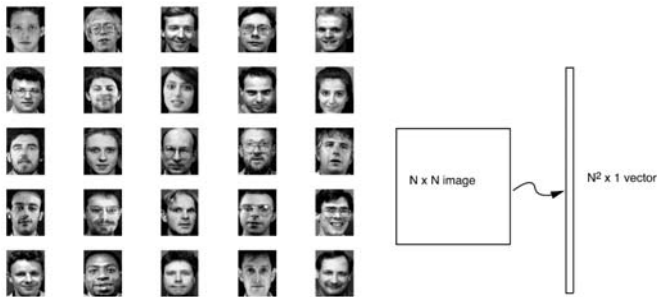


M. Turk, A. Pentland, *Eigenfaces for Recognition*. Journal of Cognitive Neuroscience, 3(1), pp. 71-86, 1991.



# PCA application: Face recognition

- Step 1. Obtain face images (training data set)
- Step 2. Preprocess the face images (centered and same size) and represent each image as a vector  $I^{(i)} \rightarrow x^{(i)}, \forall i = 1, \dots, n$



# PCA application: Face recognition

- Step 3. Compute the average face

$$\mu = \frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)}$$



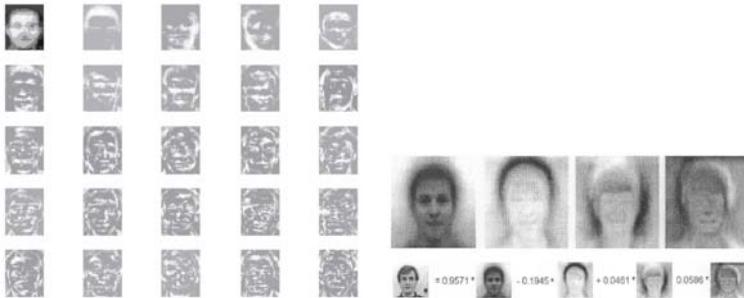
- Step 4. Compute the covariance matrix

$$\mathbf{S} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}^{(i)} - \mu)(\mathbf{x}^{(i)} - \mu)^T$$

- Step 5. Compute eigenvalues of the covariance matrix
- Step 6. Choose the  $k$  largest eigenvalues and their corresponding eigenvectors

# PCA application: Face recognition

- Now, we can represent each face as a linear combination of the best  $k$  eigenvectors
- Eigenfaces: eigenvectors of face dataset



# PCA application: Face recognition

- Step 7. For face recognition of a new image  $x$ , normalize the input image. Project on the eigenspace.  $\text{size}(x) = [N^2, 1]$

$$y_i = u_i^T x \quad \forall i = 1, \dots, k \quad \mathbf{y} = [y_1 \dots y_k]^T$$

最小距离分类

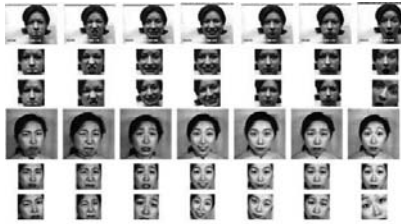
- Step 8. Compare the projected input image  $\mathbf{y}$  with face classes  $\mathbf{y}_j$ . Find the class which results in smallest (Euclidean / Mahalanobis) distance. If the minimum Euclidean distance is smaller than a threshold, the image is recognized as the face class.  
If  $\min \| \mathbf{y} - \mathbf{y}_j \| < \tau$ ,  $\mathbf{y}$  belongs to the  $j$ -th face class.

# LDA Application: Emotion recognition

- Biased linear discriminant analysis (BLDA) method that impose large penalties on interclass samples with small differences and small penalties on those samples with large differences

$$S_B = g(i,j)(\mu_1 - \mu_2)(\mu_1 - \mu_2)^T$$

类间差距小的加权大，  
类间差距大的加权小。



**Figure:** Examples of the original, well-aligned, and misaligned images of one subject from the (upper half) Cohn-Kanade and (lower half) JAFFE database. From left to right are the facial images with anger, disgust, fear, happy, neutral, sad, and surprise expressions, respectively.



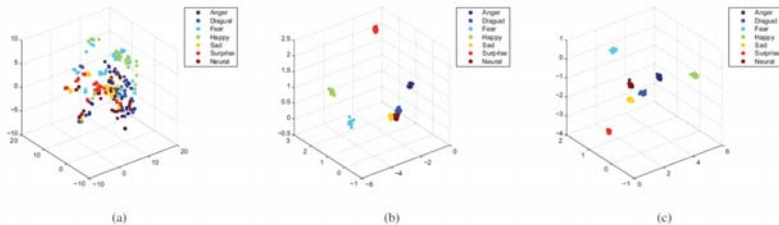
Haibin Yan, Marcelo H. Ang Jr, and Aun Neow Poo, *Weighted biased linear discriminant analysis for misalignment-robust facial expression recognition*. ICRA 2011.





# LDA Application

- LDA is a supervised subspace learning approach which searches for a set of most discriminative projections to maximize the ratio of between-class variance to within-class variance simultaneously
- When the class information is available, LDA usually outperforms PCA for classification tasks



# LDA Application

- LDA is a supervised subspace learning approach which searches for a set of most discriminative projections to maximize the ratio of between-class variance to within-class variance simultaneously
- When the class information is available, LDA usually outperforms PCA for classification tasks

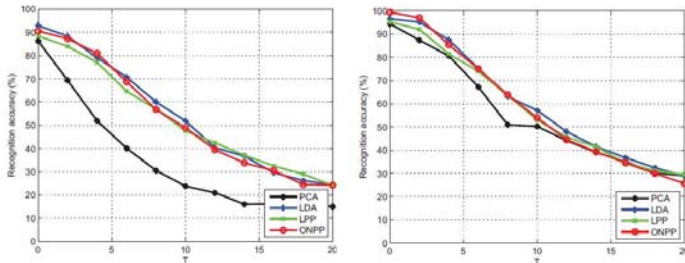
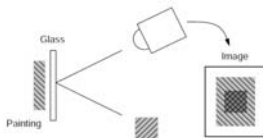


Figure: Recognition accuracy versus different amounts of spatial misalignments. (a) Results obtained on the Cohn-Kanade database. (b) Results obtained on the JAFFE database.

# Separating reflections using Independent Components Analysis

A mixed image of a painting and an object is represented as  $y_1 = aP + bR$ , by photographing through a linear polarizer, the relative strength of the reflections can be adjusted.



**Figure:** A photograph of a painting behind glass contains a superposition of the light reflected by the painting and the light reflected directly off the glass.



Farid, Hany, and Edward H. Adelson, *Separating reflections and lighting using independent components analysis*. IEEE Computer Society Conference on Computer Vision and Pattern Recognition, Vol.1, 1999.



# Announcements

- Further Reading
  - Duda: Chap 3.7, 3.8.1, 3.8.2, 10.13
  - Bishop: Chap. 12.1-12.3, 12.4.1
  - PCA: Tutorial by J. Schlens