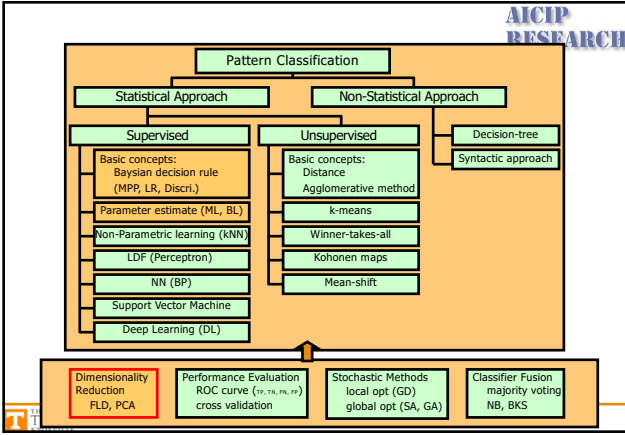


THE UNIVERSITY OF TENNESSEE KNOXVILLE **AICIP RESEARCH**

ECE471-571 – Pattern Recognition

Lecture 6 – Dimensionality Reduction – Fisher’s Linear Discriminant

Hairong Qi, Gonzalez Family Professor
 Electrical Engineering and Computer Science
 University of Tennessee, Knoxville
<http://www.eecs.utk.edu/faculty/qj>
 Email: hqi@utk.edu



Review - Bayes Decision Rule **AICIP RESEARCH**

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

Maximum Posterior Probability: For a given x , if $P(\omega_1 | x) > P(\omega_2 | x)$, then x belongs to class 1, otherwise, 2.

Likelihood Ratio: If $\frac{p(x|\omega_1)}{p(x|\omega_2)} > \frac{P(\omega_2)}{P(\omega_1)}$, then x belongs to class 1, otherwise, 2.

Discriminant Function: The classifier will assign a feature vector x to class ω_1 if $g_+(x) > g_-(x)$.

Case 1: Minimum Euclidean Distance (Linear Machine), $\Sigma_i = \sigma^2 I$
 Case 2: Minimum Mahalanobis Distance (Linear Machine), $\Sigma_i = \Sigma$
 Case 3: Quadratic classifier, $\Sigma_i = \text{arbitrary}$

THE UNIVERSITY OF TENNESSEE 3

The Curse of Dimensionality – 1st Aspect

- ◆ The number of training samples
- ◆ What would the probability density function look like if the dimensionality is very high?
 - For a 7-dimensional space, where each variable could have 20 possible values, then the 7-d histogram contains 20^7 cells. To distribute a training set of some reasonable size (1000) among this many cells is to leave virtually all the cells empty

Curse of Dimensionality – 2nd Aspect

- ◆ Accuracy and overfitting
- ◆ In theory, the higher the dimensionality, the less the error, the better the performance. However, in realistic pattern recognition problems, the opposite is often true. Why?
 - The assumption that pdf behaves like Gaussian is only approximately true
 - When increasing the dimensionality, we may be **overfitting** the training set.
 - Problem: **excellent** performance on the training set, **poor** performance on new data points which are in fact very close to the data within the training set

Curse of Dimensionality - 3rd Aspect

- ◆ Computational complexity

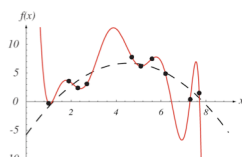


FIGURE 3.4. The "training data" (black dots) were selected from a quadratic function plus Gaussian noise, i.e., $f(x) = ax^2 + bx + c + \epsilon$ where $\epsilon \sim N(0, \sigma^2)$. The 10th-degree polynomial shown fits the data perfectly, but we desire instead the second-order function $f(x)$, because it would lead to better predictions for new samples. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

Dimensionality Reduction

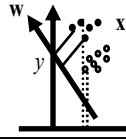
- Fisher's linear discriminant
 - Best **discriminating** the data
- Principal component analysis (PCA)
 - Best **representing** the data

Fisher's Linear Discriminant

- For two-class cases, projection of data from d-dimension onto a line
- **Principle:** We'd like to find vector w (direction of the line) such that the projected data set can be best separated

$$y = w^T x$$

$$J(w) = |\tilde{m}_1 - \tilde{m}_2|^2 = |w^T(m_1 - m_2)|^2$$



$$\tilde{m}_1 = \frac{1}{n_1} \sum_{x \in \mathcal{C}_1} y = \frac{1}{n_1} \sum_{x \in \mathcal{C}_1} w^T x = w^T m_1$$

$$m_1 = \frac{1}{n_1} \sum_{x \in \mathcal{C}_1} x$$

Projected mean

Sample mean

Other Approaches?

- ♦ Solution 1: make the projected mean as apart as possible
- ♦ Solution 2?

$$J(w) = \frac{|\tilde{m}_1 - \tilde{m}_2|^2}{\hat{s}_1^2 + \hat{s}_2^2} = \frac{|w^T(m_1 - m_2)|^2}{w^T S_1 w + w^T S_2 w} = \frac{w^T S_B w}{w^T S_W w}$$

$$\hat{s}_1^2 = \sum_{x \in \mathcal{C}_1} (y - \tilde{m}_1)^2 = \sum_{x \in \mathcal{C}_1} (w^T x - w^T m_1)^2 = \sum_{x \in \mathcal{C}_1} w^T (x - m_1)(x - m_1)^T w = w^T S_1 w$$

$$\text{Scatter matrix } S_i = \sum_{x \in \mathcal{C}_i} (x - m_i)(x - m_i)^T$$

$$\text{Between-class scatter matrix } S_B = (m_1 - m_2)(m_1 - m_2)^T$$

$$\text{Within-class scatter matrix } S_W = S_1 + S_2 = \sum_{x \in \mathcal{C}_1 \cup \mathcal{C}_2} (x - m_i)(x - m_i)^T$$

*The Generalized Rayleigh Quotient

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}$$

$$\frac{dJ(\mathbf{w})}{d\mathbf{w}} = \frac{2\mathbf{S}_B \mathbf{w} (\mathbf{w}^T \mathbf{S}_W \mathbf{w}) - 2\mathbf{S}_W \mathbf{w} (\mathbf{w}^T \mathbf{S}_B \mathbf{w})}{(\mathbf{w}^T \mathbf{S}_W \mathbf{w})^2} = 0$$

$$\mathbf{S}_B \mathbf{w} = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} \mathbf{S}_W \mathbf{w} \Rightarrow \mathbf{S}_W^{-1} \mathbf{S}_B \mathbf{w} = \lambda \mathbf{w}$$

$\mathbf{S}_B \mathbf{w}$ is always in the direction of $\mathbf{m}_1 - \mathbf{m}_2$

$$\mathbf{w} = \mathbf{S}_W^{-1} (\mathbf{m}_1 - \mathbf{m}_2) \quad \text{Canonical variate}$$

Some Math Preliminaries

- ◆ Positive definite
 - A matrix \mathbf{S} is positive definite if $y = \mathbf{x}^T \mathbf{S} \mathbf{x} > 0$ for all $\mathbf{x} \in \mathbb{R}^d$ except 0
 - $\mathbf{x}^T \mathbf{S} \mathbf{x}$ is called the quadratic form
 - The derivative of a quadratic form is particularly useful

$$\frac{d}{d\mathbf{x}} (\mathbf{x}^T \mathbf{S} \mathbf{x}) = (\mathbf{S} + \mathbf{S}^T) \mathbf{x}$$

- ◆ Eigenvalue and eigenvector
 - \mathbf{x} is called the eigenvector of \mathbf{A} iff \mathbf{x} is not zero, and $\mathbf{A}\mathbf{x} = \lambda \mathbf{x}$
 - λ is the eigenvalue of \mathbf{x}

* Multiple Discriminant Analysis

- ◆ For c-class problem, the projection is from d-dimensional space to a (c-1)-dimensional space (assume $d \geq c$)
- ◆ Sec. 3.8.3
