

Outline

- Recognition Based on Decision-Theoretic Methods
 - Optimum Statistical Classifiers
 - Neural Networks

Optimum Statistical Classifiers

- Previously, we had considered very simple classifiers without using any high order statistics of the underlying patterns
- If we have knowledge about some sort of statistics of the patterns, we can use information to design a classifier that would minimize the probability of incorrect classifications
- Let us define the probability that the pattern \mathbf{x} belongs to ω_i as $p(\omega_i/\mathbf{x})$
- If our classifier incorrectly decides on ω_j , then we define this error as L_{ij}
- Averaging this error for all classes we obtain the average error as

$$r_j(\mathbf{x}) = \sum_{k=1}^W L_{kj} p(\omega_k/\mathbf{x})$$

Optimum Statistical Classifiers

- Using the Bayes' rule we obtain

$$r_j(\mathbf{x}) = \frac{1}{p(\mathbf{x})} \sum_{k=1}^W L_{kj} p(\mathbf{x}/\omega_k) P(\omega_k)$$

- We use this average error as our cost function, and choose the class that results in the minimum error
- All r_j 's are calculated for $j = 1, 2, \dots, W$ and the class with minimum error is selected
- Now let us select a value for L_{ij} , a reasonable choice is to make it zero when $i = j$ and one otherwise
- Then we have

$$d_j(\mathbf{x}) = r_j(\mathbf{x}) = p(\mathbf{x}/\omega_j) P(\omega_j)$$

Optimum Statistical Classifiers

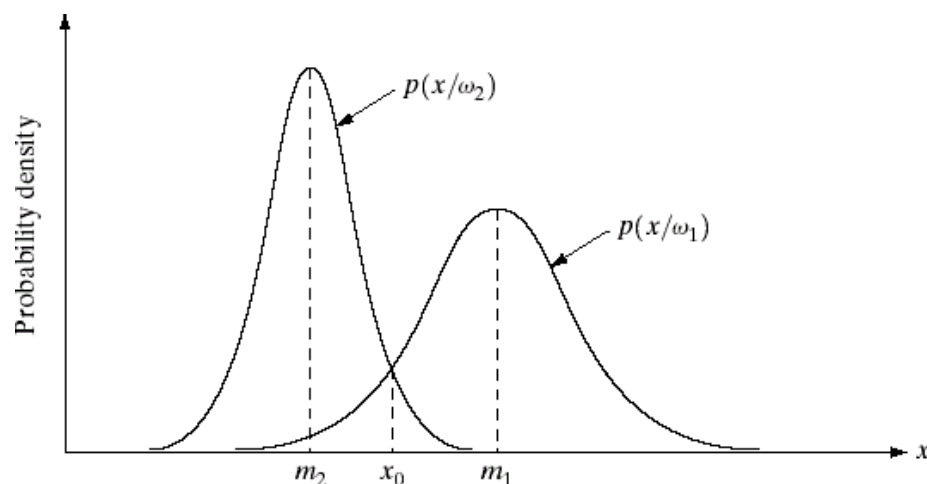
- To be able to calculate this optimum classifier we need to know the probability of each class ($P(\omega)$) and the pdf of the individual patterns must be known
- It is a reasonable assumption that we can infer the knowledge of $P(\omega)$ based on the problem at hand
- However, it is not an easy task to estimate $p(\mathbf{x}/\omega_j)$
- Therefore, instead of estimating these directly, we use a parametric model and estimate a few parameters hence obtain an approximate of $p(\mathbf{x}/\omega_j)$'s

Optimum Statistical Classifiers

- Let us use a Gaussian model and assume a scalar x for simplicity. Then the decision function is

$$d_j(x) = \frac{1}{\sqrt{2\pi}\sigma_j} e^{-\frac{(x-m_j)^2}{2\sigma_j^2}} P(\omega_j)$$

- If we assume that we have only two classes the problem of classifying is summarized in the following figure



Optimum Statistical Classifiers

- The 1-D example illustrates the idea, but in reality we will almost always have multidimensional patterns
- Then the Gaussian will be multivariate

$$p(\mathbf{x}/\omega_j) = \frac{1}{(2\pi)^{\frac{n}{2}} |C_j|^{1/2}} e^{-\frac{1}{2}(\mathbf{x}-\mathbf{m}_j)^T C_j^{-1} (\mathbf{x}-\mathbf{m}_j)}$$

where C_j is the covariance matrix and \mathbf{m}_j is the mean vector

- These quantities can be estimated based on the samples

$$\hat{\mathbf{m}}_j = \frac{1}{N_j} \sum \mathbf{x}$$

$$\hat{C}_j = \frac{1}{N_j} \sum \mathbf{x}\mathbf{x}^T - \hat{\mathbf{m}}_j \hat{\mathbf{m}}_j^T$$

Optimum Statistical Classifiers

- Now let us minimize $d_j(\mathbf{x})$ by minimizing its log (which is a monoton function)

$$\log P(\omega_j) - \frac{n}{2} \log(2\pi) - \frac{1}{2} \log |C_j| - \frac{1}{2} [(\mathbf{x} - \mathbf{m}_j)^T C_j^{-1} (\mathbf{x} - \mathbf{m}_j)]$$

- This expression has a special form when all covariance matrices are identity and all classes are equally probable

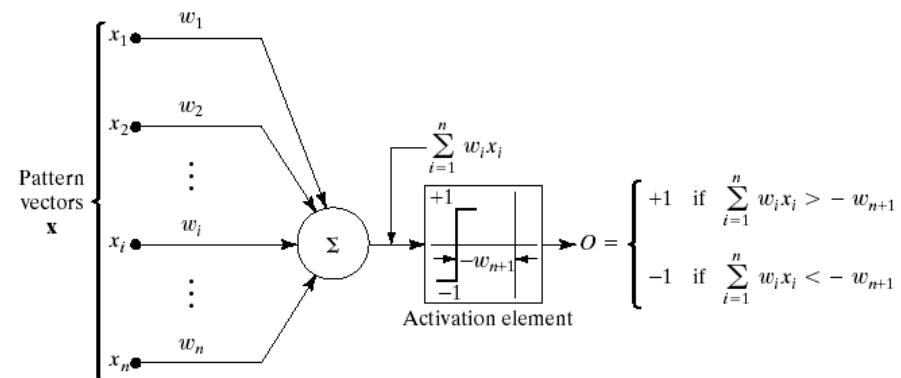
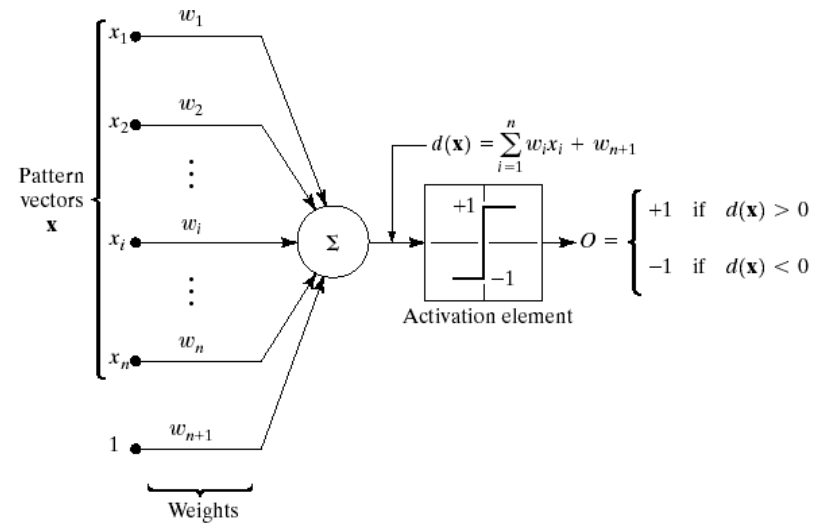
$$d_j(\mathbf{x}) = \mathbf{x}^T \mathbf{m}_j - \frac{1}{2} \mathbf{m}_j^T \mathbf{m}_j$$

- Remember this! Then the minimum distance classifier is the optimum bayesian classifier when all pattern classes are Gaussian distributed with identity covariance matrix and all classes are equally probable

Neural Networks: Introduction

- Instead of assuming a parametric model for the pdf's of patterns, we can use an approach called “training”
- Known patterns are used first for “training” and learning. In this way, the decision functions are created based on these known patterns
- This approach brings us to neural networks

Neural Networks: A simple two pattern classes example



Neural Networks: A simple two pattern classes example

- The decision is based on $d(\mathbf{x})$, for positive $d(\mathbf{x})$ the particular class is selected, for negative $d(\mathbf{x})$ it is not selected
- When $d(\mathbf{x})$ is exactly zero, we can obtain the decision boundary

$$d(\mathbf{x}) = w_1x_1 + w_2x_2 + \dots + w_{n+1}$$

- To be able to reconstruct a decision function we need to calculate the weights with training algorithms

Neural Networks: Training Algorithms - Linearly Separable Classes

- We select an initial coefficient vector $\mathbf{w}(1)$ and then using the training data at each class, update the coefficient vector as follows

$$\mathbf{w}(k+1) = \mathbf{w}(k) + c\mathbf{y}(k) \quad \text{when } \mathbf{w}^T(k)\mathbf{y}(k) < 0, \mathbf{y}(k) \in \omega_1$$

$$\mathbf{w}(k+1) = \mathbf{w}(k) - c\mathbf{y}(k) \quad \text{when } \mathbf{w}^T(k)\mathbf{y}(k) > 0, \mathbf{y}(k) \in \omega_2$$

where c is a positive constant $\mathbf{y} = [x_1, x_2, \dots, x_n, 1]$, and

$$\mathbf{w} = [w_1, w_2, \dots, w_{n+1}]$$

- This algorithm converges when the two sets are linearly separable

Neural Networks: Training Algorithms - Nonseparable classes, Widrow-Hoff Method

- This training method can be used with linearly nonseparable classes
- The algorithm minimizes the error between the desired and actual response

$$[r - \mathbf{w}^T \mathbf{y}]^2$$

where r is the desired response, 1 if the training pattern belongs to class 1, and -1 if it belongs to class 2

- This cost function can be minimized using gradient based numerical optimization methods

$$\mathbf{w}(k+1) = \mathbf{w}(k) - \alpha \left[\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} \right]_{\mathbf{w}=\mathbf{w}(k)}$$

Neural Networks: Training Algorithms - Nonseparable classes, Widrow-Hoff Method

- Substituting the derivative of the cost function we can obtain

$$\mathbf{w}(k + 1) = \mathbf{w}(k) + \alpha[r - \mathbf{w}^T \mathbf{y}] \mathbf{y}$$

- The only parameter to be selected is α , a reasonable choice is between 0.1 and 1