

Pattern Recognition Lecture

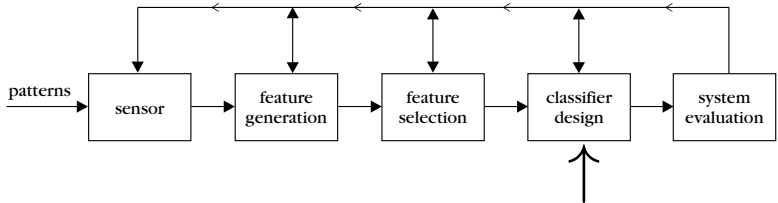
Bayes Decision Theory

Prof. Dr. Marcin Grzegorzek

Research Group for Pattern Recognition
Institute for Vision and Graphics
University of Siegen, Germany



Pattern Recognition Chain



Introduction

Bayes Decision
Theory

Discriminant
Functions and
Decision
Surfaces

Bayesian
Classification
for Normal
Distributions

Estimation of
Unknown
Probability
Density
Functions

Overview

Introduction

Bayes Decision
Theory

Discriminant
Functions and
Decision
Surfaces

Bayesian
Classification
for Normal
Distributions

Estimation of
Unknown
Probability
Density
Functions

- 1 Introduction
- 2 Bayes Decision Theory
- 3 Discriminant Functions and Decision Surfaces
- 4 Bayesian Classification for Normal Distributions
- 5 Estimation of Unknown Probability Density Functions

Overview

Introduction

Bayes Decision
Theory

Discriminant
Functions and
Decision
Surfaces

Bayesian
Classification
for Normal
Distributions

Estimation of
Unknown
Probability
Density
Functions

- 1 Introduction
- 2 Bayes Decision Theory
- 3 Discriminant Functions and Decision Surfaces
- 4 Bayesian Classification for Normal Distributions
- 5 Estimation of Unknown Probability Density Functions

Statistical Classification - Problem Statement

Introduction

Bayes Decision
Theory

Discriminant
Functions and
Decision
Surfaces

Bayesian
Classification
for Normal
Distributions

Estimation of
Unknown
Probability
Density
Functions

Classification of an unknown pattern in the most probable of the classes!

- Set of classes: $\{\omega_1, \omega_2, \dots, \omega_M\}$
- Unknown pattern represented by its feature vector \mathbf{x}
- Conditional probabilities: $P(\omega_i|\mathbf{x}), \quad i = 1, 2, \dots, M$
- Classification result: the class with the maximum conditional probability

But how to compute the conditional probability for a particular class?

Probability P vs. Density p

Introduction

Bayes Decision Theory

Discriminant Functions and Decision Surfaces

Bayesian Classification for Normal Distributions

Estimation of Unknown Probability Density Functions

Probability P

is a real number describing an event belonging to the range $< 0, 1 >$.

Density p

is a value of a function¹ $p(x)$ describing the distribution of the random variable x .

If the random variable takes only discrete values, the densities become probabilities!

¹This function is often referred as pdf - probability density function.

Overview

Introduction

**Bayes Decision
Theory**

Discriminant
Functions and
Decision
Surfaces

Bayesian
Classification
for Normal
Distributions

Estimation of
Unknown
Probability
Density
Functions

- 1 Introduction
- 2 Bayes Decision Theory**
- 3 Discriminant Functions and Decision Surfaces
- 4 Bayesian Classification for Normal Distributions
- 5 Estimation of Unknown Probability Density Functions

A Priori Probability vs. A Posteriori Probability

Introduction

Bayes Decision Theory

Discriminant Functions and Decision Surfaces

Bayesian Classification for Normal Distributions

Estimation of Unknown Probability Density Functions

A priori probability - probability before classification

- How probable is a particular class ω_i for a pattern \mathbf{x} before applying any classification algorithm?
- Answer: $P(\omega_i)$

A posteriori probability - probability after classification

- How probable is a particular class ω_i for a pattern \mathbf{x} after applying a statistical classification algorithm?
- Answer: $P(\omega_i|\mathbf{x})$

Likelihood Density Function

Introduction

Bayes Decision
Theory

Discriminant
Functions and
Decision
Surfaces

Bayesian
Classification
for Normal
Distributions

Estimation of
Unknown
Probability
Density
Functions

Likelihood Density Function

- How feature vectors \mathbf{x} are distributed in a class ω_i ?
- Answer: $p(\mathbf{x}|\omega_i)$
- $p(\mathbf{x}|\omega_i)$ is the likelihood function of ω_i with respect to \mathbf{x}
- $p(\mathbf{x}|\omega_i)$ can be trained from examples

Bayes Decision Theory for a Two-Class Problem

Known

Classes:	$\{\omega_1, \omega_2\}$
A priori probabilities:	$P(\omega_1)$ and $P(\omega_2)$
Likelihood density functions:	$p(\mathbf{x} \omega_1)$ and $p(\mathbf{x} \omega_2)$
Pattern to be classified:	$\mathbf{x} = [x_1, x_2, \dots, x_l]^T$

Assumption

The feature vectors can take any value in the l -dimensional feature space: $\mathbf{x} = [x_1, x_2, \dots, x_l]^T \in \mathbb{R}^l$

Unknown

A posteriori probabilities: $P(\omega_1|\mathbf{x})$ and $P(\omega_2|\mathbf{x})$

Introduction

Bayes Decision Theory

Discriminant Functions and Decision Surfaces

Bayesian Classification for Normal Distributions

Estimation of Unknown Probability Density Functions

Computation of the A Posteriori Probability

Introduction

Bayes Decision
Theory

Discriminant
Functions and
Decision
Surfaces

Bayesian
Classification
for Normal
Distributions

Estimation of
Unknown
Probability
Density
Functions

Using the Bayes Rule

$$P(\omega_i|\mathbf{x}) = \frac{p(\mathbf{x}|\omega_i)P(\omega_i)}{p(\mathbf{x})} \quad i = 1, 2 \quad (1)$$

$p(\mathbf{x})$ – density function for \mathbf{x}

Bayes Classification Rule (1)

Introduction

Bayes Decision
Theory

Discriminant
Functions and
Decision
Surfaces

Bayesian
Classification
for Normal
Distributions

Estimation of
Unknown
Probability
Density
Functions

Higher a posteriori probability wins

If $P(\omega_1|\mathbf{x}) > P(\omega_2|\mathbf{x})$, \mathbf{x} is classified to ω_1

If $P(\omega_1|\mathbf{x}) < P(\omega_2|\mathbf{x})$, \mathbf{x} is classified to ω_2

Bayes Classification Rule (2)

Introduction

Bayes Decision
Theory

Discriminant
Functions and
Decision
Surfaces

Bayesian
Classification
for Normal
Distributions

Estimation of
Unknown
Probability
Density
Functions

Considering the Bayes Rule (Eq. 1)

If $\frac{p(\mathbf{x}|\omega_1)P(\omega_1)}{p(\mathbf{x})} > \frac{p(\mathbf{x}|\omega_2)P(\omega_2)}{p(\mathbf{x})}$, \mathbf{x} is classified to ω_1

If $\frac{p(\mathbf{x}|\omega_1)P(\omega_1)}{p(\mathbf{x})} < \frac{p(\mathbf{x}|\omega_2)P(\omega_2)}{p(\mathbf{x})}$, \mathbf{x} is classified to ω_2

Bayes Classification Rule (3)

Introduction

Bayes Decision
Theory

Discriminant
Functions and
Decision
Surfaces

Bayesian
Classification
for Normal
Distributions

Estimation of
Unknown
Probability
Density
Functions

$p(\mathbf{x})$ can be disregarded, because it is the same for all classes

If $p(\mathbf{x}|\omega_1)P(\omega_1) > p(\mathbf{x}|\omega_2)P(\omega_2)$, \mathbf{x} is classified to ω_1

If $p(\mathbf{x}|\omega_1)P(\omega_1) < p(\mathbf{x}|\omega_2)P(\omega_2)$, \mathbf{x} is classified to ω_2

Bayes Classification Rule (4)

Introduction

Bayes Decision
Theory

Discriminant
Functions and
Decision
Surfaces

Bayesian
Classification
for Normal
Distributions

Estimation of
Unknown
Probability
Density
Functions

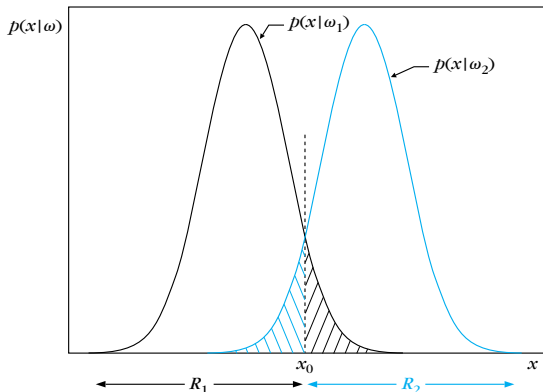
If the a priori probabilities are equal: $P(\omega_1) = P(\omega_2)$

If $p(\mathbf{x}|\omega_1) > p(\mathbf{x}|\omega_2)$, \mathbf{x} is classified to ω_1

If $p(\mathbf{x}|\omega_1) < p(\mathbf{x}|\omega_2)$, \mathbf{x} is classified to ω_2

We are done, since the likelihood density functions $p(\mathbf{x}|\omega_1)$ and $p(\mathbf{x}|\omega_2)$ are assumed to have been trained from examples!

Classification Error Probability



Error Probability:
$$P_e = \frac{1}{2} \int_{-\infty}^{x_0} p(x|\omega_2) dx + \frac{1}{2} \int_{x_0}^{\infty} p(x|\omega_1) dx$$

Introduction

Bayes Decision Theory

Discriminant Functions and Decision Surfaces

Bayesian Classification for Normal Distributions

Estimation of Unknown Probability Density Functions

Classification Error Probability in General

Introduction

Bayes Decision
Theory

Discriminant
Functions and
Decision
Surfaces

Bayesian
Classification
for Normal
Distributions

Estimation of
Unknown
Probability
Density
Functions

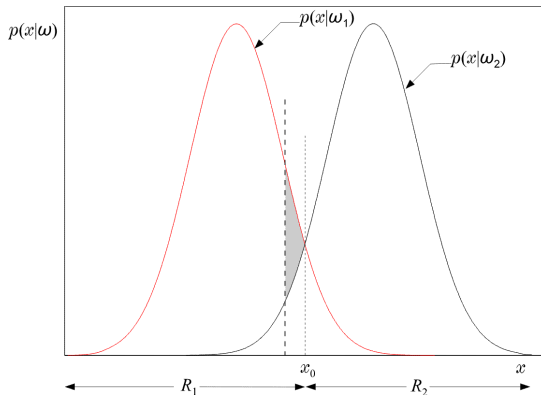
- A priori probabilities are not equal: $P(\omega_1) \neq P(\omega_2)$
- Feature vectors have more than one dimension: $l > 1$

$$\mathbf{x} = [x_1, x_2, \dots, x_l]^T$$

- General form:

$$P_e = P(\omega_1) \int_{R_2} p(\mathbf{x}|\omega_1) d\mathbf{x} + P(\omega_2) \int_{R_1} p(\mathbf{x}|\omega_2) d\mathbf{x}$$

Classification Error Probability



Bayesian Classifier is OPTIMAL with respect to minimising the classification error probability!

Introduction

Bayes Decision Theory

Discriminant Functions and Decision Surfaces

Bayesian Classification for Normal Distributions

Estimation of Unknown Probability Density Functions

Minimising Average Risk for Two Classes

Introduction

Bayes Decision Theory

Discriminant Functions and Decision Surfaces

Bayesian Classification for Normal Distributions

Estimation of Unknown Probability Density Functions

- Classification error probability assigns the same importance to all errors, which is wrong for many applications (e. g., $\omega_1 \rightarrow$ “malignant tumour”, $\omega_2 \rightarrow$ “benign tumour”).
- In such cases a penalty term is assigned to weight each error.
- A modified version of the error probability has to be minimised:

$$r = \lambda_{12}P(\omega_1) \int_{R_2} p(\mathbf{x}|\omega_1)d\mathbf{x} + \lambda_{21}P(\omega_2) \int_{R_1} p(\mathbf{x}|\omega_2)d\mathbf{x}$$

- For the tumour example λ_{12} is much greater than λ_{21} .

Modified Bayes Classification Rule

Introduction

Bayes Decision Theory

Discriminant Functions and Decision Surfaces

Bayesian Classification for Normal Distributions

Estimation of Unknown Probability Density Functions

If the a priori probabilities are equal: $P(\omega_1) = P(\omega_2)$

If $p(\mathbf{x}|\omega_2) > p(\mathbf{x}|\omega_1) \frac{\lambda_{12}}{\lambda_{21}}$, \mathbf{x} is classified to ω_2

If $p(\mathbf{x}|\omega_2) < p(\mathbf{x}|\omega_1) \frac{\lambda_{12}}{\lambda_{21}}$, \mathbf{x} is classified to ω_1

Overview

Introduction

Bayes Decision
Theory

**Discriminant
Functions and
Decision
Surfaces**

Bayesian
Classification
for Normal
Distributions

Estimation of
Unknown
Probability
Density
Functions

- 1 Introduction
- 2 Bayes Decision Theory
- 3 Discriminant Functions and Decision Surfaces**
- 4 Bayesian Classification for Normal Distributions
- 5 Estimation of Unknown Probability Density Functions

Discriminant Functions

- Sometimes it is more convenient to work with functions of probabilities instead of probabilities

$$g_i(\mathbf{x}) \equiv f(P(\omega_i|\mathbf{x}))$$

- $f(\cdot)$ is a monotonically increasing function
- $g_i(\mathbf{x})$ is known as discriminant function
- The decision test is now stated as

$$\text{classify } \mathbf{x} \text{ into } \omega_i \quad \text{if} \quad g_i(\mathbf{x}) > g_j(\mathbf{x}) \quad \forall j \neq i$$

- The decision surfaces, separating contiguous regions, are described by

$$g_{ij}(\mathbf{x}) \equiv g_i(\mathbf{x}) - g_j(\mathbf{x}) = 0, \quad i, j = 1, 2, \dots, M \quad i \neq j$$

Overview

Introduction

Bayes Decision
Theory

Discriminant
Functions and
Decision
Surfaces

**Bayesian
Classification
for Normal
Distributions**

Estimation of
Unknown
Probability
Density
Functions

- 1 Introduction
- 2 Bayes Decision Theory
- 3 Discriminant Functions and Decision Surfaces
- 4 Bayesian Classification for Normal Distributions**
- 5 Estimation of Unknown Probability Density Functions

Assumption

Introduction

Bayes Decision
Theory

Discriminant
Functions and
Decision
Surfaces

Bayesian
Classification
for Normal
Distributions

Estimation of
Unknown
Probability
Density
Functions

- The likelihood density functions describing the data in each of the classes, are multivariate Gaussian (normal) distributions

$$p(\mathbf{x}|\omega_i) = (2\pi)^{-\frac{l}{2}} |\boldsymbol{\Sigma}_i|^{-\frac{1}{2}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1}(\mathbf{x}-\boldsymbol{\mu}_i)}$$

- This “monster” will be denoted by

$$p(\mathbf{x}|\omega_i) = \mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \quad i = 1, 2, \dots, M$$

Discriminant Function $f(\cdot) = \ln(\cdot)$

- Due to the exponential form of the involved densities, the following discriminant function is applied:

$$g_i(\mathbf{x}) = \ln(p(\mathbf{x}|\omega_i)P(\omega_i)) = \ln p(\mathbf{x}|\omega_i) + \ln P(\omega_i)$$

\Updownarrow considering the “monster”

$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1}(\mathbf{x} - \boldsymbol{\mu}_i) + \ln P(\omega_i) + c_i \quad (2)$$

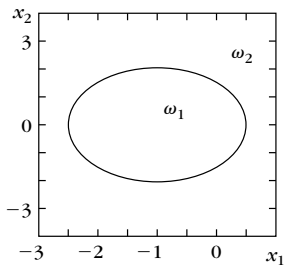
- Where: $c_i = -\frac{l}{2} \ln 2\pi - \frac{1}{2} \ln |\boldsymbol{\Sigma}_i|$

Quadrics as Decision Curves

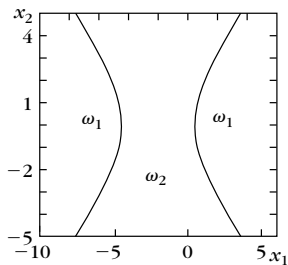
Assuming $l = 2$ and $\sigma_{1,2} = \sigma_{2,1} = 0$, the decision curves

$$g_i(\mathbf{x}) - g_j(\mathbf{x}) = 0$$

are quadrics (i. e., ellipsoids, parabolas, hyperbolas, pairs of lines)



(a)



(b)

Decision Hyperplanes

- The only quadric contribution in Equation (2) is $\mathbf{x}^T \Sigma_i^{-1} \mathbf{x}$
- Assuming that the covariance matrix is the same for all classes $\Sigma_i = \Sigma$ the quadric term will be the same for all discriminant functions
- Thus, the quadric term can be disregarded by decision surface equations. The same is true for the constant c_i
- The simplified version of the discriminant function is just a linear function

$$g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i0}$$

where

$$\mathbf{w}_i = \Sigma^{-1} \mu_i \quad \text{and} \quad w_{i0} = \ln P(\omega_i) - \frac{1}{2} \mu_i^T \Sigma^{-1} \mu_i$$

Introduction

Bayes Decision Theory

Discriminant Functions and Decision Surfaces

Bayesian Classification for Normal Distributions

Estimation of Unknown Probability Density Functions

Minimum Distance Classifiers

- Assuming equiprobable classes with the same covariance matrix and neglecting the constants Eq. 2 is simplified to:

$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_i)$$

- If $\boldsymbol{\Sigma} = \sigma^2 \mathbf{I}$ (diagonal matrix) the maximum $g_i(\mathbf{x})$ implies the minimum Euclidean distance $d_\epsilon = \|\mathbf{x} - \boldsymbol{\mu}_i\|$
- Thus, a feature vector \mathbf{x} is assigned to a class \hat{i} according to its Euclidean distance to the respective mean points $\boldsymbol{\mu}_i$

$$\hat{i} = \underset{i}{\operatorname{argmax}}(g_i(\mathbf{x})) = \underset{i}{\operatorname{argmin}}(\|\mathbf{x} - \boldsymbol{\mu}_i\|)$$

Remarks

Introduction

Bayes Decision
Theory

Discriminant
Functions and
Decision
Surfaces

Bayesian
Classification
for Normal
Distributions

Estimation of
Unknown
Probability
Density
Functions

- In practice, it is quite common to assume the Gaussian distribution of the data. In this case, the Bayesian classifier is either linear or quadratic in nature. These approaches are known as linear discriminant analysis (LDA) or quadratic discriminant analysis (QDA).
- A major problem associated with LDA and QDA is the large number of parameters to be estimated. Thus, I parameters in each mean vector and approximately $\frac{I^2}{2}$ in each covariance matrix. Moreover, a large number of training points N is needed.
- LDA and QDA perform very good for many different applications. However, in many cases the assumed normal distribution is not the right method to statistically model the data.

Overview

Introduction

Bayes Decision
Theory

Discriminant
Functions and
Decision
Surfaces

Bayesian
Classification
for Normal
Distributions

**Estimation of
Unknown
Probability
Density
Functions**

- 1 Introduction
- 2 Bayes Decision Theory
- 3 Discriminant Functions and Decision Surfaces
- 4 Bayesian Classification for Normal Distributions
- 5 Estimation of Unknown Probability Density Functions**

Problem Statement

Introduction

Bayes Decision
Theory

Discriminant
Functions and
Decision
Surfaces

Bayesian
Classification
for Normal
Distributions

Estimation of
Unknown
Probability
Density
Functions

- So far, we have assumed that the likelihood density functions $p(\mathbf{x}|\omega_i)$ for $i = 1, 2, \dots, M$ are known.
- This is not the most common case. In many problems, the likelihood density functions have to be estimated from the available training data.
- Here, two estimation methods will be considered, namely
 - Maximum Likelihood Parameter Estimation
 - Maximum a Posteriori Probability Estimation

Maximum Likelihood Parameter Estimation (1)

Introduction

Bayes Decision Theory

Discriminant Functions and Decision Surfaces

Bayesian Classification for Normal Distributions

Estimation of Unknown Probability Density Functions

- Let us consider an M -class problem with feature vectors distributed according to $p(\mathbf{x}|\omega_i)$, $i = 1, 2, \dots, M$.
- The likelihood functions are assumed to be given in a parametric form. The statistical parameters for the classes ω_i form vectors θ_i which are unknown

$$p(\mathbf{x}|\omega_i) = p(\mathbf{x}|\omega_i; \theta_i)$$

- Goal: to estimate the unknown parameters using a set of known feature vectors in each class.
- Since the estimation process is the same for all classes, the index i will be skipped for further investigations.

Maximum Likelihood Parameter Estimation (2)

- Let $X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ be a set of feature vectors describing training samples of a particular class.
- Assuming statistical independence between the different feature vectors, we can form the joint density function

$$p(X; \theta) = p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N; \theta) = \prod_{k=1}^N p(\mathbf{x}_k; \theta)$$

- The ML method estimates θ so that the likelihood function takes its maximum value

$$\hat{\theta}_{\text{ML}} = \underset{\theta}{\operatorname{argmax}} \prod_{k=1}^N p(\mathbf{x}_k; \theta)$$

Maximum Likelihood Parameter Estimation (3)

- To find a maximum, the gradient has to be zero

$$\frac{\partial \prod_{k=1}^N p(\mathbf{x}_k; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = 0$$

- Due to the monotonicity of the logarithmic function, we can use also the log-likelihood function

$$L(\boldsymbol{\theta}) = \ln \prod_{k=1}^N p(\mathbf{x}_k; \boldsymbol{\theta})$$

- Looking for the maximum here, we have

$$\frac{\partial L(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \sum_{k=1}^N \frac{\partial \ln p(\mathbf{x}_k; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \sum_{k=1}^N \frac{1}{p(\mathbf{x}_k; \boldsymbol{\theta})} \frac{\partial p(\mathbf{x}_k; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = 0$$

Maximum a Posteriori Probability Estimation

- Set of feature vectors $X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$
- θ is an unknown random vector
- The starting point is the following density function

$$p(\theta|X) = \frac{p(\theta)p(X|\theta)}{p(X)}$$

- The MAP estimate $\hat{\theta}_{\text{MAP}}$ is defined as a point where $p(\theta|X)$ becomes maximum

$$\hat{\theta}_{\text{MAP}} : \frac{\partial}{\partial \theta} p(\theta|X) = 0 \quad \text{or} \quad \frac{\partial}{\partial \theta} (p(\theta)p(X|\theta)) = 0$$