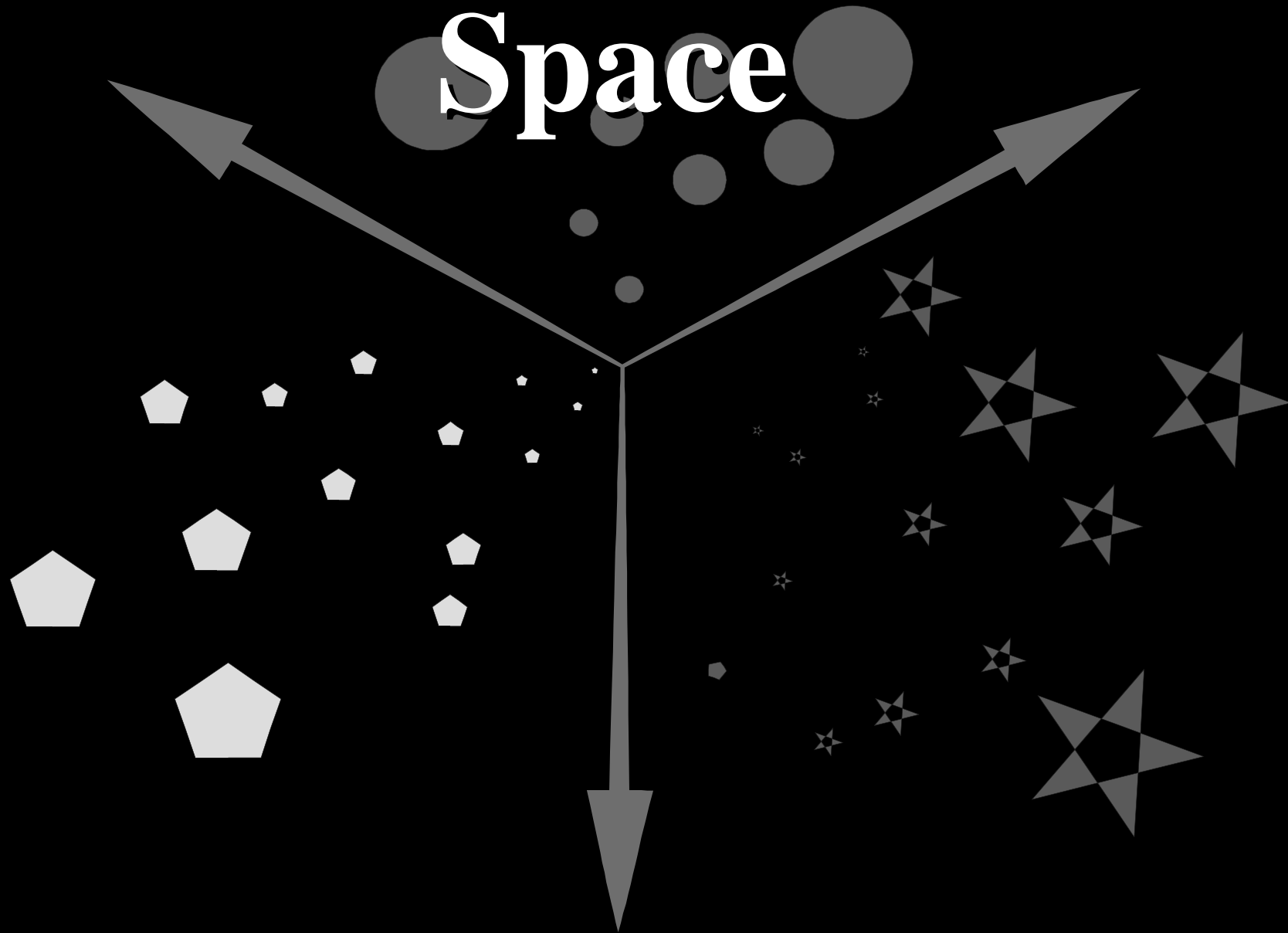


Classification in Pattern Space



Classification in Pattern

Topics

1. Introduction

1.1. Why classification ?

1.2. Different methods for

2. Formal view

2.1. Statistical model of speech creation

2.2. Classification functions

2.3. The optimal decision

3. Learning methods

3.1. Maximum Likelihood

3.2. The EM-

Algorithm

Examples

Experiments

Classification in Pattern

1. Introduction

1.1 Why classification ?

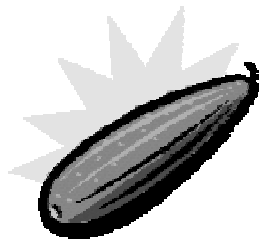
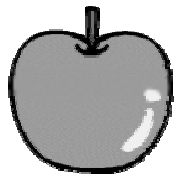
Classification in Pattern

1. Introduction

1.1 Why classification ?

- What is classification ?

Fruit or vegetable ?

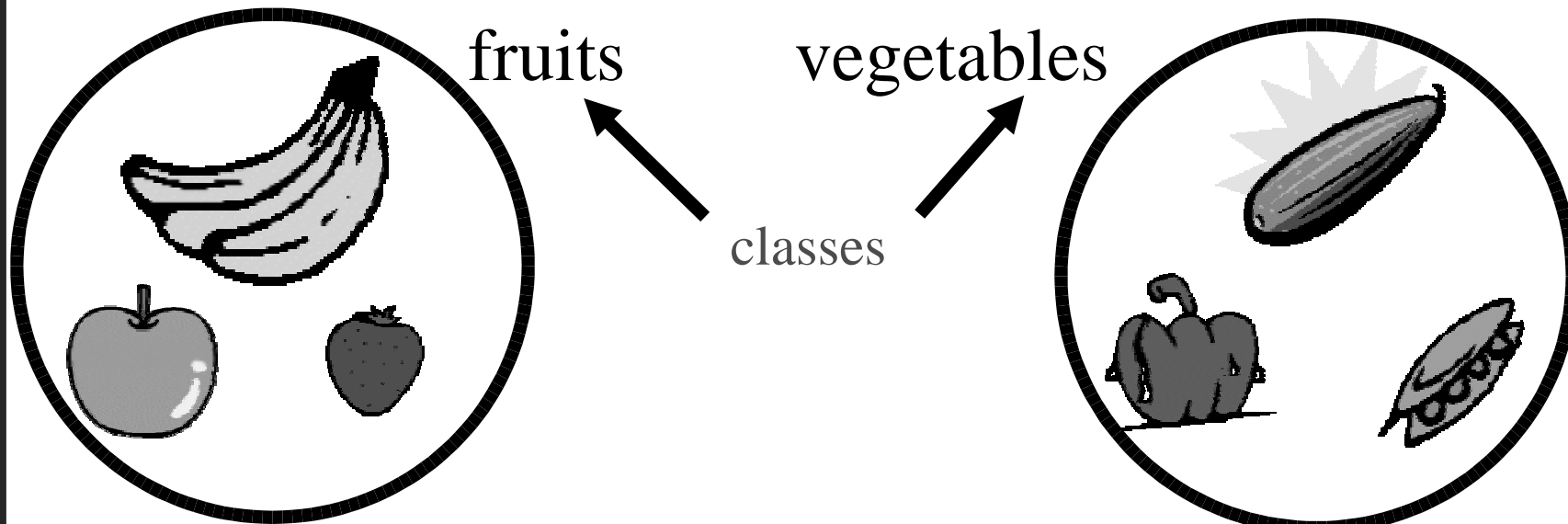


Classification in Pattern

1. Introduction

1.1 Why classification ?

- What is classification ?

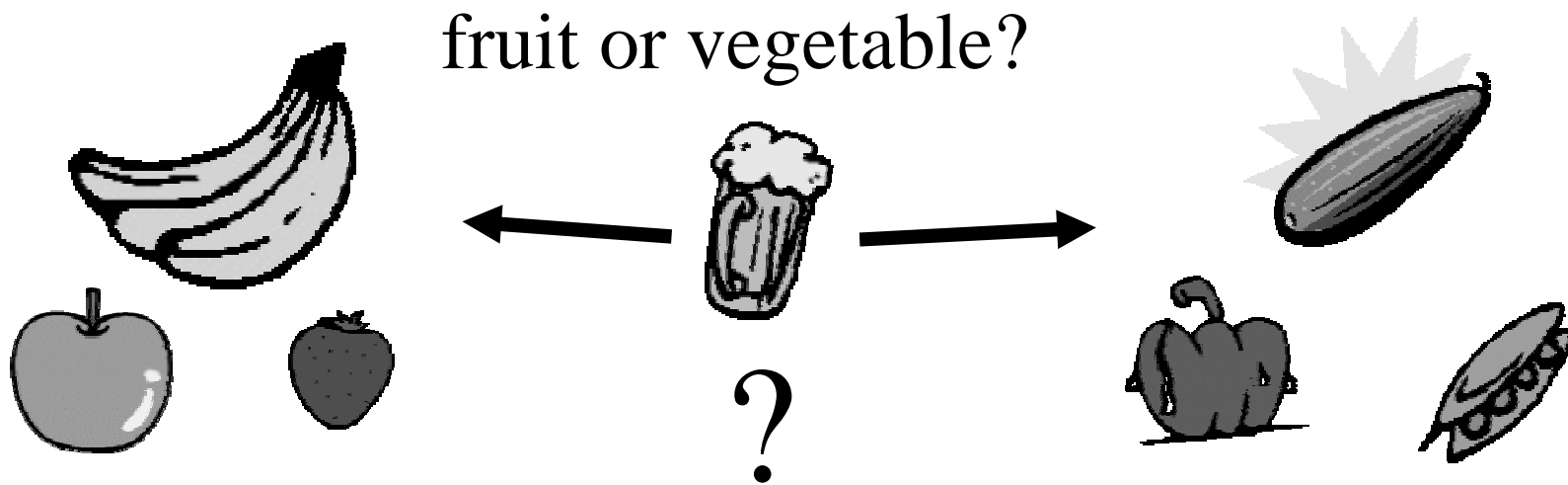


Classification in Pattern

1. Introduction

1.1 Why classification ?

- What is classification ?

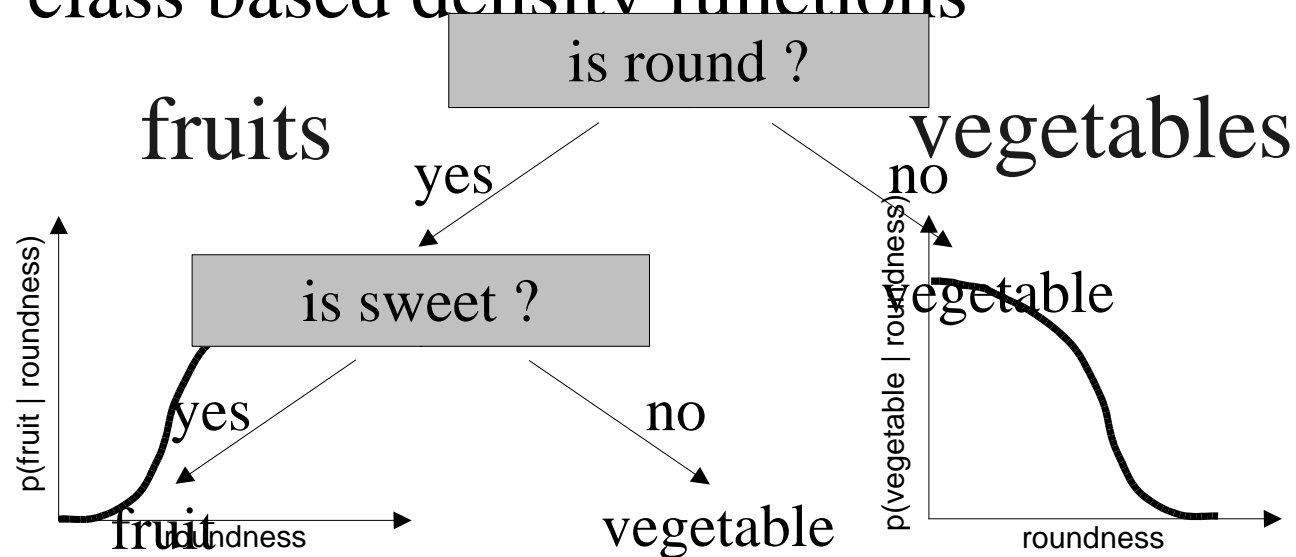


Classification in Pattern

1. Introduction

1.2 Different methods for classification

- decision trees
- artificial neural networks
- class based density functions



Classification in Pattern

2. Formal view

Situation for a typical classification problem:

- you have some data / objects you want to classify
- you have some preclassified data (optional)

Then you need to:

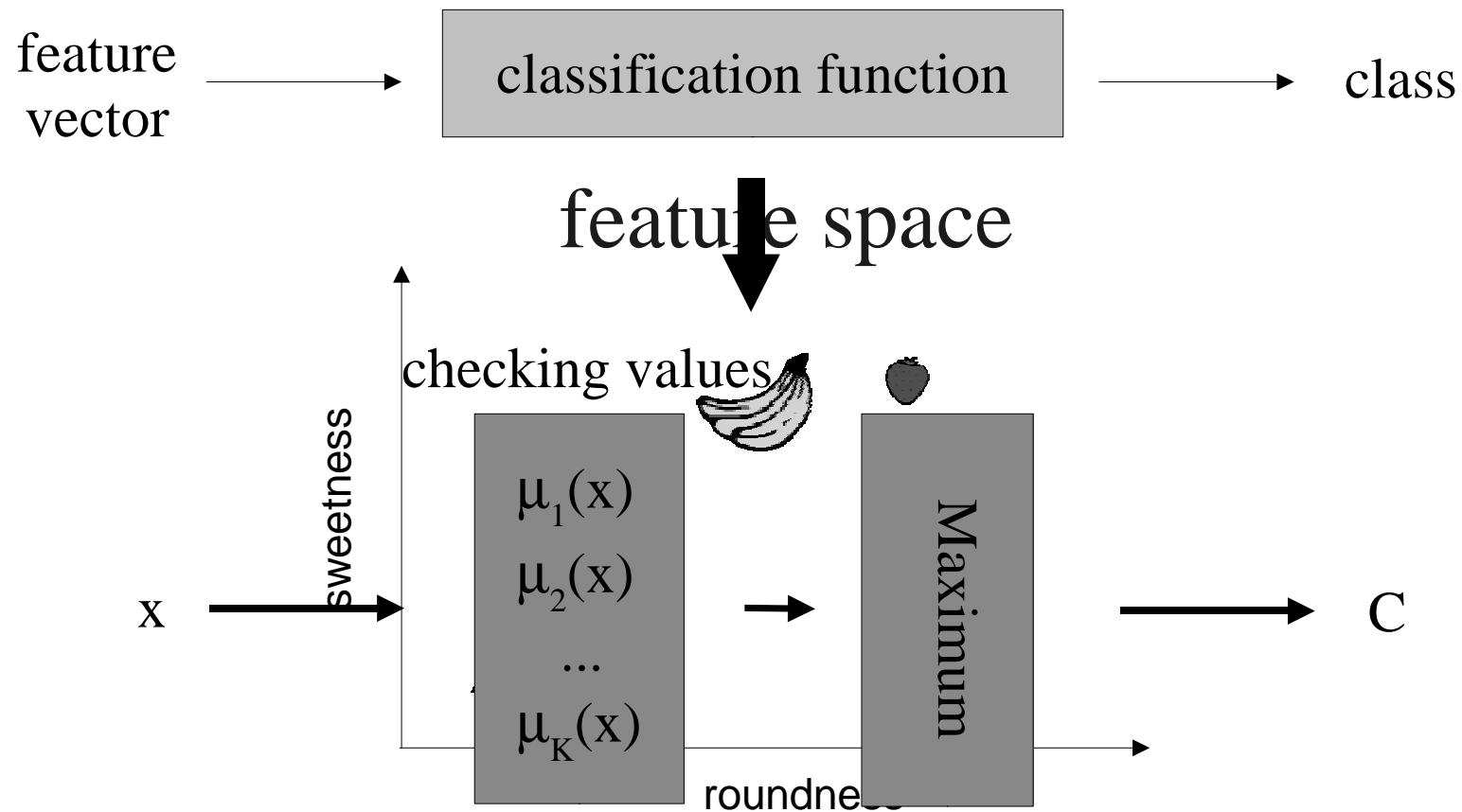
- define some features (eg. roundness, sweetness, ...)
- do a feature-extraction for your data

Result

	apple	strawberry	banana	cucumber	peas	paprica
roundness:	0.9	0.7	0.5	0.3	0.25	0.7
sweetness:	0.6	0.8	0.7	0.25	0.3	0.2

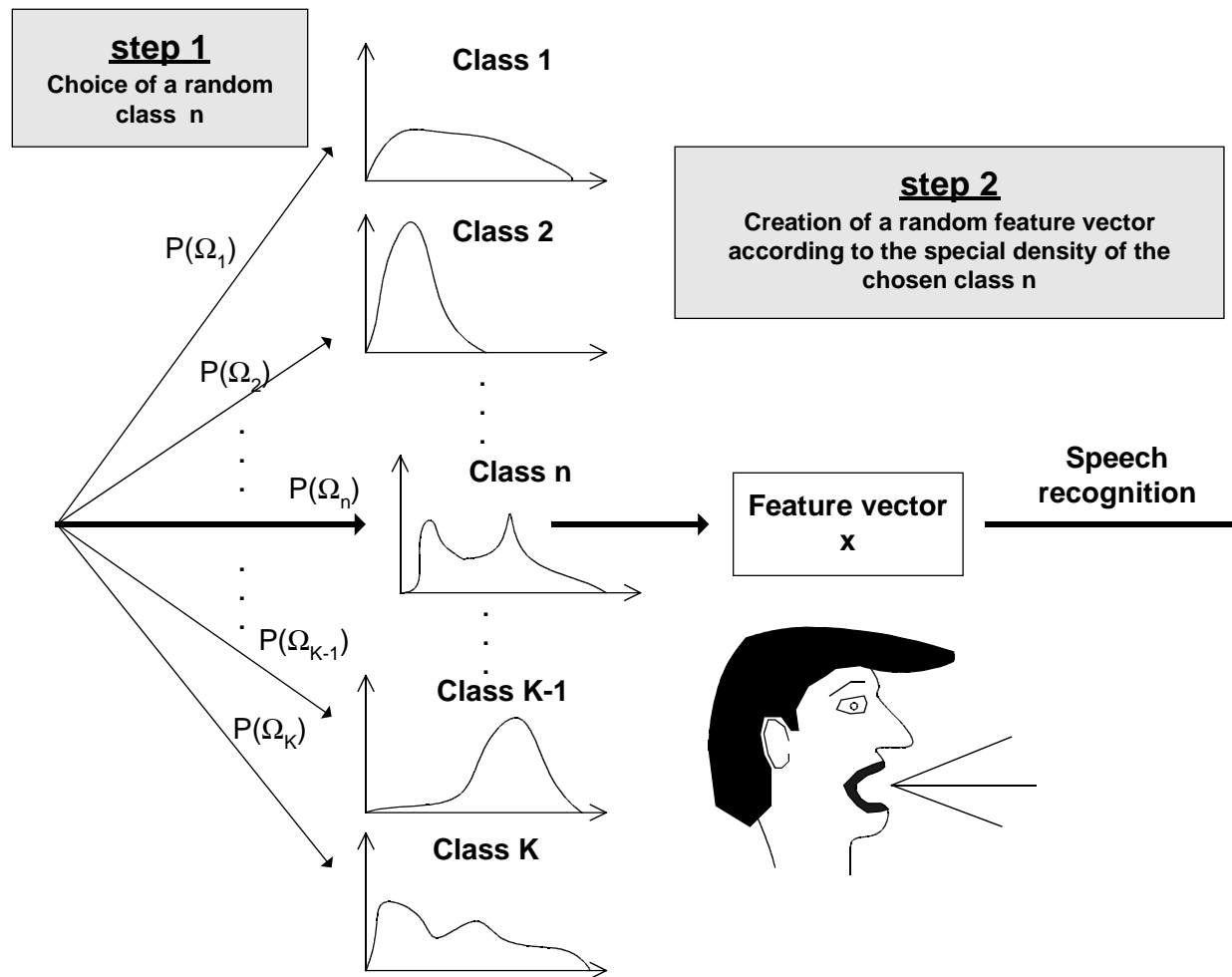
Classification in Pattern

2. Formal view



Classification in Pattern

2.1. Statistical model of speech creation



Classification in Pattern

2.1. Statistical model of speech creation

$$P(\Omega_i), i=1 \dots K \quad \text{and} \quad \sum_{i=1}^K P(\Omega_i) = 1$$

$$P(\mathbf{x} | \Omega_i) \quad \text{and} \quad \int_{\mathbb{R}^D} P(\mathbf{x} | \Omega_i) d\mathbf{x} = 1$$

but we get a temporal stream of feature vectors:

$$\begin{bmatrix} x_1 \\ \vdots \end{bmatrix}^{t=1}, \begin{bmatrix} x_2 \\ \vdots \end{bmatrix}^{t=2}, \begin{bmatrix} x_3 \\ \vdots \end{bmatrix}^{t=3}, \dots \quad P(\mathbf{x}) = \int_{-\infty}^{\infty} P(\mathbf{x}, \Omega) d\Omega$$

How to approximate $P(\mathbf{x} | \Omega)$ and $P(\Omega)$?
→ learning problem (later in this lecture)

Classification in Pattern

2.2. Classification functions

Intuitive: Maximum Likelihood (ML)

→ take the class which has the highest probability to produce \mathbf{x}

$$\Omega_i \text{ with } \max_i P(\mathbf{x} | \Omega_i)$$

less Intuitive: Maximum a-posteriori (MAP)

→ take the class that is likeliest to produce \mathbf{x}

$$\Omega_i \text{ with } \max_i P(\Omega_i | \mathbf{x})$$

$$P(\Omega_i | \mathbf{x}) = \frac{P(\mathbf{x} | \Omega_i) \cdot P(\Omega_i)}{P(\mathbf{x})}$$

↑
constant

$$P(\mathbf{x} | \Omega_1) = 0.6$$

$$P(\Omega_1) = 0.01$$

$$P(\mathbf{x} | \Omega_2) = 0.3$$

$$P(\Omega_2) = 0.99$$

Classification in Pattern

2.2. Classification functions

Empirical Risk Minimization

$$\mathbf{r} = \{r_{ik}\}$$

→ **risk of classifying a vector of class i in class j**

Expected risk:

$$R(\mathbf{x}) = \sum_{i=1}^K P(\Omega_i) \cdot P(\mathbf{x} | \Omega_i) R(\mathbf{x} | \Omega_i)$$

$R(\mathbf{x} | \Omega_i)$ – expected risk for a vector \mathbf{x} of class i

$$R(\mathbf{x} | \Omega_j) = \sum_{i=1}^K r_{i,j} \cdot \delta(\Omega_j | \mathbf{x})$$

$\delta(\Omega_j | \mathbf{x})$ – probability for decision for class j given the vector \mathbf{x}

$$R(\mathbf{x}) = \sum_j \underbrace{\delta(\Omega_j | \mathbf{x})}_{\delta_j} \underbrace{\sum_{i=1}^K r_{i,j} P(\Omega_i) \cdot P(\mathbf{x} | \Omega_i)}_{\mu_j}$$

Classification in Pattern

2.3. The optimal decision rule

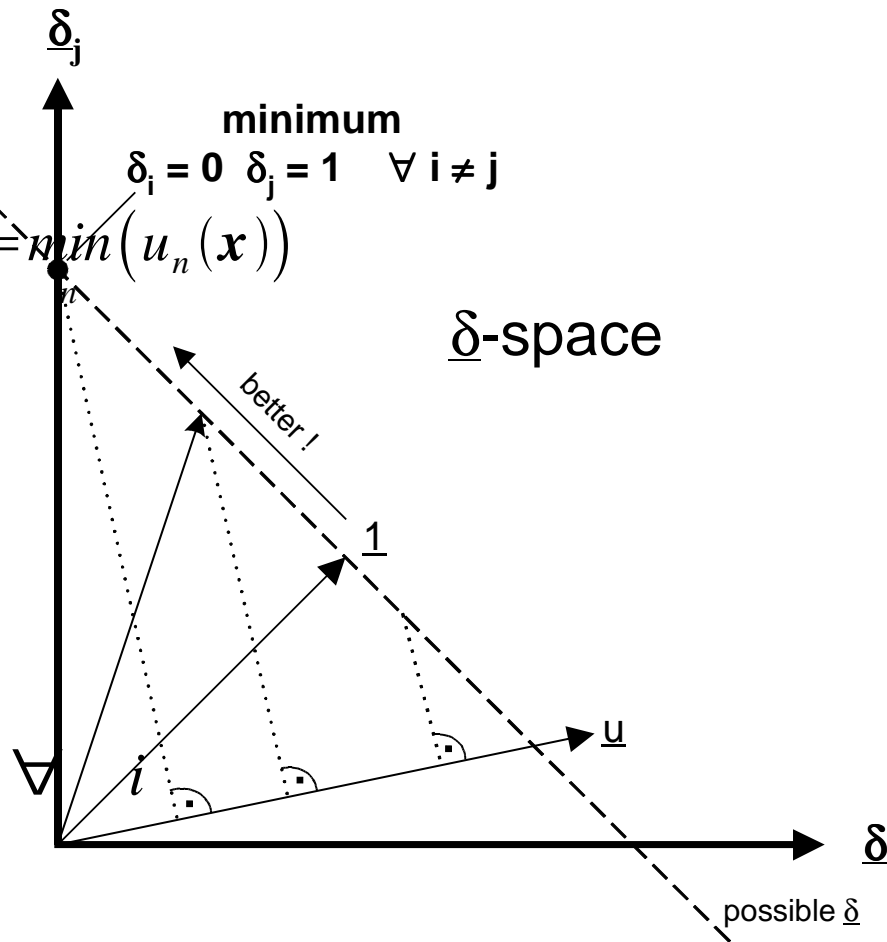
$$\min_{\delta_j} R(\mathbf{x}) = \sum_j^K \delta_j \cdot \mu_j$$

$$\delta^* = \begin{cases} \mathbf{1} & \text{falls } u_j(\mathbf{x}) = \min(u_n(\mathbf{x})) \\ \mathbf{0} & \text{sonst} \end{cases}$$

Hard decision !

Soft decision:

$$\sum_{i=1}^K \delta_i = 1, \quad \delta_i \geq 0$$



Classification in Pattern

3. Learning methods

Supervised Training

- we have got preclassified trainingsdata
- we know about the number of classes
- we want to find the position of the classes in the feature space

Unsupervised Training

- we have data, which is not preclassified
- we want to find clusters in the feature space

Classification in Pattern

3. Learning methods

3.1. Maximum likelihood

We have n training – examples:

$$\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]$$

We have got a model function for the distribution:

$P(\mathbf{x} | \phi)$ ϕ – parameter vector

$$\max_{\phi} P(\phi | \mathbf{X}) = \max_{\phi} \sum_{i=1}^N \log P(\mathbf{x}_i | \phi) + \log P(\phi) \quad \rightarrow \quad \hat{\phi} = \begin{bmatrix} \hat{\mu} \\ \hat{\sigma} \end{bmatrix}$$

$$P(\mathbf{X} | \phi) = \prod_{i=1}^N P(\mathbf{x}_i | \phi)$$

advantage : easy, just differentiate $L(\phi)$!

disadvantage : does not work most of the times !

Classification in Pattern

3. Learning methods

3.2. EM algorithm (expectation maximization)

We observe data $X = [x_1, x_2, \dots, x_N]$, which is incomplete !

eg. first step of the statistical model for speech creation

→ we don't know the classes Ω_i for each vector x_i

X → observed U → unknown data $P(x | \phi)$ → model

$Z = [X, U]$ → complete data

$$\max_{\phi} \log P(\mathbf{Z} | \phi) \quad \rightarrow P(\mathbf{Z} | \phi) = P(U | X, \phi) P(X | \phi)$$

$$\max_{\phi} E[\log P(\mathbf{Z} | \hat{\phi}) | X, \phi]$$

random variable !

