## **Fuzzy Classification**

## **Fuzzy Classification**

- Using Informal knowledge about problem domain for classification
- Example:
  - Adult salmon is oblong and light in color
  - Sea bass is stouter and dark
- Goal of fuzzy classification
  - Create fuzzy "category memberships" function
    - To convert objectively measurable parameters to "category memberships"
  - Which are then used for classification

#### "Categories"

- Does not refer to final classes
- Refer to overlapping ranges of feature values
- Example:
  - Lightness is divided into four categories
  - Dark, medium-dark, medium-light, light

Reflectivity

#### Category membership functions

Dark

**FIGURE 4.23.** "Category membership" functions, derived from the designer's prior knowledge, together with a conjunction rule lead to discriminants. In this figure, *x* might represent an objectively measurable value such as the reflectivity of a fish's skin. The designer believes there are four relevant ranges, which might be called dark, medium-dark, medium-light, and light. The designer believes there are four relevant ranges or "categories" for the reflectivity feature, which might be called dark, medium-dark, medium-light, and light. The categories for the feature, of course, are not the same as the true categories or classes for the patterns. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

- X

#### Light

### **Conjunction Rule**

- Merging several category functions corresponding to different features
- to yield a number to make the final decision
- Example: two category membership functions can be merged using

 $\mu_{\mathbf{x}}(\mathbf{x}).\mu_{\mathbf{y}}(\mathbf{y})$ specifies category Measured value of feature function for x

#### Discriminant function based on category membership functions



**FIGURE 4.24.** "Category membership" functions and a conjunction rule based on the designer's prior knowledge lead to discriminant functions. Here  $x_1$  and  $x_2$  are objectively measurable feature values. The designer believes that a particular class can be described as the conjunction of two "category memberships," here shown bold. Here the conjunction rule of Eq. 61 is used to give the discriminant function. The resulting discriminant function for the final category is indicated by the grayscale in the middle: the greater the discriminant, the darker. The designer constructs discriminant functions for other categories in a similar way (possibly also using disjunctions or other logical relations). During classification, the maximum discriminant function is chosen. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

#### Fuzzy category memberships

- Are they probabilities (or proportional to them?)
- Classical Probability
  - applies to more than *relative frequency*
  - Quantifies our notion of *uncertainty*
- Notion of *subjective* probability

# Do category membership functions represent probabilities?

- Half teaspoon of sugar placed in tea
  - Implies sweetness is 0.5
  - Not probability of sweetness is 50%
- But we can treat sweetness feature as having value 0.5

#### Limitations of fuzzy methods

- Cumbersome to use in
  - high dimensions (dozens or hundreds of features)
  - Complex problems
- Amount of information user can bring to bear is limited
  - no., positions and widths of category memberships
- Poorly suited to changing cost matrices
- Do not use training data
  - Neuro-fuzzy methods are tried
- Main contribution:
  - converting knowledge in linguistic form to discriminant functions

#### Reduced Coulomb Energy Networks

- Intermediate method to Parzen window and k-nearest neighbor estimation
  - Parzen window uses fixed window size
  - K-nn uses variable window size: increase window size until enough samples are enclosed
- Adjust window size until you encounter points of a different category
- Can be implemented as a neural network
- Gets name from electrostatics
  - Energy associated with charged particles

#### Decision regions created by RCE network

Gray: Class1 Pink: Class 2 Red: Ambiguous



**FIGURE 4.26.** During training of an RCE network, each pattern has a parameter equivalent to a radius in the *d*-dimensional space—that is adjusted to be as large as possible without enclosing any points from a different category (up to a maximum  $\lambda_m$ ). As new patterns are presented, each such radius is decreased so that no sphere encloses a pattern of a different category. In this way, each sphere can enclose only patterns having the same category label. In this figure, the regions corresponding to one category are pink, and the other category are gray. Ambiguous regions (those enclosed by spheres of both categories) are shown in dark red. The number of points is shown in each component figure. The figure at the bottom shows the final decision regions, colored by category. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

#### Training Reduced Coulomb Energy Networks

- Adjust each radius to be as large as possible (upto a maximum) without containing point from another category
- For each training sample x<sub>i</sub>, j=1,...,n set radius

```
Algorithm 4. (RCE Training)
      begin initialize j \leftarrow 0, n \leftarrow \# patterns, \epsilon \leftarrow small param, \lambda_m \leftarrow max radius
  2
                    do i \leftarrow i+1
  3
                          w_{ii} \leftarrow x_i (train weight)
 .4
                          \hat{\mathbf{x}} \leftarrow \arg\min D(\mathbf{x}, \mathbf{x}') (find nearest point not in \omega_i)
                                         \mathbf{x} \notin \omega_i
  5
                         \lambda_i \leftarrow \min[D(\hat{\mathbf{x}}, \mathbf{x}') - \epsilon, \lambda_m] (set radius)
                       if \mathbf{x} \in \omega_k then a_{ik} \leftarrow 1
  6
  7
                    until j = n
  8 end
```

#### Reduced Coulomb Energy Network



**FIGURE 4.25.** An RCE network is topologically equivalent to the PNN of Fig. 4.9. During training, normalized weights are adjusted to have the same values as the normalized pattern presented, just as in a PNN. In this way, distances can be calculated by an inner product. Pattern units in an RCE network also have a modifiable threshold corresponding to a "radius"  $\lambda$ . During training, each threshold is adjusted so that its radius is as large as possible without containing training patterns from a different category. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

#### Classification with Reduced Coulomb Energy Networks

Algorithm 5. (RCE Classification)  $1 \quad \underline{\text{begin initialize } j \leftarrow 0, k \leftarrow 0, \mathbf{x} \leftarrow \text{test pattern, } \mathcal{D}_t \leftarrow \{\} \\ 2 \qquad \underline{\text{do } j \leftarrow j+1} \\ 3 \qquad \underline{\text{if } D(\mathbf{x}, \mathbf{x}'_j) < \lambda_j \text{ then } \mathcal{D}_t \leftarrow \mathcal{D}_t \cup \mathbf{x}'_j \\ 4 \qquad \underline{\text{until } j = n} \\ 5 \qquad \underline{\text{if } \text{ label of all } \mathbf{x}'_j \in \mathcal{D}_t \text{ is the same } \underline{\text{then } \text{return } \text{ label of all } \mathbf{x}_k \in \mathcal{D}_t \\ 6 \qquad \underline{\text{else } \text{return } \text{"ambiguous" label}} \\ 7 \quad \underline{\text{end}} \end{cases}$ 

#### **Approximations by Series Expansions**

Memory requirements are severe in non-parametric methods

#### Approximate window function by series expansion

In certain circumstances the Parzen-window procedure can be modified to reduce these problems considerably. The basic idea is to approximate the window function by a finite series expansion that is acceptably accurate in the region of interest. If we are fortunate and can find two sets of functions  $\psi_j(\mathbf{x})$  and  $\chi_j(\mathbf{x})$  that allow the expansion

$$\varphi\left(\frac{\mathbf{x}-\mathbf{x}_i}{h_n}\right) = \sum_{j=1}^m a_j \psi_j(\mathbf{x}) \chi_j(\mathbf{x}_i), \tag{62}$$

then we can split the dependence upon  $\mathbf{x}$  and  $\mathbf{x}_i$  as

$$\sum_{i=1}^{n} \varphi\left(\frac{\mathbf{x} - \mathbf{x}_i}{h_n}\right) = \sum_{j=1}^{m} a_j \psi_j(\mathbf{x}) \sum_{i=1}^{n} \chi_j(\mathbf{x}_i).$$
(63)

Then from Eq. 11 we have

$$p_n(\mathbf{x}) = \sum_{j=1}^m b_j \psi_j(\mathbf{x}), \tag{64}$$

where

$$b_j = \frac{a_j}{nV_n} \sum_{i=1}^n \chi_j(\mathbf{x}_i).$$
(65)

#### Advantage of series expansion

In certain circumstances the Parzen-window procedure can be modified to reduce these problems considerably. The basic idea is to approximate the window function by a finite series expansion that is acceptably accurate in the region of interest. If we are fortunate and can find two sets of functions  $\psi_j(\mathbf{x})$  and  $\chi_j(\mathbf{x})$  that allow the expansion

$$\varphi\left(\frac{\mathbf{x}-\mathbf{x}_i}{h_n}\right) = \sum_{j=1}^m a_j \psi_j(\mathbf{x}) \chi_j(\mathbf{x}_i), \tag{62}$$

then we can split the dependence upon  $\mathbf{x}$  and  $\mathbf{x}_i$  as

$$\sum_{i=1}^{n} \varphi\left(\frac{\mathbf{x} - \mathbf{x}_i}{h_n}\right) = \sum_{j=1}^{m} a_j \psi_j(\mathbf{x}) \sum_{i=1}^{n} \chi_j(\mathbf{x}_i).$$
(63)

Then from Eq. 11 we have

$$p_n(\mathbf{x}) = \sum_{j=1}^m b_j \psi_j(\mathbf{x}), \tag{64}$$

where

$$b_j = \frac{a_j}{nV_n} \sum_{i=1}^n \chi_j(\mathbf{x}_i).$$
(65)

- Information in *n* samples is reduced to *m* coefficients  $b_i$
- Additional samples don't change no of coefficients

#### Taylor's Series Expansion of Window Function

Assuming a one-dimensional example using a Gaussian window function:

$$\sqrt{\pi} \varphi(u) = e^{-u^2}$$
$$\simeq \sum_{j=0}^{m-1} (-1)^j \frac{u^{2j}}{j!}.$$

#### Taylor's Series Expansion with quadratic terms

• If *m*=2 the window function can be approximated as

$$\sqrt{\pi}\varphi\left(\frac{x-x_i}{h}\right) \simeq 1 - \left(\frac{x-x_i}{h}\right)^2 = 1 + \frac{2}{h^2}x\,x_i - \frac{1}{h^2}\,x^2 - \frac{1}{h^2}\,x_i^2,$$

• And thus

$$\sqrt{\pi} p_n(x) = \frac{1}{nh} \sum_{i=1}^n \sqrt{\pi} \varphi\left(\frac{x-x_i}{h}\right) \simeq b_0 + b_1 x + b_2 x^2,$$

Where the coefficients are

$$b_0 = \frac{1}{h} - \frac{1}{h^3} \frac{1}{n} \sum_{i=1}^n x_i^2$$

$$b_1 = \frac{2}{h^3} \frac{1}{n} \sum_{i=1}^n x_i$$
$$b_2 = -\frac{1}{h^3}.$$

### Error in approximating $p_n(x)$

This simple expansion condenses the information in *n* samples into the values  $b_0$ ,  $b_1$ , and  $b_2$ . It is accurate if the largest value of  $|x - x_i|$  is not greater than *h*. Unfortunately, this restricts us to a very wide window that is not capable of much resolution. By taking more terms we can use a narrower window. If we let *r* be the largest value of  $|x - x_i|$  and use the fact that the error in the *m*-term expansion of  $\sqrt{\pi} \varphi((x - x_i)/h)$  is less than  $(r/h)^{2m}/m!$ , then using Stirling's approximation for *m*! we find that the error in approximating  $p_n(x)$  is less than

$$\frac{1}{\sqrt{\pi h}} \frac{\left(r/h\right)^{2m}}{m!} \simeq \frac{1}{\sqrt{\pi h}\sqrt{2\pi m}} \left[ \left(\frac{e}{m}\right) \left(\frac{r}{h}\right)^2 \right]^m.$$
(67)

Thus, the error becomes small only when  $m > e(r/h)^2$ . This implies the need for many terms if the window size *h* is small relative to the distance *r* from *x* to the most distant sample. Although this example is rudimentary, similar considerations arise in the multidimensional case even when more sophisticated expansions are used, and the procedure is most attractive when the window size is relatively large.