

Chapter 4 (Part 1): Non-Parametric Classification

- Introduction
- Density Estimation
- Parzen Windows



All materials used in this course were taken from the textbook "*Pattern Classification*" by Duda et al., John Wiley & Sons, 2001 with the permission of the authors and the publisher

■ Introduction

- All Parametric densities are unimodal (have a single local maximum), whereas many practical problems involve multi-modal densities
- Nonparametric procedures can be used with arbitrary distributions and without the assumption that the forms of the underlying densities are known
- There are two types of nonparametric methods:
 - Estimating $P(x | \omega_j)$
 - Bypass probability and go directly to a-posteriori probability estimation

■ Density Estimation

■ Basic idea:

Probability that a vector x will fall in region R is:

$$P = \int_{\mathcal{R}} p(x') dx' \quad (1)$$

P is a smoothed (or averaged) version of the density function $p(x)$ if we have a sample of size n ; therefore, the probability that k points fall in R is then:

$$P_k = \binom{n}{k} P^k (1-P)^{n-k} \quad (2)$$

and the expected value for k is:

$$E(k) = nP \quad (3)$$

ML estimation of $P = \theta$

$$\underset{\theta}{\text{Max}}(P_k | \theta) \quad \text{is reached for} \quad \hat{\theta} = \frac{k}{n} \cong P$$

Therefore, the ratio k/n is a good estimate for the probability P and hence for the density function p .

$p(x)$ is continuous and that the region R is so small that p does not vary significantly within it, we can write:

$$\int_{\mathcal{R}} p(x') dx' \cong p(x)V \quad (4)$$

where x is a point within R and V the volume enclosed by R .

Combining equation (1) , (3) and (4) yields: $p(x) \cong \frac{k/n}{V}$

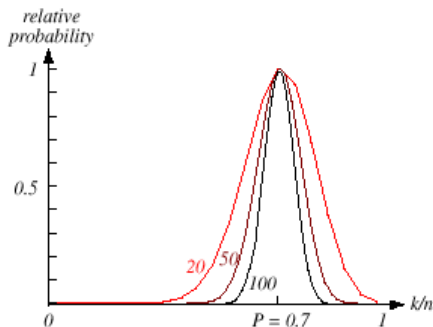


FIGURE 4.1. The relative probability an estimate given by Eq. 4 will yield a particular value for the probability density, here where the true probability was chosen to be 0.7. Each curve is labeled by the total number of patterns n sampled, and is scaled to give the same maximum (at the true probability). The form of each curve is binomial, as given by Eq. 2. For large n , such binomials peak strongly at the true probability. In the limit $n \rightarrow \infty$, the curve approaches a delta function, and we are guaranteed that our estimate will give the true probability. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

■ Density Estimation (cont'd)

■ Justification of equation (4)

$$\int_{\mathcal{R}} p(x') dx' \cong p(x)V \tag{4}$$

We assume that $p(x)$ is continuous and that region \mathcal{R} is so small that p does not vary significantly within \mathcal{R} . Since $p(x) = \text{constant}$, it is not a part of the sum.

$$\int_{\mathfrak{R}} p(\mathbf{x}') d\mathbf{x}' = p(\mathbf{x}') \int_{\mathfrak{R}} d\mathbf{x}' = p(\mathbf{x}') \int_{\mathfrak{R}} \mathbf{1}_{\mathfrak{R}}(\mathbf{x}) d\mathbf{x}' = p(\mathbf{x}') \mu(\mathfrak{R})$$

Where: $\mu(\mathfrak{R})$ is: a surface in \mathbb{R}^2
 a volume in \mathbb{R}^3
 a hypervolume in \mathbb{R}^n

Since $p(\mathbf{x}) \cong p(\mathbf{x}') = \text{constant}$, therefore in \mathbb{R}^3 :

$$\int_{\mathfrak{R}} p(\mathbf{x}') d\mathbf{x}' \cong p(\mathbf{x}) \cdot V$$

$$\text{and } p(\mathbf{x}) \cong \frac{k}{nV}$$

■ Condition for convergence

The fraction $k/(nV)$ is a space averaged value of $p(\mathbf{x})$.
 $p(\mathbf{x})$ is obtained only if V approaches zero.

$$\lim_{V \rightarrow 0, k=0} p(\mathbf{x}) = 0 \quad (\text{if } n = \text{fixed})$$

This is the case where no samples are included in \mathfrak{R} :
 it is an uninteresting case!

$$\lim_{V \rightarrow 0, k \neq 0} p(\mathbf{x}) = \infty$$

In this case, the estimate diverges: it is an
 uninteresting case!

The volume V needs to approach 0 anyway if we want to use this estimation

- Practically, V cannot be allowed to become small since the number of samples is always limited
- One will have to accept a certain amount of variance in the ratio k/n
- Theoretically, if an unlimited number of samples is available, we can circumvent this difficulty
To estimate the density of x , we form a sequence of regions R_1, R_2, \dots containing x : the first region contains one sample, the second two samples and so on.
Let V_n be the volume of R_n , k_n the number of samples falling in R_n and $p_n(x)$ be the n^{th} estimate for $p(x)$:

$$p_n(x) = (k_n/n)/V_n \quad (7)$$

Three necessary conditions should apply if we want $p_n(x)$ to converge to $p(x)$

- 1) $\lim_{n \rightarrow \infty} V_n = 0$
- 2) $\lim_{n \rightarrow \infty} k_n = \infty$
- 3) $\lim_{n \rightarrow \infty} k_n / n = 0$

There are two different ways of obtaining sequences of regions that satisfy these conditions:

- (a) Shrink an initial region where $V_n = 1/\sqrt{n}$ and show that

$$p_n(x) \xrightarrow{n \rightarrow \infty} p(x)$$

This is called “the Parzen-window estimation method”

- (b) Specify k_n as some function of n , such as $k_n = \sqrt{n}$; the volume V_n is grown until it encloses k_n neighbors of x . This is called “the k_n -nearest neighbor estimation method”

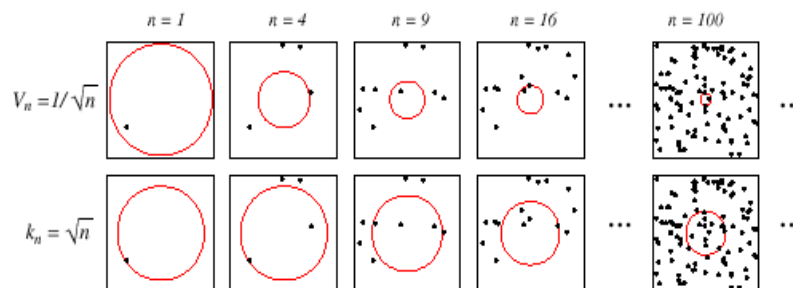


FIGURE 4.2. There are two leading methods for estimating the density at a point, here at the center of each square. The one shown in the top row is to start with a large volume centered on the test point and shrink it according to a function such as $V_n = I/\sqrt{n}$. The other method, shown in the bottom row, is to decrease the volume in a data-dependent way, for instance letting the volume enclose some number $k_n = \sqrt{n}$ of sample points. The sequences in both cases represent random variables that generally converge and allow the true density at the test point to be calculated. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

■ Parzen Windows

- Parzen-window approach to estimate densities assume that the region R_n is a d-dimensional hypercube

$$V_n = h_n^d \text{ (} h_n \text{ : length of the edge of } \mathfrak{R}_n \text{)}$$

Let $\varphi(u)$ be the following window function :

$$\varphi(u) = \begin{cases} 1 & |u_j| \leq \frac{1}{2} \quad j = 1, \dots, d \\ 0 & \text{otherwise} \end{cases}$$

- $\varphi((x-x_i)/h_n)$ is equal to unity if x_i falls within the hypercube of volume V_n centered at x and equal to zero otherwise.

- The number of samples in this hypercube is:

$$k_n = \sum_{i=1}^{i=n} \varphi\left(\frac{\mathbf{x} - \mathbf{x}_i}{\mathbf{h}_n}\right)$$

By substituting k_n in equation (7), we obtain the following estimate:

$$\mathbf{p}_n(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^{i=n} \frac{1}{V_n} \varphi\left(\frac{\mathbf{x} - \mathbf{x}_i}{\mathbf{h}_n}\right)$$

$\mathbf{P}_n(\mathbf{x})$ estimates $p(\mathbf{x})$ as an average of functions of \mathbf{x} and the samples (x_i) ($i = 1, \dots, n$). These functions φ can be general!

■ Illustration

- The behavior of the Parzen-window method

- Case where $p(\mathbf{x}) \rightarrow N(0,1)$

Let $\varphi(u) = (1/\sqrt{2\pi}) \exp(-u^2/2)$ and $h_n = h_1/\sqrt{n}$ ($n > 1$)

(h_1 : known parameter)

Thus:

$$\mathbf{p}_n(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^{i=n} \frac{1}{h_n} \varphi\left(\frac{\mathbf{x} - \mathbf{x}_i}{h_n}\right)$$

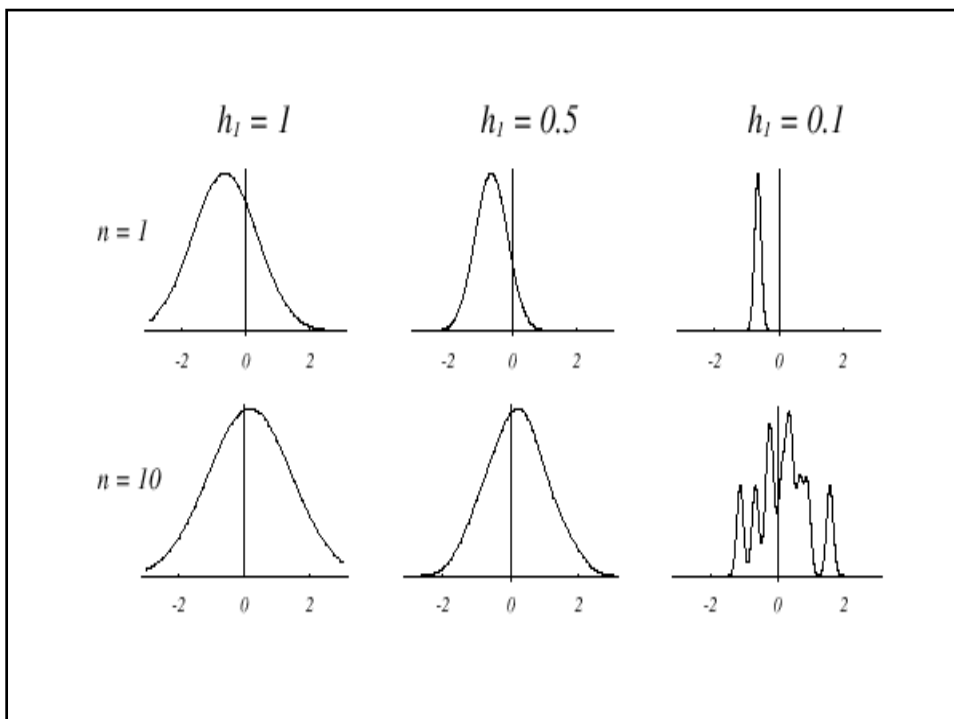
is an average of normal densities centered at the samples x_i .

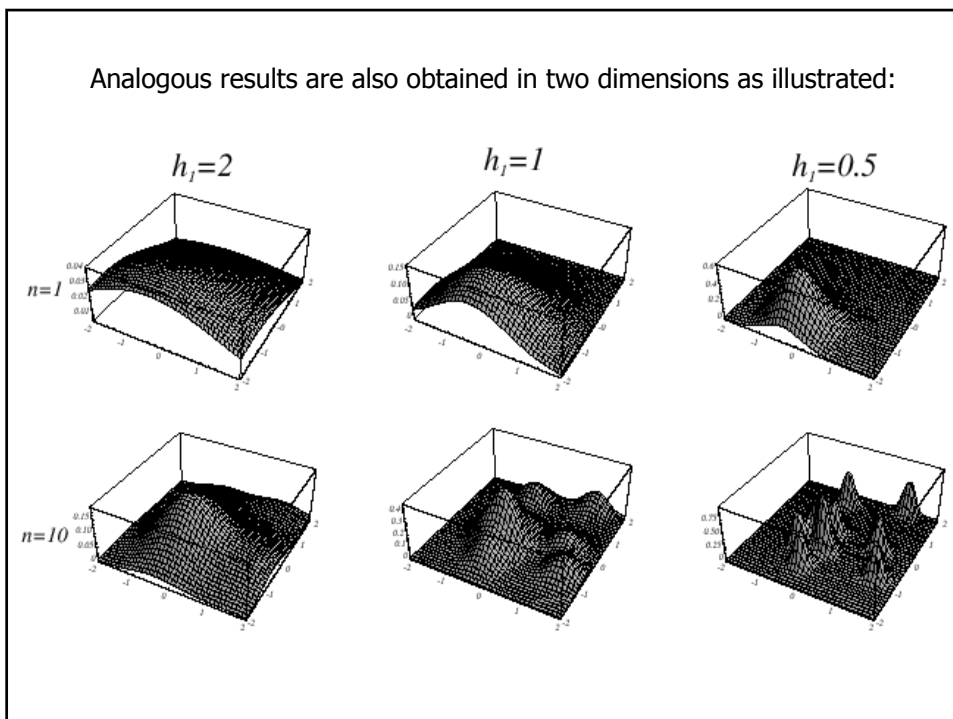
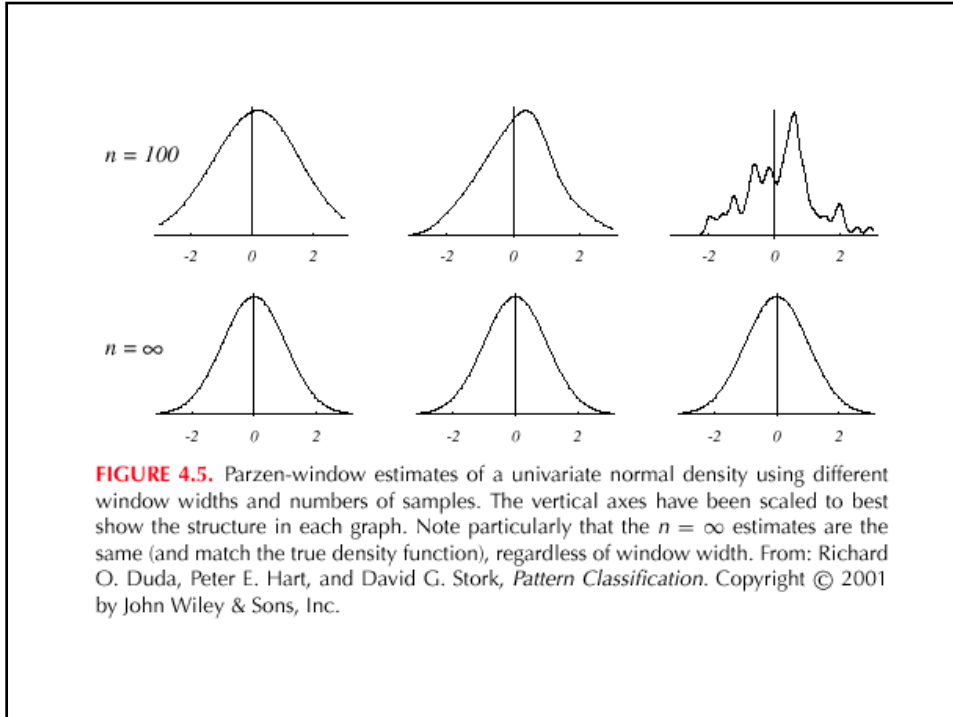
– **Numerical results:**

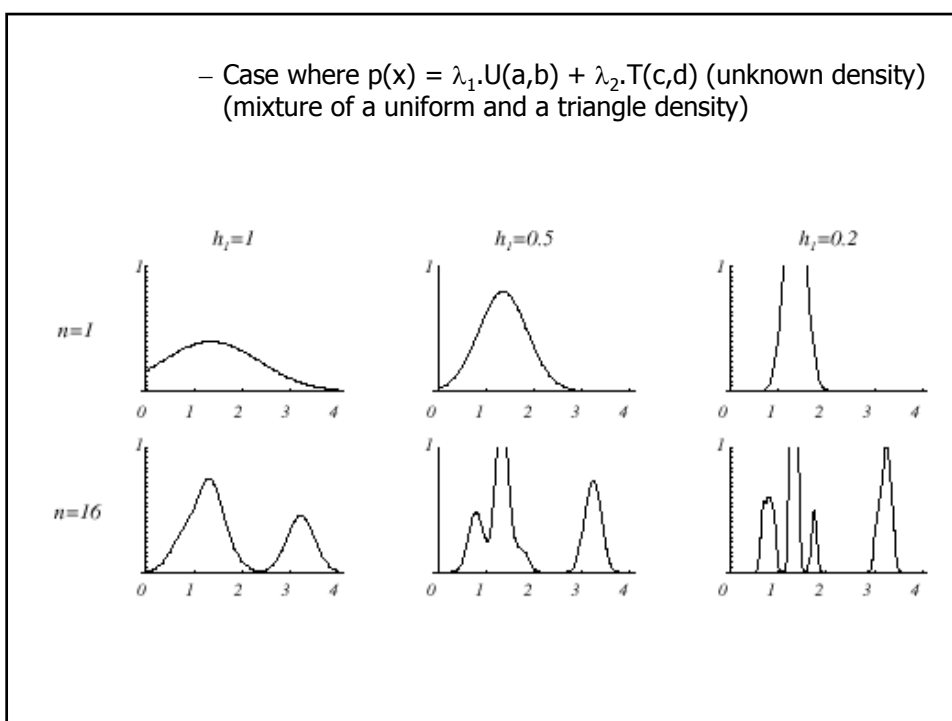
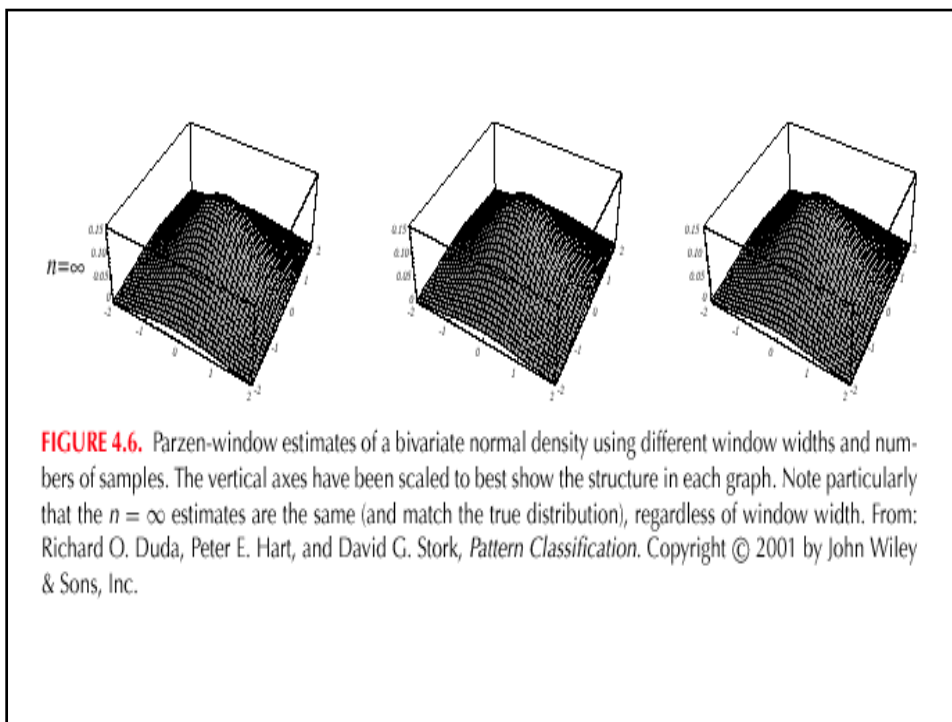
For $n = 1$ and $h_1=1$

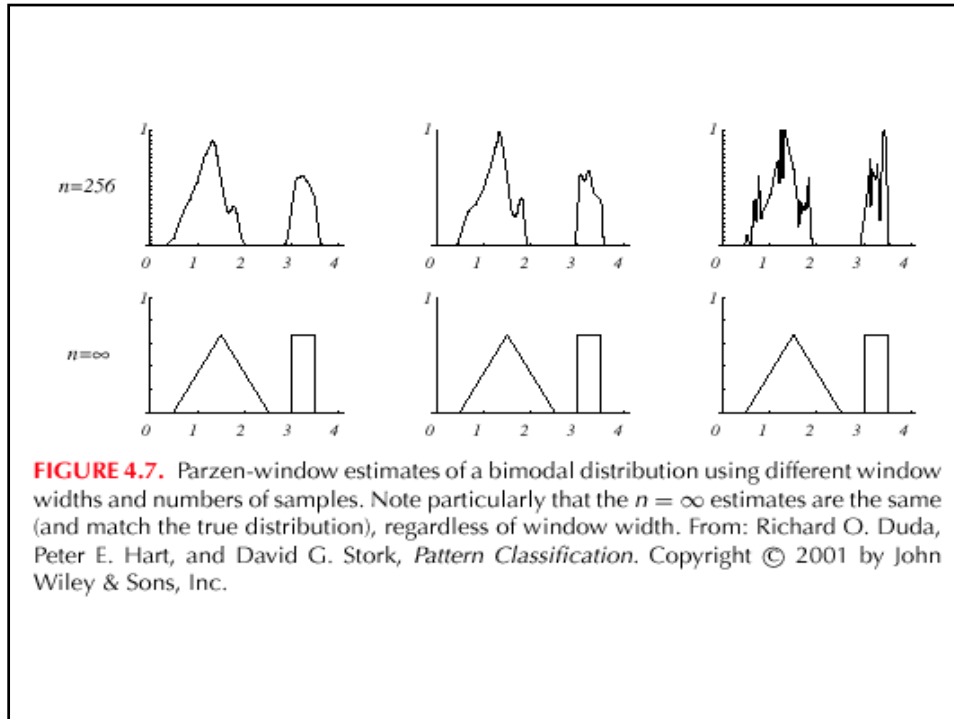
$$p_1(\mathbf{x}) = \varphi(\mathbf{x} - \mathbf{x}_1) = \frac{1}{\sqrt{2\pi}} e^{-1/2(\mathbf{x} - \mathbf{x}_1)^2} \rightarrow \mathbf{N}(\mathbf{x}_1, 1)$$

For $n = 10$ and $h = 0.1$, the contributions of the individual samples are clearly observable !









■ Classification example

In classifiers based on Parzen-window estimation:

- We estimate the densities for each category and classify a test point by the label corresponding to the maximum posterior
- The decision region for a Parzen-window classifier depends upon the choice of window function as illustrated in the following figure.

