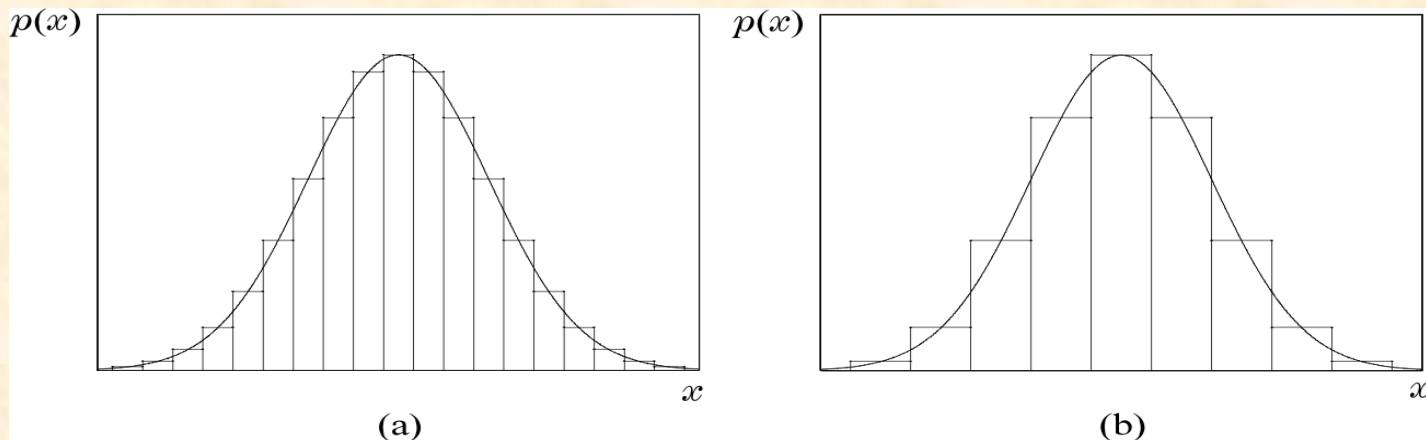


# ❖ ΑΝΑΓΝΩΡΙΣΗ ΠΡΟΤΥΠΩΝ ❖ (PATTERN RECOGNITION)

**Σέργιος Θεοδωρίδης**  
**Κωνσταντίνος Κουτρούμπας**

## ❖ Μη παραμετρική εκτίμηση



$$\triangleright P \approx \frac{k_N}{N} \begin{array}{l} \nearrow k_N \text{ στο } h \\ \searrow N \text{ συνολικά} \end{array}$$

$$\triangleright \hat{p}(x) \equiv \hat{p}(\hat{x}) = \frac{1}{h} \frac{k_N}{N}, \quad |x - \hat{x}| \leq \frac{h}{2}$$

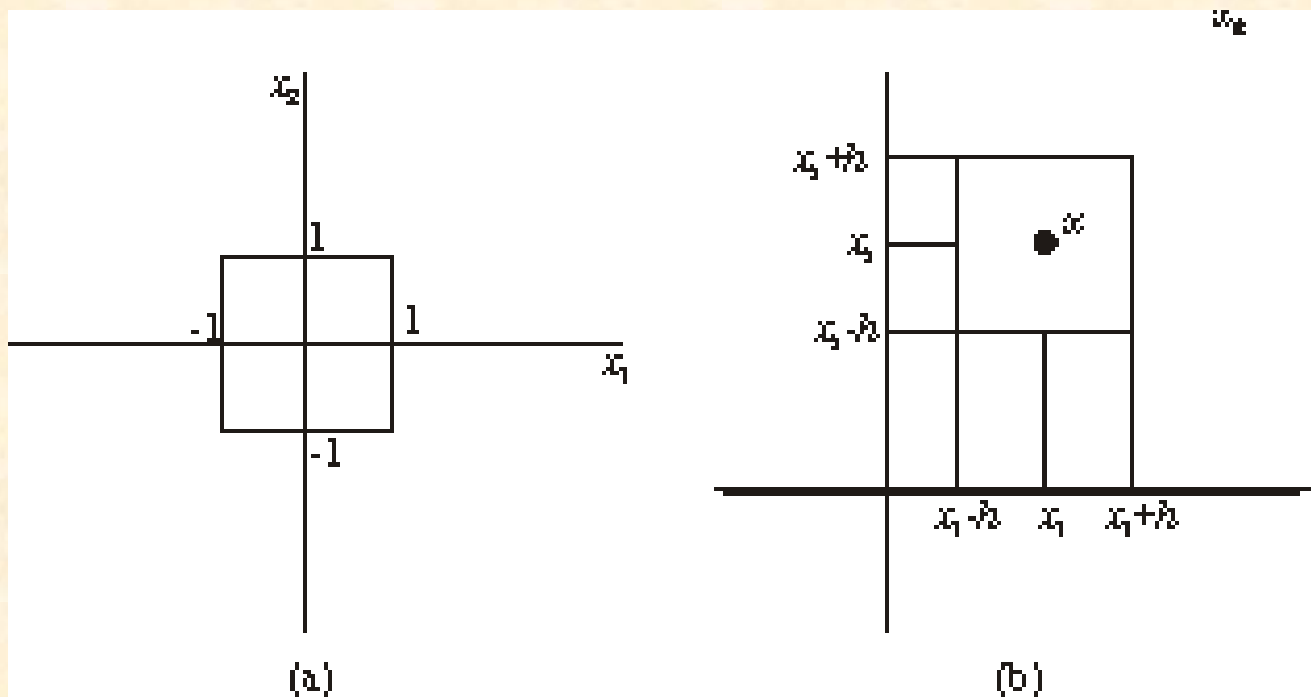
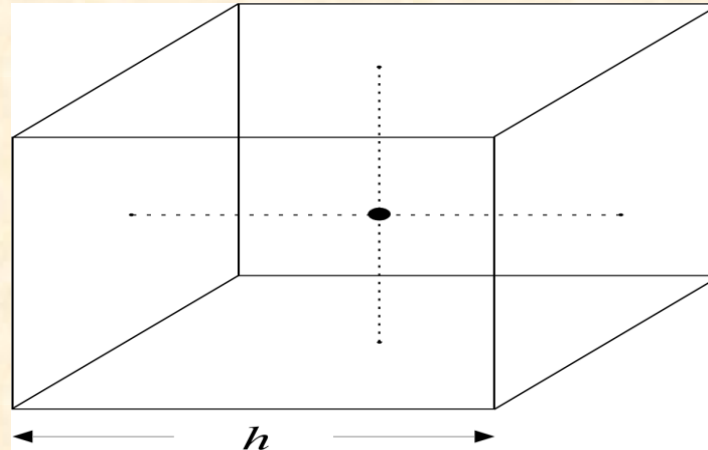
$\hat{x} - \frac{h}{2} \quad \hat{x} \quad \hat{x} + \frac{h}{2}$

$\triangleright$  Αν  $p(x)$  συνεχής,  $\hat{p}(x) \rightarrow p(x)$  καθώς  $N \rightarrow \infty$ , αν

$$h_N \rightarrow 0, \quad k_N \rightarrow \infty, \quad \frac{k_N}{N} \rightarrow 0$$

## ❖ Παράθυρα Parzen

- Διαίρεση του πολυδιάστατου χώρου σε υπερκύβους



➤ Ορίζουμε

$$\phi(\underline{x}_i) = \begin{cases} 1 & |x_{ij}| \leq 1/2 \\ 0 & \text{διαφορετικά} \end{cases}$$

- Δηλαδή, είναι 1 μέσα σε υπερκύβο μοναδιαίας πλευράς με κέντρο το 0

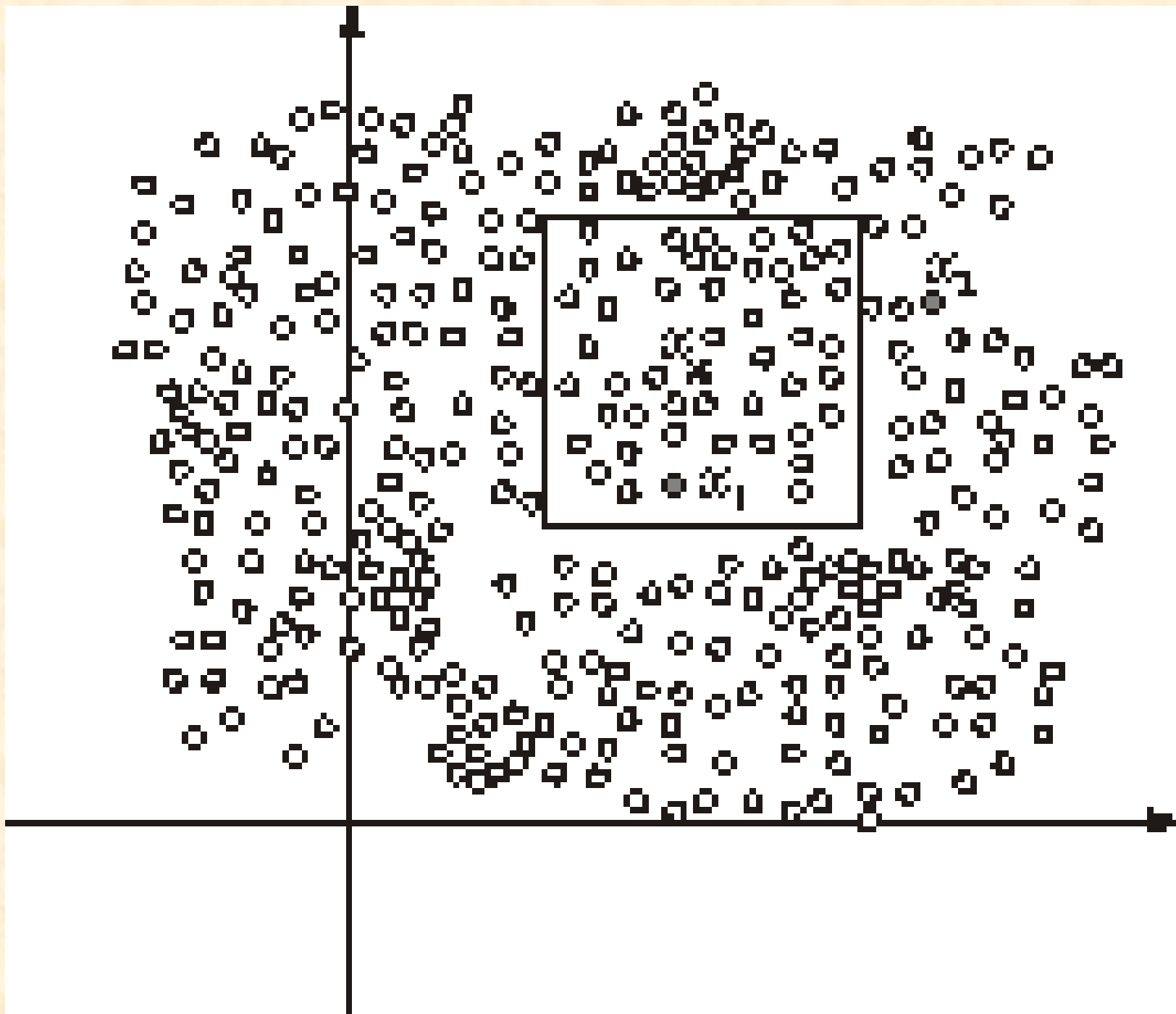
- $$\hat{p}(\underline{x}) = \frac{1}{h^l} \left( \frac{1}{N} \sum_{i=1}^N \phi\left(\frac{\underline{x}_i - \underline{x}}{h}\right) \right)$$

•  $\frac{1}{\text{όγκος}} * \frac{1}{N} *$  αριθμός σημείων εντός

ένας υπερκύβος πλευράς  $h$  κεντραρισμένος στο  $\underline{x}$

- Το πρόβλημα:  $p(\underline{x})$  συνεχής  
 $\phi(\cdot)$  ασυνεχής
- Παράθυρα Parzen-Πυρήνες-συναρτήσεις δυναμικού  
 $\phi(\underline{x})$  είναι ομαλή

$$\phi(\underline{x}) \geq 0, \quad \int_{\underline{x}} \phi(\underline{x}) d\underline{x} = 1$$



Αν  $\varphi(\mathbf{x})=N(0,I)$  ΤΟΤΕ

$$\hat{p}(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \frac{1}{(2\pi)^{\frac{l}{2}} h^l} \exp\left(-\frac{(\mathbf{x}-\mathbf{x}_i)^T(\mathbf{x}-\mathbf{x}_i)}{2h^2}\right)$$

➤ Μέση τιμή

$$E[\hat{p}(\underline{x})] = \frac{1}{h^l} \left( \frac{1}{N} \sum_{i=1}^N E[\varphi(\frac{\underline{x}_i - \underline{x}}{h})] \right) = \int_{\underline{x}'} \frac{1}{h^l} \varphi(\frac{\underline{x}' - \underline{x}}{h}) p(\underline{x}') d\underline{x}'$$

- $h \rightarrow 0, \frac{1}{h^l} \rightarrow \infty$

- $h \rightarrow 0$  το εύρος της  $\varphi(\frac{\underline{x}' - \underline{x}}{h}) \rightarrow 0$

- $\int \frac{1}{h^l} \varphi(\frac{\underline{x}' - \underline{x}}{h}) d\underline{x} = 1$

- $h \rightarrow 0 \frac{1}{h^l} \varphi(\frac{\underline{x}}{h}) \rightarrow \delta(\underline{x})$

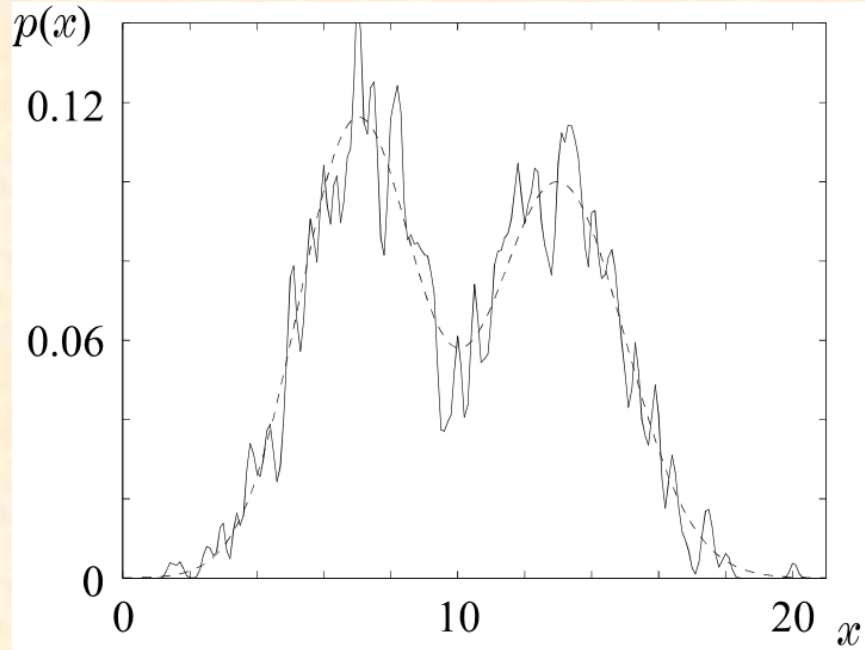
$$E[\hat{p}(\underline{x})] = \int_{\underline{x}'} \delta(\underline{x}' - \underline{x}) p(\underline{x}') d\underline{x}' = p(\underline{x})$$

Συνεπώς αμερόληπτος (unbiased) στο όριο

## ➤ Διασπορά

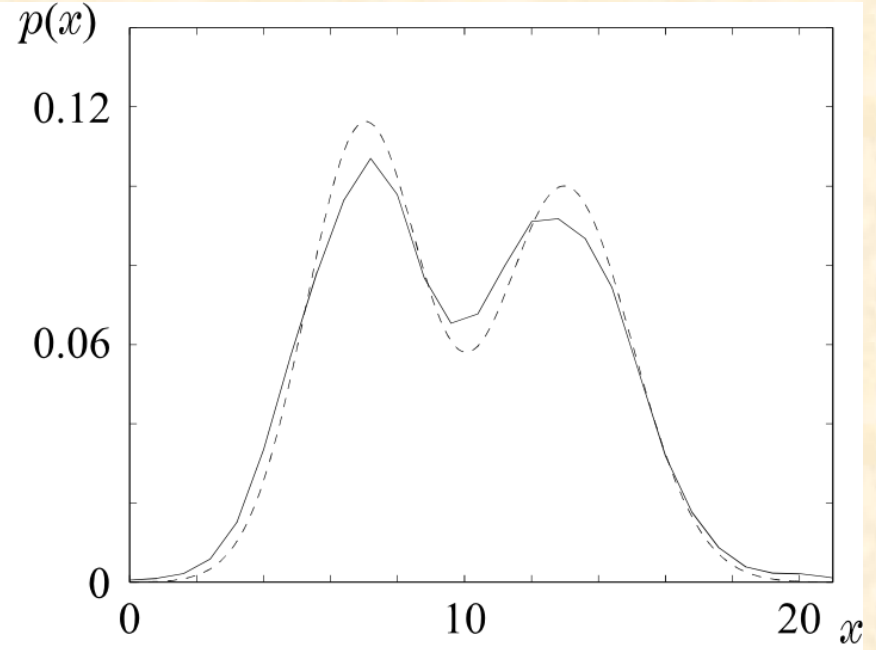
- Όσο **μικρότερο** το  $h$ , τόσο **μεγαλύτερη** η διασπορά

$h=0.1, N=1000$



(a)

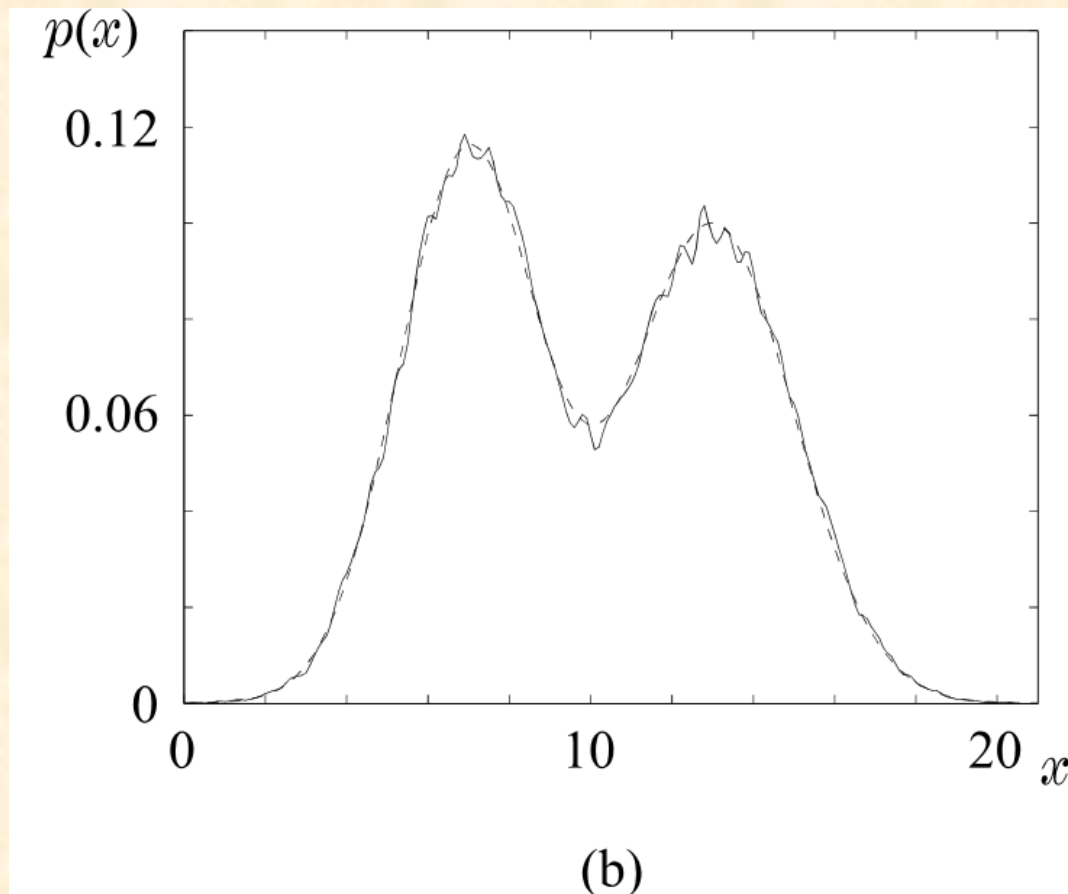
$h=0.8, N=1000$



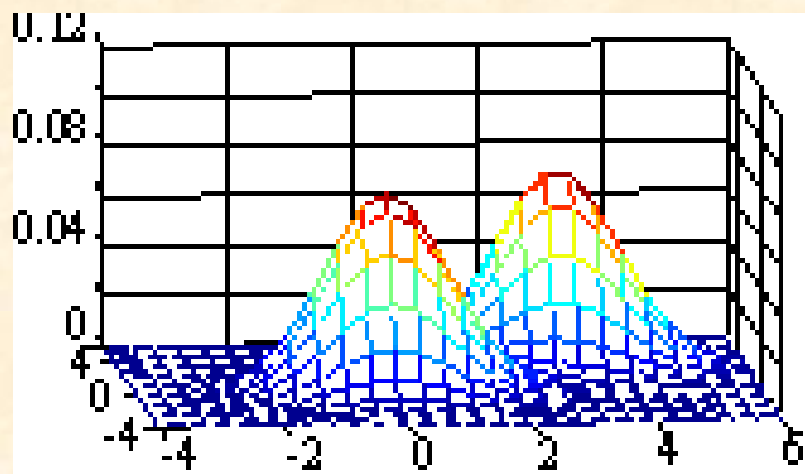
(b)



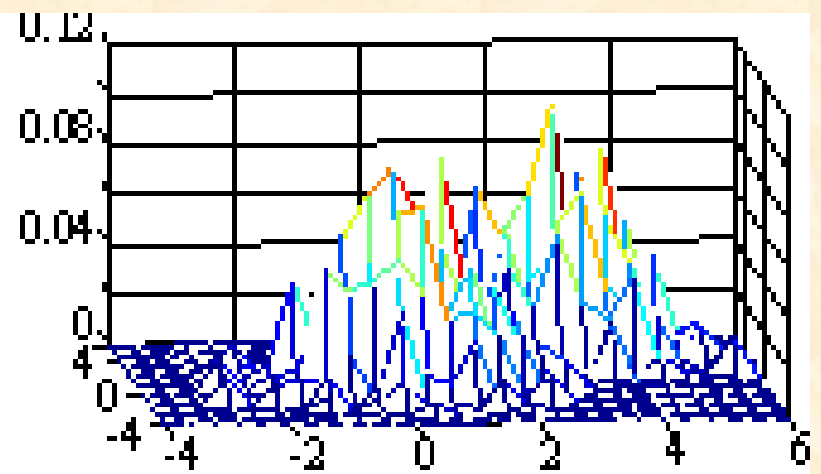
$h=0.1, N=10000$



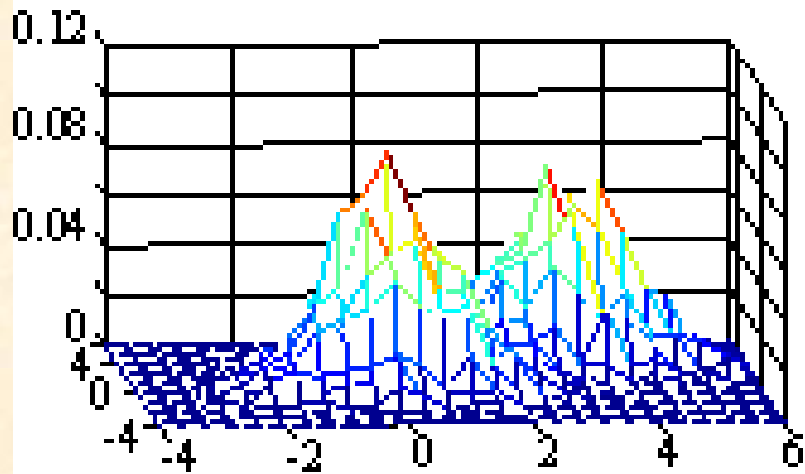
➤ Όσο μεγαλύτερο το  $N$ , τόσο καλύτερη η ακρίβεια



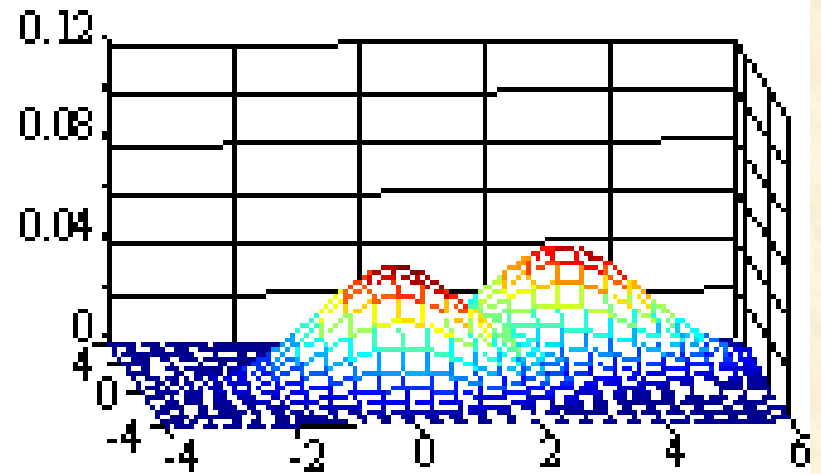
(a)



(b)



(c)



(d)

(a) αρχ. κατανομή, (b)  $h=0.05$ ,  $N=1000$ , (c)  $h=0.05$ ,  $N=20000$ , (d)  $h=0.8$ ,  $N=20000$

➤  $Av$

- $h \rightarrow 0$
- $N \rightarrow \infty$
- $hN \rightarrow \infty$

Ασυμπτωτικά απροκατάληπτος

## ❖ Εκτίμηση πυκνότητας με βάση τους K-εγγύτερους γείτονες (K Nearest Neighbor Density Estimation)

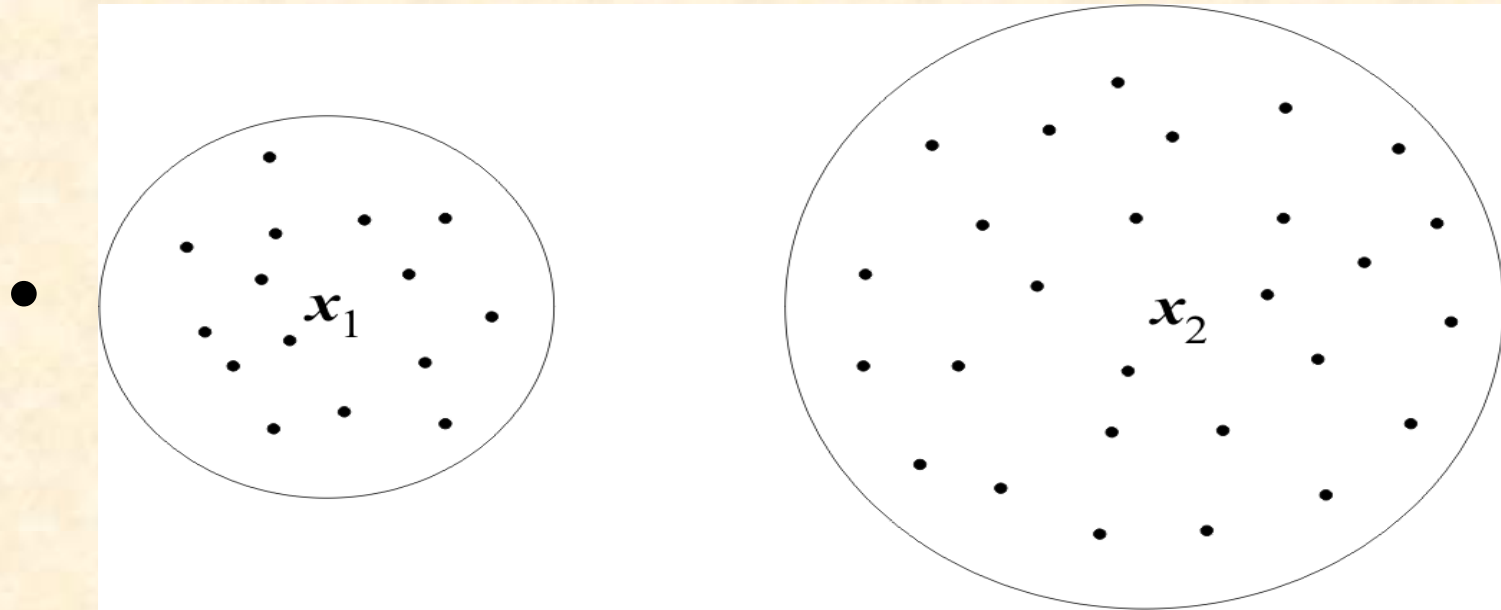
### ➤ Στα παράθυρα Parzen:

- Ο όγκος είναι σταθερός
- Ο αριθμός των σημείων στον όγκο μεταβάλλεται

### ➤ Τώρα:

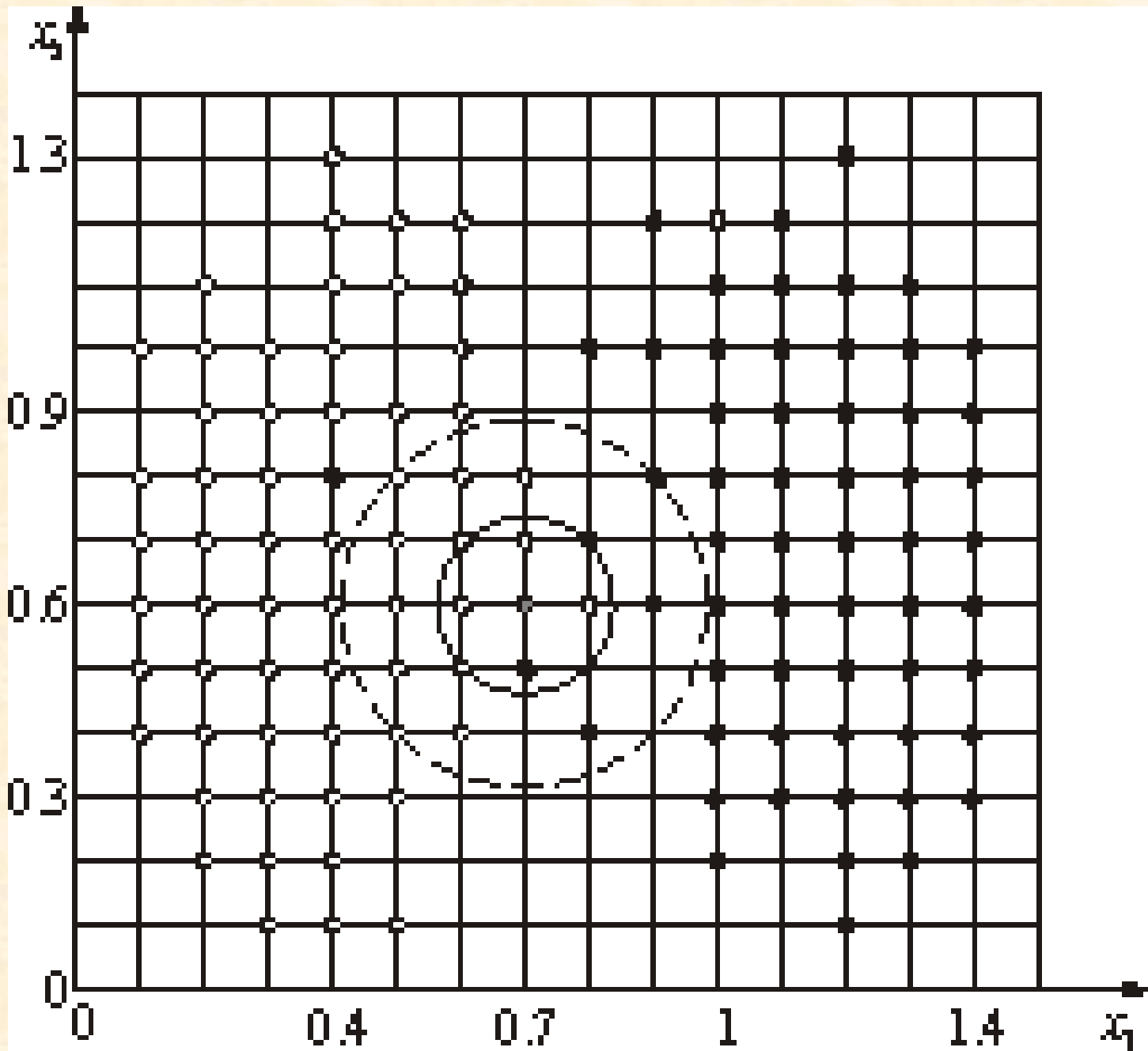
- Κρατάμε τον αριθμό των σημείων **σταθερό**  $k_N = k$
- Επιτρέπουμε στον όγκο να **μεταβάλλεται**

- $$\hat{p}(\underline{x}) = \frac{k}{NV(\underline{x})}$$



Κανόνας του Bayes:  
 $(\theta = P(\omega_2) / P(\omega_1))$

$$\frac{\frac{k}{N_1 V_1}}{\frac{k}{N_2 V_2}} = \frac{N_2 V_2}{N_1 V_1} (><) \theta$$



## ❖ ΚΑΤΑΡΑ ΤΗΣ ΔΙΑΣΤΑΤΙΚΟΤΗΤΑΣ (CURSE OF DIMENSIONALITY)

- Σε όλες τις μεθόδους, μέχρι τώρα, είδαμε ότι όσο **μεγαλύτερος** είναι ο αριθμός των σημείων,  $N$ , τόσο **καλύτερη** είναι η προκύπτουσα εκτίμηση.
- Αν στο μονοδιάστατο χώρο ένα διάστημα, που περιέχει  $N$  σημεία, είναι **αρκετό** (για καλή εκτίμηση), στον δισδιάστατο χώρο το αντίστοιχο τετράγωνο θα απαιτεί  $N^2$  και στον  $\ell$ -διάστατο χώρο ο  $\ell$ -διάστατος κύβος θα απαιτεί  $N^\ell$  σημεία.
- Η εκθετικά αύξηση του αριθμού των αναγκαίων σημείων είναι γνωστή ως **κατάρρα της διαστατικότητας** (**curse of dimensionality**). Πρόκειται για σημαντικό πρόβλημα που αντιμετωπίζει κανείς σε χώρους υψηλής διάστασης.

## ❖ ΑΠΛΟΙΚΟΣ ΤΑΞΙΝΟΜΗΤΗΣ BAYES (NAIVE – BAYES CLASSIFIER)

➤ Έστω  $\underline{x} \in \mathcal{R}^\ell$  και ότι στόχος είναι η εκτίμηση της  $p(\underline{x} | \omega_i)$   $i = 1, 2, \dots, M$ . Για μία “καλή” εκτίμηση της pdf χρειάζονται, ας πούμε,  $N^\ell$  σημεία.

➤ Έστω ότι  $x_1, x_2, \dots, x_\ell$  είναι **αμοιβαίως ανεξάρτητα**. Τότε:

$$p(\underline{x} | \omega_i) = \prod_{j=1}^{\ell} p(x_j | \omega_i)$$

➤ Σ’ αυτή την περίπτωση, κάποιος θα χρειαζόταν, περίπου,  $N$  σημεία για κάθε pdf.

➤ Ο απλοϊκός ταξινομητής Bayes δουλεύει καλά ακόμα και σε περιπτώσεις όπου παραβιάζεται η υπόθεση της ανεξαρτησίας



## ❖ Ο κανόνας του εγγύτερου γείτονα (The Nearest Neighbor Rule)

- **Για δεδομένο  $\underline{x}$**
- Από δεδομένο σύνολο  $N$  διανυσμάτων εκπαίδευσης, προσδιόρισε τα  $k$  πλησιέστερα στο  $\underline{x}$
- Από αυτά τα  $k$  προσδιόρισε τα  $k_i$  που ανήκουν στην κλάση  $\omega_i$   
Καταχώρησε  $\underline{x} \rightarrow \omega_i : k_i > k_j \quad \forall i \neq j$
- Η απλούστερη περίπτωση  

$k=1 !!!$
- Για μεγάλα  $N$  παρουσιάζει καλή συμπεριφορά. Αποδεικνύεται ότι:  
Αν  $P_B$  είναι η βέλτιστη πιθανότητα λάθους κατά Bayes, τότε:

$$P_B \leq P_{NN} \leq P_B \left( 2 - \frac{M}{M-1} P_B \right) \leq 2P_B$$

➤  $P_B \leq P_{kNN} \leq P_B + \sqrt{\frac{2P_{NN}}{k}}$       Για δύο κλάσεις

➤  $k \rightarrow \infty, P_{kNN} \rightarrow P_B$

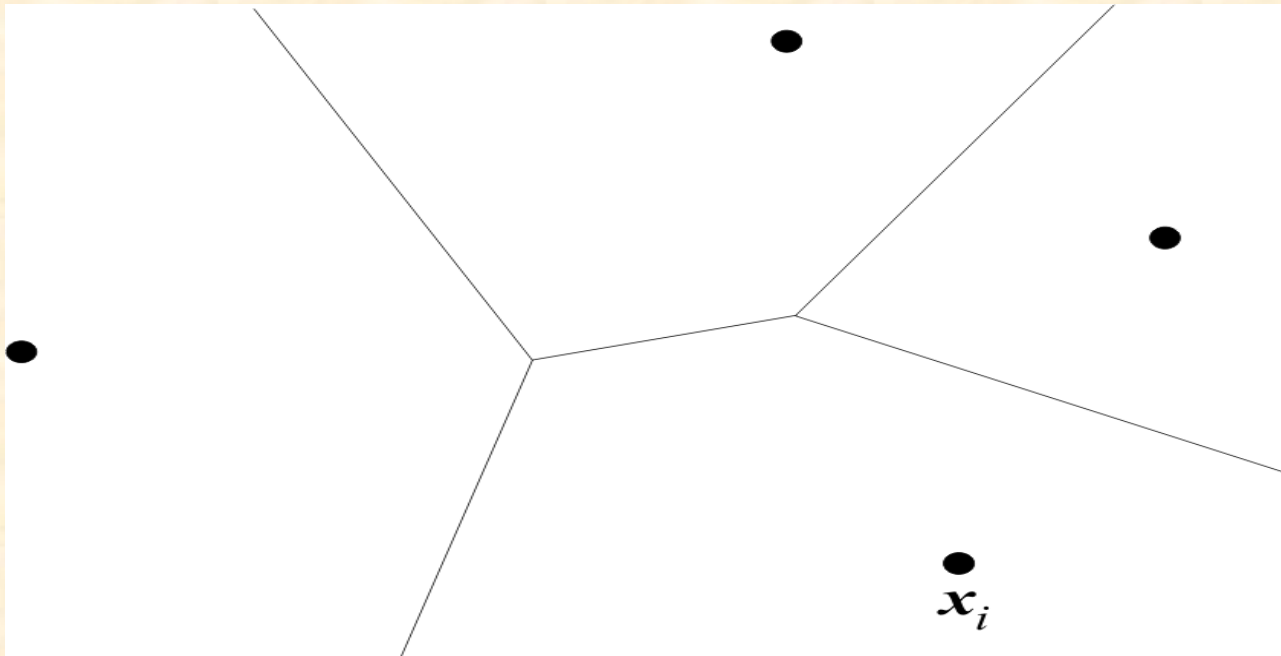
➤ Για μικρό  $P_B$ :

$$P_{NN} \cong 2P_B$$

$$P_{3NN} \cong P_B + 3(P_B)^2$$

**Πρόβλημα:** Μεγάλο υπολογιστικό κόστος ( $O(kN)$  ανά δείγμα)

❖ Ψηφοθέτηση Voronoi (Voronoi tessellation) (1-NN)



$$R_i = \{ \underline{x} : d(\underline{x}, \underline{x}_i) < d(\underline{x}, \underline{x}_j) \ i \neq j \}$$