

# Non-Parametric Techniques

Generative and Discriminative Methods

Density Estimation

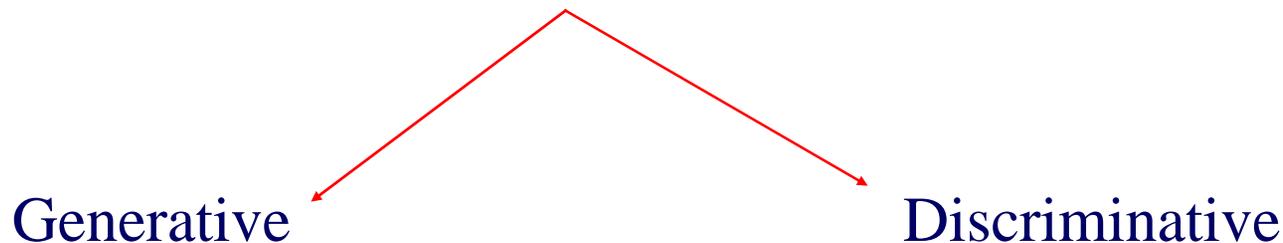
# Parametric versus Non-parametric

## 1. Parametric

- densities are uni-modal
  - have a single local maximum
- practical problems involve multi-modal densities

## 2. Nonparametric

- arbitrary distributions
- without assuming forms of the underlying densities



# Two types of nonparametric methods

## 1. *Generative*:

Estimate class-conditional density  $p(x / \omega_j)$

1. Parzen Windows
2.  $k_n$ -nearest neighbor estimation

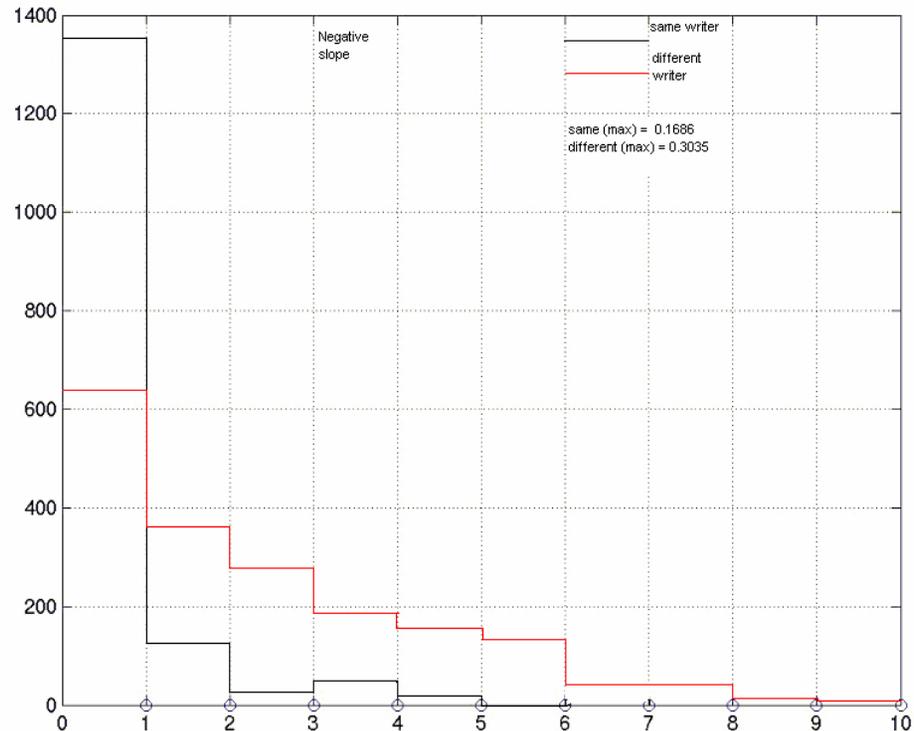
## 2. *Discriminative*:

Bypass density estimation and go directly to compute *a-posteriori* probability  $P(\omega_j/x)$

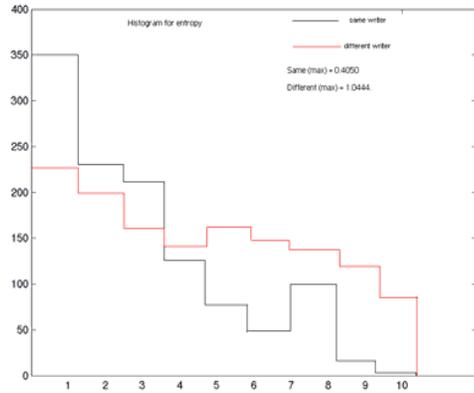
1. nearest neighbor rule
2. k-nearest neighbor rule

# Histograms are the heart of density estimation

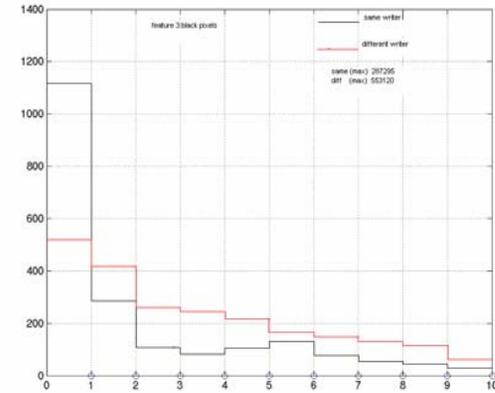
An example of a histogram (feature called negative slope)



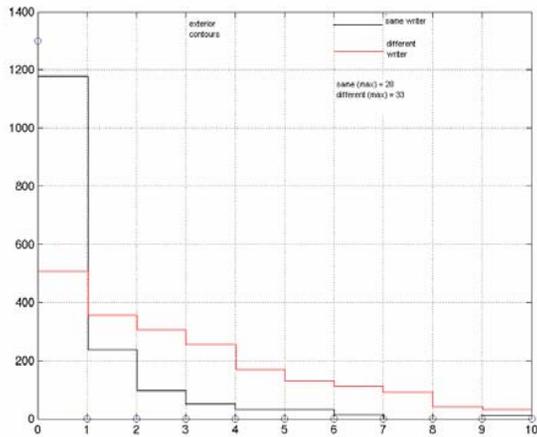
# Entropy



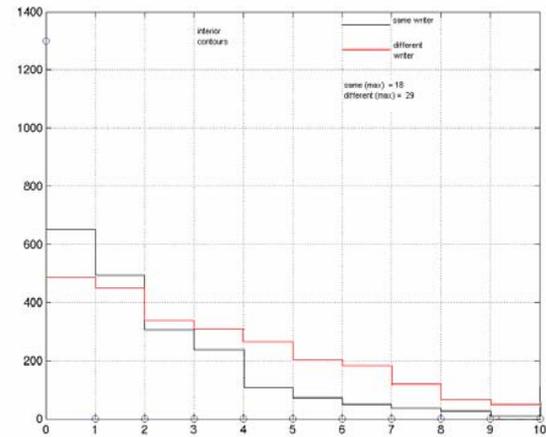
# No of black pixels



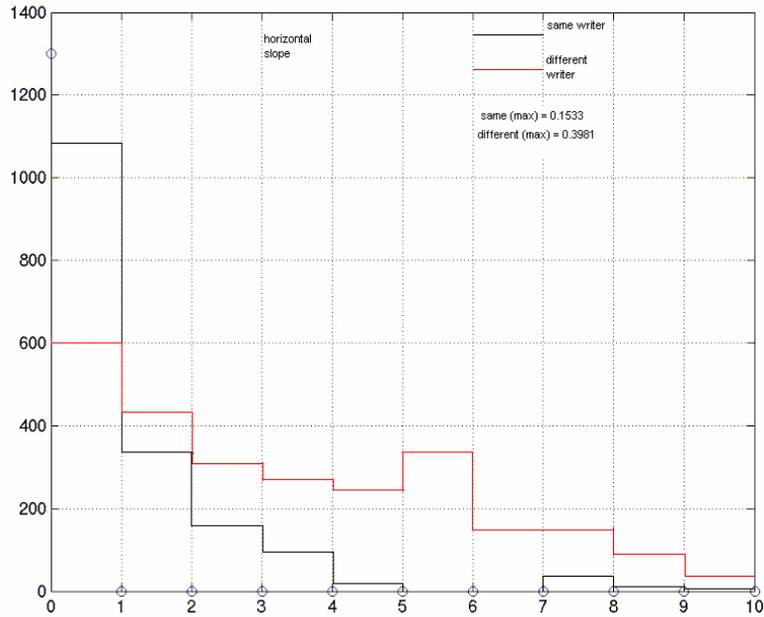
# Exterior contours



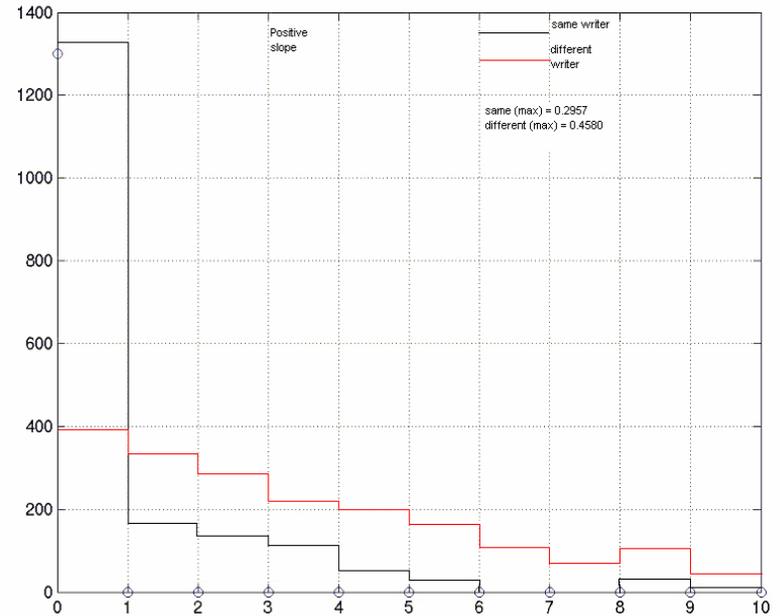
# Interior contours



# Histogram for Horizontal slope



# Histogram for Positive slope

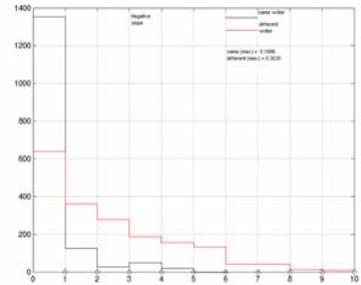


# Density Estimation

- Basic idea of estimating an unknown pdf:
- Probability  $P$  that a vector  $x$  will fall in region  $R$  is:

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

- $P$  is a smoothed (or averaged) version of the density function  $p(x)$
- We can estimate the smoothed value of  $p$  by estimating the probability  $P$



# Relating $P$ to histogram value

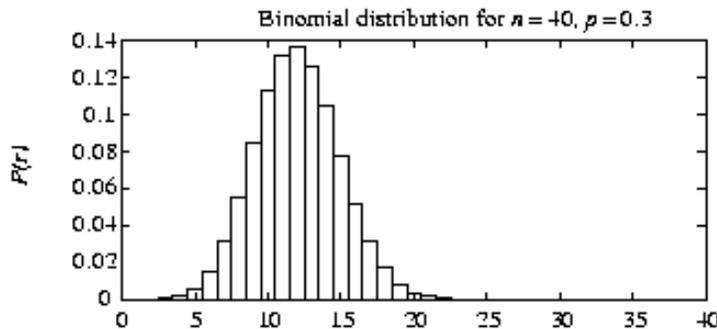
- If we have a sample of size  $n$ ,
  - the probability that  $k$  points fall in  $R$  is given by the binomial law:

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

- and the expected value (mean) for  $k$  is:

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

# Binomial Distribution



$$P(r) = \frac{n!}{r!(n-r)!} p^r (1-p)^{n-r}$$

Probability  $P(r)$  of  $r$  heads in  $n$  coin flips, if  $p = \text{Pr}(\text{heads})$

- Expected, or mean value of  $X$ ,  $E[X]$ , is

$$E[X] \equiv \sum_{i=0}^n iP(i) = np$$

- Variance of  $X$  is

$$\text{Var}(X) \equiv E[(X - E[X])^2] = np(1-p)$$

- Standard deviation of  $X$ ,  $\sigma_X$ , is

$$\sigma_X \equiv \sqrt{E[(X - E[X])^2]} = \sqrt{np(1-p)}$$

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

- This binomial distribution for  $k$  peaks very sharply about the mean

$$k \approx nP$$

- Therefore, the ratio  $k/n$  is a good estimate for the probability  $P$  and hence for the density function  $p$
- Especially accurate when  $n$  is very large

# Space average of $p(x)$

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

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

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

# Expression for $p(x)$ from histogram

Combining equations

(1), (3) and (4)



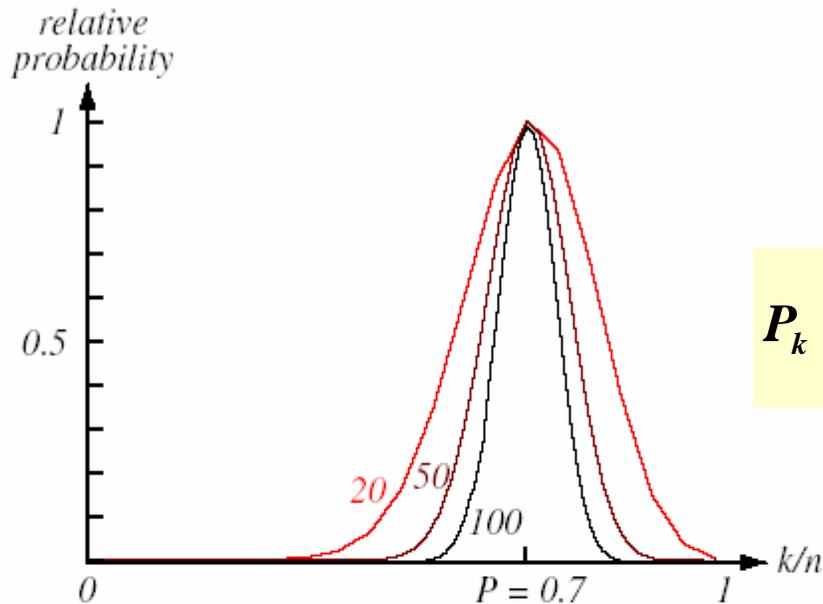
$$P = \int_{\mathfrak{R}} p(x') dx' \quad (1)$$

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

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

yields

$$p(x) \cong \frac{k/n}{V}$$



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

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.

# Practical and Theoretical Issues

- The fraction  $p(\mathbf{x}) \cong \frac{k/n}{V}$  is a space averaged value of  $p(\mathbf{x})$
- $p(\mathbf{x})$  is obtained only if  $V$  approaches zero  
This is the case where no samples are included in  $\mathcal{R}$  and our estimate  $p(\mathbf{x}) = 0$
- 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  $k/n$

# Formulation when unlimited samples are available

Theoretically, if an unlimited number of samples is available, we can circumvent this difficulty

To estimate the density at  $x$ , 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  $n^{\text{th}}$  estimate for  $p(x)$ :

$$p_n(x) \cong \frac{k_n / n}{V_n} \quad (7)$$

# Necessary conditions for convergence: $p_n(x) \rightarrow p(x)$

$$p_n(x) \cong \frac{k_n / n}{V_n}$$

$$1) \lim_{n \rightarrow \infty} V_n = 0$$

$$2) \lim_{n \rightarrow \infty} k_n = \infty$$

$$3) \lim_{n \rightarrow \infty} k_n / n = 0$$

Two different ways of obtaining sequences of regions that satisfy these conditions:

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

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

This is called “the Parzen-window estimation method”

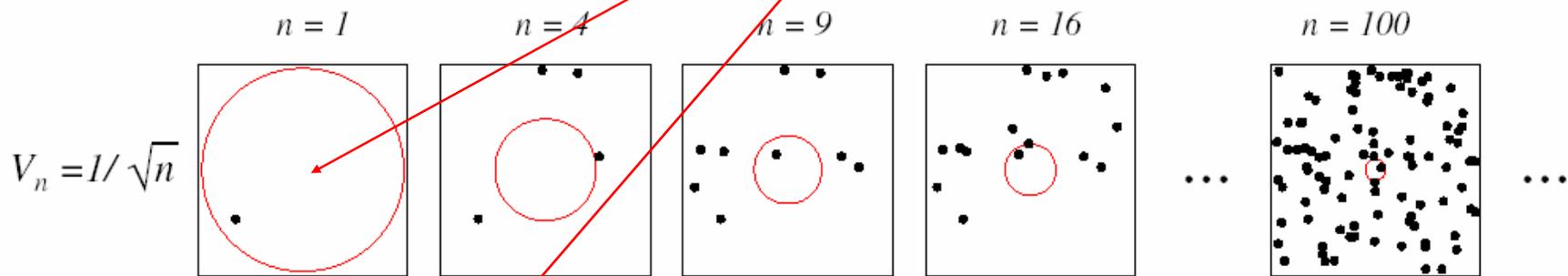
(b) Specify  $k_n$  as some function of  $n$ , such as  $k_n = \sqrt[n]{n}$ ; the volume  $V_n$  is grown until it encloses  $k_n$  neighbors of  $x$ .

This is called “the  $k_n$ -nearest neighbor estimation method”

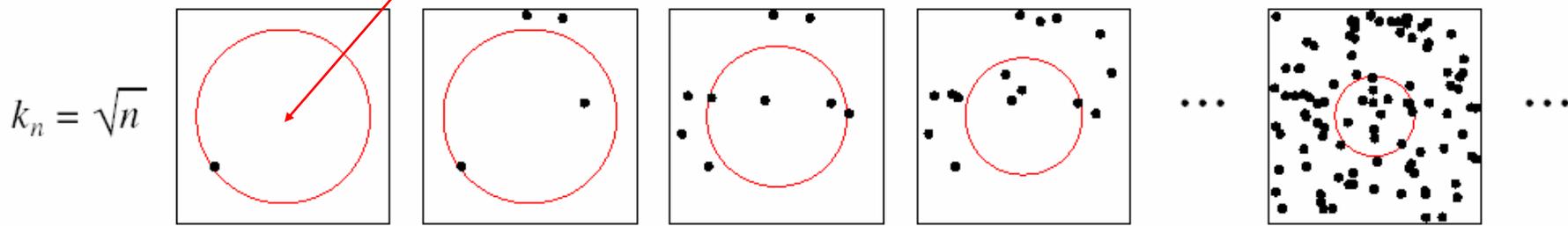
# Two Methods for Estimating Density

We are estimating density at center of square

$$V_n = \frac{1}{\sqrt{n}}$$



$$k_n = \sqrt{n}$$

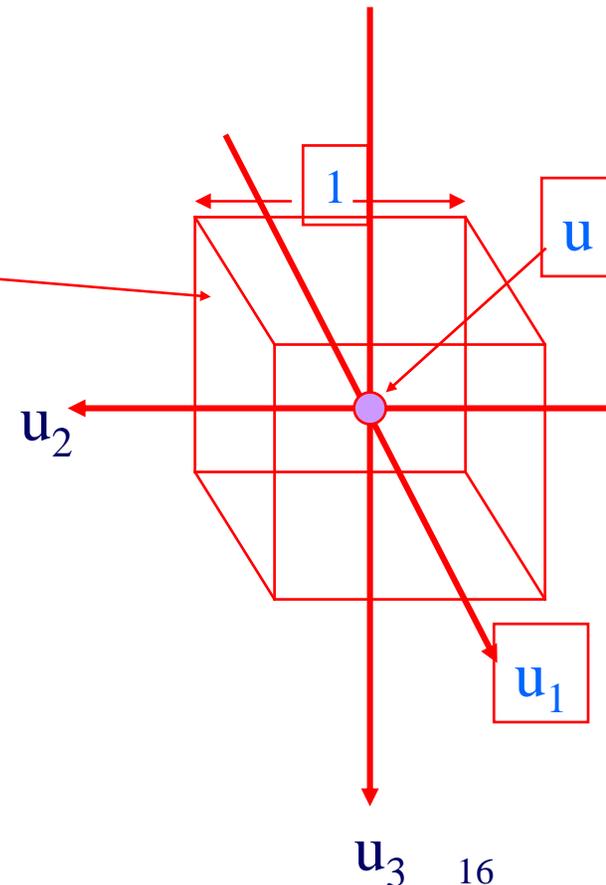


# Parzen Window Method

- Unit hypercube function

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

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



# Parzen Window using hypercubes

- Estimate density assuming that region  $\mathcal{R}_n$  is a d-dimensional hypercube

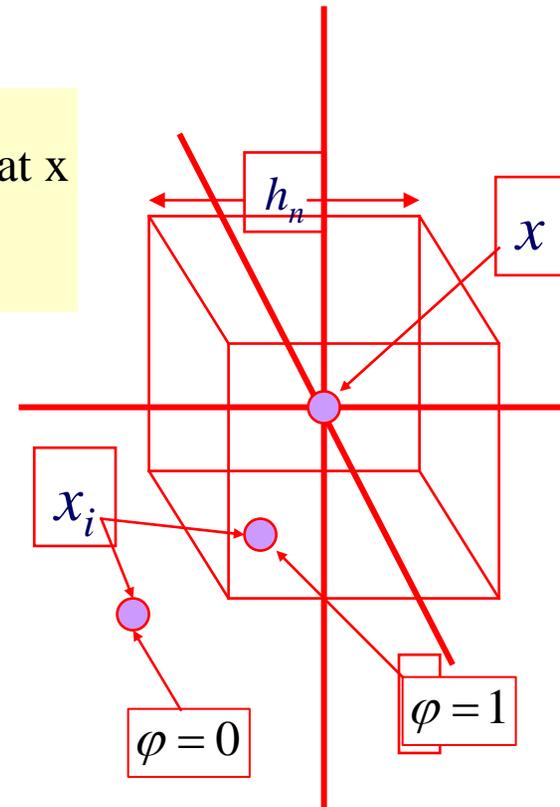
$$V_n = h_n^d \quad (h_n : \text{length of the edge of } \mathcal{R}_n)$$

$$\varphi\left(\frac{x - x_i}{h_n}\right) = 1 \text{ if } x_i \text{ falls within hypercube of volume } V_n \text{ centered at } x \\ = 0 \text{ otherwise}$$

$$\text{No. of samples in hypercube: } k_n = \sum_{i=1}^n \varphi\left(\frac{x - x_i}{h_n}\right)$$

$$p_n(x) = \frac{k_n / n}{V_n} \\ = \frac{1}{n} \sum_{i=1}^n \frac{1}{V_n} \varphi\left(\frac{x - x_i}{h_n}\right)$$

CSE 555: Srihari



# Properties of Window Function

$$p_n(x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{V_n} \varphi\left(\frac{x - x_i}{h_n}\right)$$

$p_n(x)$  estimates  $p(x)$  as an average of functions of  $x$  and the samples  $x_i$

Window function is an *interpolation function*

- Functions  $\varphi$  can be general as long as  $p_n(x)$  is a legitimate density function, i.e., require that

$$\varphi(\mathbf{x}) \geq 0$$

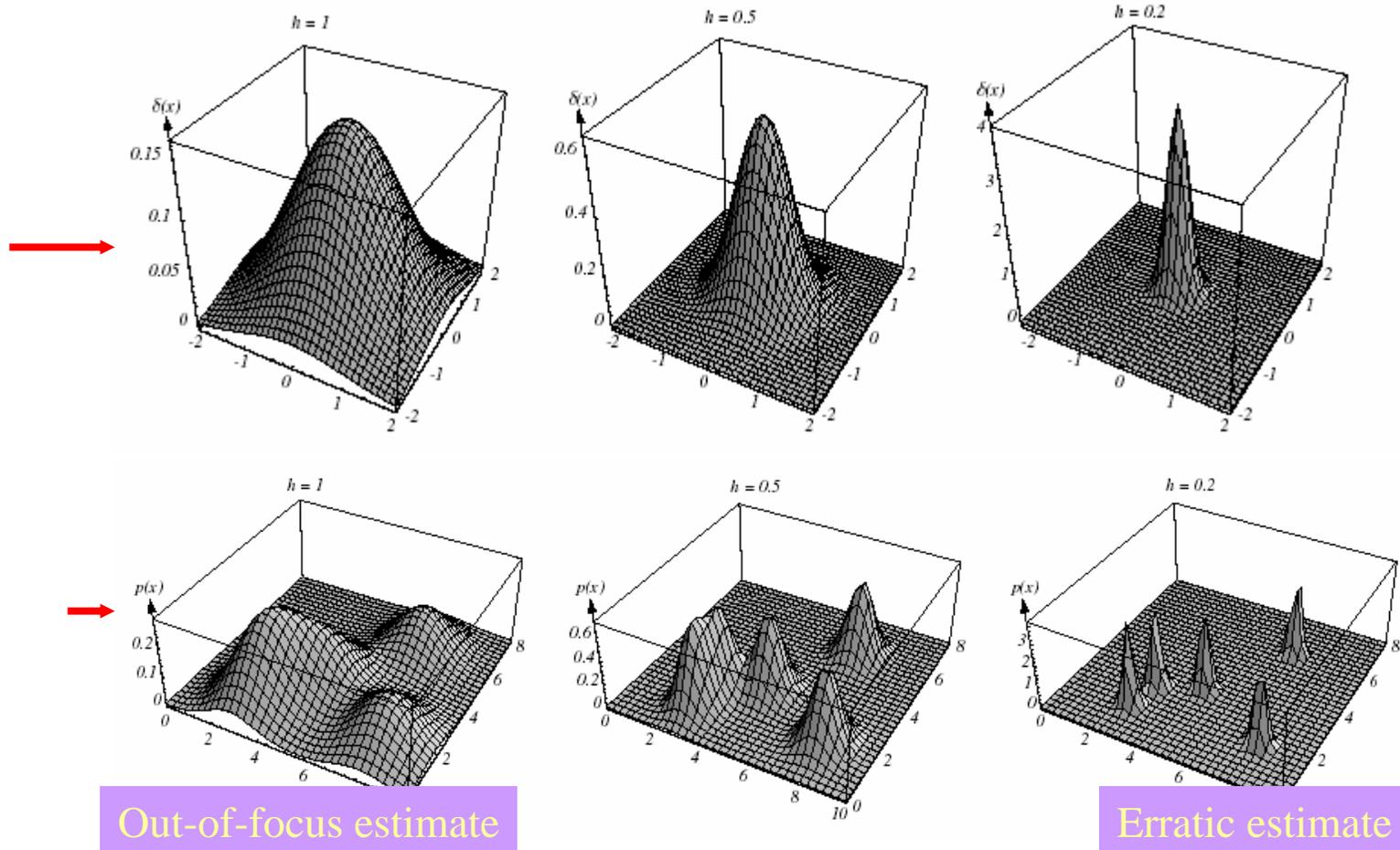
$$\int \varphi(\mathbf{u}) d\mathbf{u} = 1,$$

Can use a Gaussian that is circularly symmetric with variance  $h_n$

# Effect of window-width on estimate $p_n(x)$

$$\text{Let } \delta_n(x) = \frac{1}{V_n} \varphi\left(\frac{x}{h_n}\right) \text{ then } p_n(x) = \frac{1}{n} \sum_{i=1}^n \delta_n(x - x_i)$$

Gaussian windows with decreasing widths



For any  $h_n$ , distribution is normalized, i.e.,

$$\int \delta_n(\mathbf{x} - \mathbf{x}_i) d\mathbf{x} = \int \frac{1}{V_n} \varphi\left(\frac{\mathbf{x} - \mathbf{x}_i}{h_n}\right) d\mathbf{x} = \int \varphi(\mathbf{u}) d\mathbf{u} = 1.$$

# Conditions for Convergence of Estimate

Since we speak of convergence of sequence of random variables,  $p_n(\mathbf{x})$  has a mean and a variance

$$\lim_{n \rightarrow \infty} \bar{p}_n(\mathbf{x}) = p(\mathbf{x})$$

$$\lim_{n \rightarrow \infty} \sigma_n^2(\mathbf{x}) = 0.$$

Additional criteria for

- (1) well-behaved ness
- (2) volume approaches zero
- (3) but at rate slower than  $1/n$

$$\sup_{\mathbf{u}} \varphi(\mathbf{u}) < \infty$$

$$\lim_{\|\mathbf{u}\| \rightarrow \infty} \varphi(\mathbf{u}) \prod_{i=1}^d u_i = 0$$

$$\lim_{n \rightarrow \infty} V_n = 0$$

$$\lim_{n \rightarrow \infty} nV_n = \infty.$$

# Convergence of Mean and Variance

## Mean

$$\begin{aligned}\bar{p}_n(\mathbf{x}) &= \mathcal{E}[p_n(\mathbf{x})] \\ &= \frac{1}{n} \sum_{i=1}^n \mathcal{E}\left[\frac{1}{V_n} \varphi\left(\frac{\mathbf{x} - \mathbf{x}_i}{h_n}\right)\right] \\ &= \int \frac{1}{V_n} \varphi\left(\frac{\mathbf{x} - \mathbf{v}}{h_n}\right) p(\mathbf{v}) d\mathbf{v} \\ &= \int \delta_n(\mathbf{x} - \mathbf{v}) p(\mathbf{v}) d\mathbf{v}.\end{aligned}$$

$$\begin{aligned}\sigma_n^2(\mathbf{x}) &= \sum_{i=1}^n \mathcal{E}\left[\left(\frac{1}{nV_n} \varphi\left(\frac{\mathbf{x} - \mathbf{x}_i}{h_n}\right) - \frac{1}{n} \bar{p}_n(\mathbf{x})\right)^2\right] \\ &= n \mathcal{E}\left[\frac{1}{n^2 V_n^2} \varphi^2\left(\frac{\mathbf{x} - \mathbf{x}_i}{h_n}\right)\right] - \frac{1}{n} \bar{p}_n^2(\mathbf{x}) \\ &= \frac{1}{nV_n} \int \frac{1}{V_n} \varphi^2\left(\frac{\mathbf{x} - \mathbf{v}}{h_n}\right) p(\mathbf{v}) d\mathbf{v} - \frac{1}{n} \bar{p}_n^2(\mathbf{x}).\end{aligned}$$

$$\sigma_n^2(\mathbf{x}) \leq \frac{\sup(\varphi(\cdot)) \bar{p}_n(\mathbf{x})}{nV_n}.$$

# Example of Parzen Window on simple cases

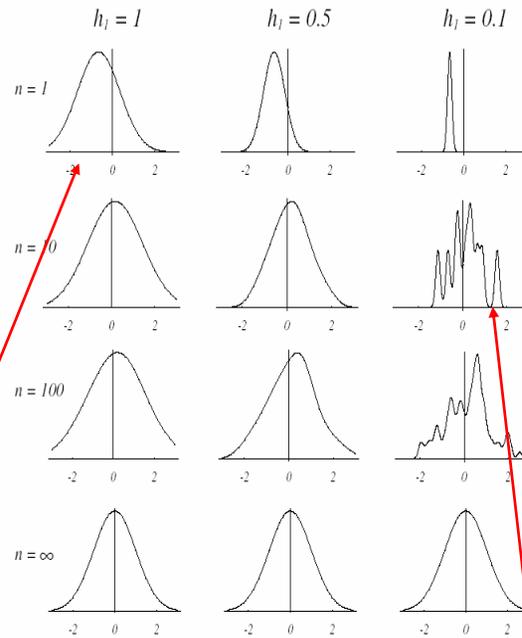
$$\varphi(u) = \frac{1}{\sqrt{2\pi}} e^{-u^2/2} \quad \text{and} \quad h_n = h_1 / \sqrt{n}$$

Case where  $p(x) \rightarrow N(0,1)$

$$p_n(x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h_n} \varphi\left(\frac{x - x_i}{h_n}\right)$$

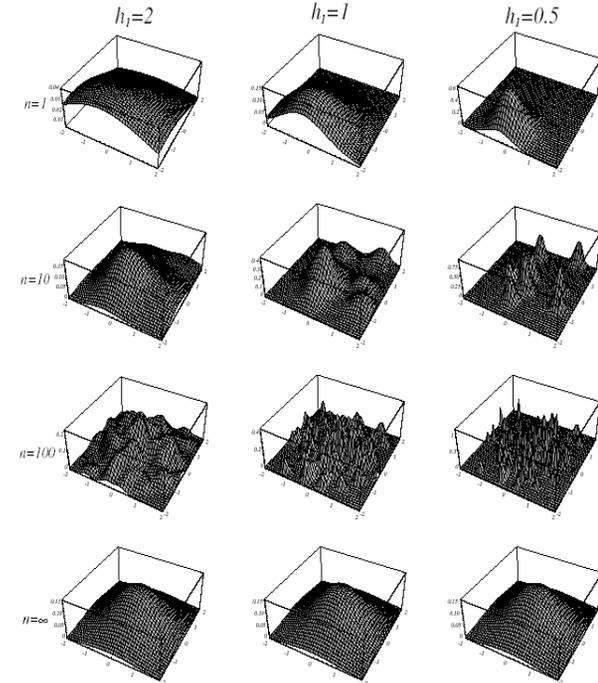
is an average of normal densities centered at the samples  $x_i$ .

## Univariate normal



Parzen-window estimates of a univariate normal density using different window widths and numbers of samples. The vertical axes have been scaled to best show the structure in each graph. Note particularly that the  $n = \infty$  estimates are the same (and match the true density function), regardless of window width

## Bivariate normal



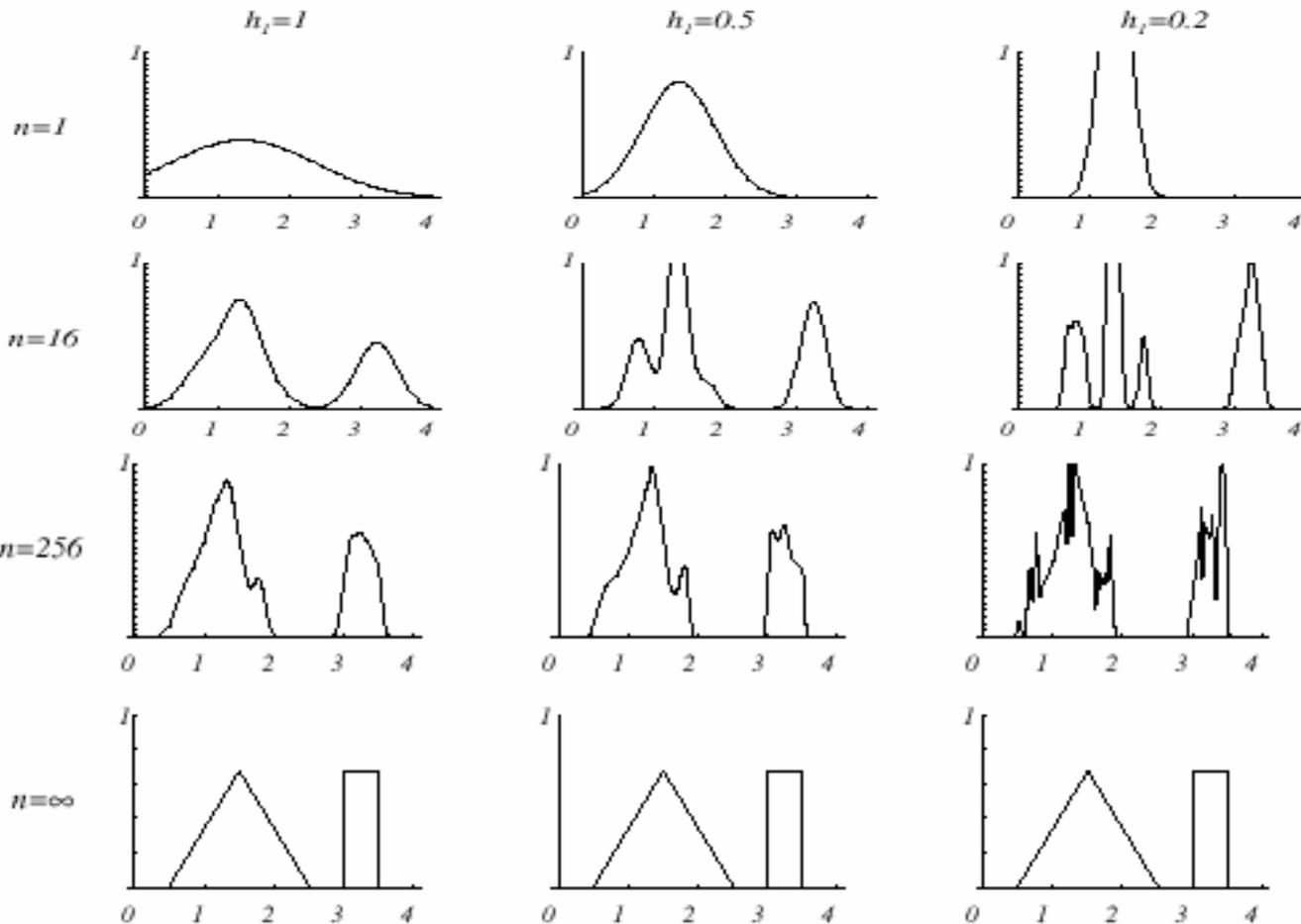
Parzen-window estimates of a bivariate normal density using different window widths and numbers of samples. The vertical axes have been scaled to best show the structure in each graph. Note particularly that the  $n = \infty$  estimates are the same (and match the true distribution), regardless of window width

$$p_1(x) = \varphi(x - x_1) = \frac{1}{\sqrt{2\pi}} e^{-1/2(x - x_1)^2} \rightarrow N(x_1, 1)$$

Contributions of samples clearly observable

# Parzen Window Estimates of Bimodal Distribution

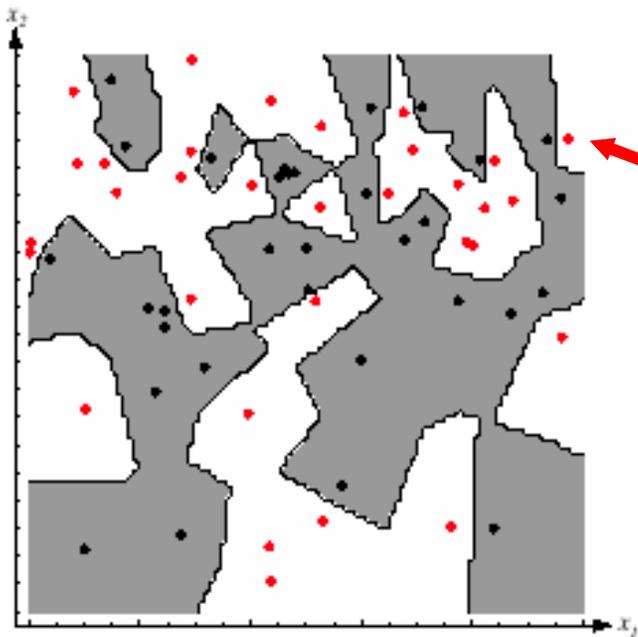
- Case where  $p(x) = \lambda_1 \cdot U(a,b) + \lambda_2 \cdot T(c,d)$
- (mixture of a uniform and a triangle density)



# Classification based on Parzen Windows

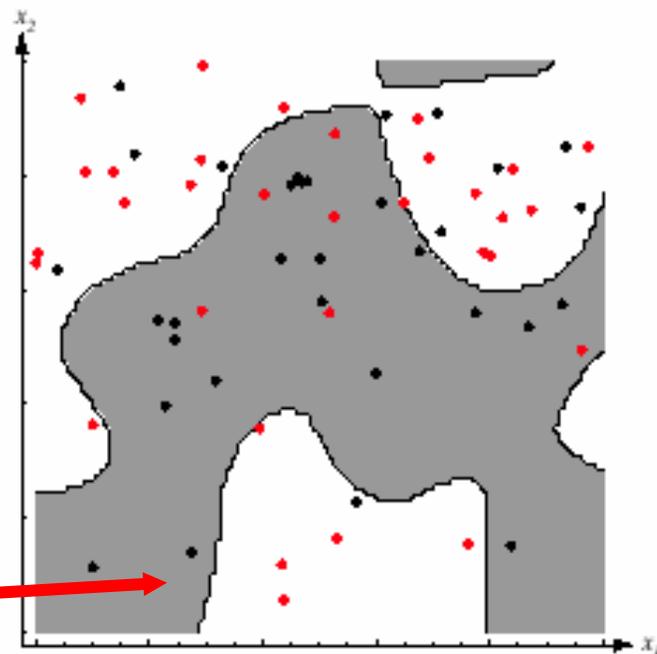
We estimate the densities for each category and classify a test point by the label corresponding to the maximum posterior

## Decision Boundaries of Two-dimensional dichotomizer



Small  $h$   
is better  
in upper  
region.

Large  $h$   
in lower  
Region.

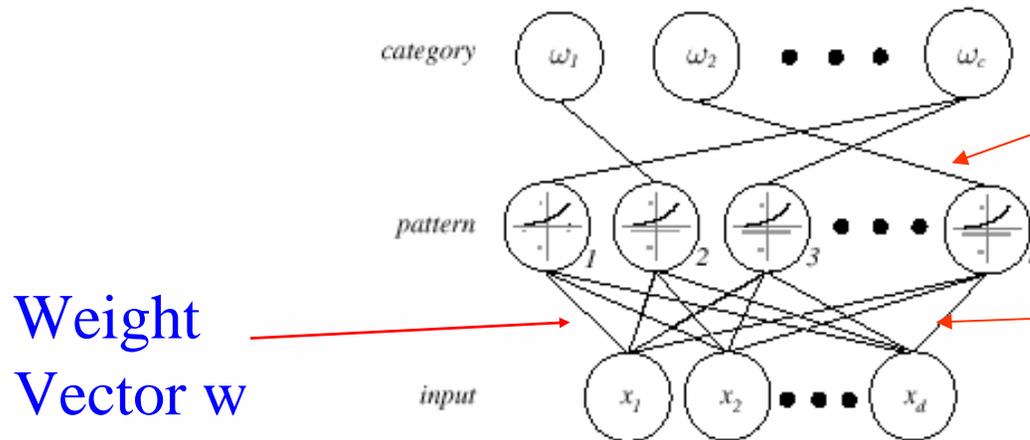


# Parzen Window Conclusion

- No assumptions made ahead of time regarding distributions
- Same procedure was used for unimodal and bimodal cases
- With enough samples assured of convergence
  - May need large no of samples however
  - Also, no data reduction provided which leads to severe computation and storage requirements
  - Demand for no of samples grows exponentially with dimensionality

# Probabilistic Neural Networks

Parzen Window Method can be implemented as an artificial neural network



Weight  
Vector  $w$

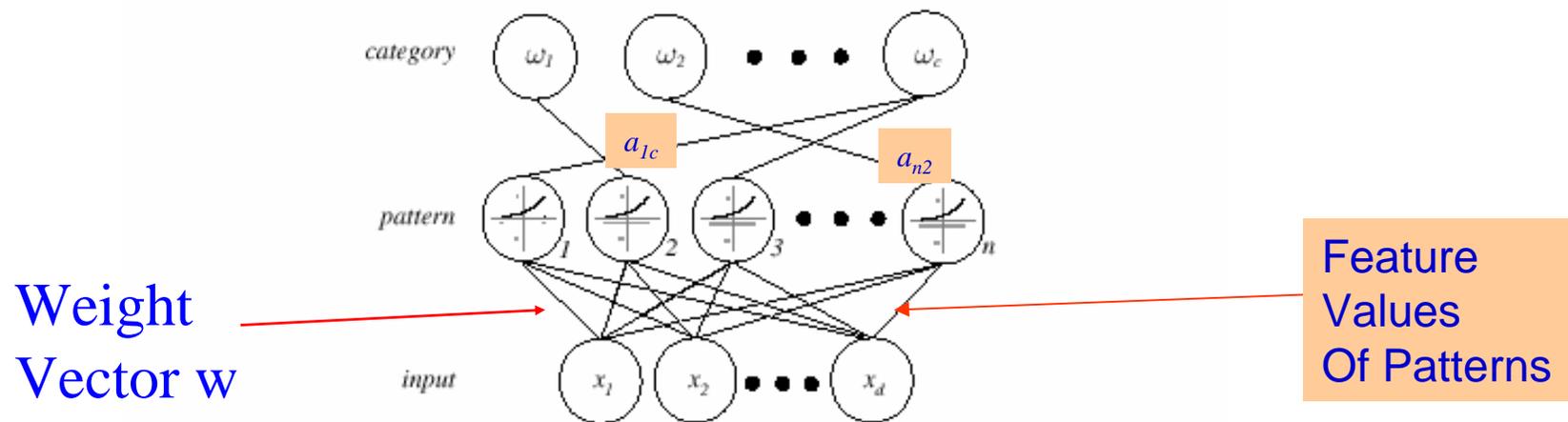
Connections  
correspond  
to class labels  
of  
patterns

Feature  
Values  
Of Patterns

A probabilistic neural network (PNN) consists of  $d$  input units,  $n$  pattern units, and  $c$  category units. Each pattern unit forms the inner product of its weight vector and the normalized pattern vector  $\mathbf{x}$  to form  $z = \mathbf{w}^t \mathbf{x}$ , and then it emits  $\exp[(z - 1)/\sigma^2]$ . Each category unit sums such contributions from the pattern unit connected to it. This ensures that the activity in each of the category units represents the Parzen-window density estimate using a circularly symmetric Gaussian window of covariance  $\sigma^2 \mathbf{I}$ , where  $\mathbf{I}$  is the  $d \times d$  identity matrix. |

# Probabilistic Neural Networks

Parzen Window Method can be implemented as an artificial neural network



A probabilistic neural network (PNN) consists of  $d$  input units,  $n$  pattern units, and  $c$  category units. Each pattern unit forms the inner product of its weight vector and the normalized pattern vector  $\mathbf{x}$  to form  $z = \mathbf{w}^t \mathbf{x}$ , and then it emits  $\exp[(z - 1)/\sigma^2]$ . Each category unit sums such contributions from the pattern unit connected to it. This ensures that the activity in each of the category units represents the Parzen-window density estimate using a circularly symmetric Gaussian window of covariance  $\sigma^2 \mathbf{I}$ , where  $\mathbf{I}$  is the  $d \times d$  identity matrix.

# Normalization

- Patterns are normalized (or scaled) to have unit length, or

$$\sum_{i=1}^d \mathbf{x}_i^2 = 1$$

- This is done by replacing each feature value by

$$\mathbf{x}_j \leftarrow \frac{\mathbf{x}_j}{\left( \sum_{i=1}^d \mathbf{x}_i^2 \right)^{1/2}}$$

- Effect of normalization

$$\mathbf{x}^t \mathbf{x} = 1$$

# Normalization Example

- Normalize  $\mathbf{x} = [3 \ 4]^t$
- $(9+16)^{1/2} = 5$
- Normalized vector is  $= [3/5 \ 4/5]^t = [0.6 \ 0.8]^t$

- Effect of normalization

$$\sum_{i=1}^d \mathbf{x}_i^2 = 0.36 + 0.64 = 1$$

$$\mathbf{x}^t \mathbf{x} = \begin{bmatrix} 0.6 \\ 0.8 \end{bmatrix} [0.6, 0.8] = [1]$$

# Activation Function

$$\varphi\left(\frac{\mathbf{x} - \mathbf{w}_k}{\mathbf{h}_n}\right) \alpha e^{-(\mathbf{x} - \mathbf{w}_k)^t (\mathbf{x} - \mathbf{w}_k) / 2\sigma^2}$$

Contribution of a sample at  $\mathbf{x}$   
Here we let weights correspond to feature values of sample

$$\begin{aligned} & e^{-(\mathbf{x} - \mathbf{w}_k)^t (\mathbf{x} - \mathbf{w}_k) / 2\sigma^2} \\ &= e^{-(\mathbf{x}^t \mathbf{x} + \mathbf{w}_k^t \mathbf{w}_k - 2\mathbf{x}^t \mathbf{w}_k) / 2\sigma^2} \\ &= e^{(\text{net}_k - 1) / \sigma^2} \end{aligned}$$

$$\text{net}_k = \mathbf{w}_k^t \mathbf{x}$$

Simplified form due to normalization

$$\mathbf{x}^t \mathbf{x} = \mathbf{w}_k^t \mathbf{w}_k = 1$$

# PNN Training Algorithm

## ■ Algorithm 1. (PNN training)

```
1 begin initialize  $j \leftarrow 0, n, a_{ji} \leftarrow 0$  for  $j = 1, \dots, n; i = 1, \dots$   
2   do  $j \leftarrow j + 1$  Over all samples j  
3      $x_{jk} \leftarrow x_{jk} / \left( \sum_i^d x_{ji}^2 \right)^{1/2}$  (normalize)  
4   Over all features k  $w_{jk} \leftarrow x_{jk}$  (train)  
5     if  $\mathbf{x} \in \omega_i$  then  $a_{ji} \leftarrow 1$   
6   until  $j = n$   
7 end
```

Normalize  
Feature values

Set weights equal  
to feature values

Only connections  
Corresponding to class  
labels are set to 1

The  $a_{ji}$  correspond to making connections between labeled samples and the corresponding classes

# PNN Classification Algorithm

## ■ Algorithm 2. (PNN Classification)

```
1 begin initialize  $k \leftarrow 0$ ,  $\mathbf{x} \leftarrow$  test pattern
2       do  $k \leftarrow k + 1$ 
3          $net_k \leftarrow \mathbf{w}_k^t \mathbf{x}$ 
4         if  $a_{ki} = 1$  then  $g_i \leftarrow g_i + \exp[(net_k - 1)/\sigma^2]$ 
5       until  $k = n$ 
6 return  $class \leftarrow \arg \max_i \bar{g}_i(\mathbf{x})$ 
7 end
```

Made possible  
by appropriate  
choice of  
activation  
function

Output is the sum of the activation functions corresponding to the labeled samples of that class

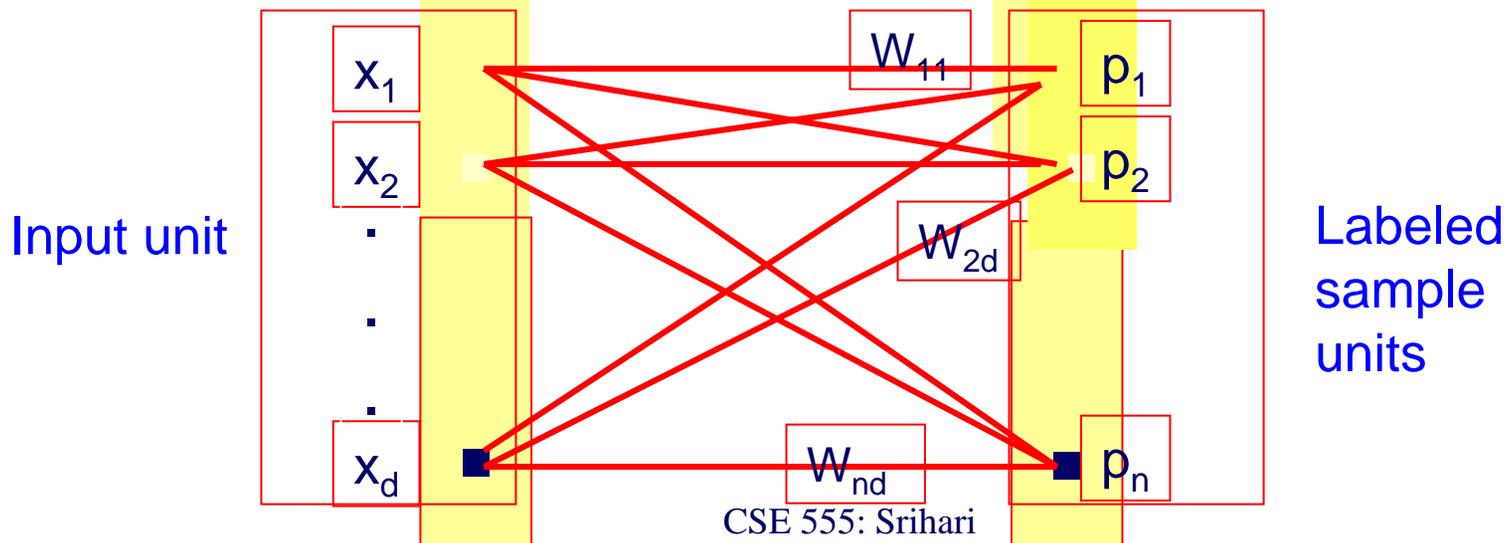
# Parzen/PNN Classifier Stage 1

- Inputs:  $d$  normalized features
- Weights correspond to feature values of  $n$  labeled samples

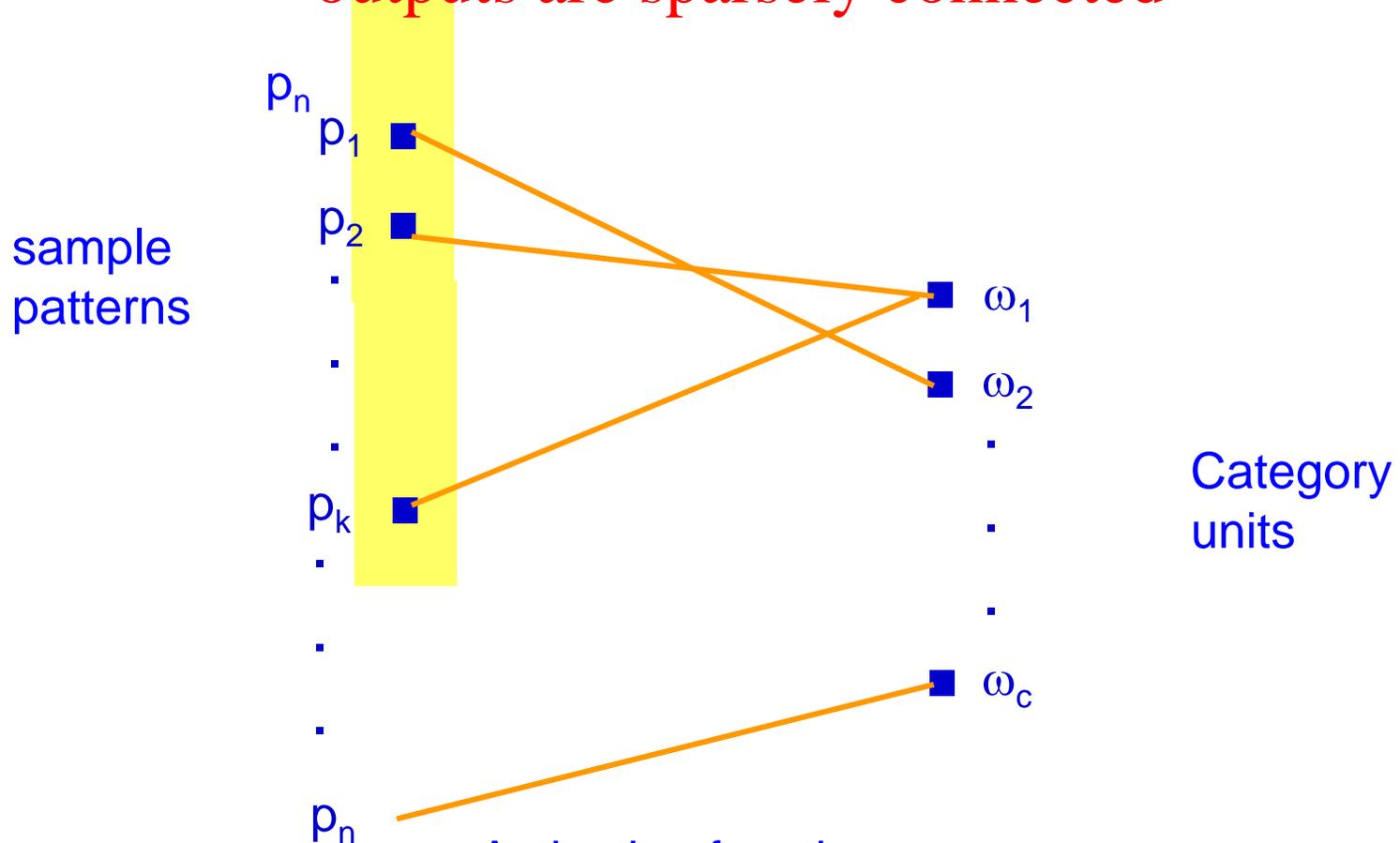
$$w_{ij} = x_{ij}, i=1, \dots, n; j=1, \dots, d$$

- Input units are fully connected to  $n$  sample units

$$\text{net}_k = \mathbf{w}_k^t \mathbf{x}; k = 1, \dots, n$$



# PNN Stage 2: outputs are sparsely connected



Activation functions  
(nonlinear functions)

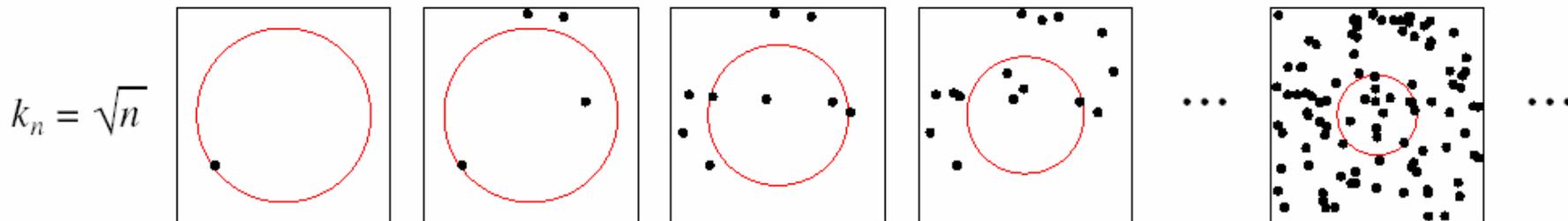
$$e^{(\text{net}_k - 1) / \sigma^2}$$

ri

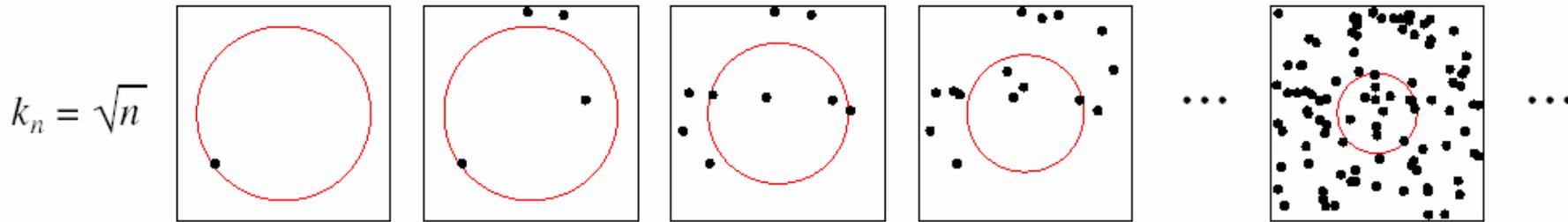
# $k_n$ -Nearest neighbor estimation

Parzen disadvantage: unknown “best” window function

- Let cell volume be a function of the training data
  - Center a cell about  $x$  and grow it until it captures  $k_n$  samples  
 *$k_n$  is a function of  $n$*
  - Or  $k_n$  nearest-neighbors of  $x$



# Choice of $k_n$



- If density is high near  $x$ 
  - cell will be small -- which provides good resolution
- If density is low
  - cell will grow large and stop until higher density regions are reached

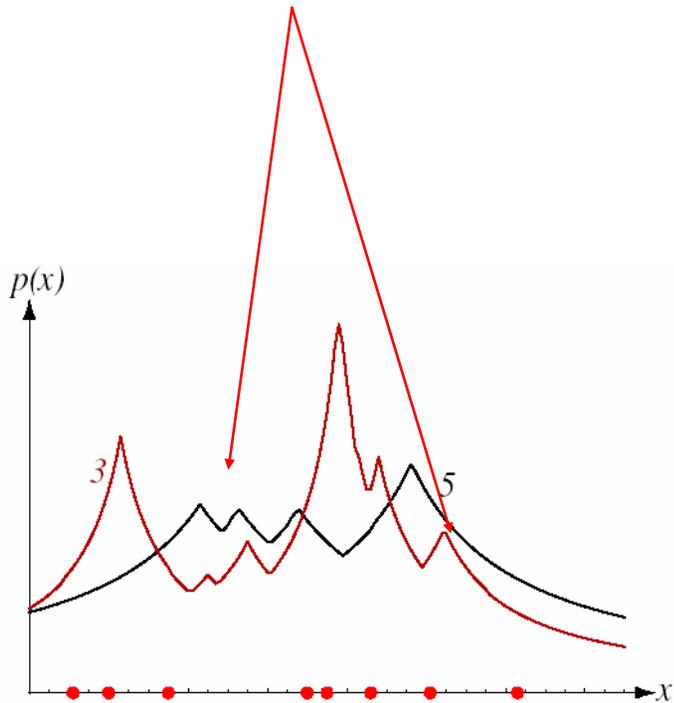
- Choose

$$p_n(\mathbf{x}) = \frac{k_n/n}{V_n}$$

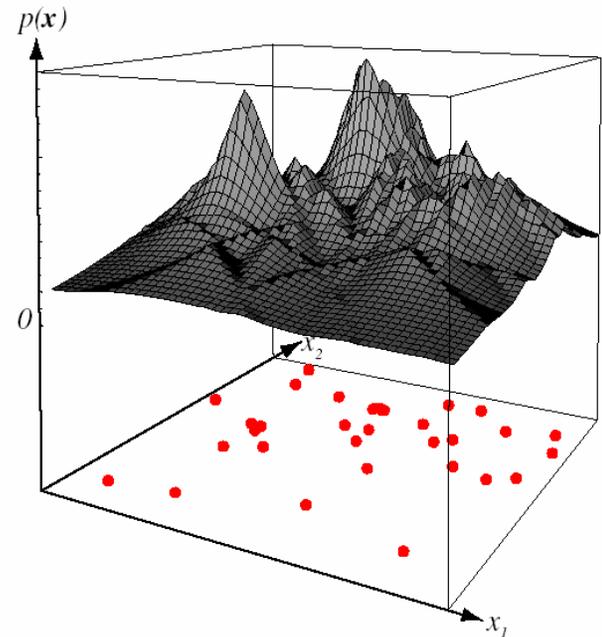
- Requirement for convergence:  $\lim k_n = \text{infinity}$  and  $\lim k_n/n = 0$
- We can obtain a family of estimates by setting  $k_n = k_1 / \sqrt{n}$  and choosing different values for  $k_1$

# Examples of estimated pdfs using k nearest neighbor estimate

Discontinuities in the estimated densities take place away from the points



Two dimensions:  $k = 5$

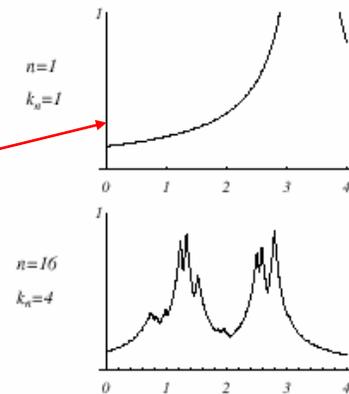


# Comparison of Parzen and k nn estimates in one-dimension

For k-nearest neighbor estimation using  $k_n = \sqrt{n} = 1$  the estimate becomes

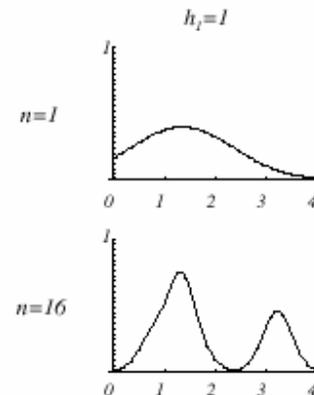
$$\begin{aligned} p_n(x) &= k_n / n \cdot V_n \\ &= 1 / V_1 \\ &= 1 / 2|x-x_1| \end{aligned}$$

Clearly a poor estimate



k-nn estimate

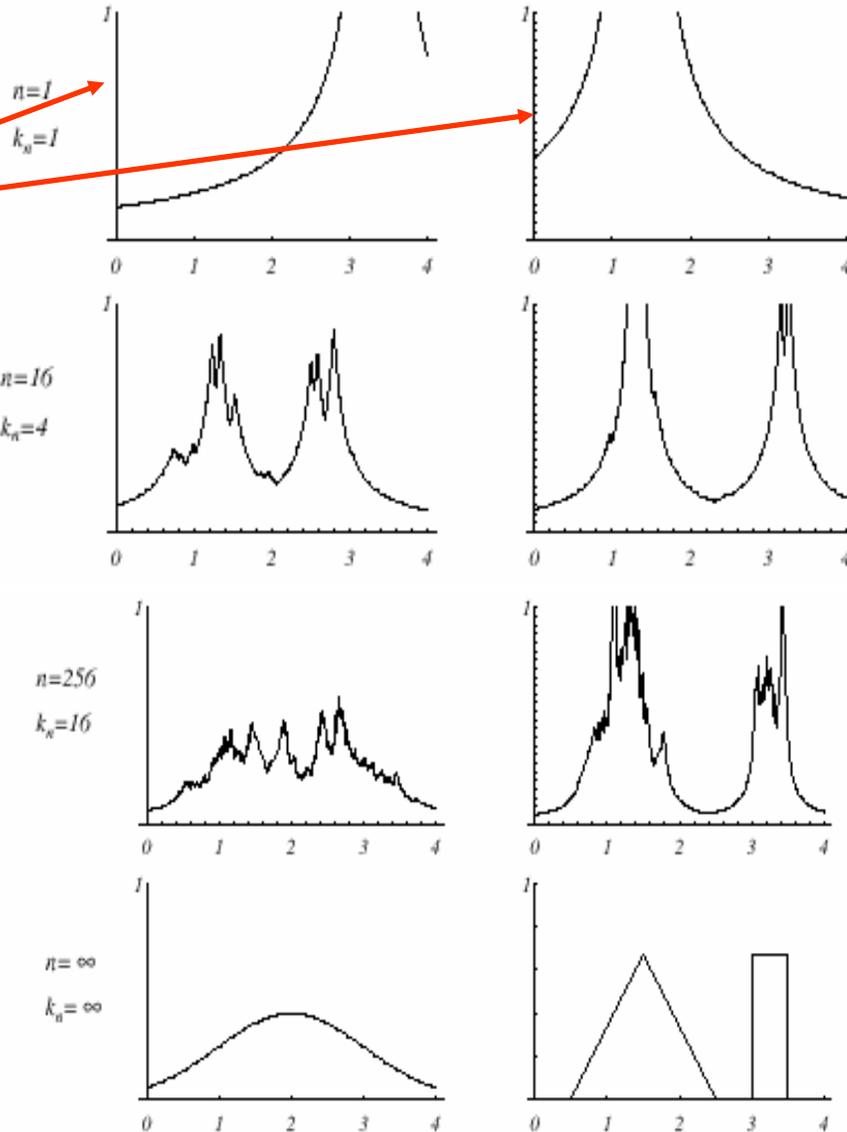
Whereas Parzen yields:



Parzen estimate

# Examples of $k_n$ -Nearest-Neighbor Estimation

$$p_n(x) = \frac{1}{2|x - x_1|}$$



Spiky estimates  
when there are  
few samples

# Estimation of a-posteriori probabilities

Product of class-conditional and a posteriori is represented among samples

$$p_n(x, \omega_i) = \frac{k_i / n}{V}$$

$$P_n(\omega_i / x) = \frac{p_n(x, \omega_i)}{\sum_{j=1}^c p_n(x, \omega_j)} = \frac{k_i}{k}$$

- $k_i/k$  is the fraction of the samples within the cell that are labeled  $\omega_i$
- For minimum error rate, the most frequently represented category within the cell is selected
- If  $k$  is large and the cell sufficiently small, the performance will approach the best possible
- Size of the cell can be chosen using either Parzen or  $k_n$  nearest neighbor