

Density Estimation

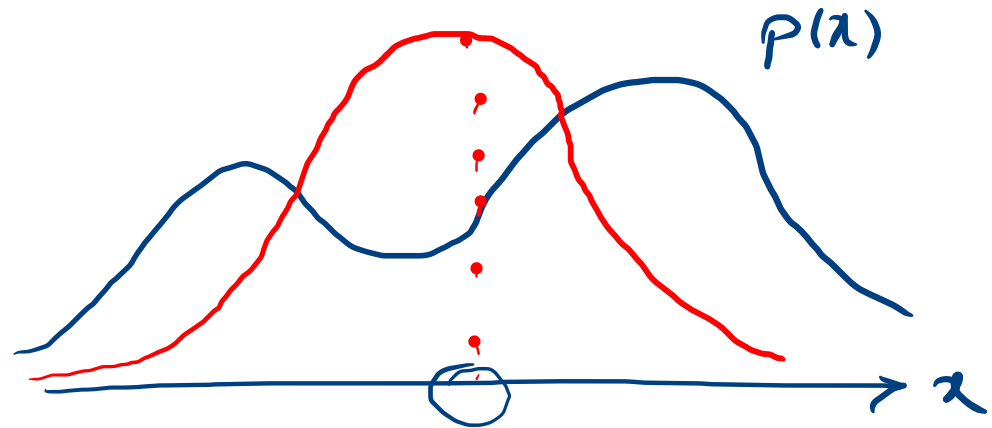
Parametric

$$x \sim P(x|\theta)$$

Estimate θ

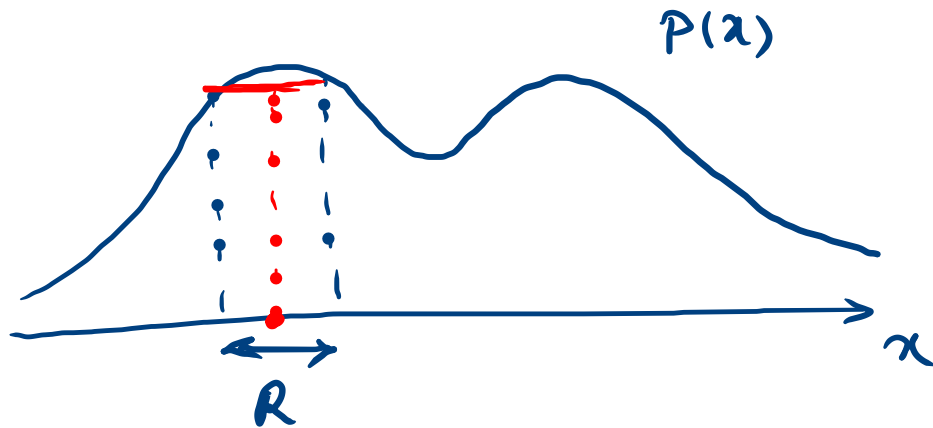
non-parametric

$P(x)$



$$\hat{\mu}_{ML} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$D = \{x_1, \dots, x_n\} \quad p(x) = ?$$



$$P(x \in R) = \int_R p(x) dx \approx \frac{K}{n}$$



$$V p(x) = \frac{K}{n} \Rightarrow$$

$$p(x) = \frac{K}{nV}$$

المعروف

$$p_n(x) = \frac{k_n}{nV_n}$$

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

$$\textcircled{1} \quad \lim_{n \rightarrow \infty} V_n \rightarrow 0$$

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

$$\textcircled{2} \quad \lim_{n \rightarrow \infty} k_n \rightarrow \infty$$

$$k_n = \sqrt{n}$$

$$\textcircled{3} \quad \lim_{n \rightarrow \infty} \frac{k_n}{n} \rightarrow 0$$

$$\text{Or} \quad \lim_{n \rightarrow \infty} \underline{nV_n} \rightarrow \infty$$

✓

Conditions for convergence



- Let n be the **number of samples** used, R_n be the region used with n samples, V_n be the volume of R_n , k_n be the number of samples falling in R_n , and the estimate for $p(x)$ be

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

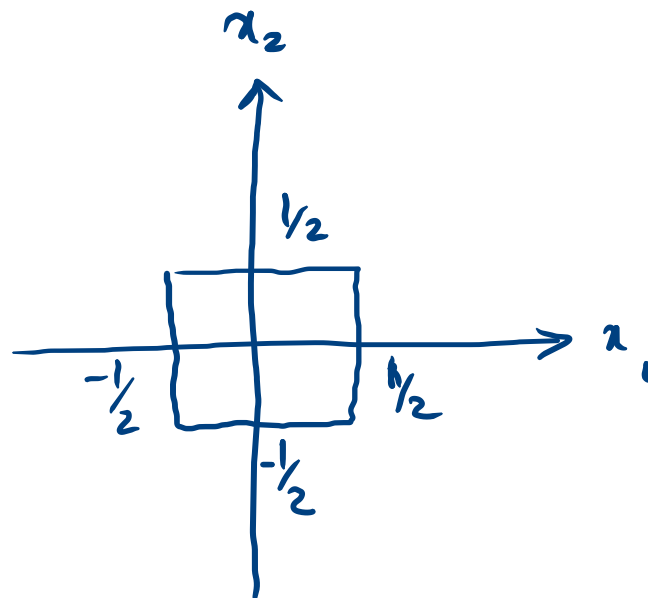
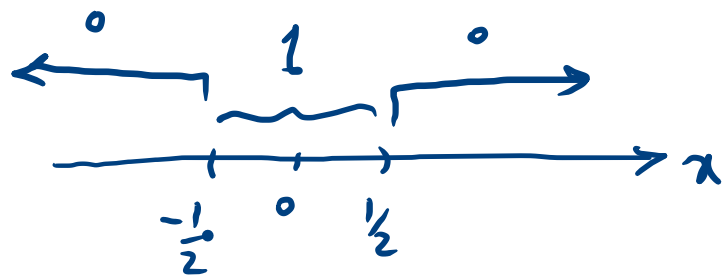
- If $p_n(x)$ is to **converge** to $p(x)$, three **conditions** are required:

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

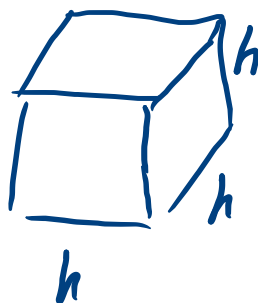
$$\lim_{n \rightarrow \infty} k_n = \infty$$

$$\lim_{n \rightarrow \infty} \frac{k_n}{n} = 0.$$

$$\phi(x) = \begin{cases} 1 & |x_i| \leq \frac{1}{2} \\ 0 & \text{o.w.} \end{cases} \quad i=1, \dots, d \quad x \in \mathbb{R}^d$$



$$\phi\left(\frac{x-x_0}{h}\right)$$



$$\int \phi\left(\frac{x-x_0}{h}\right) dx = h^d$$

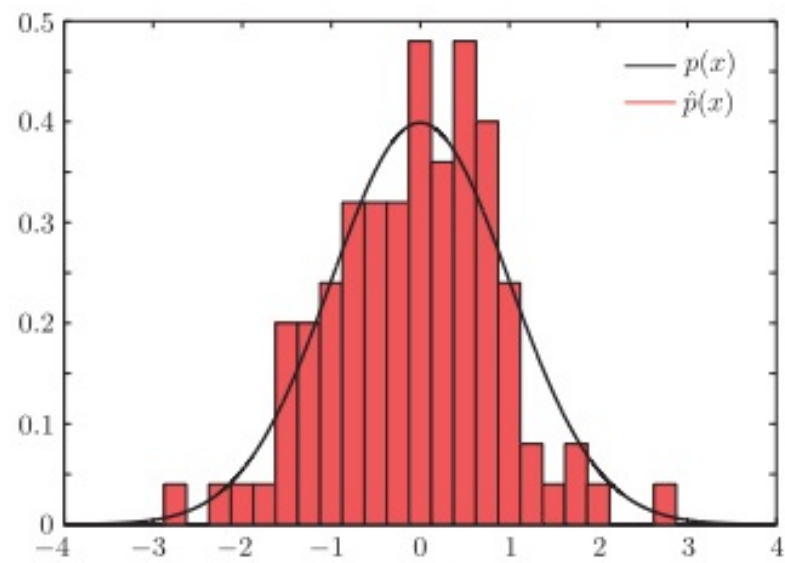
$$p_n(x) = \frac{K_n}{nV_n} = \frac{1}{nV_n} \underbrace{\sum_{i=1}^n \phi\left(\frac{x-x_i}{h_n}\right)}_{K_n}$$

$$V_n = h_n^d$$

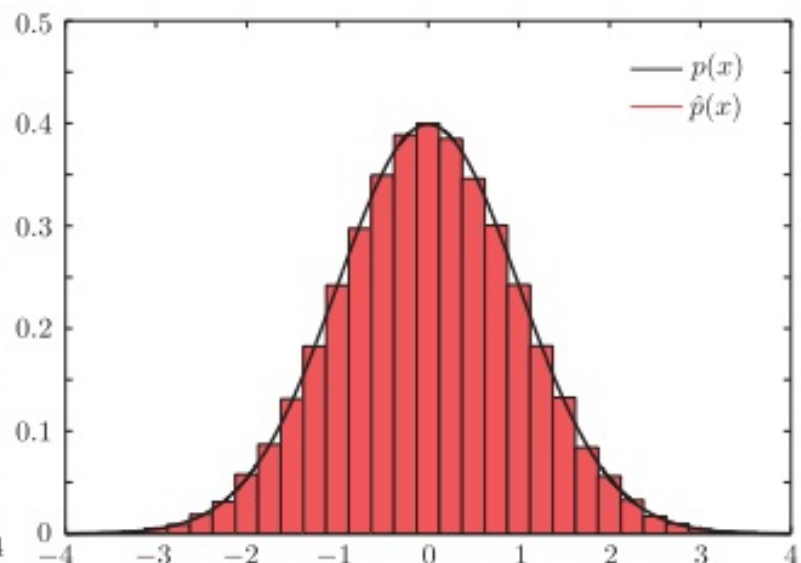
$$p(x)$$

$$\left\{ \begin{array}{l} \int \phi(x) dx = 1 \\ \phi(x) \geq 0 \end{array} \right.$$

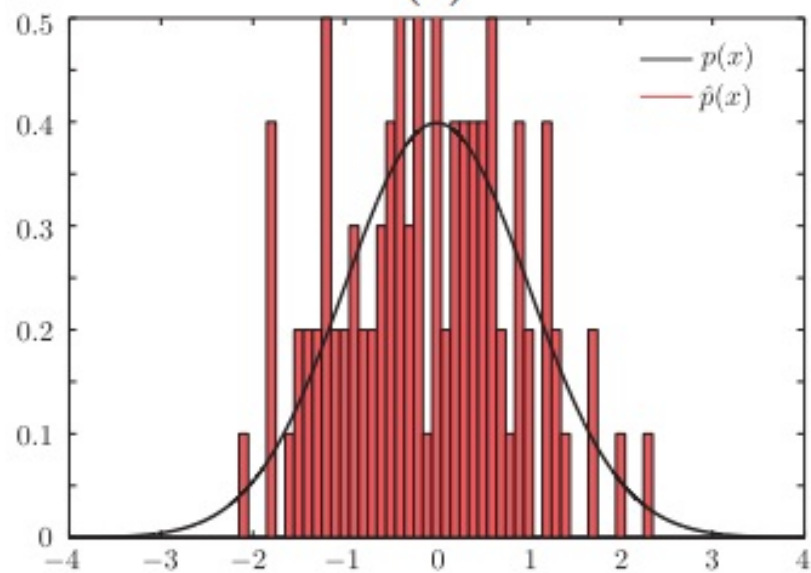
$$\boxed{\frac{1}{V_n} \phi\left(\frac{x-x_0}{h}\right)}$$



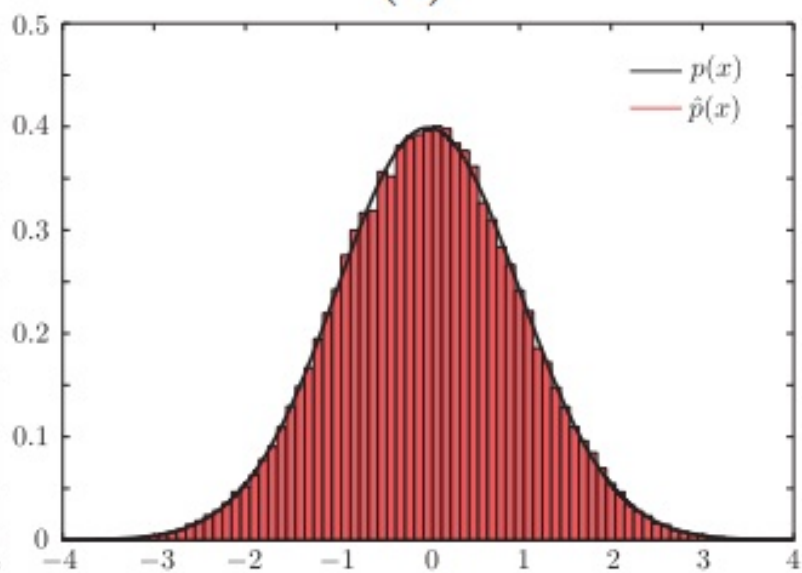
(a)



(b)



(c)



(d)

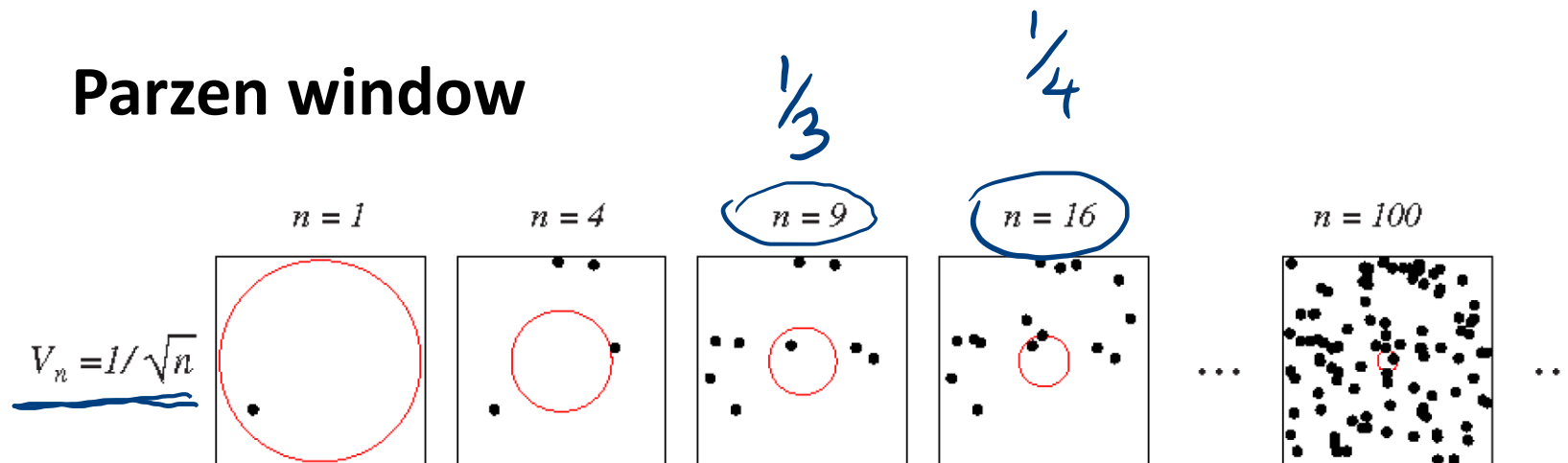
10^5

10^5

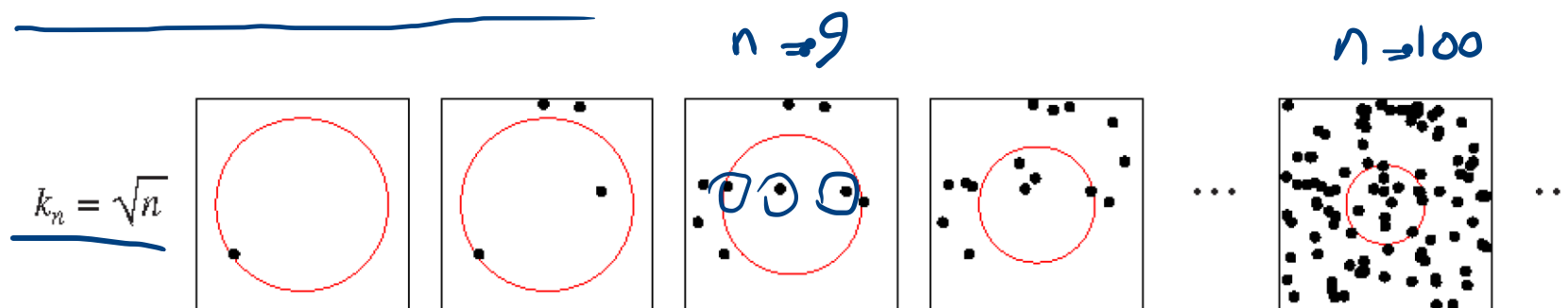
Two methods for estimating the density at a point \mathbf{x} (at the center of each square)



- Parzen window



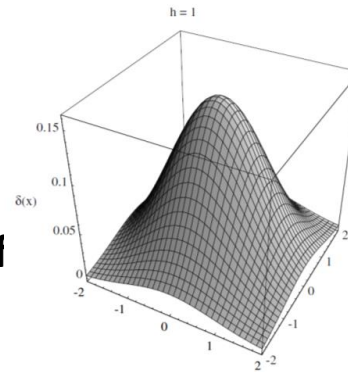
- k-nearest neighbor



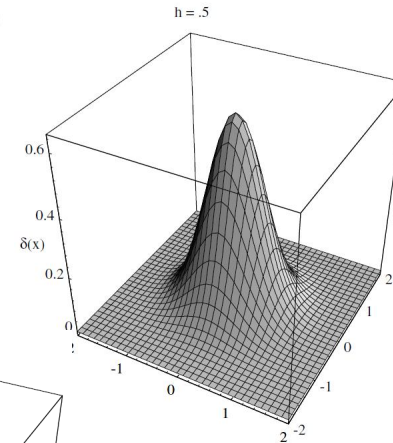
The role of h_n



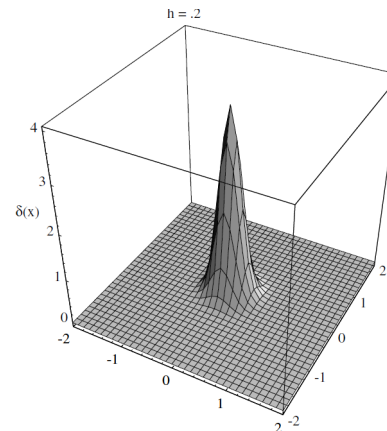
- If h_n is very **large**, $p_n(x)$ is the superposition of n **broad functions**, and is a smooth “**out-of-focus**” estimate of $p(x)$.
- If h_n is very **small**, $p_n(x)$ is the superposition of n **sharp pulses centered at the samples**, and is a “**noisy**” estimate of $p(x)$.
- As h_n approaches zero, $d_n(x - x_i)$ approaches a **Dirac delta** function centered at x_i , and $p_n(x)$ is a superposition of delta functions.



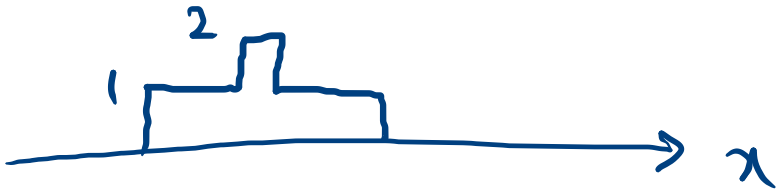
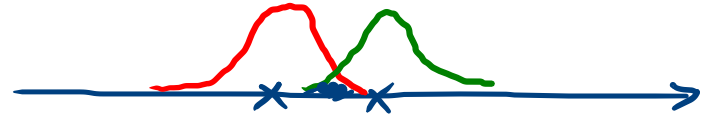
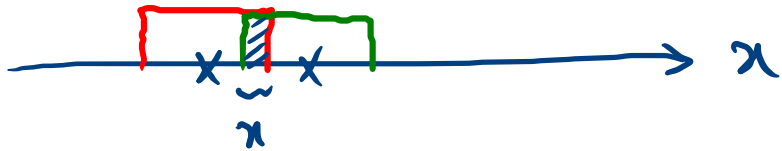
$$\phi(x) = \mathcal{N}(0,1)$$



$$\phi\left(\frac{x-x_i}{h}\right)$$



$f(u)$ as a function of h_n

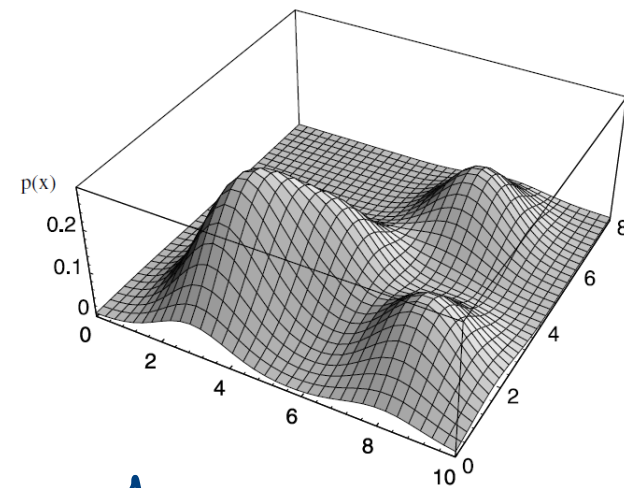


$$\textcircled{P_n(\lambda)} = \frac{1}{n v_n} \sum \overbrace{\Phi\left(\frac{\lambda - \lambda_i}{h}\right)}$$

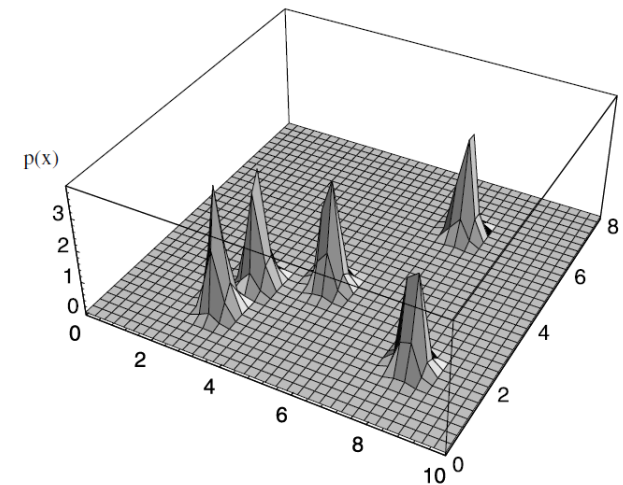
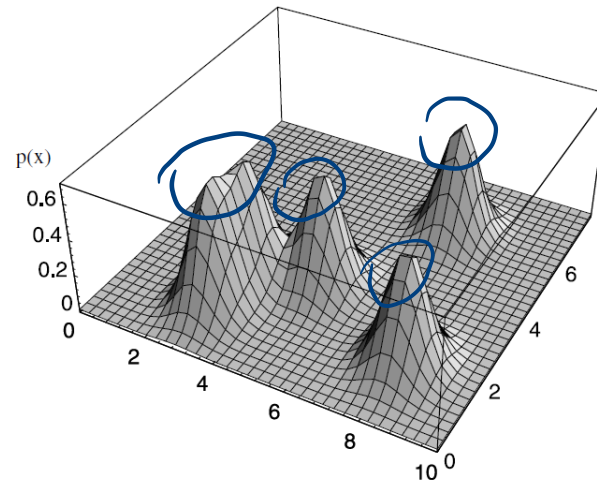
$p_n(x)$ as a function of $\underline{h_n}$



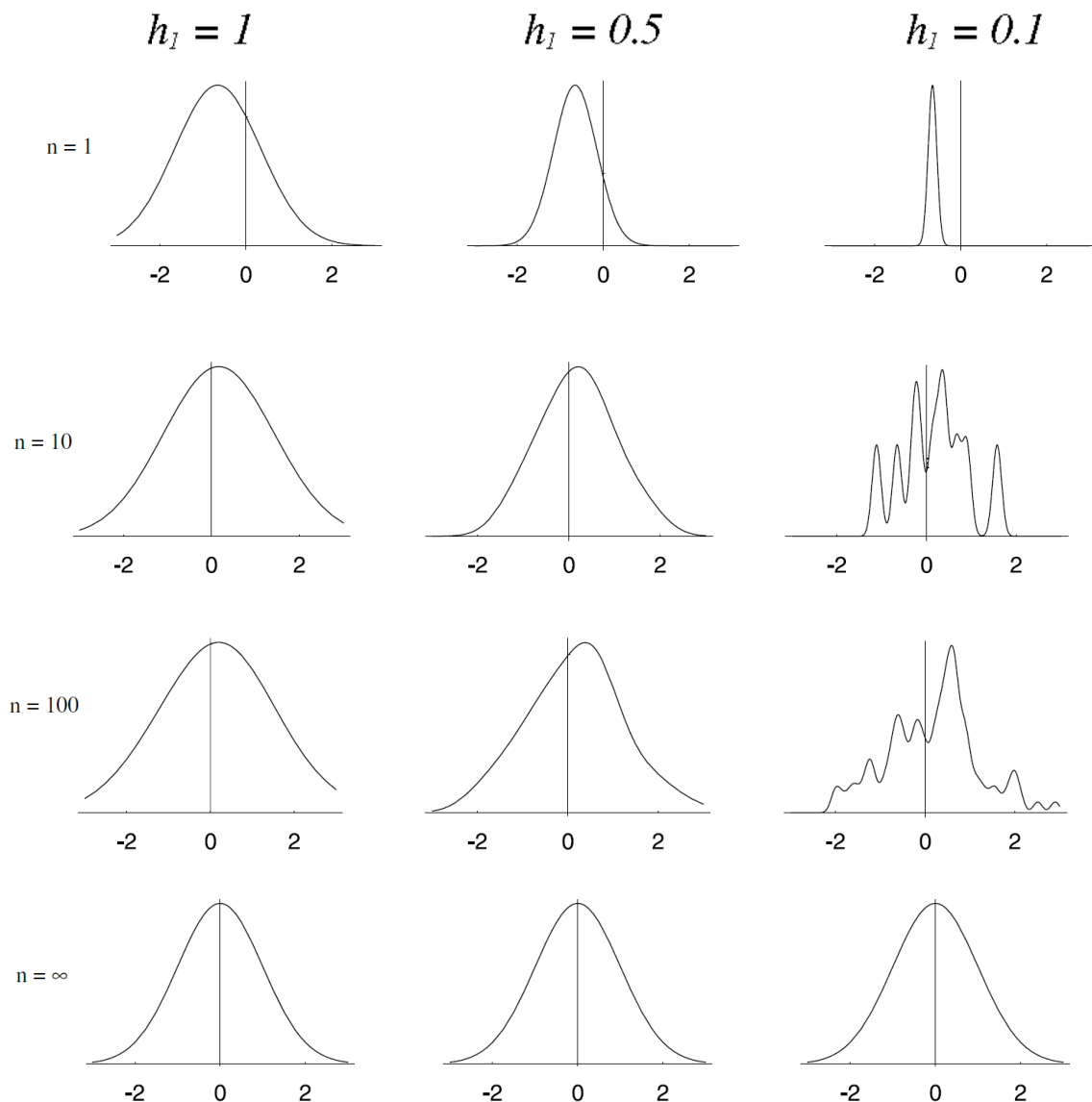
Parzen window density estimates based on the same set of five samples using the window functions in the previous figure



h_n large



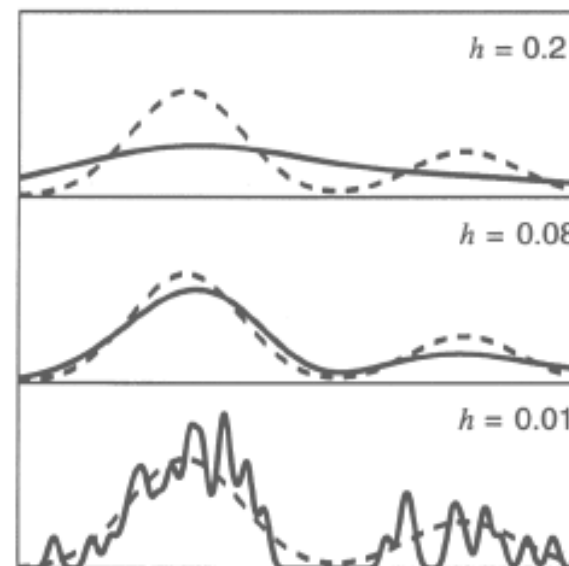
both $p(\mathbf{x})$ and $f(u)$ are Gaussian



$p_n(x)$ is an average of normal densities centered at the samples

$$\varphi(u) = \frac{1}{\sqrt{2\pi}} e^{-u^2/2}.$$

$$h_n = h_1 / \sqrt{n},$$

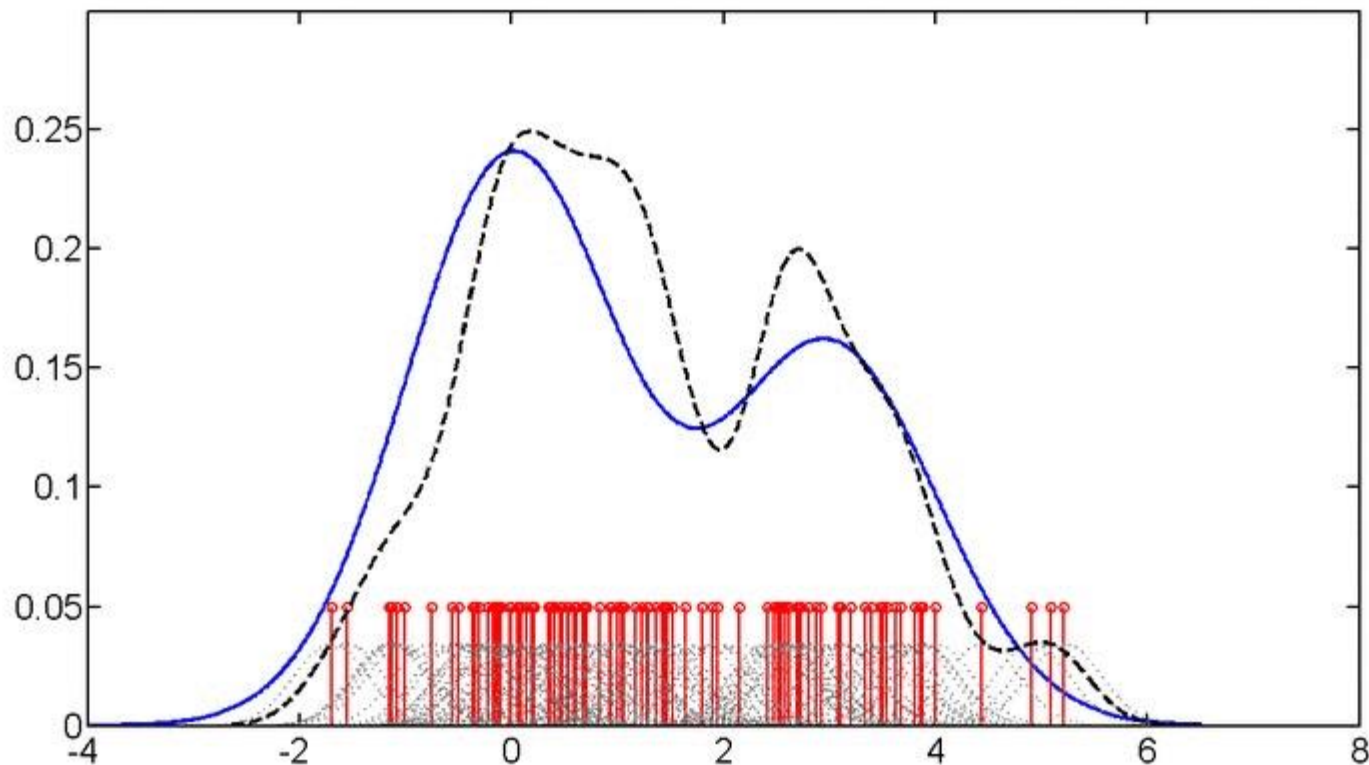


Parzen window

Blue curve: true density is mixture of two Gaussians centered around 0 and 3

In each frame, 100 samples are generated from the distribution, shown in **red**

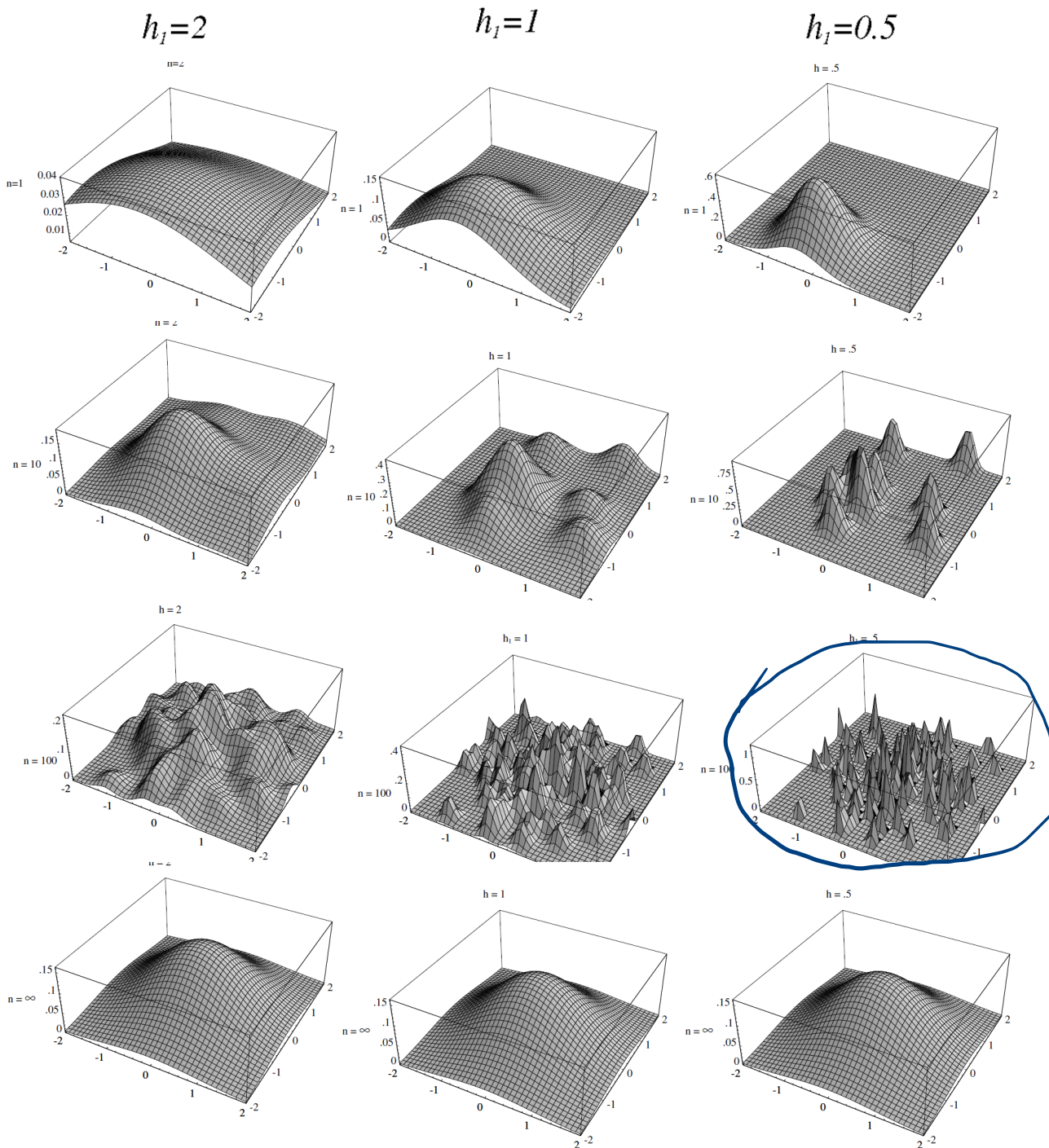
Dashed black curve: averaging the Gaussians yields the density estimate



Parzen-window estimates of a bivariate normal

$$\varphi(\mathbf{u}) = N(\mathbf{0}, \mathbf{I})$$

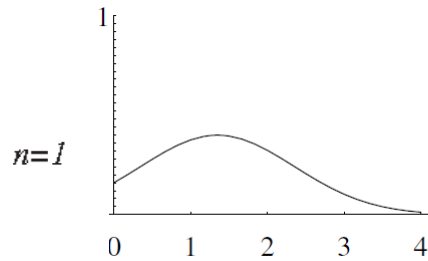
$$h_n = h_1 / \sqrt{n}.$$



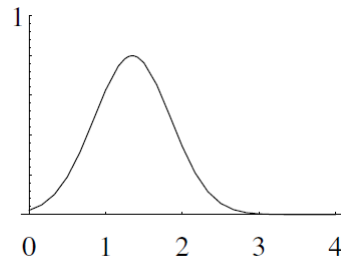
$p(x)$ consists of a uniform and triangular density and $f(x)$ is Gaussian



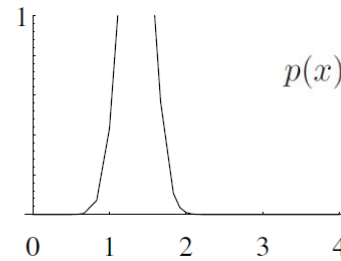
$h_1 = 1$



$h_1 = 0.5$

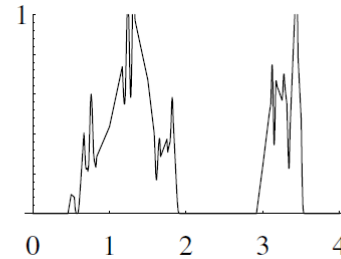
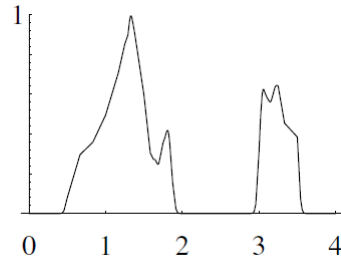
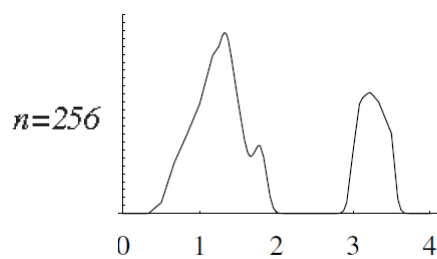
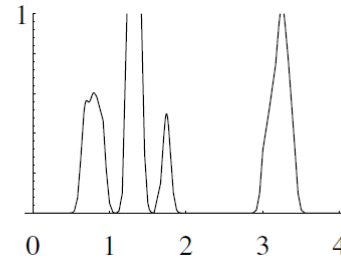
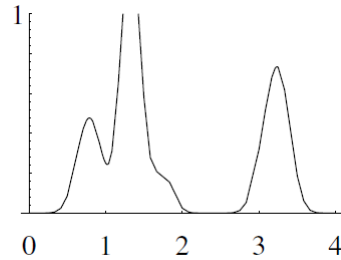
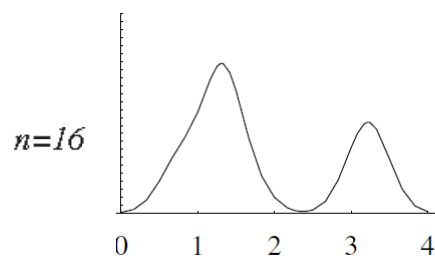


$h_1 = 0.1$

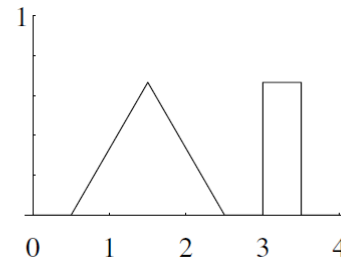
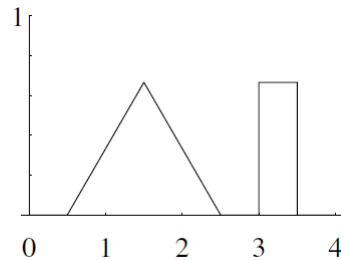
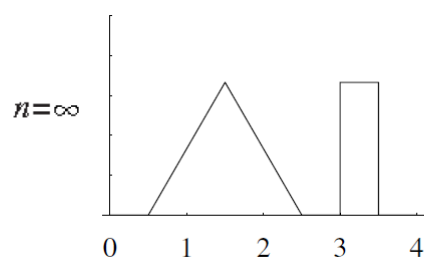


$$p(x) = \begin{cases} 1 & -2.5 < x < -2 \\ 1/4 & 0 < x < 2 \\ 0 & \text{otherwise.} \end{cases}$$

$$\varphi(u) = \frac{1}{\sqrt{2\pi}} e^{-u^2/2}$$



$$h_n = h_1 / \sqrt{n},$$





Classification using kernel-based density estimation (Bayesian decision rule)

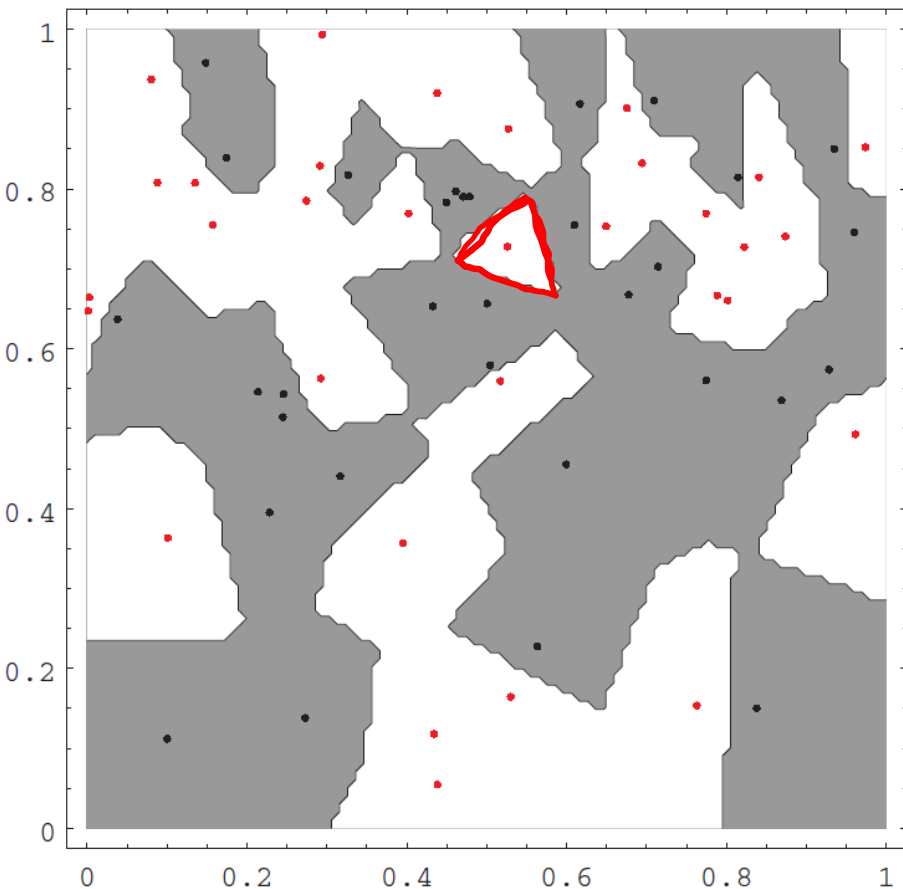
- Estimate density for **each class**.
- **Classify** a test point by computing the **posterior probabilities** and **picking the max**.
- The **decision regions** depend on the choice of the **kernel function** and h_n .
- The training error can be made **arbitrarily low** by making the window width **sufficiently small**.
- However, the **goal is to classify novel** patterns so the window width **cannot be made** too small.

$$\frac{p(y=1|x)}{p(y=2|x)} = \frac{p(x|y=1)p(y=1)}{p(x|y=2)p(y=2)}$$

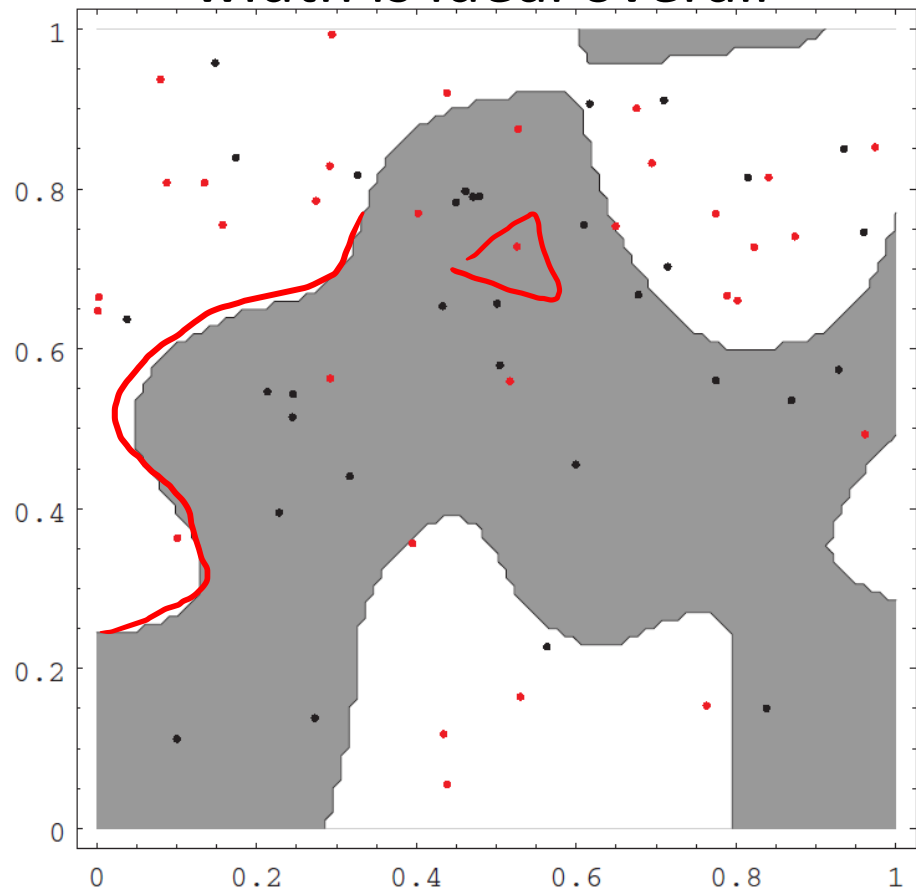
dimensional Parzen-window dichotomizer



- small h_n
- Boundaries are more complicated



- large h_n
- No single window width is ideal overall



Drawbacks of kernel-based methods



- Require a **large number** of samples.
- Require all the samples to **be stored**.
- Evaluation of the density could be **very slow** if the number of data points is large.
- Possible **solution**:
 - use **fewer kernels** and **adapt the positions** and widths in response to the data (e.g., mixtures of Gaussians!)

$$\lim_{n \rightarrow \infty} \underline{\underline{P_n(x)}} \rightarrow P(x)$$

$$x \sim P(x)$$

$$\text{Var}(x)$$

$$\begin{cases} E[P_n(x)] = P(x) \\ \text{Var}(\underline{\underline{P_n(x)}}) = 0 \end{cases}$$

$$\underline{E[P_n(x)]} = E\left[\frac{1}{n} \sum_{i=1}^n \frac{1}{v_n} \phi\left(\frac{x-x_i}{h_n}\right)\right] = \frac{1}{n} \sum_{i=1}^n E\left[\underbrace{\frac{1}{v_n} \phi\left(\frac{x-\overset{\text{R.V.}}{x_i}}{h_n}\right)}_{g(x_i)}\right]$$

$$= E\left[\frac{1}{v_n} \phi\left(\frac{x-x_i}{h_n}\right)\right] \stackrel{v_n \rightarrow 0}{=} \underbrace{\int \frac{1}{v_n} \phi\left(\frac{x-x_i}{h_n}\right) P(x_i) dx_i}_{\delta(x-x_i)} = \int \delta(x-x_i) P(x_i) dx_i = \underline{P(x)}$$

$$\text{var}(X_1 + X_2) = \text{var}(X_1) + \text{var}(X_2)$$

$$\text{var}(X) = E[X^2] - E^2[X]$$

$$\begin{aligned} \text{var}(P_n(x)) &= \text{var}\left(\frac{1}{n} \sum_{i=1}^n \frac{1}{v_n} \phi\left(\frac{x-x_i}{h_n}\right)\right) \\ &= \frac{1}{n^2} \sum_{i=1}^n \text{var}\left(\frac{1}{v_n} \phi\left(\frac{x-x_i}{h_n}\right)\right) \end{aligned}$$

$$= \frac{1}{n} \text{var}\left(\frac{1}{v_n} \phi\left(\frac{x-x_i}{h_n}\right)\right) = \frac{1}{n} E\left[\frac{1}{v_n^2} \phi^2\left(\frac{x-x_i}{h_n}\right)\right] - \frac{1}{n} E^2\left[\frac{1}{v_n} \phi\left(\frac{x-x_i}{h_n}\right)\right]$$

$\bar{P}_n(x)$

$$= \frac{1}{n} \int \underbrace{\frac{1}{v_n^2} \phi^2\left(\frac{x-x_i}{h_n}\right)}_{\Phi_{\max} \phi\left(\frac{x-x_i}{h_n}\right)} p(x_i) dx_i - \frac{1}{n} \underbrace{\bar{P}_n^2(x)}_{\geq 0}$$

$$\leq \frac{1}{n} \frac{\Phi_{\max}}{v_n} \int \frac{1}{v_n} \phi\left(\frac{x-x_i}{h_n}\right) p(x_i) dx_i$$

$$\lim_{n \rightarrow \infty} \text{var}(P_n(x)) \leq \frac{\Phi_{\max} P(x)}{n V_n} \rightarrow 0$$

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

$$\boxed{\Phi_{\max} < \infty}$$