# Non-parametric Density Estimation
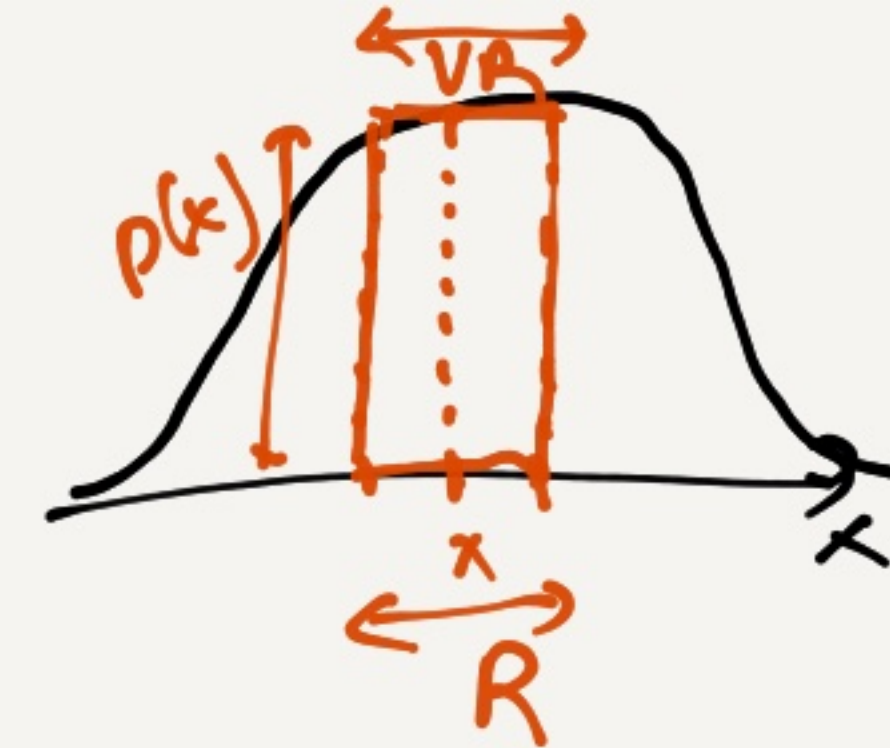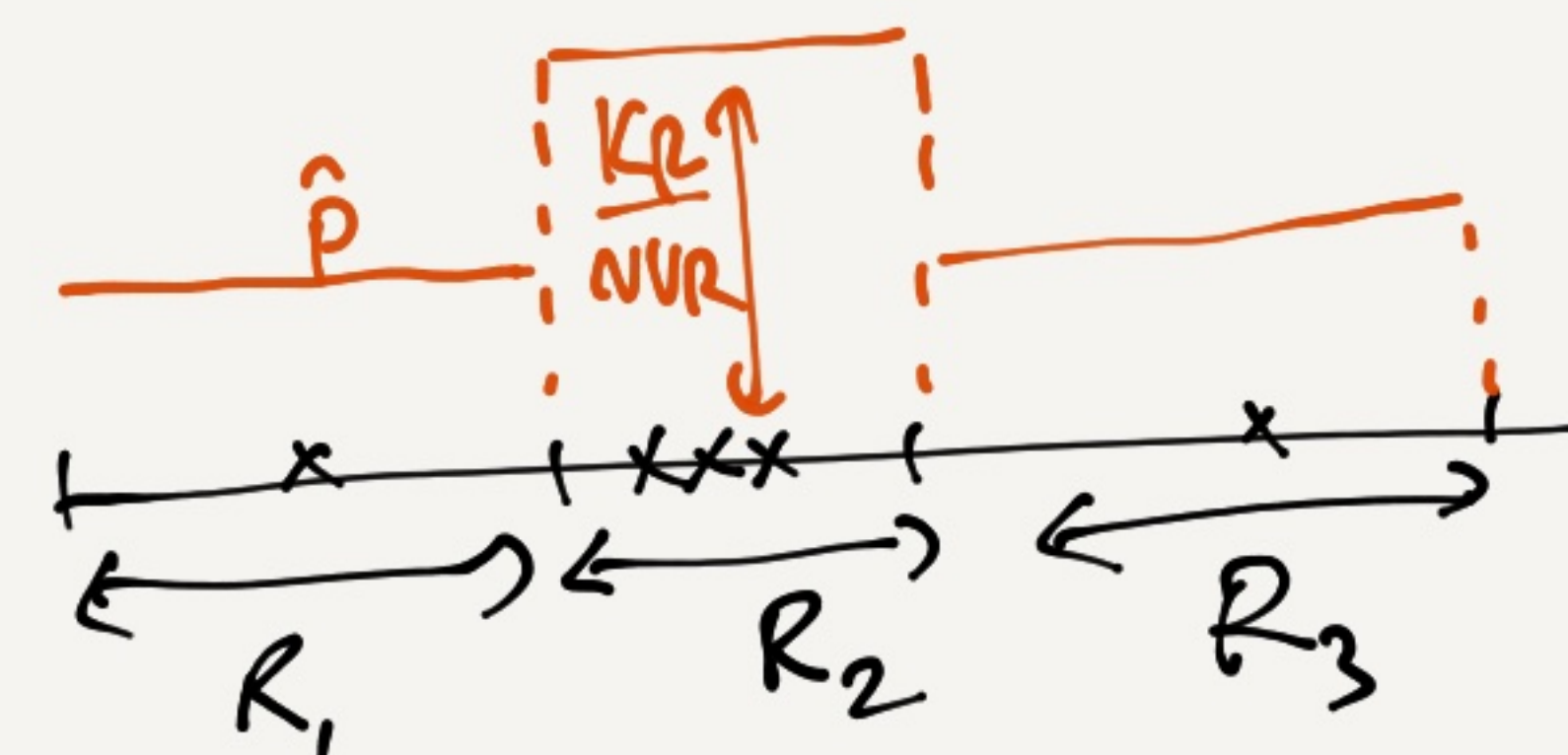
- So far we have looked at <u>parametric</u> models ⇒ make assumptions about the form of the density (Gaussian, Expo, GMM, etc)

- <u>Non parametric Estimation</u> - estimate $p(x)$ w/o assuming a form (has some parameters)

# Histogram

- Samples $\{x_1, ..., x_N\} = D$

- consider a region $R$
  - define $p = p(x \in R) = \int_R p(x) dx$
  - define $K_R$ = # of points in $D$ that are inside $R$.
  - ML estimate of $p$:
    $$\hat{p} = \frac{K_R}{N}$$
    (this is just a bin in a histogram)

$K_R = 3, \ \hat{p} = \frac{3}{5}$

---

- Assume $R$ is small enough, then
  $$\hat{p} \approx p(x) V_R \approx \int_R p(x) dx$$
  ↑ volume of $R$

- Solve for $p(x)$:
  $$\hat{p}(x) = \frac{\hat{p}}{V_R} = \frac{K_R}{N V_R}$$
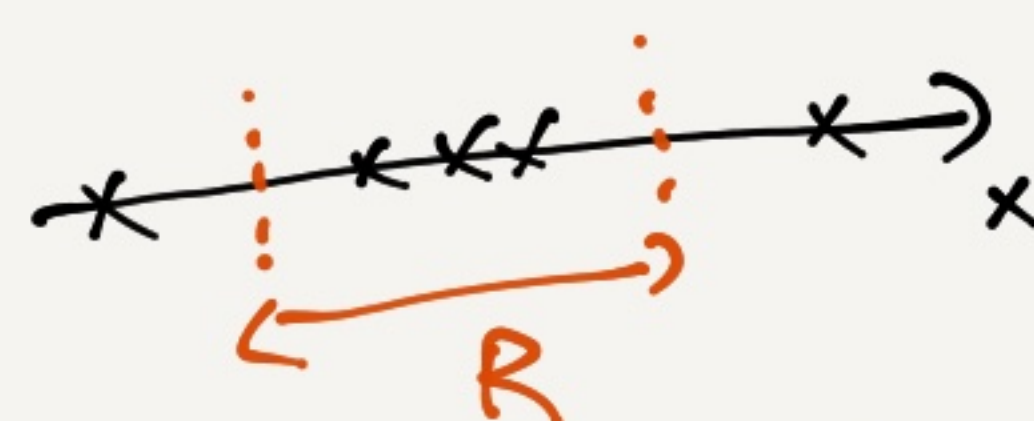  ← # points in $R$
  ← volume of $R$

We can extend this simple histogram:
<u>How to choose $R$?</u>

1) Keep $V_R$ fixed, & let $K_R$ vary. ✓✓
   ⇒ Kernel density estimator (KDE); Parzen windows

2) Keep $K_R$ fixed, let $V_R$ vary.
   ⇒ k-NN estimator (very bad) PRML
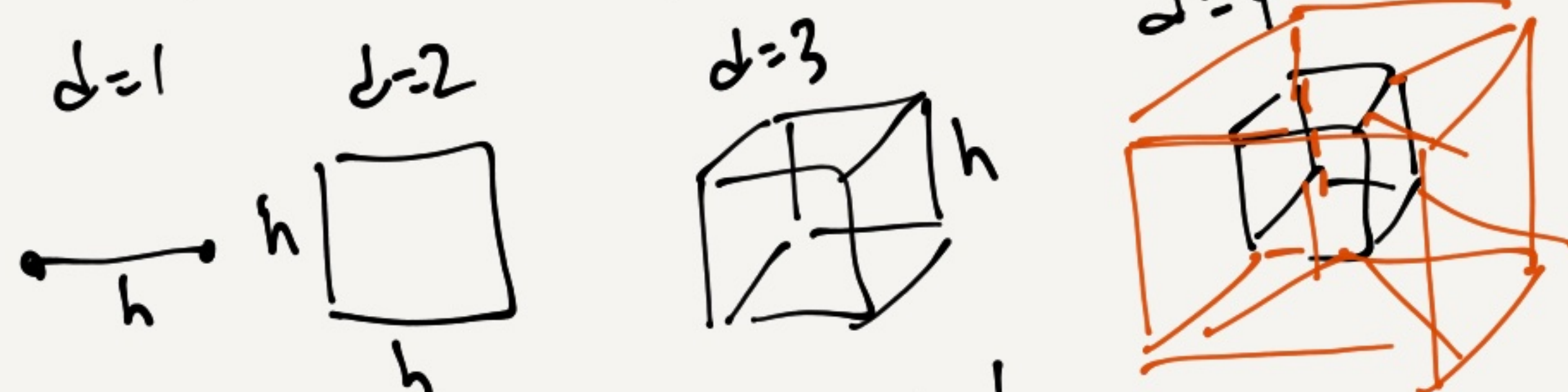
Assignment Project Exam Help
https://powcoder.com
Add WeChat powcoder

# Kernel Density Estimation (KDE)
- Parzen window
- Parzen Estimators.

- let $R$ be a $d$-dim <u>hypercube</u> w/ side length $h$.

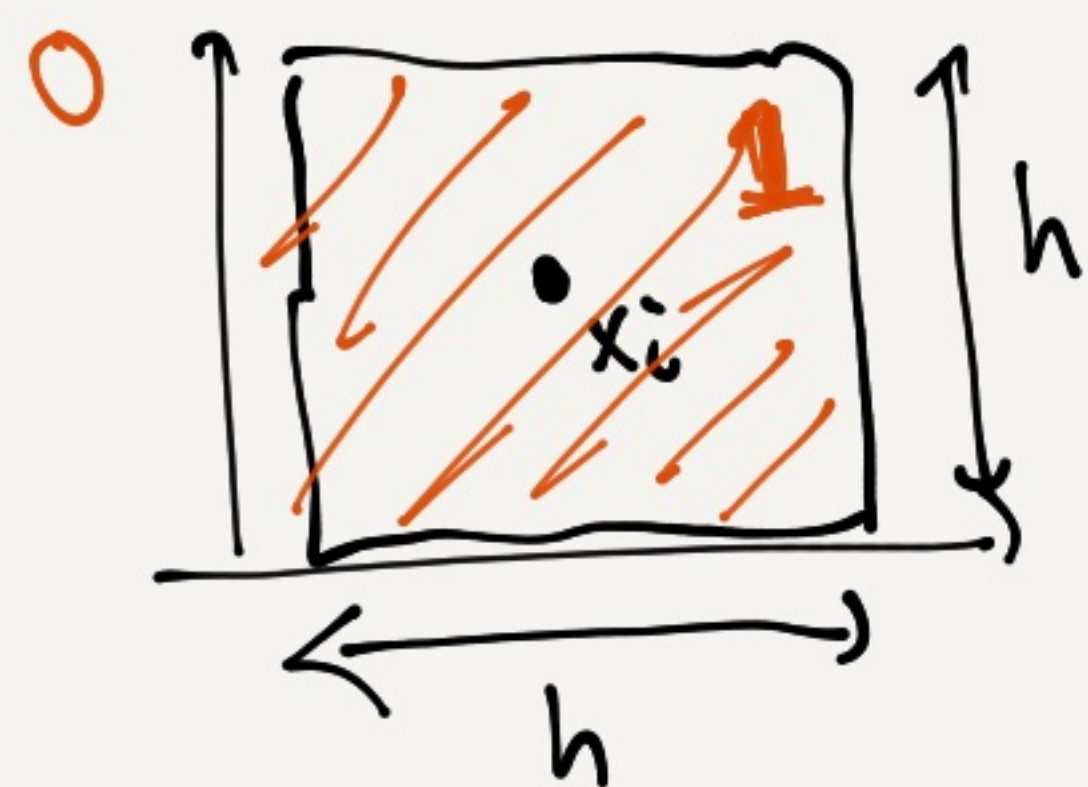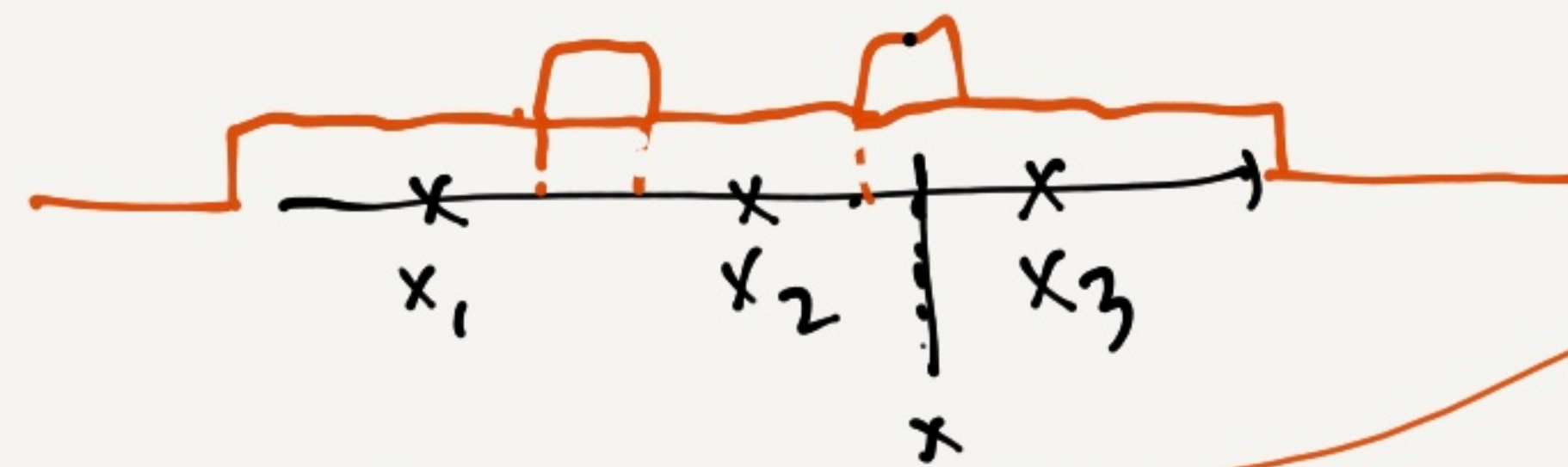$d=1$    $d=2$      $d=3$      $d=4$

volume of hypercube $= h^d$

- Introduce a window (kernel) (unit box)

$$K(x) = \begin{cases} 1, & |x_i| \leq \frac{1}{2}, \quad \forall i \in \{1, \cdots, d\} \\ 0, & \text{otherwise} \end{cases}$$

<u>Note</u>: $K\left(\dfrac{x-x_i}{h}\right) = \begin{cases} 1, & \text{if } x \text{ falls inside cube} \\ & \text{w/ side } h, \text{ centered at } x_i \\ 0, & \text{otherwise} \end{cases}$
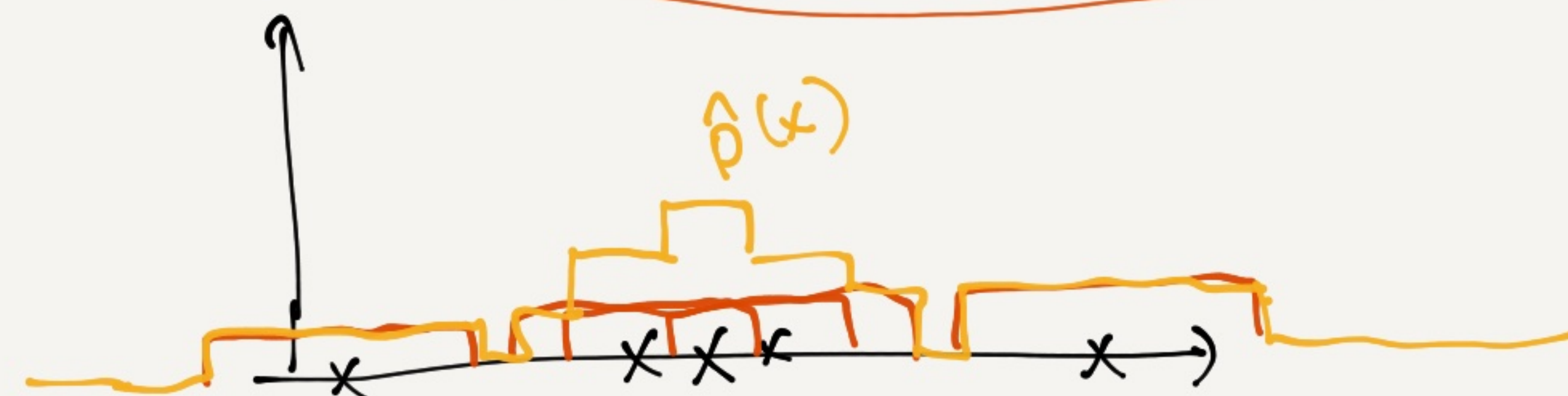
$\Rightarrow$ # of points near to $x$ $= K = \sum\limits_{i=1}^{N} K\left(\dfrac{x-x_i}{h}\right)$

<u>Thus</u>

$$\hat{p}(x) = \frac{1}{N}\frac{K_R}{V_R} = \frac{1}{Nh^d}\sum_{i=1}^{N} K\left(\frac{x-x_i}{h}\right)$$

estimate of $p(x)$ from $\{x_1, \cdots, x_N\}$

$\hat{p}(x)$

estimation using interpolation between samples $x_i$.

$\rightarrow$ each $x_i$ contributes to a local region.

# Kernel Functions

**Constraints:** $K(x) \geq 0$ — $\left. \begin{array}{c} \\ \\ \end{array} \right\}$ it must be a valid pdf.

$\int K(x)\,dx = 1$

## Example:

**uniform box:** $K(x) = \begin{cases} 1, & |x_i| \leq \frac{1}{2} \quad \forall_i = \{1,\ldots,d\} \\ 0, & \text{otherwise} \end{cases}$

**unit sphere:** $K(x) = \begin{cases} \frac{1}{c}, & \|x\|^2 \leq 1 \\ 0, & \text{otherwise} \end{cases}$

$c = $ volume of sphere.

**Gaussian:** $K(x) = \frac{1}{(2\pi)^{d/2}} e^{-\frac{1}{2}\|x\|^2}$

$$\hat{p}(x) = \frac{1}{Nh^d} \sum_i K\left(\frac{x - x_i}{h}\right)$$

$$= \frac{1}{Nh^d} \sum_i \frac{1}{(2\pi)^{d/2}} e^{-\frac{1}{2}\left\|\frac{x - x_i}{h}\right\|^2}$$

$\underbrace{|h^2 I|^{1/2}}$

$$= \frac{1}{N} \sum_i N(x \mid x_i, h^2 I)$$

$\underbrace{\quad}_{\pi_i = \frac{1}{N}}$

Gaussian component
mean $= x_i$
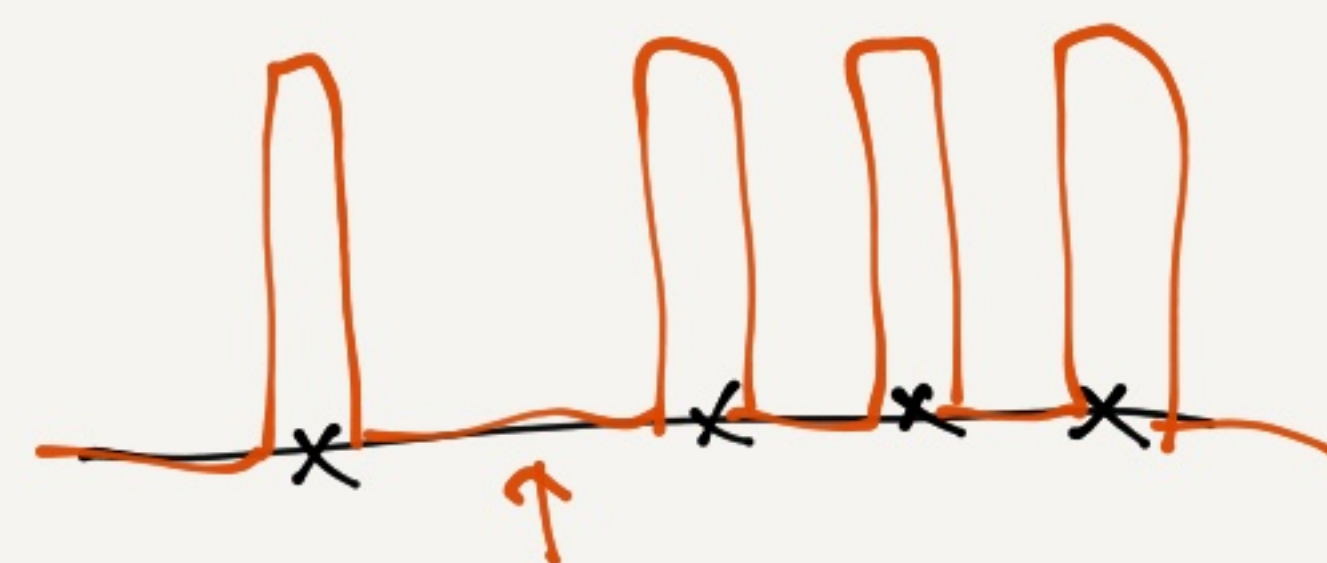cov $= h^2 I$

$\underbrace{\qquad\qquad}$ GMM w/ N components

# Bandwidth Parameter

$h$ controls the size of the region.
( covariance of the Gaussian )

Intuitively,

$h$ small:



might be noisy if not enough samples.

$h$ large:

blurry estimate if too many points.

$h$ is important and controls the quality of $\hat{p}(x)$.
Is there an optimal setting to recover the true $p(x)$?

# Convergence Analysis

☆ Will $\hat{p}(x)$ converge to the true pdf $p(x)$?

$\hat{p}(x)$ depends on samples $\{x_i\}$, which are r.v. $\to$ bias / variance.

- We say $\hat{p}(x)$ converges to $p(x)$ if:

  1) $\lim\limits_{N\to\infty} E[\hat{p}(x)] = p(x)$

  2) $\lim\limits_{N\to\infty} var(\hat{p}(x)) = 0$

- Define: $\tilde{K}(x) = \dfrac{1}{h^d} K\left(\dfrac{x}{h}\right)$

  <span style="color:orange">↑ scale amplitude. ↑ scale width</span>

  Then: $\hat{p}(x) = \dfrac{1}{N} \sum\limits_{i=1}^{N} \tilde{K}(x - x_i)$

Mean: $E[\hat{p}(x)] = \ldots\ldots$    (Tutorial S.1)

$\qquad = \int p(\mu)\, \tilde{K}(x-\mu)\, d\mu \quad \leftarrow$ defn of convolution

$\qquad = p(x) * \tilde{K}(x) \quad \leftarrow$ conv. of $p(x)$ with the kernel $\tilde{K}(x)$.

e.g.   (waveform) $*$ (box) $\quad \tilde{K}(x)$.

blurry $p(x)$, where the blur comes from the kernel.

---

to have unbiased $\hat{p}(x)$, we want

$$E[\hat{p}(x)] = p(x) * \underline{\tilde{K}(x)} = p(x)$$

$$\Rightarrow \tilde{K}(x) = \delta(x) = \lim\limits_{h\to 0} \tilde{K}(x)$$

$\qquad\qquad$ Dirac delta

to be unbiased, we want $h = 0$.   or   $\tilde{K}(x) = \delta(x)$

---

## Variance        (Tut. Sol)

$var(\hat{p}(x)) = \ldots\ldots$

$$var(\hat{p}(x)) \leq \dfrac{1}{Nh^d}\left[\max\limits_{x} K(x)\right] E[\hat{p}(x)]$$

For small variance, we need:
- <span style="color:orange">$h$</span> to be large.
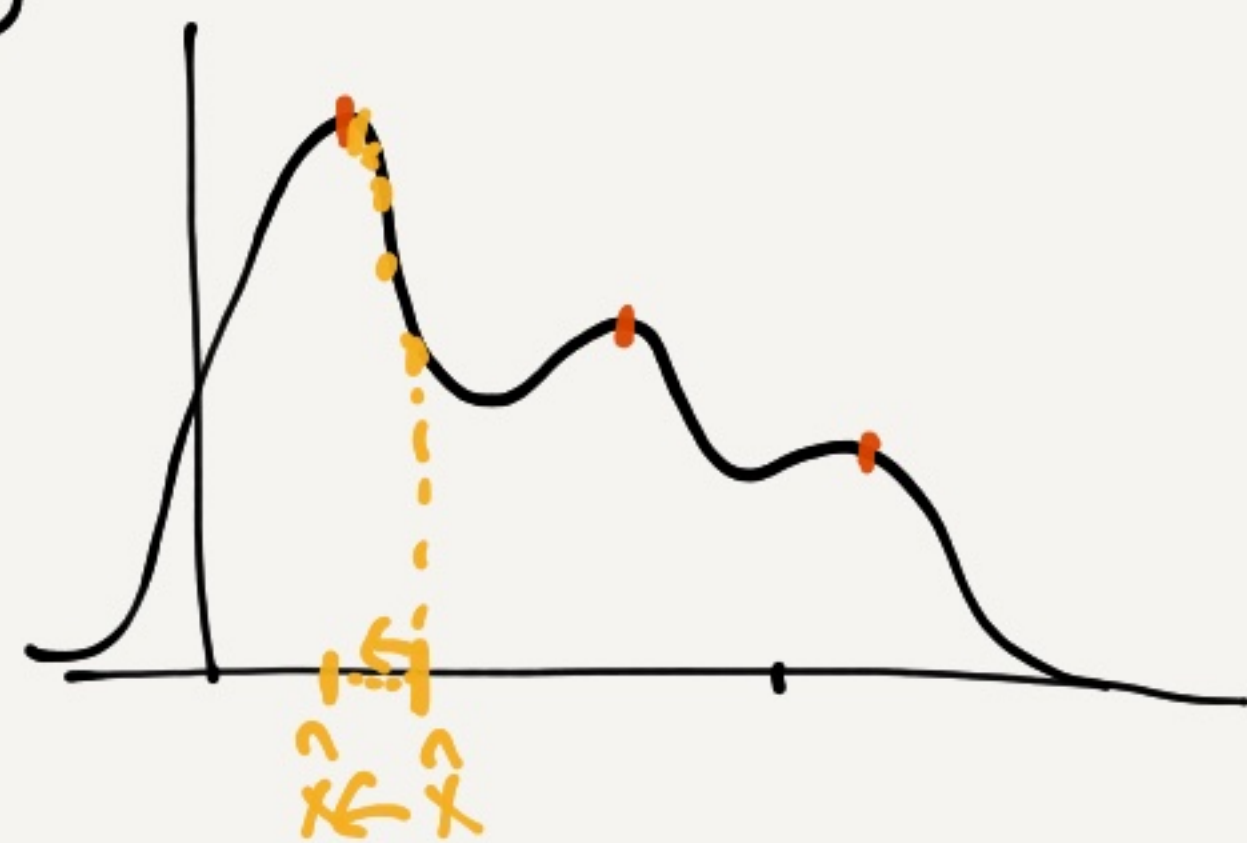
or
- $N$ to be large.

---

☆ $h$ controls the tradeoff btwn bias & variance:

$$\begin{cases} h \to 0 \Rightarrow \text{bias} = 0, \ var = \infty \\ h \to \infty \Rightarrow \text{bias} > 0, \ var = 0 \end{cases}$$

No theoretical optimal choice $\Rightarrow$ choose $h$ based on problem.

# Mean-Shift Algorithm (Comaniciu & Meer)

- Find the modes of $\hat{p}(x)$



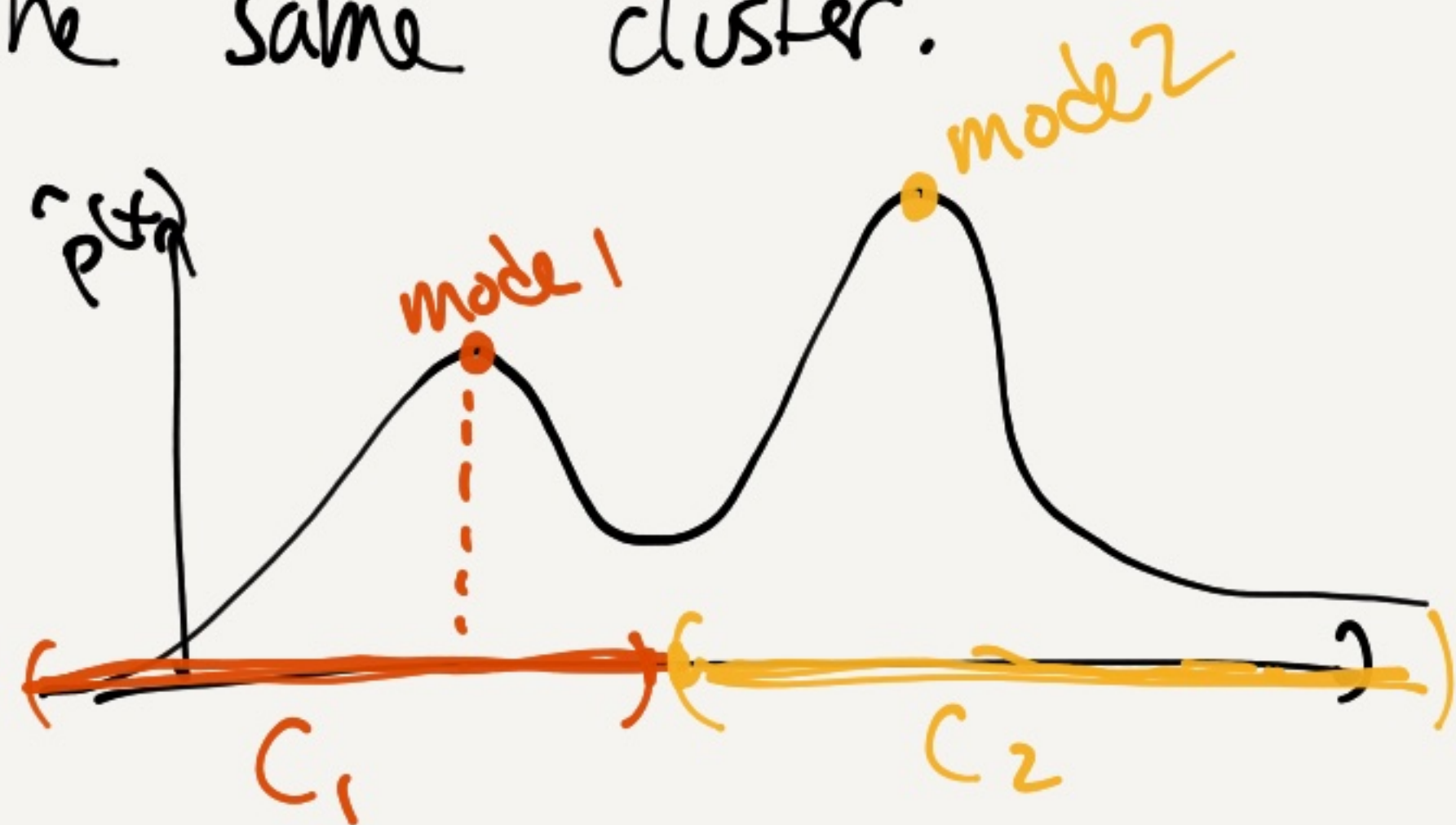- Idea:
  1) start at a point $\hat{x}$.
  2) use gradient ascent to move uphill
     $\left( \hat{x} \leftarrow \hat{x} + \lambda \nabla \hat{p}(x) \right)$
  3) eventually $\hat{x}$ will converge to the mode.

Modes: Repeat w/ many initial $\hat{x}$'s.
Remove duplicate converged $\hat{x}_s \rightarrow$ modes.

Clustering: given some $x_i$, all the $x_i$ that yield the same mode are in the same cluster.



## Consider only radially symmetric kernels

$$K(x) = \alpha \, \bar{K}(\|x\|^2)$$

constant    kernel profile

e.g. Gaussian: $\bar{K}(r) = e^{-\frac{1}{2}r}$, $\alpha = (2\pi)^{-d/2}$, $r \geq 0$

## Density Estimate

$$\hat{p}(x) = \frac{1}{Nh^d} \alpha \sum_{i=1}^{N} \bar{K}\left(\left\| \frac{x - x_i}{h} \right\|^2\right)$$

## Gradient

Define: $\bar{g}(r) = -\bar{K}'(r)$    ,    Gaussian: $\bar{g}(r) = \frac{1}{2} e^{-\frac{1}{2}r}$
$$= \frac{1}{2}\bar{K}(r)$$

$$\nabla \hat{p}(x) = \frac{\alpha}{Nh^{d+2}} \underbrace{\left( \sum_{i=1}^{N} \bar{g}\left(\left\| \frac{x-x_i}{h} \right\|^2\right) \right)}_{\approx \text{ density estimate using } \bar{g}(r) \text{ instead of } \bar{K}(r) \;=\; \hat{g}(x)} \left[ \underbrace{\frac{\sum_i x_i \bar{g}\left(\left\| \frac{x-x_i}{h} \right\|^2\right)}{\sum_i \bar{g}\left(\left\| \frac{x-x_i}{h} \right\|^2\right)}}_{\text{weighted mean of samples } x_i. \text{ weights depend on distance to a point } x. \text{ (closer samples have higher weights)}} - x \right]$$

"mean-shift vector" - difference btwn weighted mean and window center
$$= m(x)$$

## Gradient Ascent

$$\hat{x}^{(k+1)} = \hat{x}^{(k)} + \lambda \nabla \hat{p}(\hat{x}^{(k)})$$

↑ stepsize (important for convergence)

## Use an adaptive stepsize.

$$\lambda = \frac{1}{\hat{g}(x)}$$ 

⟸ $g(x)$ is small (low density region) → stepsize is large.

$g(x)$ is large (high density region) → stepsize is small

$$\Rightarrow \hat{x}^{(k+1)} = \hat{x}^{(k)} + \frac{1}{g(x^{(k)})} g(x^{(k)}) \hat{m}(x^{(k)})$$

$$\Rightarrow \hat{x}^{(k+1)} = \frac{\sum_i x_i \, \bar{g}\left(\left\|\frac{\hat{x}^{(k)} - x_i}{h}\right\|^2\right)}{\sum_i \bar{g}\left(\left\|\frac{\hat{x}^{(k)} - x_i}{h}\right\|^2\right)}$$
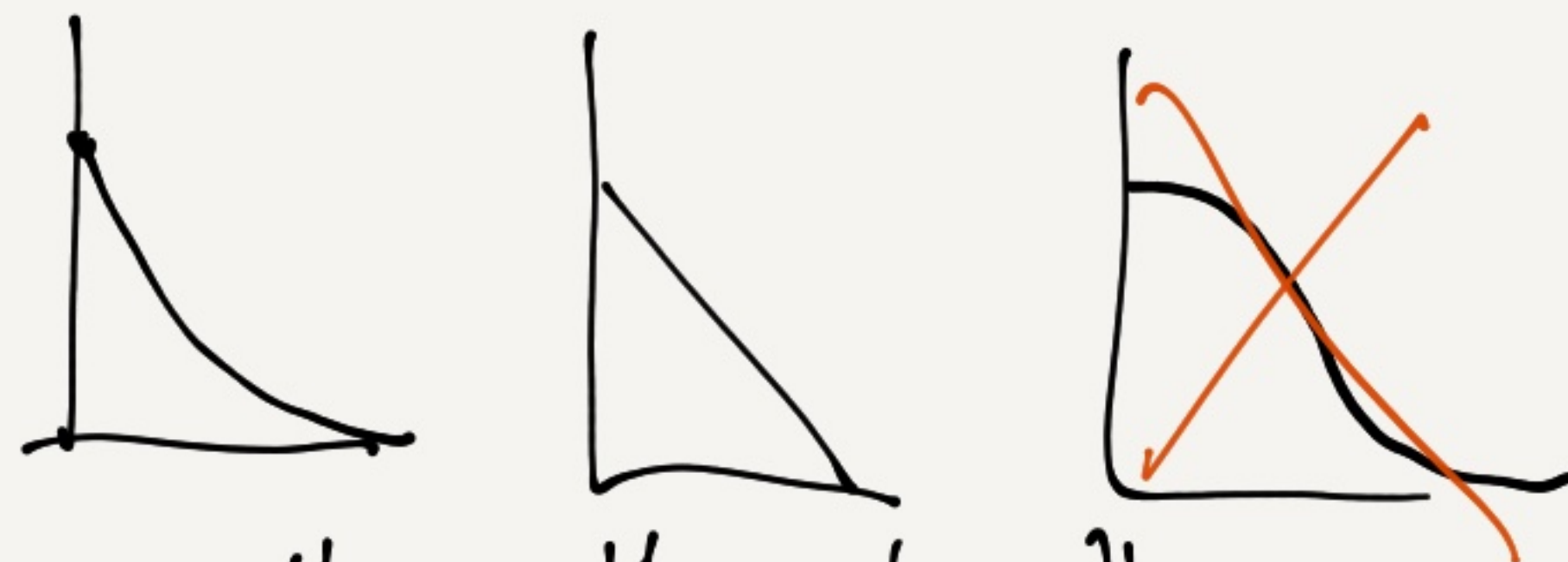
mean-shift procedure

in each iteration, shift to the weighted mean of the nearby points.

• The profile should be monotonically decreasing & convex.



if so, then the algorithm is guaranteed to converge to a stationary point.