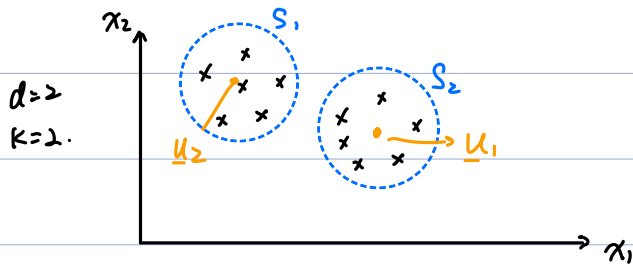


## Lec 20 Unsupervised learning

Recall:

↳ clustering:

$$\mathcal{D} = \{ \underline{x}_1, \underline{x}_2, \dots, \underline{x}_n \}$$

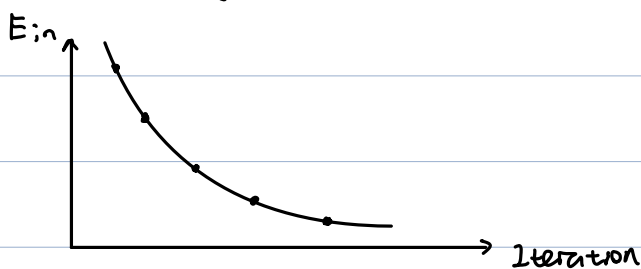


↳ k-means algo:

- ① Initialize  $\underline{\mu}_1, \dots, \underline{\mu}_k$  randomly
- ② construct  $S_1, \dots, S_k$  by solving subproblem #2. ("nearest cluster center")
- ③ Update  $\underline{\mu}_1, \dots, \underline{\mu}_k$  by solving subproblem #1 ("centroid")
- ④ Repeat step ② and ③ until converges

• Notes:

(a) converge is guaranteed.



⇒ decreasing and lower-bounded.

(Monotone convergence theorem)

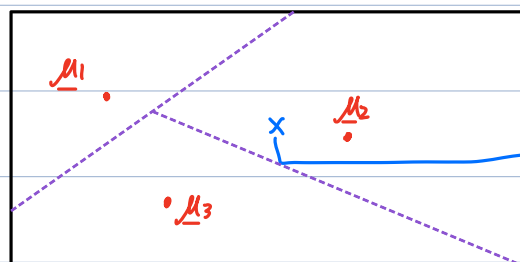
(b) It's locally optimal. Not global optimal in general.

(c) End result:

k-means create a

Voronoi diagram

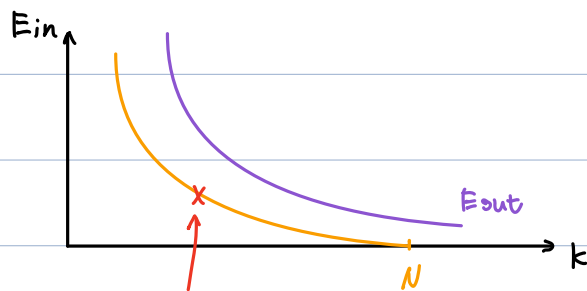
"tessellation"



→ test data.  
we find the  
nearest cluster

$k=3$

(d) How to choose  $k$ ? (Most tough question)



$N$ : Total #. of class

"Elbow point"

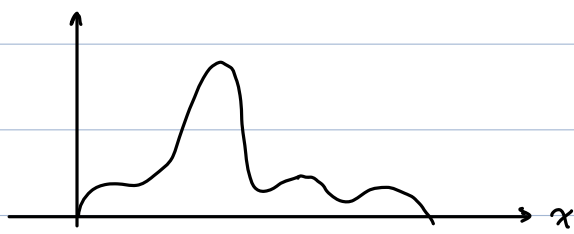
Not too much  $k$

Not too much  $E_{in}$ .

## 2. Density Estimation:

$$\hookrightarrow \mathcal{D} = \{ \underline{x}_1, \dots, \underline{x}_N \}, \quad \underline{x}_i \stackrel{\text{iid}}{\sim} p(\underline{x})$$

output  $\hat{p}(\underline{x}) \approx p(\underline{x})$



$d=1$

## 3. Review on Prob:

$\hookrightarrow$  Discrete RV:

$$\bullet X \sim \text{uniform in } \{1, 2, \dots, 10\}$$

$$\text{PMF} \quad \Pr \{X = k\} = \frac{1}{10} \quad \text{for } 1 \leq k \leq 10$$

$\hookrightarrow$  Continuous RV:

$$\bullet X \sim \text{uniform in } (a, b)$$

$$\Pr \{X = \frac{a+b}{2}\} = 0 \quad \therefore \text{we don't talk about PMF, PDF instead}$$

$$\text{PDF: } f_X(x) = \begin{cases} \frac{1}{b-a} & \text{if } a \leq x \leq b \\ 0 & \text{otherwise} \end{cases}$$

- $X \sim \mathcal{N}(\mu, \sigma^2)$

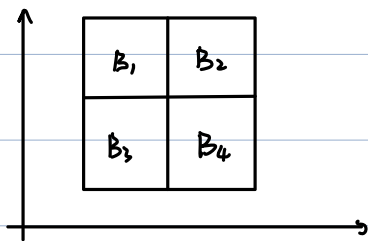
$$f_X(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

↳ Histogram Method:

- Assume  $P(\cdot)$  is non-zero only over a bounded-area

Cover the region with uniform bins (hyperparameters)

$B_1, B_2 \dots B_k$ , each with volume  $V$



- Let  $N_i$  be the nb. of samples in Bin  $B_i$

Hypothesis:  $\Pr\{X \in B_i\} = \frac{N_i}{N}$

- Assume the distribution within each bin is a uniform

$$\Rightarrow \text{PDF} = \frac{1}{V}$$

According to Total prob. Thrm:

$$\begin{aligned} \hat{P}(x) &= \sum_{i=1}^K \hat{P}(x | x \in B_i) \frac{N_i}{N} \\ &= \sum_{i=1}^K \frac{1}{V} \cdot \mathbb{1}(x \in B_i) \cdot \frac{N_i}{N} \\ &= \begin{cases} \frac{1}{V} \frac{N_i}{N} & \text{if } x \in B_i \\ \frac{1}{V} \frac{N_k}{N} & \text{if } x \in B_k \end{cases} \end{aligned}$$

e.g.

$d=1$

