# Lec 16: SVM (contd.)

## Support vector machines (max margin classifier)



**Recap:**

$$\min \ \frac{1}{2} \|W\|^2 \qquad \text{PRIMAL}$$

**Hard margin SVM**

$$\text{s.t.} \quad y_i \left( \underline{W}^T x_i + \underline{b} \right) \geq 1$$

$$\forall \ i = 1, \dots, n$$

$$\boxed{x_i \in \mathbb{R}^d}$$

$$d \gg n$$

$$w^* = \sum_{i=1}^{n} \lambda_i^* y_i x_i$$

**DUAL:**

$$\max_{\lambda \geq 0} \ \sum_{i=1}^{n} \lambda_i - \frac{1}{2} \sum_j \sum_i \lambda_i \lambda_j \underbrace{y_i y_j}_{} \underbrace{x_i^T x_j}_{}$$

$$\phi(x_i)^T \phi(x_j)$$

$$\phi(x_i) \quad \phi(x_j)$$

**Why dual?**

① less costly to solve.

② Kernelization

**easier problem**

$$\text{s.t.} \quad \sum_{i=1}^{n} \lambda_i y_i = 0$$

Given: linearly inseparable data

Goal: project the data to a high dimensional space
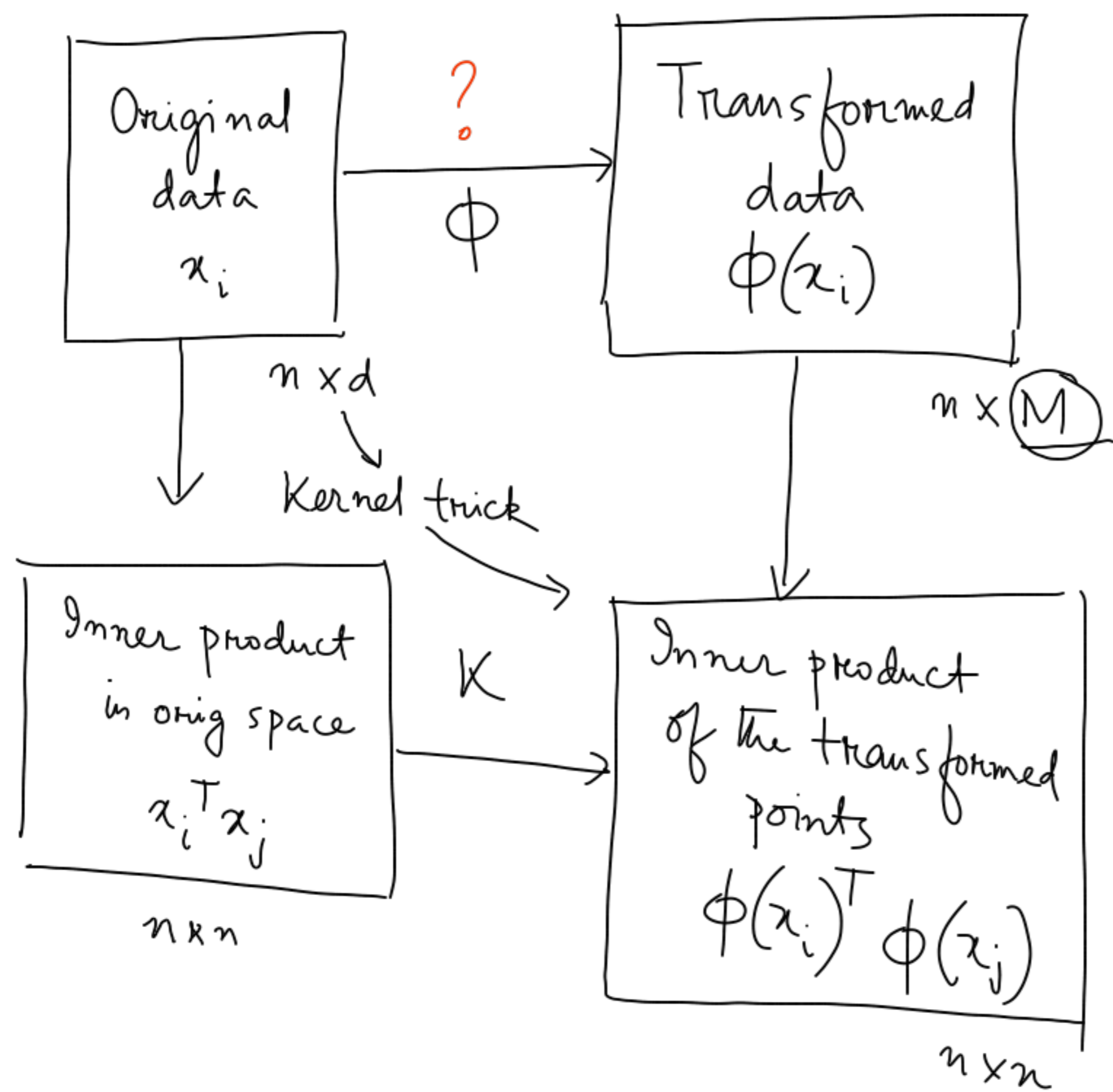
→ solve SVM → find $w, b$ in the new dimension.

E.g.

$$x_i = \begin{bmatrix} x_i^{(1)} \\ x_i^{(2)} \end{bmatrix} \xrightarrow{\phi} \begin{bmatrix} 1 \\ x_i^{(1)} \\ x_i^{(2)} \\ x_i^{(1)} x_i^{(2)} \\ x_i^{(1)\,2} \\ x_i^{(2)\,2} \end{bmatrix} = \phi(x_i)$$

Step 1:

$$\phi(x_i)^T \phi(x_j)$$

Step 2: Dual SVM

obj function has an inner product of the data
in the higher dimension

Soft margin SVMs are also not very good.

→ $\phi$ → basis transform

$$x_i \mapsto \phi(x_i) \in \mathbb{R}^M$$

high dimension

$M \gg d$

**Original data** $x_i$ — $? \atop \phi$ → **Transformed data** $\phi(x_i)$

$n \times d$

$n \times \boxed{M}$

Kernel trick

**Inner product in orig space** $x_i^T x_j$

$n \times n$

— $K$ → **Inner product of the transformed points** $\phi(x_i)^T \phi(x_j)$

$n \times n$

Q: Do there exist functions that calculate the inner product in the transformed space <u>without explicitly computing the transformations?</u>

A: <u>Yes, via the kernel function</u>

E.g. previous example, terms

$$= \{ 1, \; x_i^{(1)} x_j^{(1)}, \; x_i^{(2)} x_j^{(2)}, \; x_i^{(1)} x_i^{(2)} x_j^{(1)} x_j^{(2)}, \; x_i^{(1)^2} x_j^{(1)^2}, \; x_i^{(2)^2} x_j^{(2)^2} \}$$

$$K(x_i, x_j) = \left(1 + x_i^\top x_j\right)^2 \rightarrow \text{Same terms}$$

$$x_i^{(1)} x_j^{(1)} + x_i^{(2)} x_j^{(2)}$$

## Kernel trick

transformations are equivalent as long as we are calculating the dual of SVM.

## Kernel Regression

Different Kernels

① Linear: $K(x, z) = x^\top z$

② Polynomial: $K(x, z) = \left(1 + x^\top z\right)^m$

③ Gaussian: $K(x, z) = e^{-\|x - z\|^2 / 2\sigma^2}$

④ Laplace/Radial: $K(x, z) = e^{-\|x - z\| / 2\sigma}$

Use cases of SVM: Handwriting recognition, protein structure, medical image classification.

A set of necessary and sufficient conditions govern the kernel functions

Which Kernel to pick?

grid seach →

→ Mercer's theorem

see note on webpage.

## Limitations of SVM

1. Binary classification

→ One-vs-rest classifier

2. No probabilistic interpretation

3. Does not work very well when data is noisy.

Story so far:

Supervised learning

$$D = \left\{ (x_i, y_i) \right\}_{i=1, \cdots n}$$

↑

costly

Unsupervised learning

$$D = \left\{ (x_i) \right\}_{i=1, \cdots, n}$$
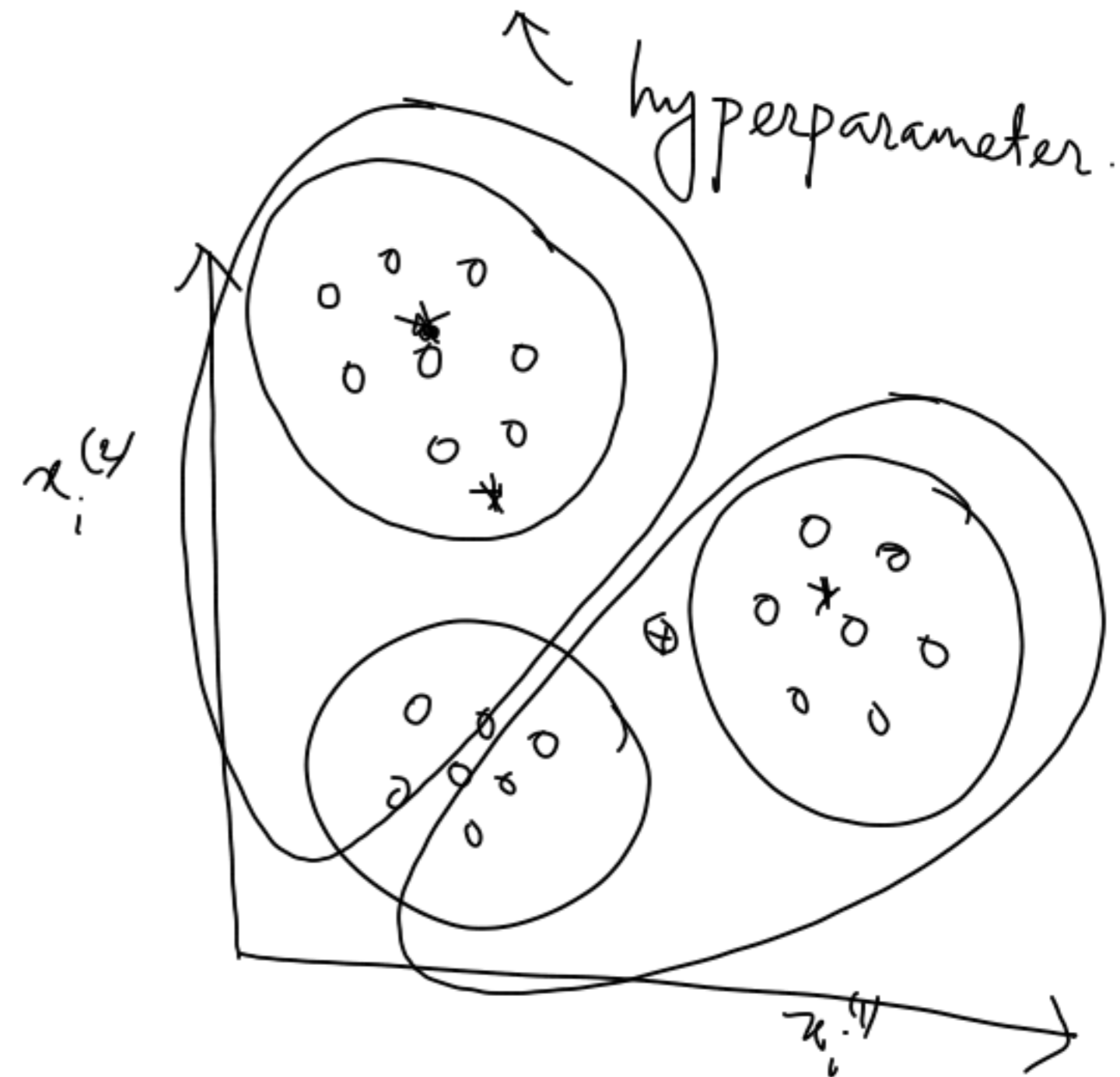
Clustering : Unsupervised learning

$$D = \left\{ x_1, x_2, \cdots, x_n \right\}, \quad x_i \in \mathbb{R}^d$$

Goal: find a "well-separated" partition of the data

$$D = D_1 \cup D_2 \cup \cdots \cup D_k$$

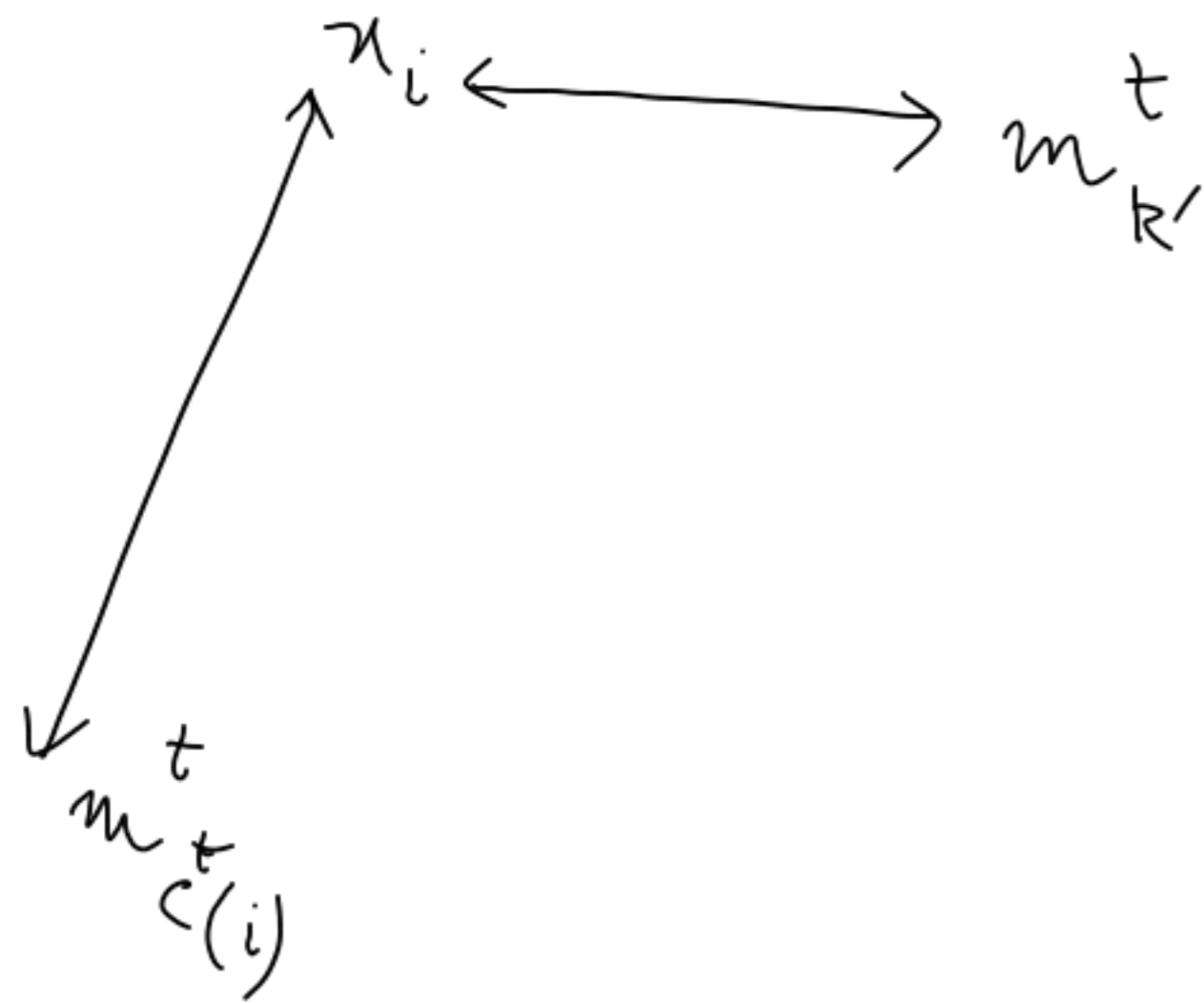hard cluster    $D_i \cap D_j = \phi$

k-means clustering

hyperparameter

$x_i^{(2)}$

$x_i^{(1)}$

$$C(i) = \tilde{k}$$

Clustering function

$i^{Th}$ data point

$$C : [n] \rightarrow [k]$$

$$[n] = \{1, 2, \cdots, n\}$$

$$C^t(i) \neq k'$$

$x_i \longleftrightarrow m_{k'}^t$

$m_{C(i)}^t$

## $k$-means

Input $= D = \{x_1, \cdots, x_n\}$

Initialize: $k$ cluster means $m_1, \cdots, m_R$

some arbitrary $C^0$

Repeat until convergence

(until assignments do not change)

for every $i = 1, \cdots, n$

if $\exists k' \neq C^t(i)$
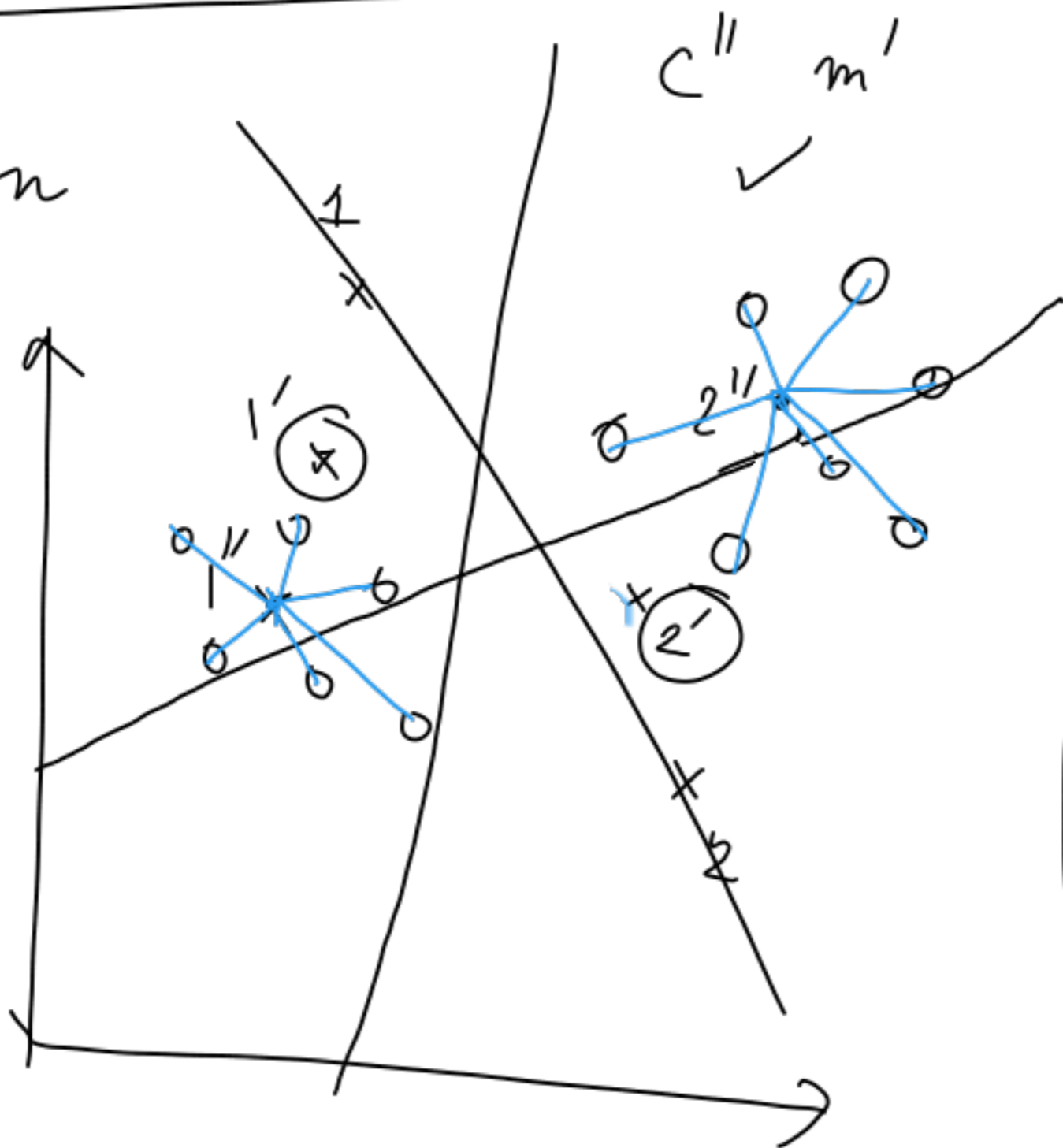
$$\| m_{k'}^t - x_i \| < \| m_{C^t(i)}^t - x_i \|$$

$$C^{t+1}(i) \leftarrow k'$$

Update cluster centers by average

of its current associated
data points

$m_{k'}^{t+1}$

---

Illustration



$C''$  $m'$

$$\min_{C \in \mathcal{C}} \sum_{i=1}^{n} \| x_i - m_{C(i)} \|^2$$

non-convex

NP-hard

Squared error  SE

$|\mathcal{C}| = \boxed{k^n}$

$C \in \mathcal{C}$

set of all
possible
clusters

$C(i) \rightarrow$

k values

1  2          $n$  $\rightarrow$ point

k means is a
"reasonable"  $\rightarrow$ converges. local optima.
approach

# Convergence

## Lemma:

$$\operatorname{argmin}_{x} \sum_{i=1}^{n_k} \|x_i - x\|^2 = \overline{x}$$

$$= \frac{1}{n_k} \sum_{i=1}^{n} x_i$$

fixed clustering



C

Why is this sufficient?

- no clustering are repeated
- finite number of clusterings

## Theorem:   k-means converges to a local minima

- consider $t \to t+1$
- $SE\left(c^{t+1}, m^{t+1}\right) < SE\left(c^t, m^t\right)$

Proof:

$$m^t, c^t \to m^t, c^{t+1} \to$$

$$m^{t+1}, c^{t+1}$$

① $SE\left(c^{t+1}, m^t\right) < SE\left(c^t, m^t\right)$
from the algorithm itself

② $SE\left(c^{t+1}, m^{t+1}\right) < SE\left(c^{t+1}, m^t\right)$
Lemma

Combining the two
claims, get the
theorem.

---

$$C^k(i) = prob\left(i \text{ belongs to } k\right)$$

↑

Soft clustering