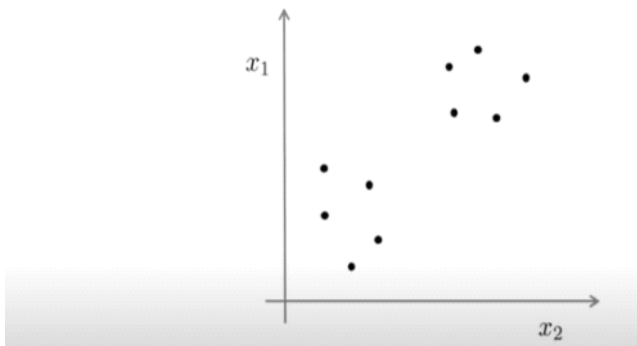


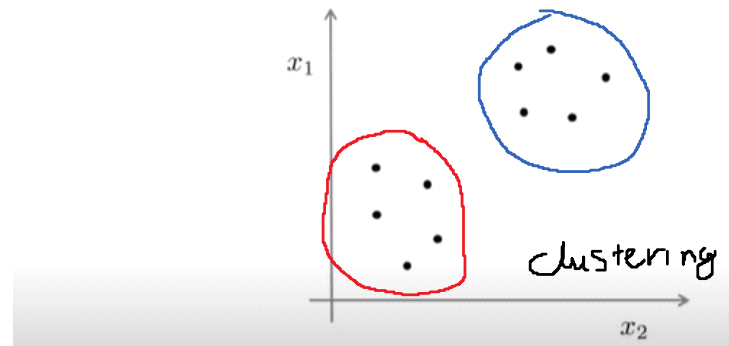
Clustering

Unsupervised Learning

Unsupervised learning

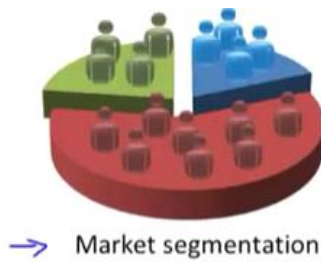


Unsupervised learning

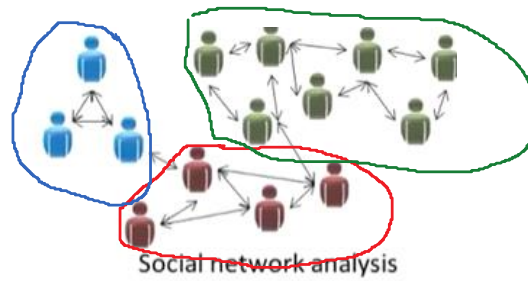


Phân cụm các dữ liệu có tính chất giống nhau

Ứng dụng



→ Market segmentation

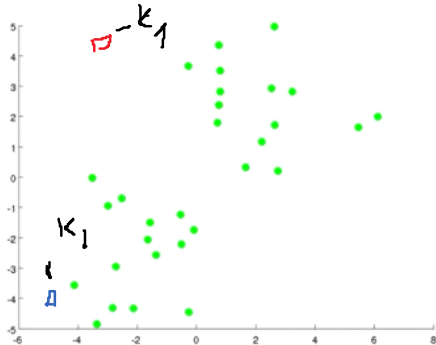


Social network analysis

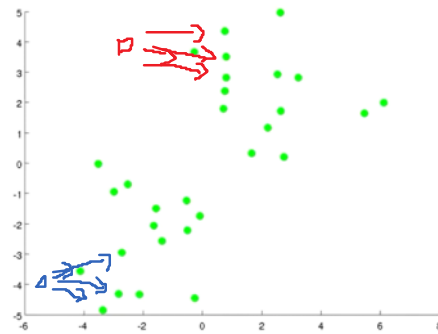


thuật toán cơ bản trong unsupervised Learning
các bước của thuật toán K-mean

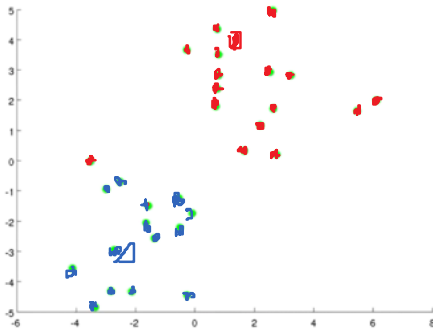
B₁: Random K center



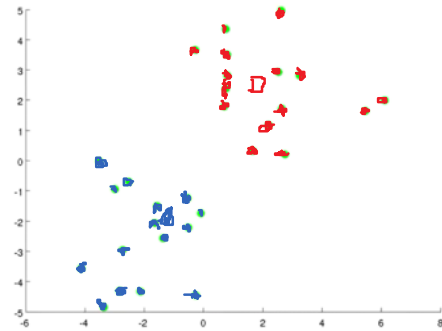
B₂: tính K center tới những điểm gần nhất



B₃: di chuyển center đến trung tâm của liên



B₄: lặp lại tính K và dịch chuyển



thuật toán

- Input: K

* training Set

mô tả

sinh ngẫu nhiên K center $u_1, u_2, u_3, \dots, u_K$

$\left\{ \begin{array}{l} \text{for } i=1 \text{ to } m \quad (m \text{ là số example data}) \\ \text{c}^{(i)} = \text{ch' số của center gần } x^{(i)} \\ \text{for } j=1 \text{ to } K \\ u_j = \text{trung tâm (trung bình) của tất cả các điểm có nhãn } j \end{array} \right.$

Optimization

$c^{(i)}$ = index của center mà $x^{(i)}$ gần nhất trong

μ_k = trung tâm (center) của nhóm k

$u_{c^{(i)}}$ = center mà $x^{(i)}$ gần nhất

$$J(c^1, c^2, \dots, c^m, \mu_1, \mu_2, \dots, \mu_K) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c^{(i)}}\|^2$$

min J
 c^1, c^2, \dots, c^m
 $\mu_1, \mu_2, \dots, \mu_K$

Repeat {

for $i = 1$ to m

$c^{(i)} :=$ index (from 1 to K) of cluster centroid
 closest to $x^{(i)}$

} gán nguyên μ_k
 và cập nhật

for $k = 1$ to K

$\mu_k :=$ average (mean) of points assigned to cluster k

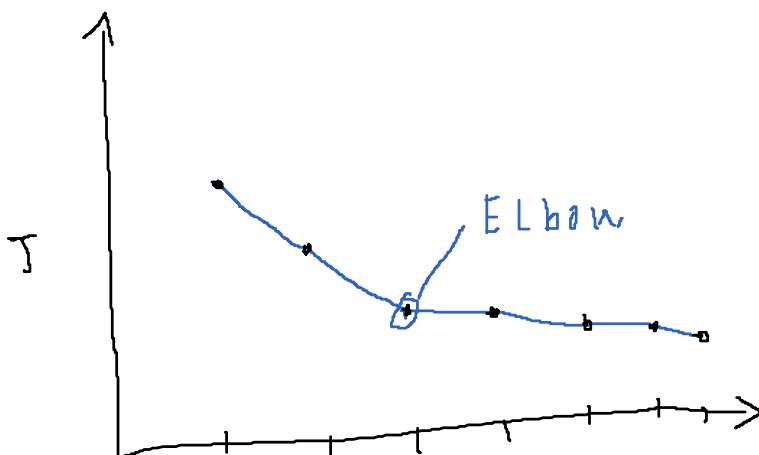
} cập nhật
 μ_k

choosing value of K

Elbow method

hoặc là theo mục đích

và có thể size của $K=6$



K

Cho tập dữ liệu $X = [x_1, x_2, \dots, x_N] \in \mathbb{R}^{d \times N}$ và $k < N$
 và center $m_1, m_2, \dots, m_k \in \mathbb{R}^{1 \times k}$

Với mỗi điểm x_i đặt 1 vector $y_i = [y_{i1}, y_{i2}, y_{i3}, \dots, y_{ik}]$
 y là vector one-hot $[0, 1, 0, 0, \dots, 0]$

điểm khoảng cách tới center nhỏ nhất

$$\|x_i - m_k\|_2^2 \quad (\text{norm } L)$$

$$\text{vì } x_i \in \text{cluster (center } k) \Rightarrow \begin{cases} y_{ik} = 1 & \forall i \neq k \\ y_{ij} = 0 \end{cases}$$

$$\Rightarrow y_i \|x_i - m_k\|_2^2 = \sum_{j=1}^k y_{ij} \cdot \|x_i - m_j\|_2^2 - k = \text{số cluster}$$

one-hot tính k từ $x_i \rightarrow$ center
 nên chỉ cần k ngẫu nhiên

áp dụng cho toàn bộ data set

$$J(M, Y) = \sum_{i=1}^m \sum_{j=1}^k y_{ij} \|x_i - m_j\|_2^2$$

$$\Rightarrow \min_{M, Y} J$$

$$M, Y = \arg \min_{M, Y} \sum_{i=1}^m \sum_{j=1}^k y_{ij} \|x_i - m_j\|_2^2$$

thỏa mãn $y_{ij} \in \{0, 1\}$ và $\sum_{j=1}^k y_{ij} = 1 \Leftrightarrow y_i$ one hot

* khái niệm $\arg \min$ là giá trị của biến để hàm số min

$$\forall n \quad f(x) = x^2 - 2x + 1 = (x-1)^2$$

$$\min f(x) = 0$$

$$\arg \min f(x) = 1 \quad \text{vì } \min f(x) \text{ tại } x=1$$

$$\min_x f(x) = 0$$

$$\arg \min_x f(x) = 1 \quad \text{vì } \min_x f(x) \text{ tại } x=1$$

giải bài toán. xem kẽ giữa Y và M khi cần cần lại cố định
 cố định M tìm Y: Sau khi tìm đc hoặc di chuyển các center

$$x_i = \arg \min_j \sum_{i=1}^K y_{ij} \|x_i - m_j\|_2^2$$

$$\Leftrightarrow i = \arg \min_j \sum_{j=1}^K \|x_i - m_j\|_2^2$$

Cố định Y tìm M. sau khi đư ợc nh ậ n cho tất cả data

$$m_j = \arg \min_{m_j} \sum_{i=1}^N y_{ij} \|x_i - m_j\|_2^2 \quad N \text{ tập dữ liệu}$$

$$\text{Đặt } L(m_j) = \sum_{i=1}^N y_{ij} \|x_i - m_j\|_2^2$$

$$\Rightarrow \frac{\partial L(m_j)}{\partial m_j} = \sum_{i=1}^N 2 y_{ij} (m_j - x_i) \quad (I)$$

$$\text{giải } (I) = 0$$

$$\Rightarrow m_j \sum_{i=1}^N 2 y_{ij} = \sum_{i=1}^N 2 y_{ij} \cdot x_i$$

$$\Rightarrow m_j = \frac{\sum_{i=1}^N y_{ij} x_i}{\sum_{i=1}^N y_{ij}} = \frac{\text{tổng các } x_i \in \text{cluster } j}{\text{phép đếm số } x_i \in \text{cluster } j}$$

$$\Rightarrow m_j = \frac{\text{tổng}}{\text{SL}} \Leftrightarrow \text{trung bình cộng của cluster } j$$