

Học Máy (IT 4862)

Nguyễn Nhật Quang

quangnn-fit@mail.hut.edu.vn

Trường Đại học Bách Khoa Hà Nội
Viện Công nghệ thông tin và truyền thông
Năm học 2011-2012

Nội dung môn học:

- Giới thiệu chung
- Đánh giá hiệu năng hệ thống học máy
- Các phương pháp học dựa trên xác suất
- Các phương pháp học có giám sát
- **Các phương pháp học không giám sát**
 - **Giới thiệu về phân cụm**
 - **Phân cụm dựa trên phân tách: k-Means**
- Lọc cộng tác
- Học tăng cường

Học có vs. không có giám sát

■ Học có giám sát (Supervised learning)

- Tập dữ liệu (dataset) bao gồm các ví dụ, mà mỗi ví dụ được *gắn kèm với một nhãn lớp/giá trị đầu ra mong muốn*
- Mục đích là học (xấp xỉ) một giả thiết (vd: một phân lớp, một hàm mục tiêu,...) phù hợp với tập dữ liệu hiện có
- Giả thiết học được (learned hypothesis) sau đó sẽ được dùng để phân lớp/dự đoán đối với các ví dụ mới

■ Học không có giám sát (Unsupervised learning)

- Tập dữ liệu (dataset) bao gồm các ví dụ, mà mỗi ví dụ *không có thông tin về nhãn lớp/giá trị đầu ra mong muốn*
- Mục đích là tìm ra (học) các cụm/các cấu trúc/các quan hệ tồn tại trong tập dữ liệu hiện có

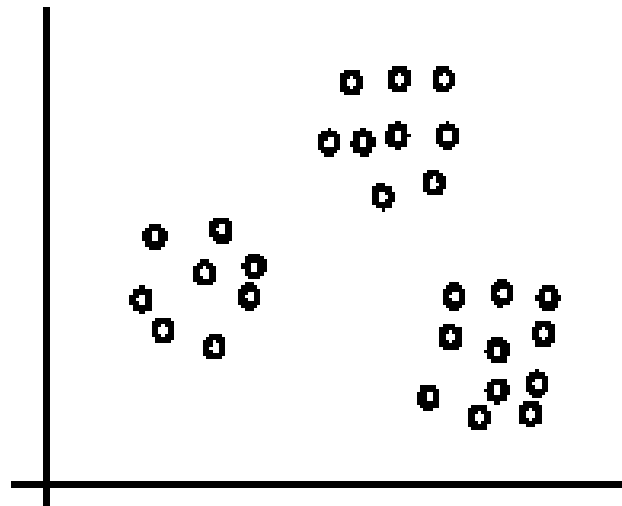
Phân cụm

- Phân cụm/nhóm (Clustering) là phương pháp học không có giám sát được sử dụng phổ biến nhất
 - Tồn tại các phương pháp học không có giám sát khác, ví dụ: Lọc cộng tác (Collaborative filtering), Khai phá luật kết hợp (Association rule mining), ...
- Học phân cụm
 - Đầu vào: một tập dữ liệu không có nhãn (các ví dụ không có nhãn lớp/giá trị đầu ra mong muốn)
 - Đầu ra: các cụm (nhóm) của các ví dụ
- Một **cụm (cluster)** là một tập các ví dụ
 - Tương tự với nhau (theo một ý nghĩa, đánh giá nào đó)
 - Khác biệt với các ví dụ thuộc các cụm khác

Phân cụm – Ví dụ

Một ví dụ về phân cụm:

Các ví dụ được phân chia thành 3 cụm



[Liu, 2006]

Phân cụm – Các thành phần

- Hàm tính khoảng cách (độ tương tự, độ khác biệt)
- Giải thuật phân cụm
 - **Dựa trên phân tách (Partition-based clustering)**
 - **Dựa trên tích tụ phân cấp (Hierarchical clustering)**
 - Bản đồ tự tổ chức (Self-organizing map – SOM)
 - Các mô hình hỗn hợp (Mixture models)
 - ...
- Đánh giá chất lượng phân cụm (Clustering quality)
 - Khoảng cách/sự khác biệt *giữa các cụm* → Cần được *cực đại* hóa
 - Khoảng cách/sự khác biệt *bên trong một cụm* → Cần được *cực tiểu* hóa

Phân cụm k-Means

- Là phương pháp phổ biến nhất trong các phương pháp phân cụm dựa trên chia cắt (partition-based clustering)
- Tập dữ liệu $D = \{x_1, x_2, \dots, x_r\}$
 - x_i là một ví dụ (một vector trong một không gian n chiều)
- Giải thuật k -means phân chia (partitions) tập dữ liệu thành k cụm
 - Mỗi cụm (cluster) có một điểm trung tâm, được gọi là **centroid**
 - k (tổng số các cụm thu được) là một giá trị được xác định trước (vd: được chỉ định bởi người thiết kế hệ thống phân cụm)

k-Means – Các bước chính

Với một giá trị k được xác định trước

- Bước 1. Chọn ngẫu nhiên k ví dụ (được gọi là **các hạt nhân – seeds**) để sử dụng làm *các điểm trung tâm ban đầu (initial centroids)* của k cụm
- Bước 2. Đối với mỗi ví dụ, *gán nó vào cụm* (trong số k cụm) có điểm trung tâm (centroid) gần ví dụ đó nhất
- Bước 3. Đối với mỗi cụm, *tính toán lại điểm trung tâm (centroid) của nó* dựa trên tất cả các ví dụ thuộc vào cụm đó
- Bước 4. Dừng lại nếu *điều kiện hội tụ (convergence criterion)* được thỏa mãn; nếu không, quay lại Bước 2

***k*-means(*D*, *k*)**

D: Tập ví dụ học

k: Số lượng cụm kết quả (thu được)

Lựa chọn ngẫu nhiên *k* ví dụ trong tập *D* để làm các điểm trung tâm ban đầu (initial centroids)

while not CONVERGENCE

for each ví dụ $x \in D$

Tính các khoảng cách từ x đến các điểm trung tâm (centroid)

Gán x vào cụm có điểm trung tâm (centroid) gần x nhất

end for

for each cụm

Tính (xác định) lại điểm trung tâm (centroid) dựa trên các ví dụ hiện thời đang thuộc vào cụm này

end while

return {*k* cụm kết quả}

Điều kiện hội tụ

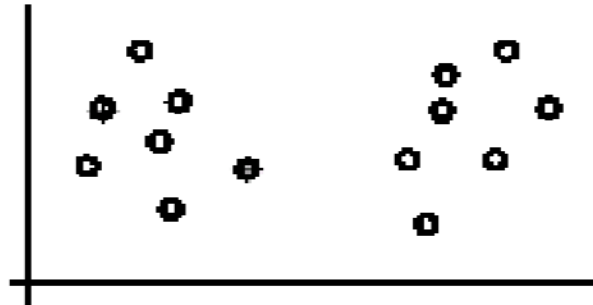
Quá trình phân cụm kết thúc, nếu:

- Không có (hoặc có không đáng kể) việc gán lại các ví dụ vào các cụm khác, *hoặc*
- Không có (hoặc có không đáng kể) thay đổi về các điểm trung tâm (centroids) của các cụm, *hoặc*
- Giảm không đáng kể về tổng lỗi phân cụm:

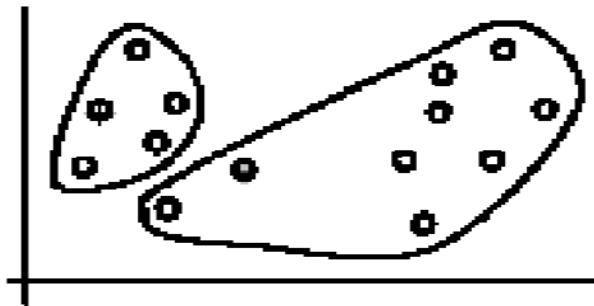
$$Error = \sum_{i=1}^k \sum_{\mathbf{x} \in C_i} d(\mathbf{x}, \mathbf{m}_i)^2$$

- C_i : Cụm thứ i
- \mathbf{m}_i : Điểm trung tâm (centroid) của cụm C_i
- $d(\mathbf{x}, \mathbf{m}_i)$: Khoảng cách (khác biệt) giữa ví dụ \mathbf{x} và điểm trung tâm \mathbf{m}_i

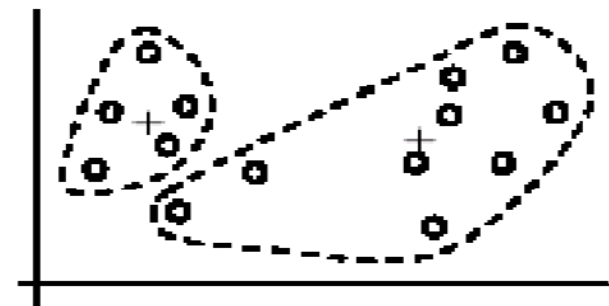
k-Means – Minh họa (1)



(A). Random selection of k centers



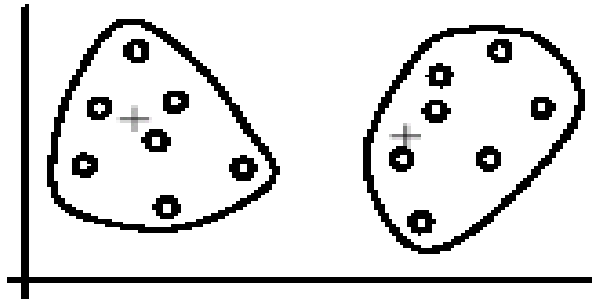
Iteration 1: (B). Cluster assignment



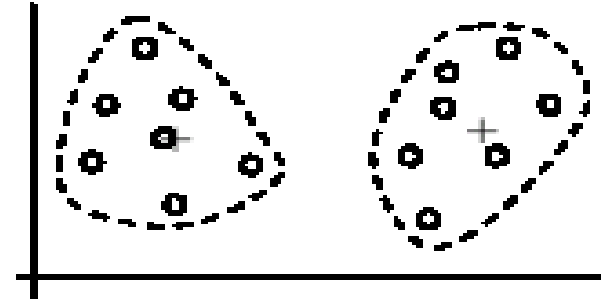
(C). Re-compute centroids

[Liu, 2006]

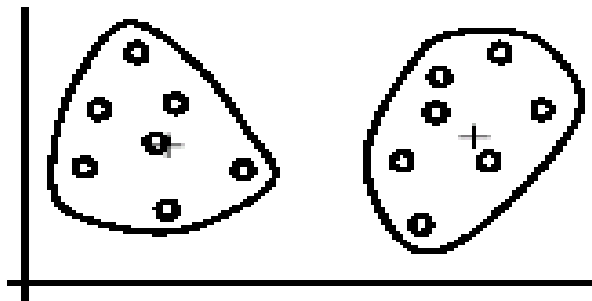
k-Means – Minh họa (2)



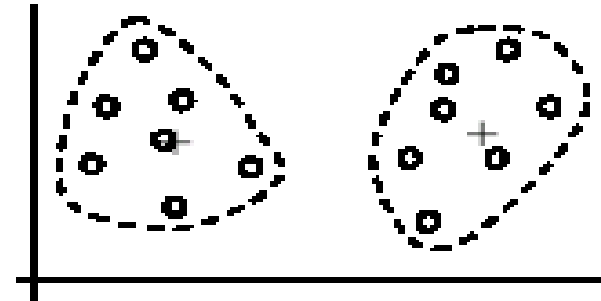
Iteration 2: (D). Cluster assignment



(E). Re-compute centroids



Iteration 3: (F). Cluster assignment



(G). Re-compute centroids

[Liu, 2006]

Điểm trung tâm, Hàm khoảng cách

- Xác định điểm trung tâm: Điểm trung bình (*Mean centroid*)

$$\mathbf{m}_i = \frac{1}{|C_i|} \sum_{\mathbf{x} \in C_i} \mathbf{x}$$

- (vector) \mathbf{m}_i là điểm trung tâm (centroid) của cụm C_i
- $|C_i|$ kích thước của cụm C_i (tổng số ví dụ trong C_i)

- Hàm khoảng cách: *Euclidean distance*

$$d(\mathbf{x}, \mathbf{m}_i) = \|\mathbf{x} - \mathbf{m}_i\| = \sqrt{(x_1 - m_{i1})^2 + (x_2 - m_{i2})^2 + \dots + (x_n - m_{in})^2}$$

- (vector) \mathbf{m}_i là điểm trung tâm (centroid) của cụm C_i
- $d(\mathbf{x}, \mathbf{m}_i)$ là khoảng cách giữa ví dụ \mathbf{x} và điểm trung tâm \mathbf{m}_i

k-Means – Các ưu điểm

■ Đơn giản

- Rất dễ cài đặt
- Rất dễ hiểu

■ Hiệu quả

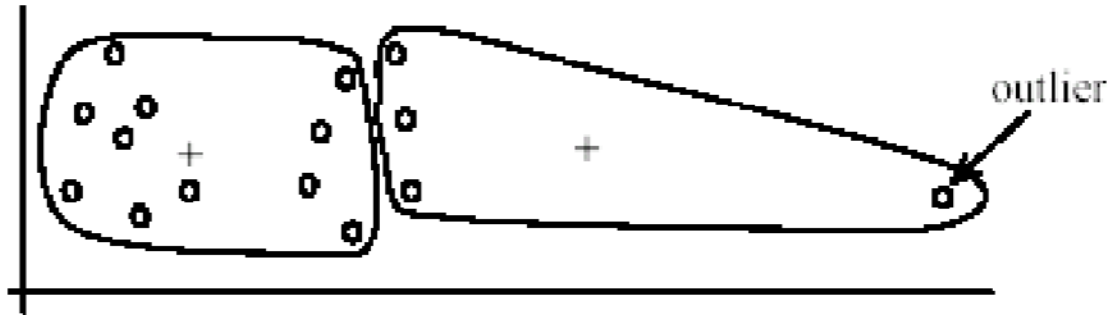
- Độ phức tạp về thời gian $\sim O(r \cdot k \cdot t)$
 - r : Tổng số các ví dụ (kích thước của tập dữ liệu)
 - k : Tổng số cụm thu được
 - t : Tổng số bước lặp (của quá trình phân cụm)
- Nếu cả 2 giá trị k và t đều nhỏ, thì giải thuật k -means được xem như là có độ phức tạp ở mức tuyến tính

■ k -means là giải thuật phân cụm được dùng phổ biến nhất

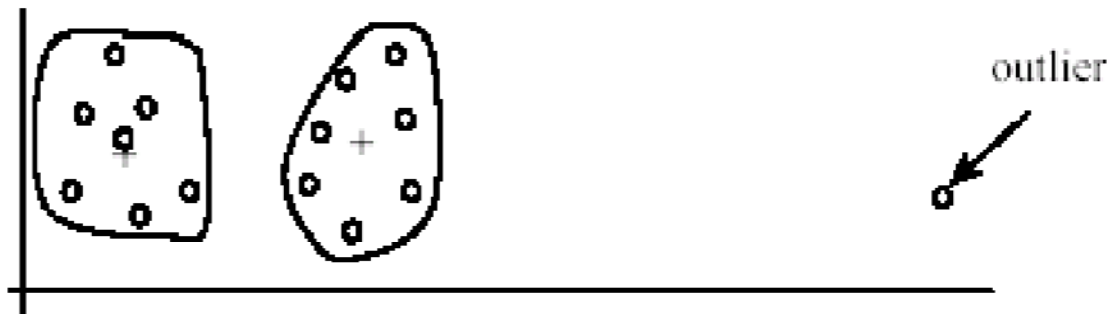
k-Means – Các nhược điểm (1)

- Giá trị k (số cụm thu được) phải được xác định trước
- Giải thuật k -means cần xác định cách tính điểm trung bình (centroid) của một cụm
 - Đối với các thuộc tính định danh (nominal attributes), giá trị trung bình có thể được xác định là giá trị phổ biến nhất
- Giải thuật k -means nhạy cảm (gặp lỗi) với ***các ví dụ ngoại lai (outliers)***
 - Các ví dụ ngoại lai là các ví dụ (rất) khác biệt với tất các ví dụ khác
 - Các ví dụ ngoại lai có thể do lỗi trong quá trình thu thập/lưu dữ liệu
 - Các ví dụ ngoại lai có các giá trị thuộc tính (rất) khác biệt với các giá trị thuộc tính của các ví dụ khác

k-Means – Các ví dụ ngoại lai



(A): Undesirable clusters



(B): Ideal clusters

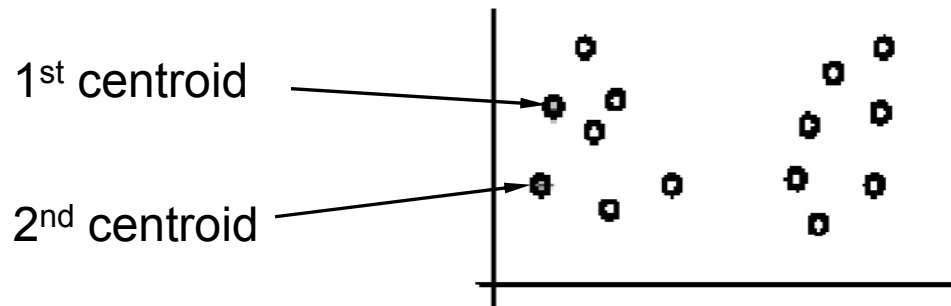
[Liu, 2006]

Giải quyết vấn đề ngoại lai

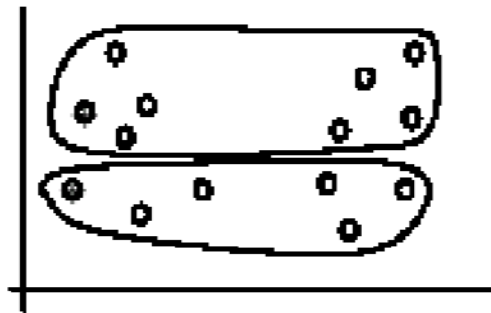
- Giải pháp 1. Trong quá trình phân cụm, cần loại bỏ một số các ví dụ quá khác biệt với (cách xa) các điểm trung tâm (centroids) so với các ví dụ khác
 - Để chắc chắn (không loại nhầm), theo dõi các ví dụ ngoại lai (outliers) qua một vài (thay vì chỉ 1) bước lặp phân cụm, trước khi quyết định loại bỏ
- Giải pháp 2. Thực hiện việc lấy mẫu ngẫu nhiên (a random sampling)
 - Do quá trình lấy mẫu chỉ lựa chọn một tập con nhỏ của tập dữ liệu ban đầu, nên khả năng một ngoại lai (outlier) được chọn là rất nhỏ
 - Gán các ví dụ còn lại của tập dữ liệu vào các cụm tùy theo đánh giá về khoảng cách (hoặc độ tương tự)

k-Means – Các nhược điểm (2)

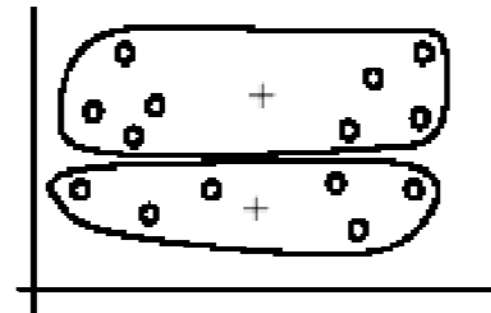
- Giải thuật k -means phụ thuộc vào việc chọn các điểm trung tâm ban đầu (initial centroids)



(A). Random selection of seeds (centroids)



(B). Iteration 1



(C). Iteration 2

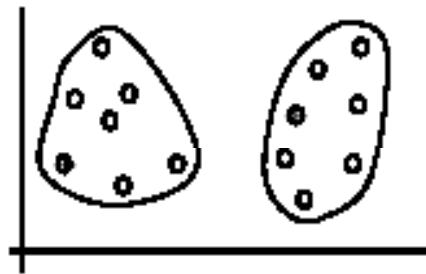
[Liu, 2006]

k-Means – Các hạt nhân ban đầu (1)

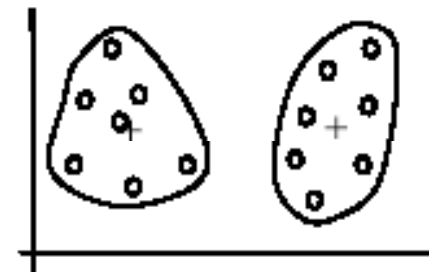
- Sử dụng các hạt nhân (seeds) khác nhau → Kết quả tốt hơn!
 - Thực hiện giải thuật k -means nhiều lần, mỗi lần bắt đầu với một tập (khác lần trước) các hạt nhân được chọn ngẫu nhiên



(A). Random selection of k seeds (centroids)



(B). Iteration 1



(C). Iteration 2

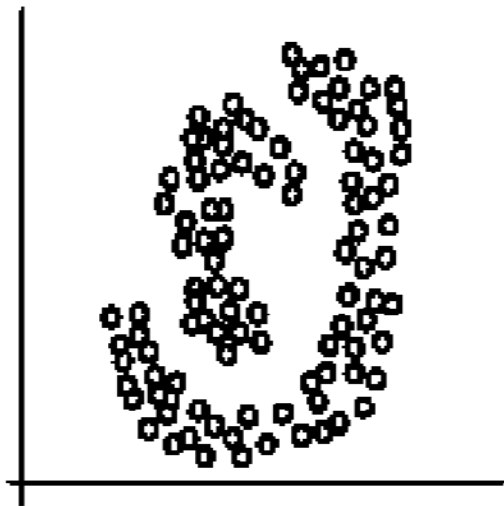
[Liu, 2006]

k-Means – Các hạt nhân ban đầu (2)

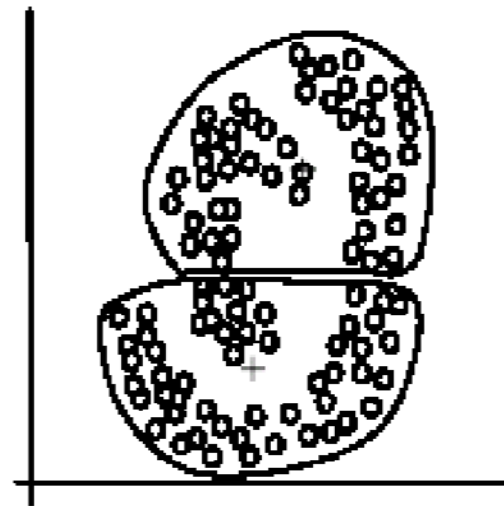
- Lựa chọn ngẫu nhiên hạt nhân thứ 1 (\mathbf{m}_1)
- Lựa chọn hạt nhân thứ 2 (\mathbf{m}_2) càng xa càng tốt so với hạt nhân thứ 1
- ...
- Lựa chọn hạt nhân thứ i (\mathbf{m}_i) càng xa càng tốt so với hạt nhân gần nhất trong số $\{\mathbf{m}_1, \mathbf{m}_2, \dots, \mathbf{m}_{i-1}\}$
- ...

k-Means – Các nhược điểm (3)

- Giải thuật k -means không phù hợp để phát hiện các cụm (nhóm) không có dạng hình elip hoặc hình cầu



(A): Two natural clusters



(B): k -means clusters

[Liu, 2006]

k-Means – Tổng kết

- Mặc dù có những nhược điểm như trên, k -means vẫn là giải thuật phổ biến nhất được dùng để giải quyết các bài toán phân cụm – do tính đơn giản và hiệu quả
 - Các giải thuật phân cụm khác cũng có các nhược điểm riêng
- Về tổng quát, không có lý thuyết nào chứng minh rằng một giải thuật phân cụm khác hiệu quả hơn k -means
 - Một số giải thuật phân cụm có thể phù hợp hơn một số giải thuật khác đối với một số kiểu tập dữ liệu nhất định, hoặc đối với một số bài toán ứng dụng nhất định
- So sánh hiệu năng của các giải thuật phân cụm là một nhiệm vụ khó khăn (thách thức)
 - Làm sao để biết được các cụm kết quả thu được là chính xác?

Tài liệu tham khảo

- B. Liu. *Web Data Mining: Exploring Hyperlinks, Contents, and Usage Data*. Springer, 2006.