

Personal Notebook: Unsupervised Machine Learning

Comprehensive Study Guide - IBM Machine Learning Course

Generated by AI Research Agent

October 29, 2025

Contents

1	Module 1: Introduction to Unsupervised Learning and K-Means Clustering	3
1.1	Introduction to Unsupervised Learning	3
1.2	K-Means Clustering	3
1.2.1	Algorithm Steps	3
1.3	Initialization and Elbow Method	3
1.4	Silhouette Score	4
1.5	Gaussian Mixture Models	4
2	Module 2: Distance Metrics and Curse of Dimensionality	5
2.1	Distance Metrics	5
2.2	Curse of Dimensionality	5
3	Module 3: Advanced Clustering Algorithms	7
3.1	Hierarchical Agglomerative Clustering	7
3.2	DBSCAN	7
3.3	Mean Shift	7
3.4	So sanh cac thuat toan clustering	9
4	Module 4: Principal Component Analysis (PCA)	10
5	Module 5: Advanced Dimensionality Reduction Techniques	11
5.1	Kernel PCA	11

5.2	Multidimensional Scaling (MDS)	11
6	Module 6: Non-Negative Matrix Factorization (NMF)	12
7	Module 7: Summary and Best Practices	13

1 Module 1: Introduction to Unsupervised Learning and K-Means Clustering

1.1 Introduction to Unsupervised Learning

Unsupervised Learning là hình thức học máy mà dữ liệu không có nhãn, mục tiêu là tìm cấu trúc ẩn trong dữ liệu. Các loại bao gồm:

- Clustering (Phân cụm)
- Dimensionality Reduction (Giảm chiều dữ liệu)
- Association Rule Learning
- Anomaly Detection

1.2 K-Means Clustering

Mục tiêu: Phân cụm dữ liệu thành k cụm, tối thiểu hóa tổng bình phương khoảng cách đến centroid.

$$J = \sum_{i=1}^k \sum_{x \in C_i} \|x - \mu_i\|^2$$

Công thức centroid:

$$\mu_i = \frac{1}{|C_i|} \sum_{x \in C_i} x$$

Khoảng cách Euclidean:

$$d(x, \mu) = \sqrt{\sum_{j=1}^n (x_j - \mu_j)^2}$$

1.2.1 Algorithm Steps

1. Chọn số cụm k
2. Khởi tạo k centroid
3. Gán mỗi điểm dữ liệu cho centroid gần nhất
4. Cập nhật lại centroid
5. Lặp lại đến khi hội tụ

1.3 Initialization and Elbow Method

K-Means++ chọn centroid đầu tiên ngẫu nhiên, tiếp theo với xác suất tỷ lệ bình phương khoảng cách. Elbow method chọn k tại điểm uốn đồ thị WCSS.

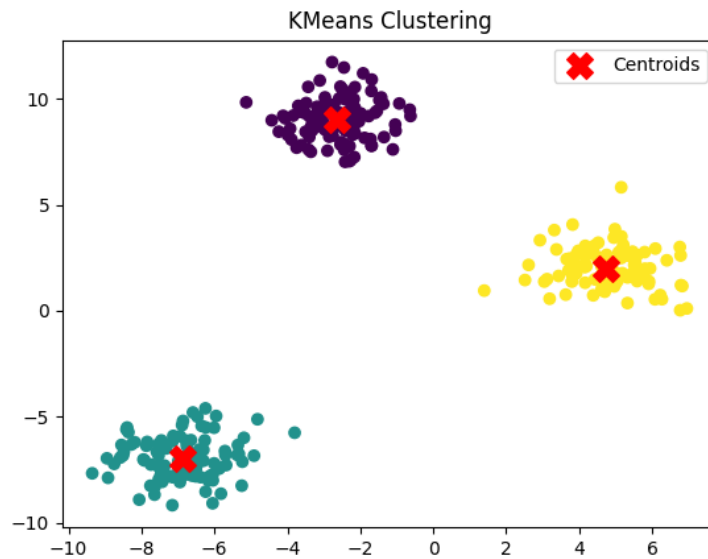


Figure 1: Minh hoa KMeans Clustering

1.4 Silhouette Score

Silhouette score $s(i)$ đánh giá chất lượng clustering:

$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}$$

$a(i)$ là khoảng cách trung bình trong cụm, $b(i)$ là khoảng cách trung bình tới cụm gần nhất.

1.5 Gaussian Mixture Models

Mô hình hóa dữ liệu bằng hỗn hợp các Gaussian:

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x | \mu_k, \Sigma_k)$$

Thuật toán EM lặp giữa E-step (tính responsibility) và M-step (cập nhật tham số).

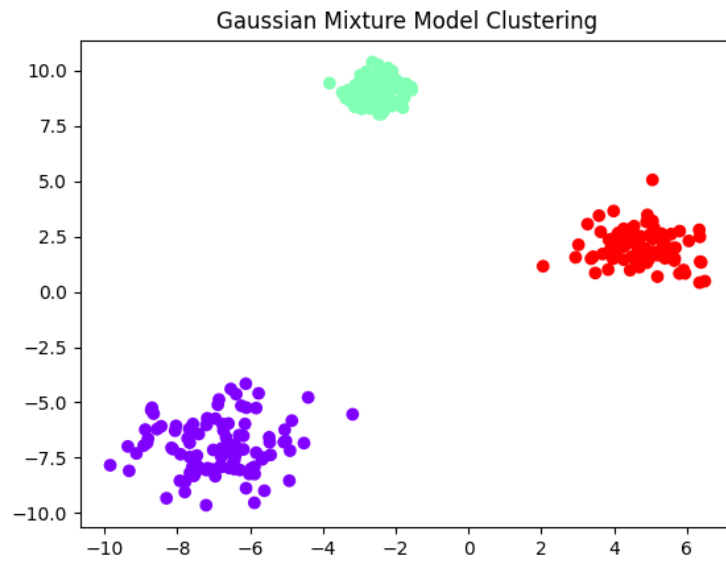


Figure 2: Minh hoa Gaussian Mixture Model

2 Module 2: Distance Metrics and Curse of Dimensionality

2.1 Distance Metrics

Euclidean Distance:

$$d_E(p, q) = \sqrt{\sum_{i=1}^n (p_i - q_i)^2}$$

Manhattan Distance:

$$d_M(p, q) = \sum_{i=1}^n |p_i - q_i|$$

Cosine Similarity & Distance:

$$\text{cosine similarity} = \frac{A \cdot B}{\|A\| \|B\|}, \quad \text{cosine distance} = 1 - \text{cosine similarity}$$

Jaccard Similarity & Distance:

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}, \quad d_J = 1 - J(A, B)$$

2.2 Curse of Dimensionality

Không gian đa chiều làm khoảng cách đồng đều và dữ liệu thừa thớt, gây khó khăn cho thuật toán.

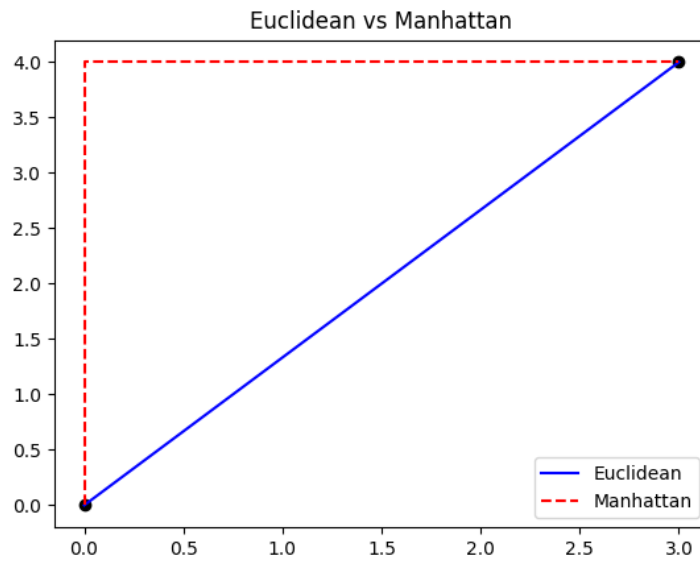



Figure 3: So sánh Euclidean và Manhattan distance

distance_metrics_comparison.png

Figure 4: Các loại distance metrics phổ biến



curse_dimensionality.png

Figure 5: Curse of Dimensionality

3 Module 3: Advanced Clustering Algorithms

3.1 Hierarchical Agglomerative Clustering

Cac phuong phap linkage: Single, Complete, Average, Ward.

3.2 DBSCAN

Chia diem thanh diem loi, bien va nhieu dua tren mat do lan can.

3.3 Mean Shift

Tim mode cua phan phoi mat do kernel.



Figure 6: Hierarchical Dendrogram

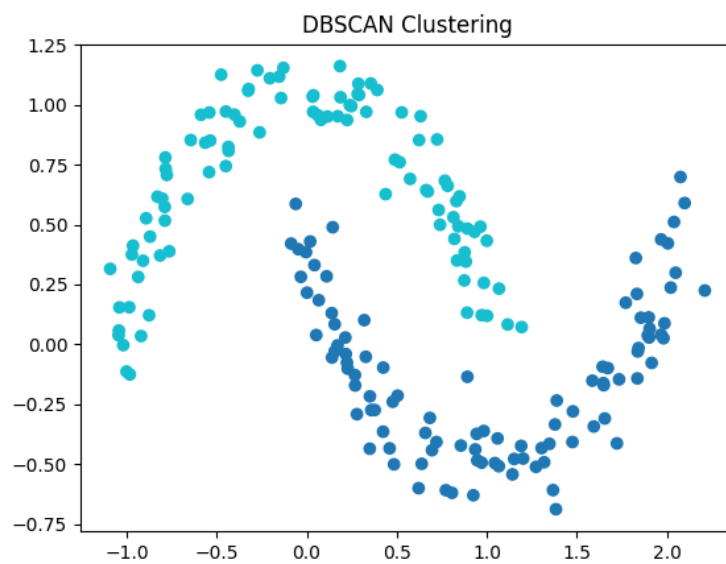


Figure 7: Minh hoa DBSCAN



Figure 8: Mean Shift Clustering

3.4 So sanh cac thuat toan clustering

Aspect	K-Means	HAC	DBSCAN	Mean Shift	GMM
Number of Clusters	Need k	Dendrogram	Auto	Auto	Need k
Cluster Shape	Spherical	Depends	Arbitrary	Arbitrary	Elliptical
Scalability	Good	Poor	Medium	Poor	Medium
Noise Handling	Poor	Depends	Good	Good	Moderate
Parameters	k , init	Linkage	ϵ , MinPts	Bandwidth	k , covar
Probabilistic	No	No	No	No	Yes

4 Module 4: Principal Component Analysis (PCA)

PCA tìm các thành phần chính orthogonal với phương sai lớn nhất.

$$C = \frac{1}{n-1} X^T X, \quad Cv = \lambda v, \quad Z = XW$$

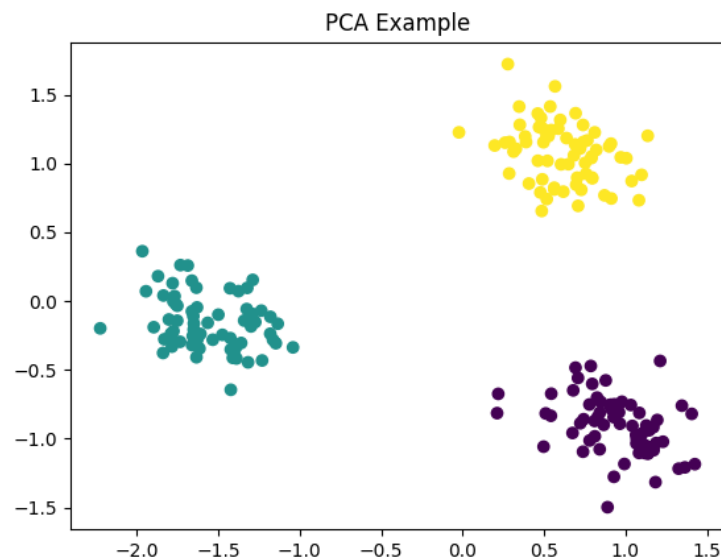


Figure 9: PCA Visualization

5 Module 5: Advanced Dimensionality Reduction Techniques

5.1 Kernel PCA

Chieu du lieu phi tuyen vao space cao chieu bang kernel, lam PCA.

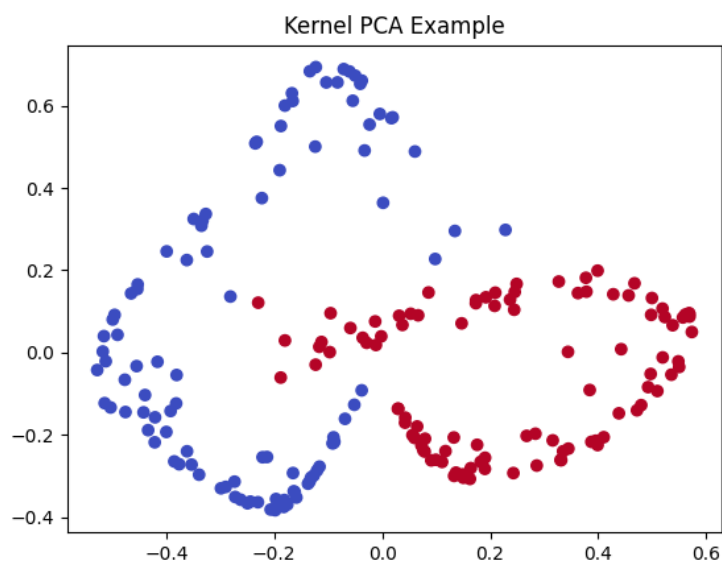


Figure 10: Kernel PCA Visualization

5.2 Multidimensional Scaling (MDS)

Giam chieu dua tren ma tran khoang cach.

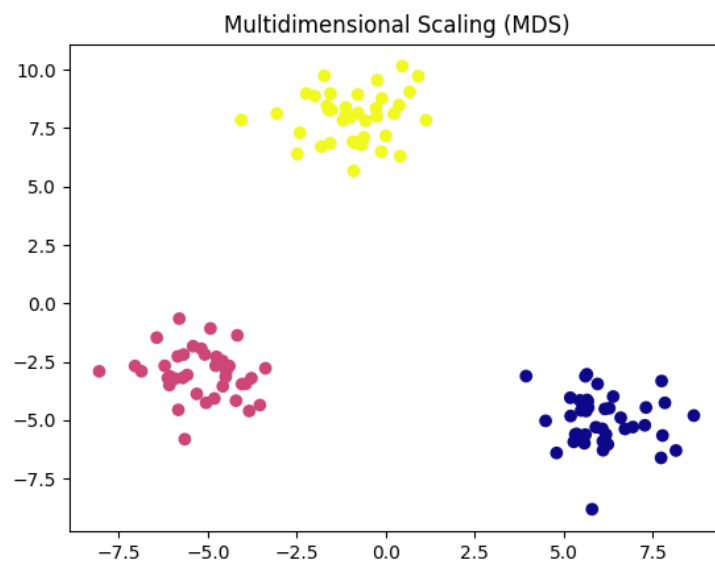


Figure 11: MDS Visualization

6 Module 6: Non-Negative Matrix Factorization (NMF)

Phân tích ma trận không âm V thành W, H cũng không âm:

$$V \approx WH$$

Quy tắc cập nhật multiplicative:

$$W \leftarrow W \odot \frac{VH^T}{WHH^T}, \quad H \leftarrow H \odot \frac{W^TV}{W^TWH}$$

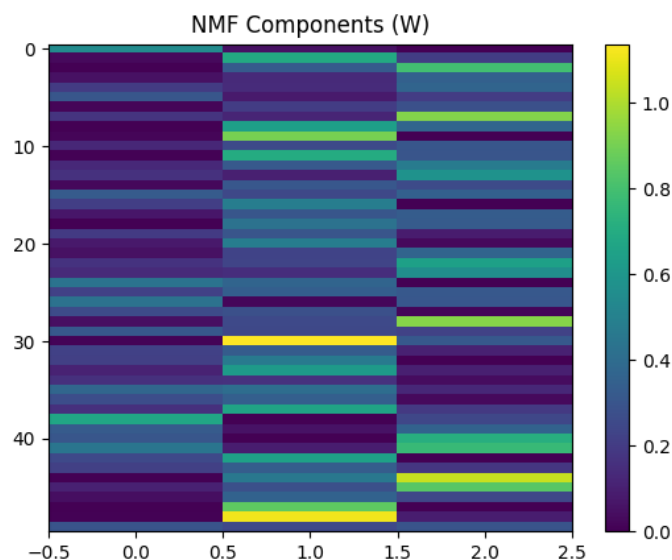


Figure 12: NMF Topic Modeling

7 Module 7: Summary and Best Practices

Tom tat cac ky thuat hoc khong giam sat, cac metric, thuat toan clustering, giam chieu:

- Clustering: K-Means, HAC, DBSCAN, Mean Shift, GMM
- Do khoang cach: Euclidean, Manhattan, Cosine, Jaccard
- Giam chieu: PCA, Kernel PCA, MDS, NMF
- Cac diem quan trong: Curse of dimensionality, chon so cum, danh gia clustering

References to images

- K-Means clustering visualization `kmeans_clustering.png`
- Distance metrics comparison `euclidean_vs_manhattan.png`, `distance_metrics_comparison.png`
- PCA visualizations `pca_visualization.png`
- Curse of dimensionality `curse_dimensionality.png`
- Cluster comparison plots
- Mean Shift clustering `mean_shift_clustering.png`
- MDS visualization `mds_voting_patterns.png`
- Kernel PCA `kernel_pca.png`
- GMM clustering `gmm_clustering.png`

- Silhouette score
- Covariance matrix