Unsupervised Learning
ooooo

Hierarchical Clustering
oooooooooooo

K Means Clustering
oooooooooooo

# CS540 Introduction to Artificial Intelligence
## Lecture 15

Young Wu

Based on lecture slides by Jerry Zhu, Yingyu Liang, and Charles Dyer

July 18, 2022

**Unsupervised Learning**
●○○○

Hierarchical Clustering
○○○○○○○○○○○

*K* Means Clustering
○○○○○○○○○○○○○

# Midterm

## Admin

- The midterm is:
- *A* : Too Easy
- *B* : Easy
- *C* : Just right
- *D* : Hard
- *E* : Too Hard

Q1

Socrative

Room CS540E

**Unsupervised Learning**
○●○○

Hierarchical Clustering
○○○○○○○○○○○

*K* Means Clustering
○○○○○○○○○○○○

# Midterm Discussion
## Admin

- Go over some new questions at the end of the lecture.

- Post the stats later in the week.

- If you are planning to take the make-up midterm, there is no need to notify me.

- Same format, join by Zoom, Q6 questions still on the exam (with different randomization).

- You can start the exam and not submit it, but if you submit, your current grade will be replaced.

**Unsupervised Learning**
○○○●○

Hierarchical Clustering
○○○○○○○○○○○

K Means Clustering
○○○○○○○○○○○○

# Unsupervised Learning
## Motivation

- Supervised learning: $(x_1, y_1), (x_2, y_2), ..., (x_n, y_n)$.
- Unsupervised learning: $x_1, x_2, ..., x_n$.
- There are a few common tasks without labels.

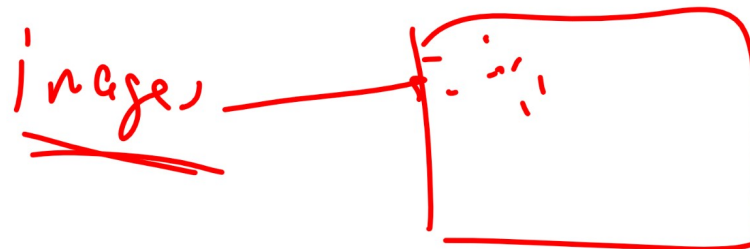1. Clustering: separate instances into groups.
2. Novelty (outlier) detection: find instances that are different.
3. Dimensionality reduction: represent each instance with a lower dimensional feature vector while maintaining key characteristics.
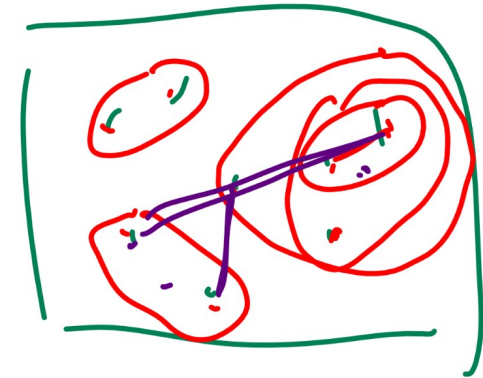
*Handwritten annotations:*

image → label, features

$y_i \approx f(x_i)$

$\Pr\{y_i | x_i\} \leftarrow \Pr\{x_i | y_i\}$

Bayes Rule

images

**Unsupervised Learning**
○○○●

Hierarchical Clustering
○○○○○○○○○○○

K Means Clustering
○○○○○○○○○○○

# Unsupervised Learning Applications
## Motivation

1. Google News
2. Google Photo
3. Image Segmentation
4. Text Processing

Unsupervised Learning
○○○○

Hierarchical Clustering
●○○○○○○○○○○

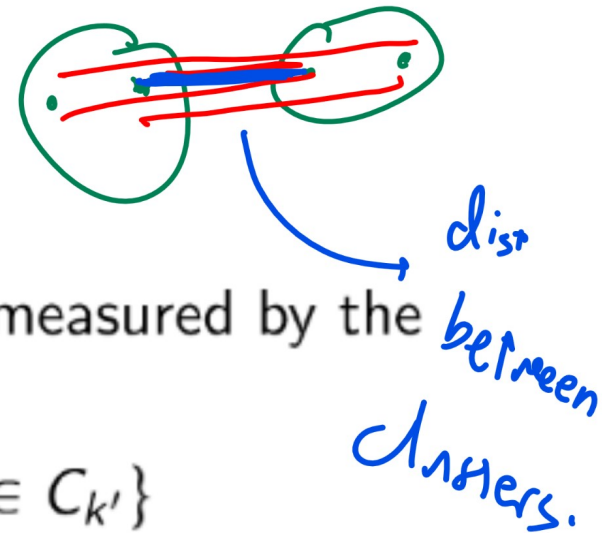K Means Clustering
○○○○○○○○○○○○

# Hierarchical Clustering
## Description

- Start with each instance as a cluster.
- Merge clusters that are closest to each other.
- Result in a binary tree with close clusters as children.

# Hierarchical Clustering Diagram

## Description

Unsupervised Learning
◦◦◦◦

Hierarchical Clustering
◦◦●◦◦◦◦◦◦◦◦◦

K Means Clustering
◦◦◦◦◦◦◦◦◦◦◦◦

# Single Linkage Distance

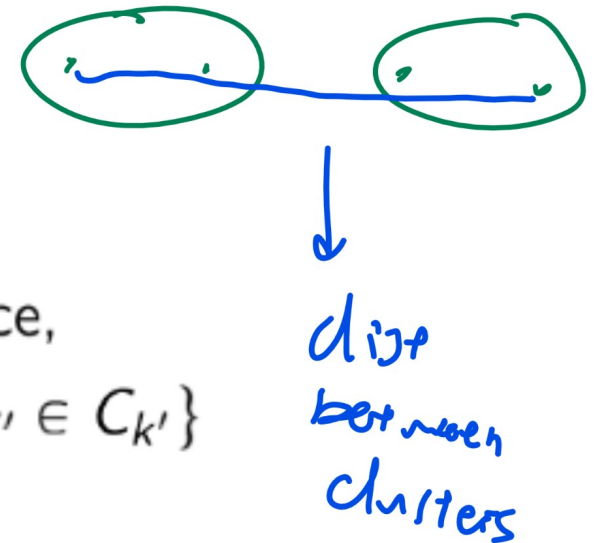### Definition

dist between clusters.

- Usually, the distance between two clusters is measured by the single-linkage distance.

$$d\left(C_k, C_{k'}\right) = \min \left\{ d\left(x_i, x_{i'}\right) : x_i \in C_k, x_{i'} \in C_{k'} \right\}$$

- It is the shortest distance from any instance in one cluster to any instance in the other cluster.

Unsupervised Learning
○○○○

Hierarchical Clustering
○○○●○○○○○○○

K Means Clustering
○○○○○○○○○○○○

# Complete Linkage Distance

## Definition

- Another measure is complete-linkage distance,

$$d\left(C_k, C_{k'}\right) = \max\left\{d\left(x_i, x_{i'}\right) : x_i \in C_k, x_{i'} \in C_{k'}\right\}$$

*dist between clusters*

- It is the longest distance from any instance in one cluster to any instance in the other cluster.

Unsupervised Learning
oooo

Hierarchical Clustering
oooo●ooooooo

K Means Clustering
oooooooooooo

# Average Linkage Distance Diagram
## Definition

- Another measure is average-linkage distance.

$$d\left(C_k, C_{k'}\right) = \frac{1}{|C_k||C_{k'}|} \sum_{x_i \in C_k, x_{i'} \in C_{k'}} d\left(x_i, x_{i'}\right)$$

*dist between clusters*

- It is the average distance from any instance in one cluster to any instance in the other cluster.

Unsupervised Learning
○○○○

Hierarchical Clustering
○○○○○●○○○○○○

K Means Clustering
○○○○○○○○○○○○

# Hierarchical Clustering 1

Quiz

- Given three clusters $A = \{0, 2, 6\}$, $B = \{3, 9\}$, $C = \{11\}$. What is the next iteration of hierarchical clustering with Euclidean distance and single and complete linkage?

single
$$\text{dist } A, B = 1$$
$$AC = 5$$
$$BC = 2$$

complete
$$\text{dist } A, B = 9$$
$$A, C = 11$$
$$B, C = 8$$

next step merge $\boxed{A, B}$, C

next merge $\boxed{B, C}$, A

Unsupervised Learning
○○○○

Hierarchical Clustering
○○○○○○●○○○○○

K Means Clustering
○○○○○○○○○○○○○

# Hierarchical Clustering 2

## Quiz

- Given three clusters $A = \{0, 1\}$, $B = \{4, 6\}$, $C = \{8\}$. What is the next iteration of hierarchical clustering with Euclidean distance and complete linkage?

  max dist.

  - $A$ : Merge $A$ and $B$.
  - $B$ : Merge $A$ and $C$.
  - $C$ : Merge $B$ and $C$.
  - $D$ : I don't understand.

Q2

dist $A, B = 6$

$A, C = 8$

$B, C = 4$

Unsupervised Learning

○○○○

Hierarchical Clustering
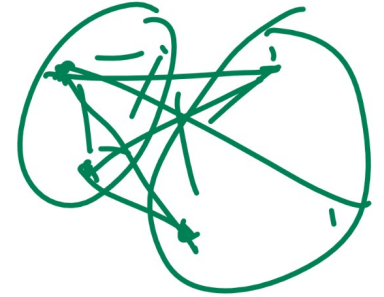
○○○○○○○●○○○○

K Means Clustering

○○○○○○○○○○○○○

# Hierarchical Clustering 3

## Quiz

- Spring 2018 Midterm $Q5$
- Given three clusters $A = \{0, 1\}$, $B = \{4, 6\}$, $C = \{8\}$. What is the next iteration of hierarchical clustering with Euclidean distance and single linkage?
- $A$ : Merge $A$ and $B$.
- $B$ : Merge $A$ and $C$.
- $C$ : Merge $B$ and $C$.
- $D$ : I don't understand.

Unsupervised Learning
ooooo

Hierarchical Clustering
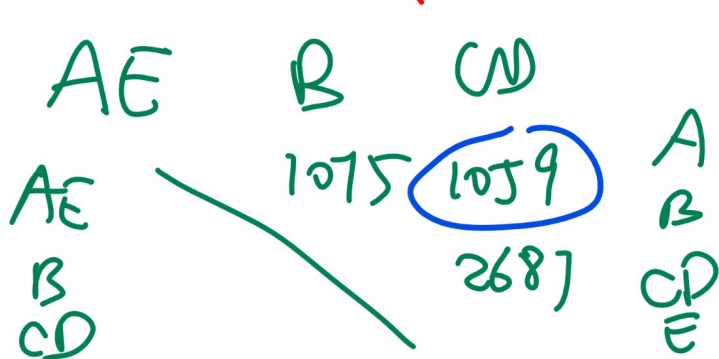ooooooooo●ooo

K Means Clustering
oooooooooooooo

# Hierarchical Clustering 4

## Quiz

- Given the distance between the clusters so far. Which pair of clusters will be merged using single linkage.

| —  | A    | B    | C    | D    | E    |
|----|------|------|------|------|------|
| A  | 0    | 1075 | 2013 | 2054 | 996  |
| B  | 1075 | 0    | 3272 | 2687 | 2037 |
| C  | 2013 | 3272 | 0    | 808  | 1307 |
| D  | 2054 | 2687 | 808  | 0    | 1059 |
| E  | 996  | 2037 | 1307 | 1059 | 0    |

*(handwritten annotations)* merge CD

AE    B    CD
AE         1075   1059
B                 2687
CD

A    B    CD    E
     1075  2013  996
           2687  2037
                 1059

merge AE

A
B
CD
E

Unsupervised Learning
ooooo

Hierarchical Clustering
ooooooooo●oo

K Means Clustering
oooooooooooooo

# Hierarchical Clustering 4, Diagram

## Quiz

| $-$ | A | B | C | D | E |
|---|---|---|---|---|---|
| A | 0 | 1075 | 2013 | 2054 | 996 |
| B | 1075 | 0 | 3272 | 2687 | 2037 |
| C | 2013 | 3272 | 0 | 808 | 1307 |
| D | 2054 | 2687 | 808 | 0 | 1059 |
| E | 996 | 2037 | 1307 | 1059 | 0 |

Unsupervised Learning
○○○○

Hierarchical Clustering
○○○○○○○○○●○

K Means Clustering
○○○○○○○○○○○○

# Hierarchical Clustering 5

## Quiz

Q3

- Given the distance between the clusters so far. Which pair of clusters will be merged using complete linkage.

merge
CD

| — | A | B | C | D |
|---|---|---|---|---|
| A | 0 | 1075 | 2013 | 2054 |
| B | 1075 | 0 | 3272 | 2687 |
| C | 2013 | 3272 | 0 | 808 |
| D | 2054 | 2687 | 808 | 0 |

$dist(A, CD)$

$\max d(A,C), d(A,D)$

- E : I don't understand.

Q4

A
B
CD → C

merge A, B

|  | A | B | CD |
|---|---|---|---|
| A | 0 | 1075 | 2054 |
| B | 1075 | 0 | 3272 |
| CD | 2054 | 3272 | 0 |

Unsupervised Learning
○○○○

Hierarchical Clustering
○○○○○○○○○○○●

K Means Clustering
○○○○○○○○○○○○○

# Number of Clusters

## Discussion

- $K$ can be chosen using prior knowledge about $X$.

- The algorithm can stop merging as soon as all the between-cluster distances are larger than some fixed $R$.

- The binary tree generated in the process is often called dendrogram, or taxonomy, or a hierarchy of data points.

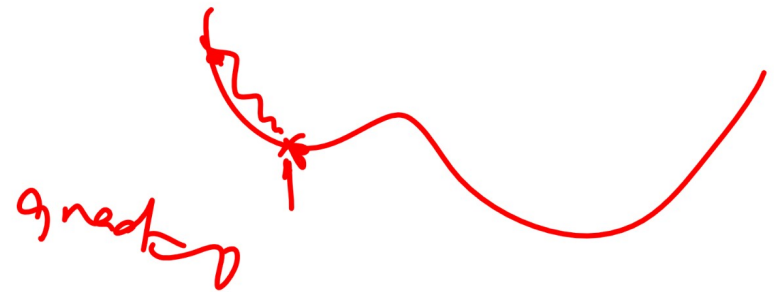- An example of a dendrogram is the tree of life in biology.

Unsupervised Learning
oooo

Hierarchical Clustering
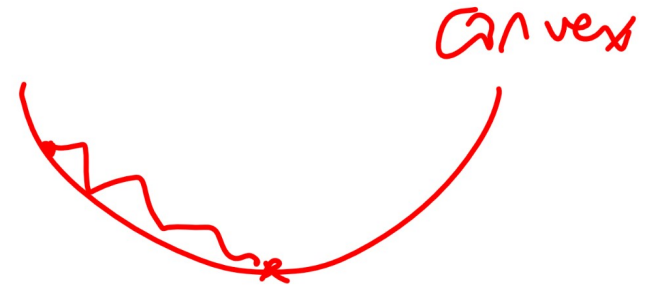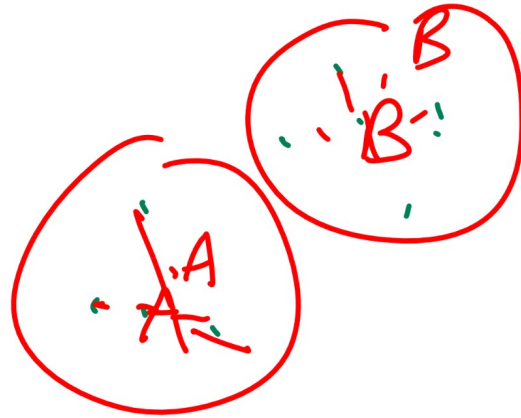oooooooooooo

K Means Clustering
●ooooooooooo

# K Means Clustering

## Description

- This is not $K$ Nearest Neighbor.
- Start with random cluster centers.
- Assign each point to its closest center.
- Update all cluster centers as the center of its points.

Unsupervised Learning
○○○○

Hierarchical Clustering
○○○○○○○○○○○

K Means Clustering
○●○○○○○○○○○○

# K Means Clustering Demo
## Description

Unsupervised Learning
○○○○

Hierarchical Clustering
○○○○○○○○○○○○

K Means Clustering
○○●○○○○○○○○○○

# Distortion

## Distortion

- Distortion for a point is the distance from the point to its cluster center.

- Total distortion is the sum of distortion for all points.

$$D_K = \sum_{i=1}^{n} d\left(x_i, c_{k^\star(x_i)}(x_i)\right)$$

$$k^\star(x) = \underset{k=1,2,\ldots K}{\operatorname{argmin}} \; d(x, c_k)$$

Cost

Unsupervised Learning

○○○○

Hierarchical Clustering

○○○○○○○○○○○○○

*K* Means Clustering

○○○●○○○○○○○○

# Objective Function Counterexample

## Definition

Unsupervised Learning
○○○○

Hierarchical Clustering
○○○○○○○○○○○○

K Means Clustering
○○○○●○○○○○○○

# Gradient Descent

## Definition

- When $d$ is the Euclidean distance. $K$ Means algorithm is the gradient descent when distortion is the objective (cost) function.

$$\frac{\partial}{\partial c_k} \sum_{k=1}^{K} \sum_{x \in C_k} \|x - c_k\|_2^2 = 0$$

$$\Rightarrow -2 \sum_{x \in C_k} (x - c_k) = 0$$

$$\Rightarrow c_k = \frac{1}{|C_k|} \sum_{x \in C_k} x$$

$$w \ = \ w - \frac{\partial C}{\partial w}$$

Unsupervised Learning
○○○○

Hierarchical Clustering
○○○○○○○○○○○○

K Means Clustering
○○○○○●○○○○○○

# K Means Clustering 1

## Quiz

- Given data $x = \{-1, 0, 2\}$ and initial cluster centers $c_1 = 0, c_2 = 1$, what is the initial clusters and what is the initial total distortion (sum of squares without square root)?

$9\sqrt{3}$

| | dist to $C_1$ | dist to $C_2$ | cluster | distortion |
|---|---|---|---|---|
| $-1$ | $1$ | $2$ | $C_1$ | $1^2$ |
| $0$ | $0$ | $1$ | $C_1$ | $1^2$ |
| $2$ | $2$ | $1$ | $C_2$ | $\dfrac{1^2}{3}$ |

distortion

$\left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2 + 0^2$

$= \frac{1}{2}$

$\boxed{\{-1, 0\}}$   $\boxed{\{2\}}$

$C_1' = -\frac{1}{2}$

$C_2' = 2$

Unsupervised Learning
○○○○

Hierarchical Clustering
○○○○○○○○○○○○

K Means Clustering
○○○○○○●○○○○○○

# K Means Clustering 2

## Quiz

$$\begin{pmatrix} c_{11} \\ c_{12} \\ c_{13} \end{pmatrix} \qquad \begin{pmatrix} x_{i1} \\ x_{i2} \\ x_{i3} \end{pmatrix}$$

Q5

- Given data $x = \{-1, 0, 2\}$ and initial cluster centers $c_1 = 0, c_2 = 5$, what is the initial clusters?

- $A : \{\varnothing\}$ and $\{-1, 0, 2\}$

- $B : \{-1\}$ and $\{0, 2\}$

- $C : \{-1, 0\}$ and $\{2\}$

- $D : \{-1, 0, 2\}$ and $\{\varnothing\}$

- $E : I$ don't understand.

$\text{dist}_1 \qquad \text{dist}_2 \qquad ?$

$\begin{array}{cccc} -1 & & & c_1 \\ 0 & & & c_1 \\ 2 & \boxed{2} & 3 & c_1 \end{array}$

$P4 \rightarrow$ k-means, try diff initial cluster

Unsupervised Learning
○○○○

Hierarchical Clustering
○○○○○○○○○○○○

K Means Clustering
○○○○○○○●○○○○

# Total Distortion 2

## Quiz

$$\frac{1}{7}$$

- Given data $x = \{-1, 0, 2\}$ and initial cluster centers $c_1 = 0, c_2 = 5$, what is the initial total distortion (sum of squares without square root)?

Q6

$$1^2 + 0^2 + 2^2 = 5$$

- A : 2
- B : 5
- C : 10
- D : 50
- E : I don't understand.

Unsupervised Learning
○○○○

Hierarchical Clustering
○○○○○○○○○○○○

*K* Means Clustering
○○○○○○○○○●○○○

# Number of Clusters

### Discussion

- There are a few ways to pick the number of clusters $K$.

1. $K$ can be chosen using prior knowledge about $X$.
2. $K$ can be the one that minimizes distortion? No, when $K = n$, distortion $= 0$.
3. $K$ can be the one that minimizes distortion $+$ regularizer.

$$K^\star = \operatorname*{argmin}_{k} \left( D_k + \lambda \cdot m \cdot k \cdot \log n \right)$$

- $\lambda$ is a fixed constant chosen arbitrarily.

Unsupervised Learning
○○○○

Hierarchical Clustering
○○○○○○○○○○○○

K Means Clustering
○○○○○○○○○●○○

# Initial Clusters

## Discussion

- There are a few ways to initialize the clusters.

1. $K$ uniform random points in $\{x_i\}_{i=1}^n$.

2. 1 uniform random point in $\{x_i\}_{i=1}^n$ as $c_1^{(0)}$, then find the farthest point in $\{x_i\}_{i=1}^n$ from $c_1^{(0)}$ as $c_2^{(0)}$, and find the farthest point in $\{x_i\}_{i=1}^n$ from the closer of $c_1^{(0)}$ and $c_2^{(0)}$ as $c_3^{(0)}$, and repeat this $K$ times.

Unsupervised Learning
OOOO

Hierarchical Clustering
OOOOOOOOOOOOO

K Means Clustering
OOOOOOOOOO●O

# Gaussian Mixture Model
## Discussion

- In $K$ means, each instance belong to one cluster with certainty.

- One continuous version is called the Gaussian mixture model: each instance belongs to one of the clusters with a positive probability.

- The model can be trained using Expectation Maximization Algorithm (EM Algorithm).

back 7:15

Unsupervised Learning
○○○○

Hierarchical Clustering
○○○○○○○○○○○○

K Means Clustering
○○○○○○○○○○○○●

# Gaussian Mixture Model Demo

## Discussion