



# Clustering

Sayan Ranu

Nick Mckeown Chair Associate Professor

Department of Computer Science and Engineer

School of AI

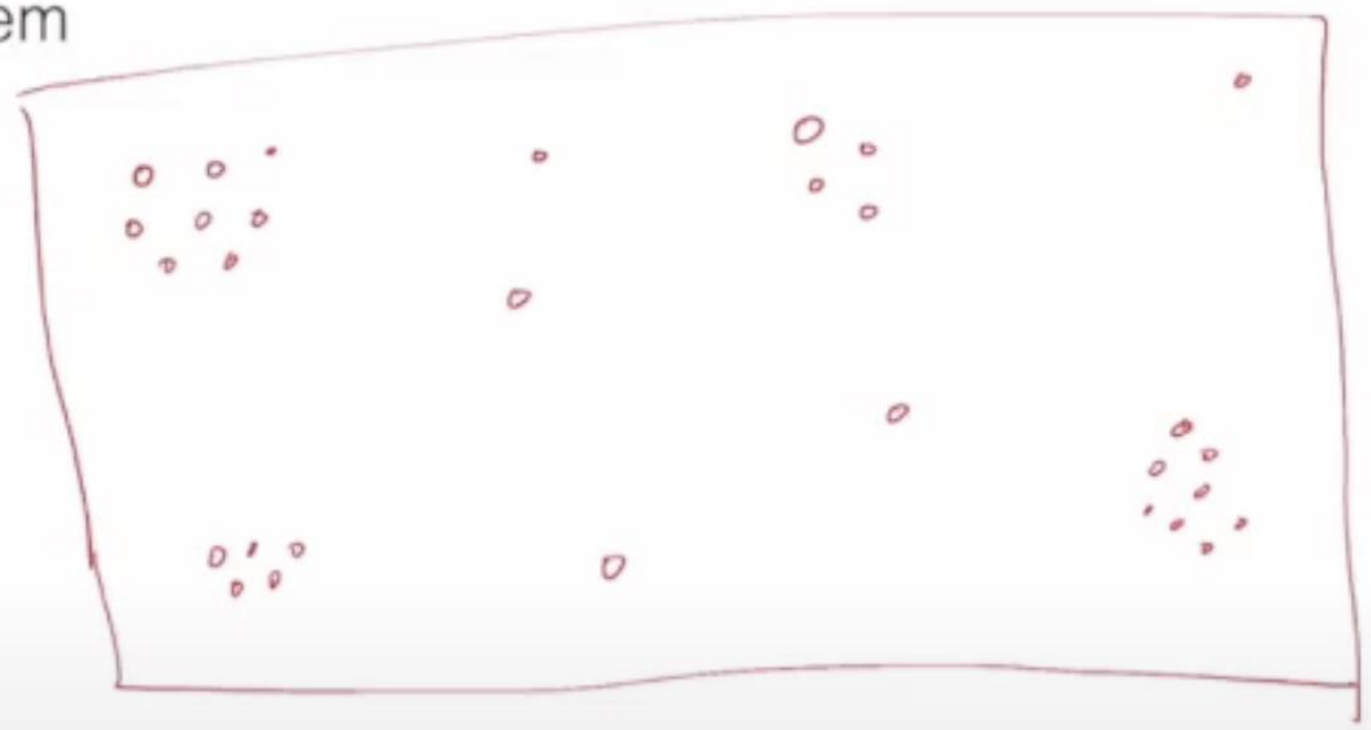
IIT Delhi

Please ask questions!



All the airports in that spa

The problem



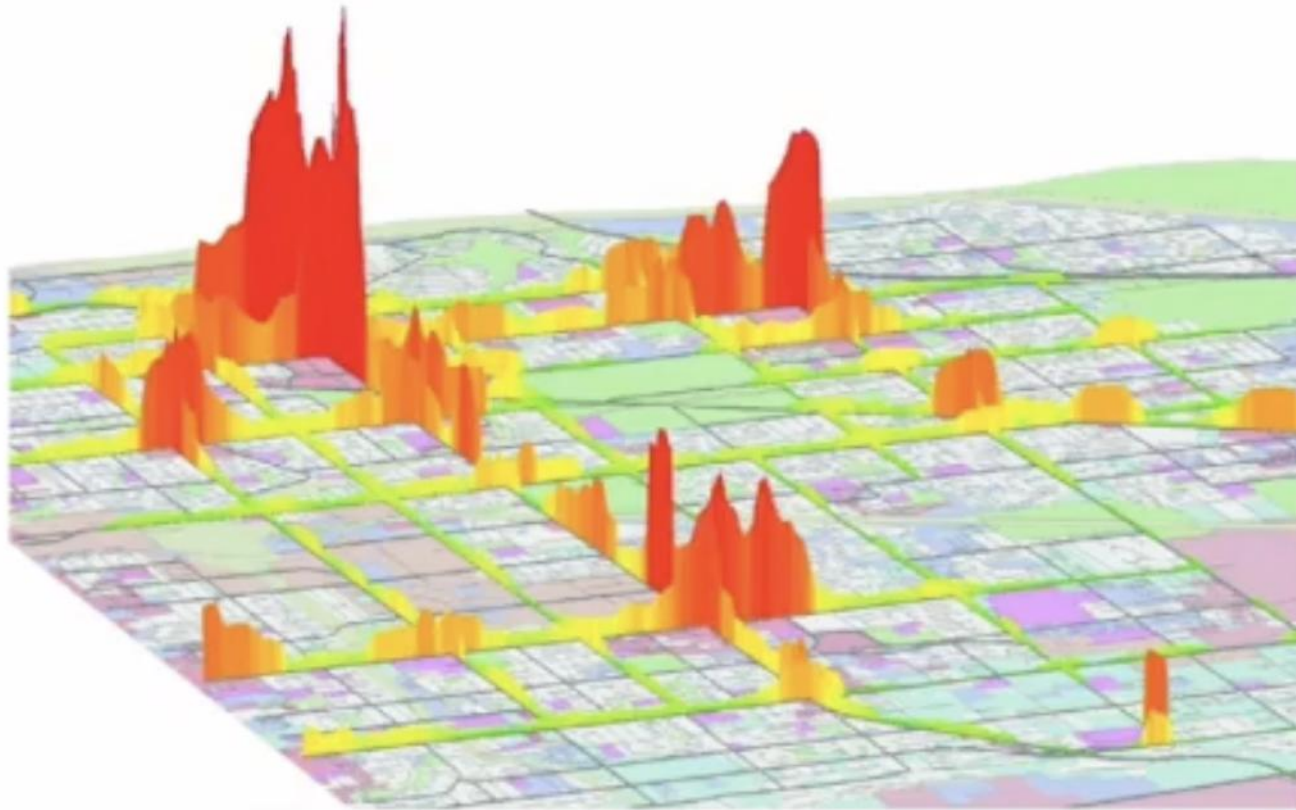
The problem

All the airports in that space  
crime event  
earthquakes

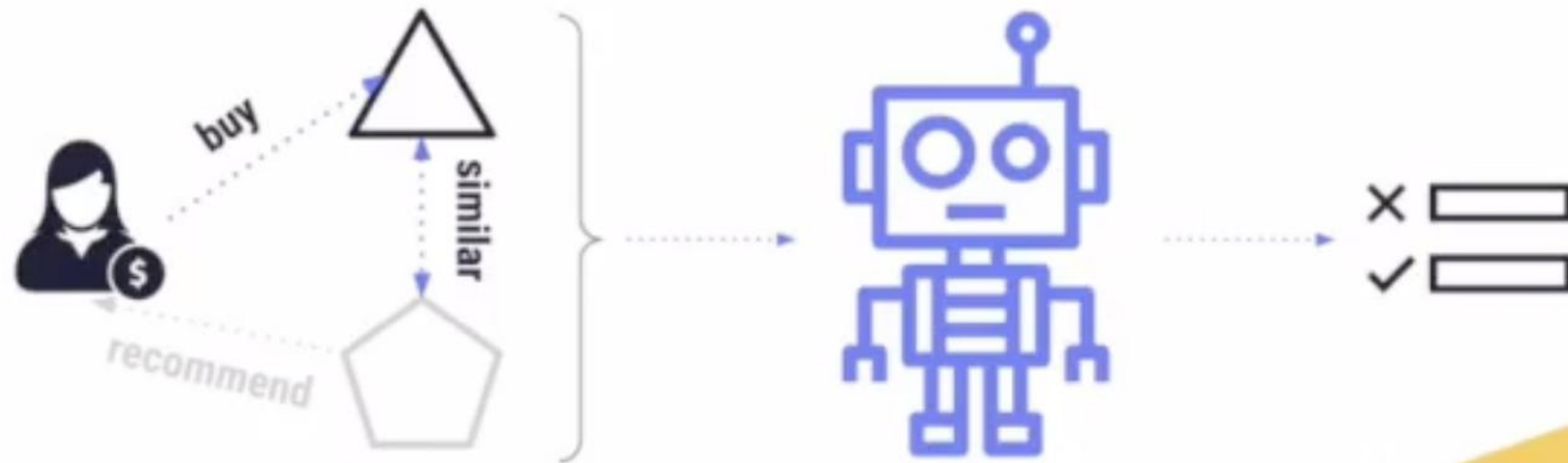


point  $(x, y)$   
(lat, long)

# Crime Hotspots



## Recommendation Engines



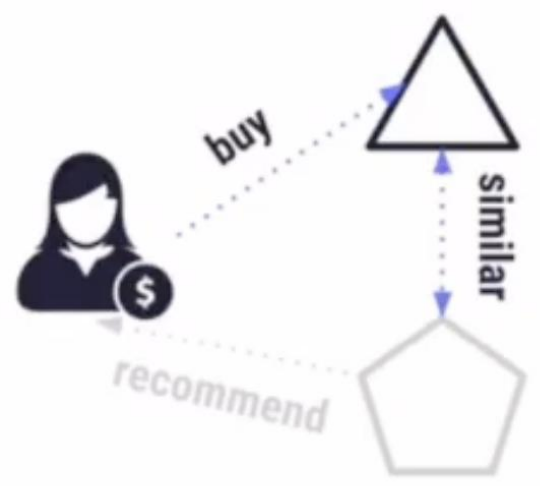
Consumption History  
Unstructured Text  
Purchasing Groups

Recommended  
Product

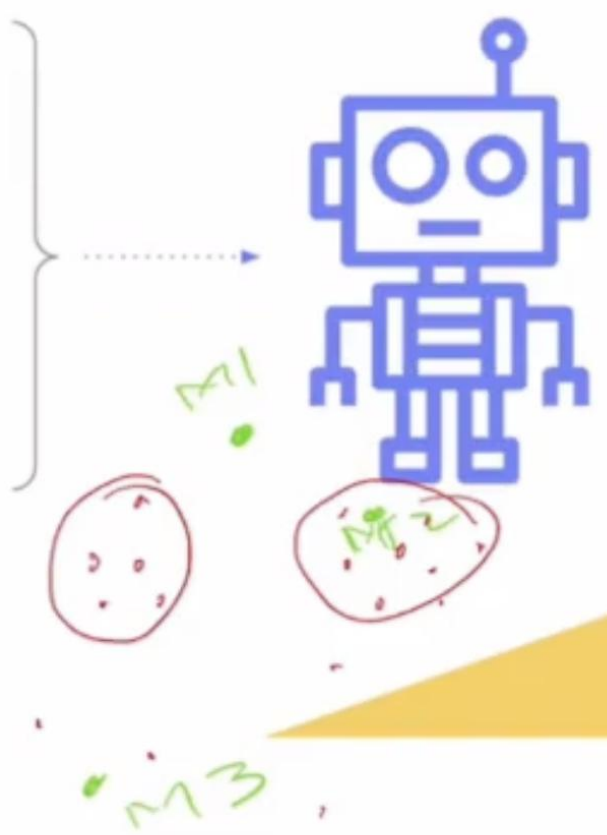
Youtube, Netflix, Amazon Prime

Length  
Language

# Recommendation Engines



Consumption History  
Unstructured Text  
Purchasing Groups



video → [ ]  
high-dimensional point

X [ ]  
✓ [ ]

Recommended Product

# Market Segmentation

Supervised  
Car →

BMW<sup>a</sup>  
customer →



[ ↓ ↓ ↓ ↓ ]  
Age Salary Education



## What do we need to cluster data?

1. Dataset of Points ✓
2. Distance Function
  - Typical to assume a *metric* distance function

close to each other



e

r

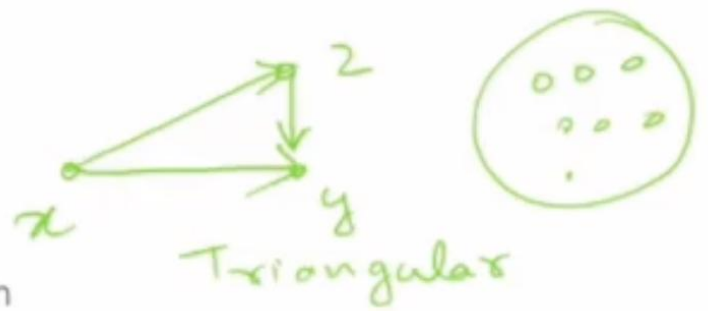
o

# What do we need to cluster data?

close to each other  
distance in km / time

1. Dataset of Points ✓✓
2. Distance Function

- Typical to assume a metric distance function



①  $d(x, y) = 0$  if  $x = y$   
 $> 0$  if  $x \neq y$

②  $d(x, y) = d(y, x)$

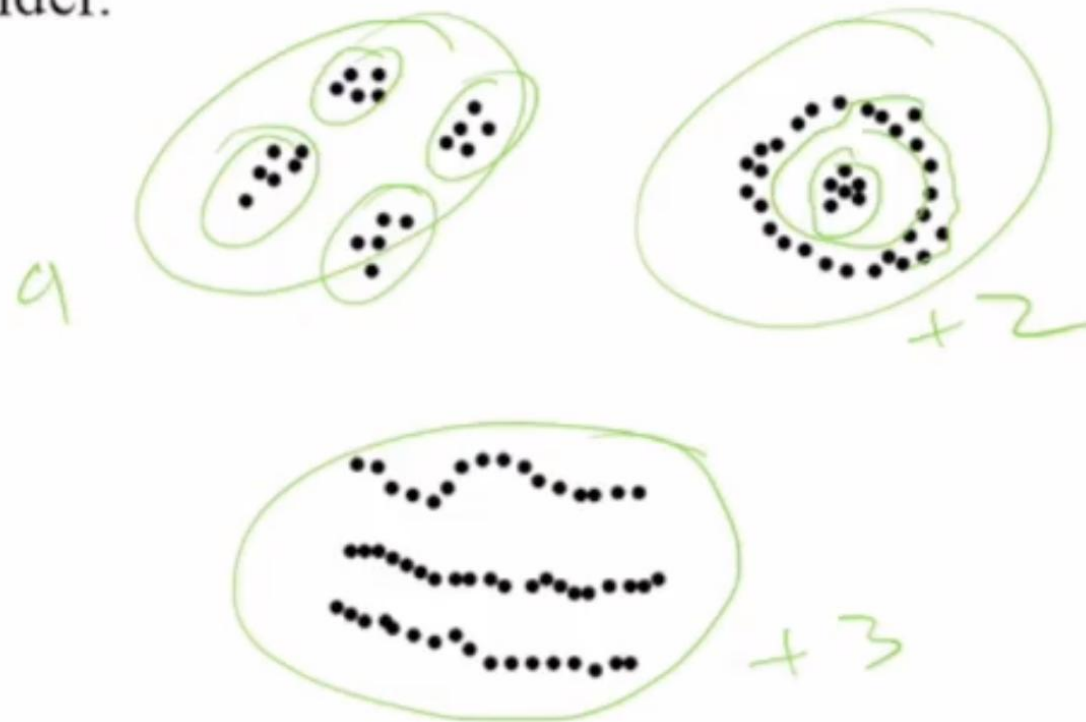
③  $d(x, y) \leq d(x, z) + d(z, y)$

Euclidean distance

$P_1 = (x_1, y_1)$   
 $P_2 = (x_2, y_2)$

$$d(P_1, P_2) = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

Clustering is an inherently ill-defined problem since the correct clusters depend upon context and are in the eye of the beholder.



How many clusters are there?

## The *K-Means* Clustering Method

- Given  $k$ , the *k-means* algorithm is implemented in four steps:
  - Partition objects into  $k$  nonempty subsets
  - Compute seed points as the centroids of the clusters of the current partition (the centroid is the center, i.e., *mean point*, of the cluster)
  - Assign each object to the cluster with the nearest seed point
  - Go back to Step 2, stop when no more new assignment

# The K-Means Clustering Method

← Input Parameters

Can you tell me the  
no of clusters in the  
dataset

• Given  $k$ , the  $k$ -means algorithm is implemented in four steps:

• Partition points into  $k$  nonempty subsets

→ avg of the points

• Compute seed points as the centroids of the clusters of the current partition (the centroid is the center, i.e., *mean point*, of the cluster)

• Assign each object to the cluster with the nearest seed point

• Go back to Step 2, stop when no more new assignment

↻ Cente

$$x_1, y_1$$

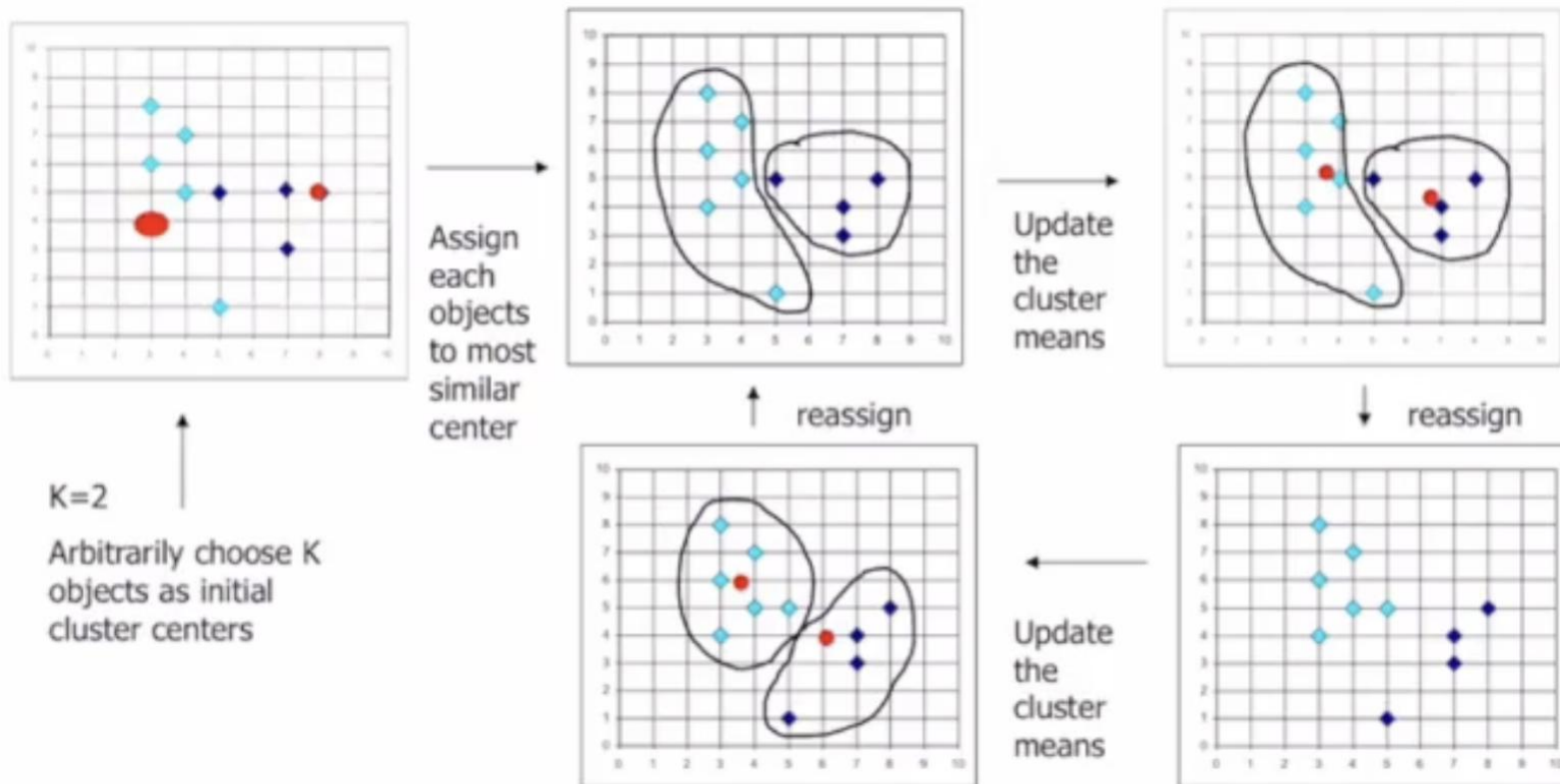
$$x_2, y_2$$

$$x_3, y_3$$

$$\text{Centroid} \left( \frac{x_1 + x_2 + x_3}{3}, \frac{y_1 + y_2 + y_3}{3} \right)$$

# The *K-Means* Clustering Method

- Example



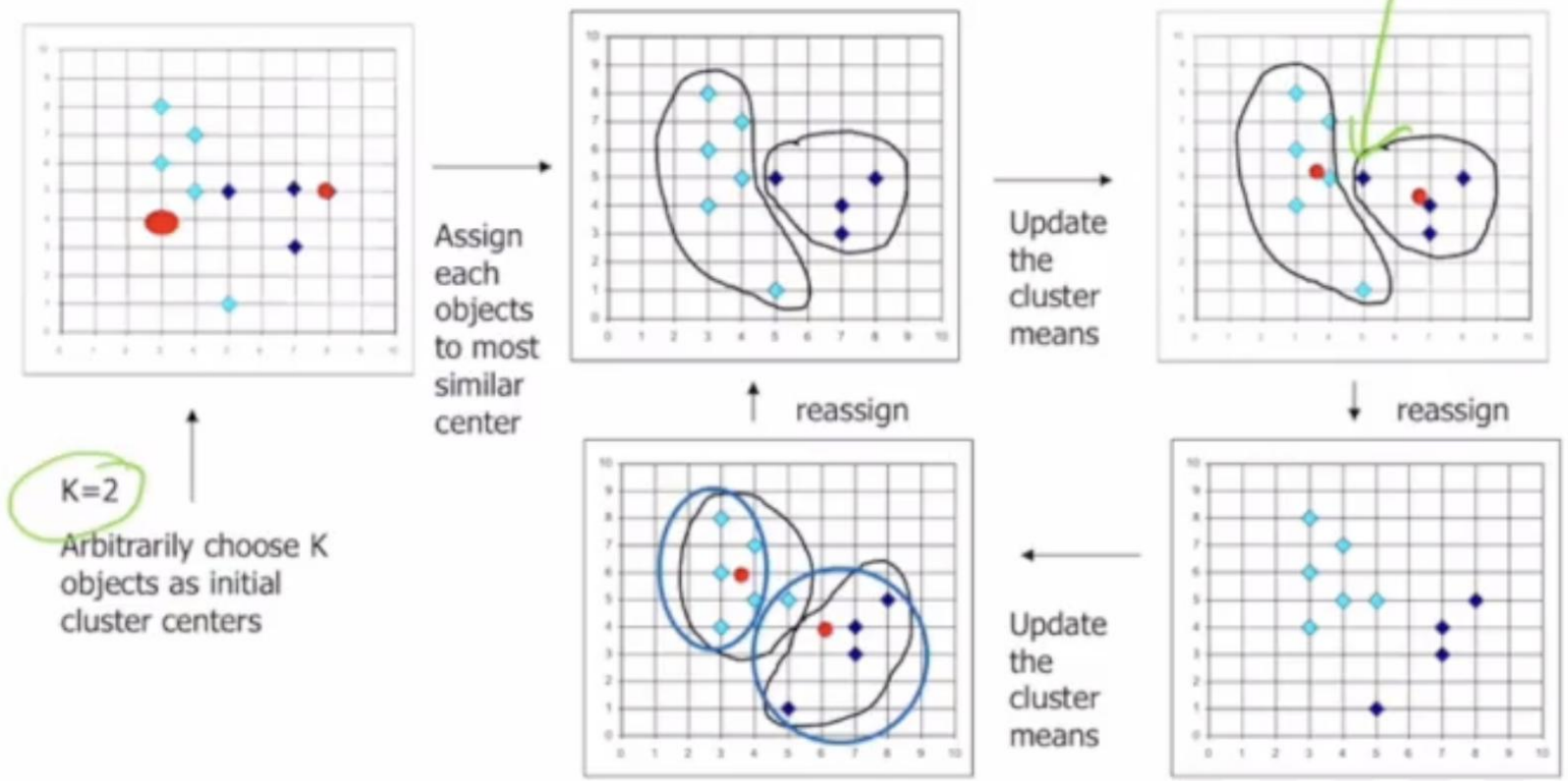
5

Time: Linear to the no. of points in the dataset

Saturation /  
Convergence

### The K-Means Clustering Method

- Example



S



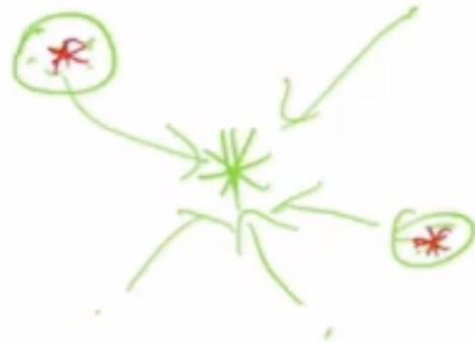
## Properties of k-means

- How do you select the initial k clusters?
  - This is a black art.
  - Just apply your heuristic and hope it works
    - Choose diverse centers
- How do you find out the value of k?
  - Identify the value of k where cluster quality is best
  - Heuristic: Elbow Plot

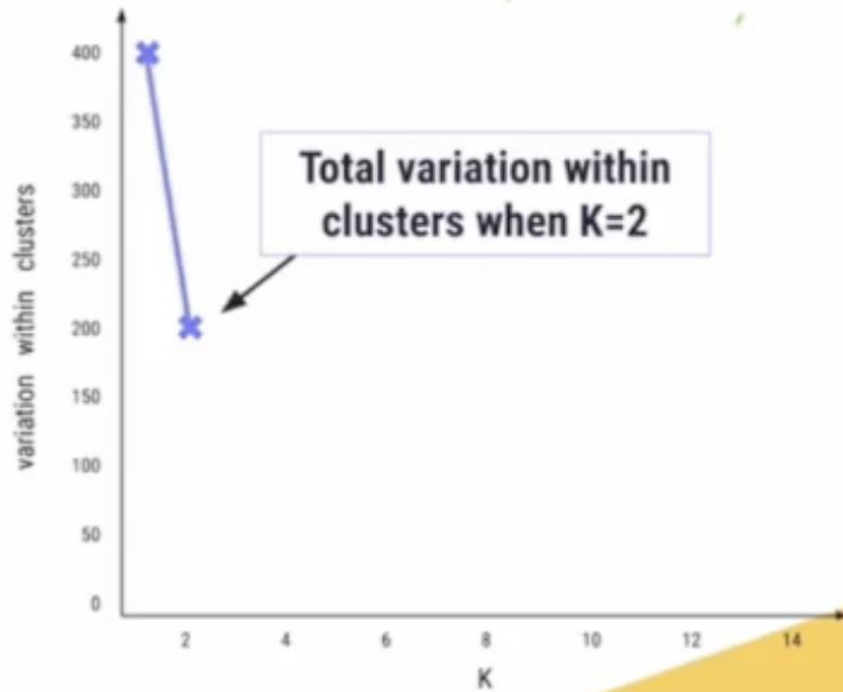




# Elbow plot



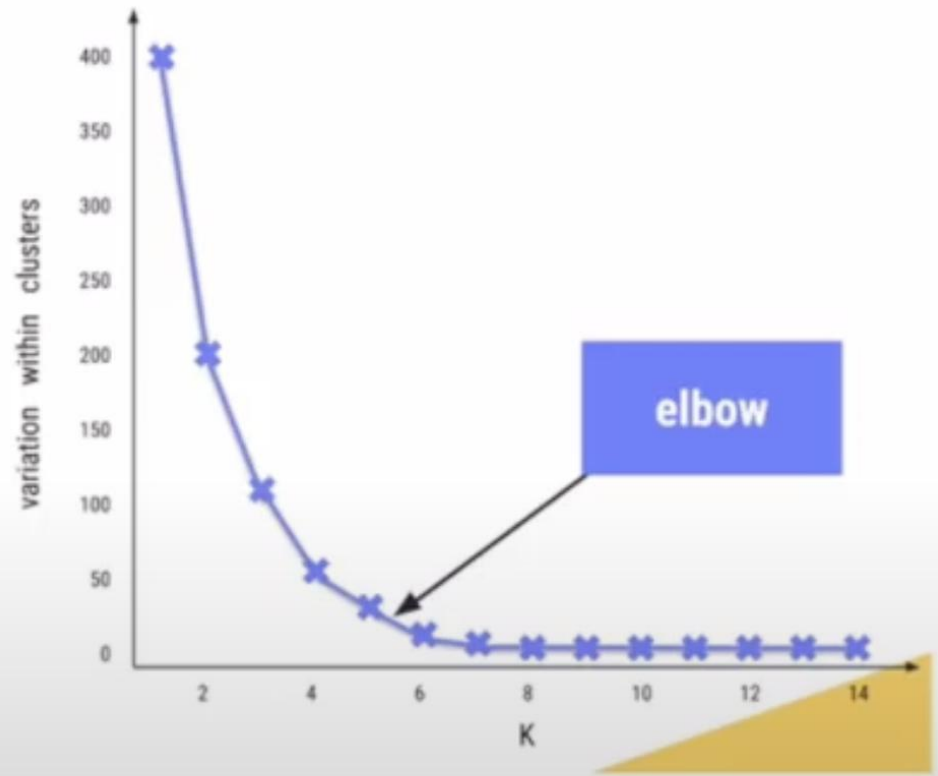
$K=1 \Rightarrow$  Variation within a cluster  
 $K=2$   
⋮  
 $K = \text{very large number}$



Variation: Avg dist from the centroid

# Elbow plot

4, 5, 6



## Weakness of k-means

- Can it identify clusters of all shapes?
  - Limited to convex shapes
- Prone to noise
  - Outliers generate spurious centroids
- How do you cluster non-vector data?
  - Text
  - Graphs
  - Time-series

