

→ Class start 9:05 PM

ML 1.2 Unsupervised ML

- Clustering (K Means, K means++, Hierarchical, DBSCAN, GMM)
- Anomaly Detection like Isolation forest etc.
- High Dimension Visualization (PCA, t-SNE, UMAP)

Unsupervised ML

→ Supervised ML → $\left\{ \begin{array}{l} \text{Features} \\ \text{target or labels, or ground truth} \end{array} \right.$

→ Unsupervised ML → Features

→ No target or label

Classification

Binary
 $\left\{ (x_i, y_i)_{i=1}^n, x_i \in \mathbb{R}^d, y_i \in (0, 1) \right\}$

Multi-class

$y_i \in S \rightarrow \text{Set of classes}$

$y_i \in \mathbb{R}$

Regression

$$\left\{ x_i^m, x_i \in \mathbb{R}^d \right\}_{i=1}^m$$

Examples of Unsupervised ML

- Anomaly Detection / Fraud Detection
- Clustering problem
- Dimensionality Reduction like PCA
- Recom. System like MF
- Word-2-vec (NLP)
- Autoencoders (CV) ✓

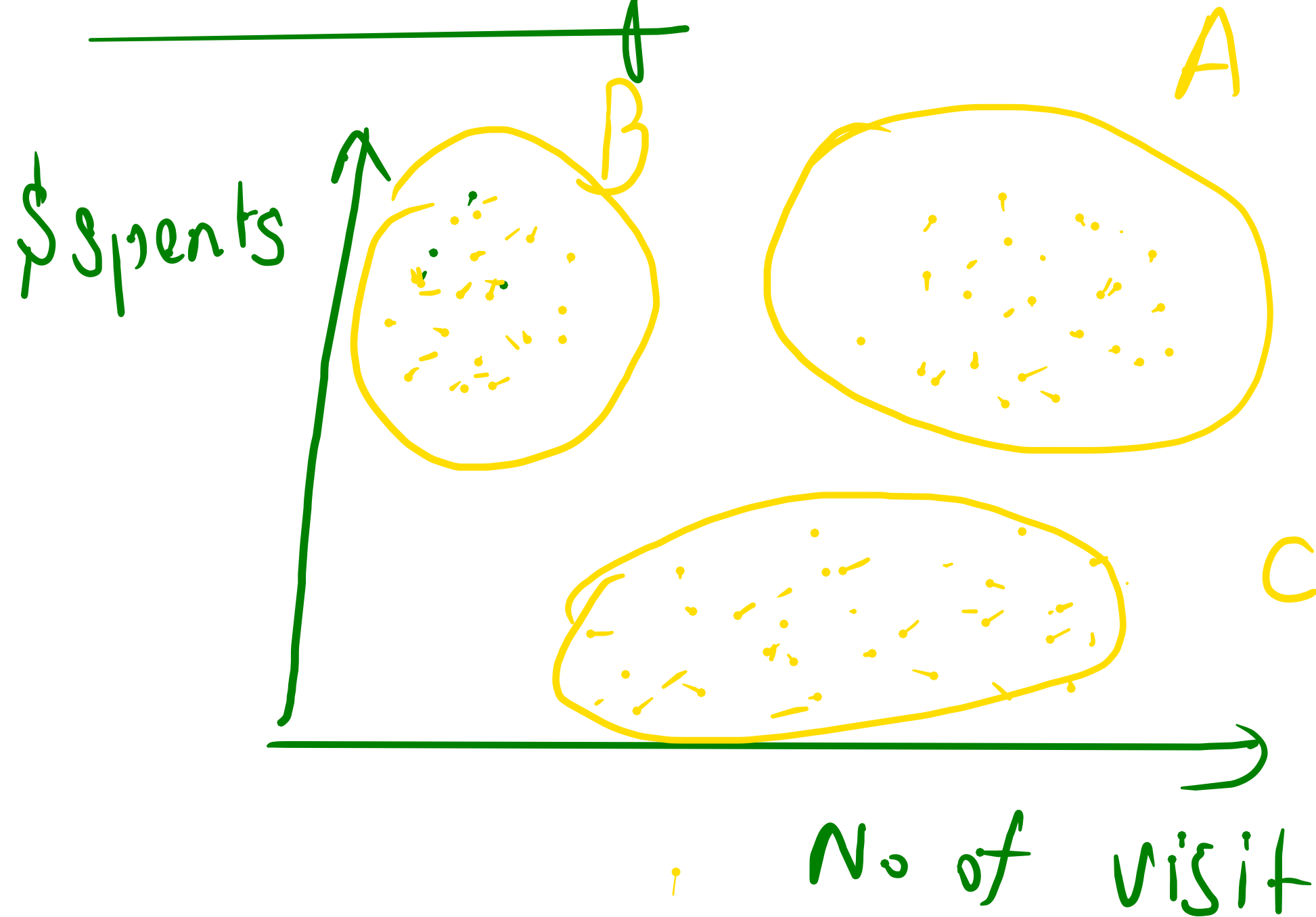
Clustering

→ Process of grouping any kind of data based on similarity of these features

Eg

- Customer Segmentation / Product Segmentation
- Detecting similar stock
- Google Photos groups similar in gallery

Clustering



A: Heavy shopper,
vist a lot, spent
a lot

B: Rich people /
Impulse buyer

C: Window shopper

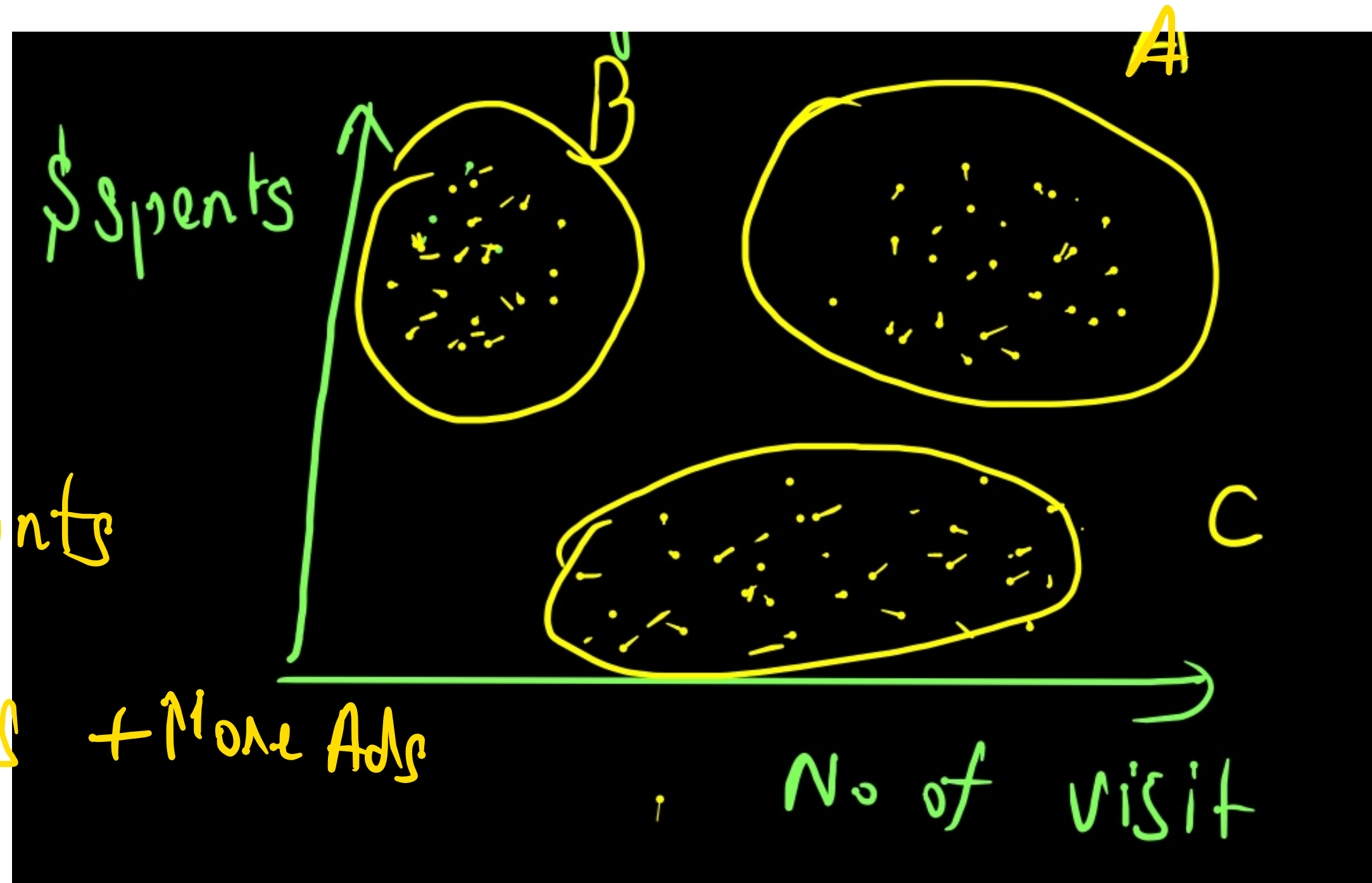
→ No of Discount coupon

→ No of Ads

B → More Ads

A → Discounts

C → More Discounts + More Ads

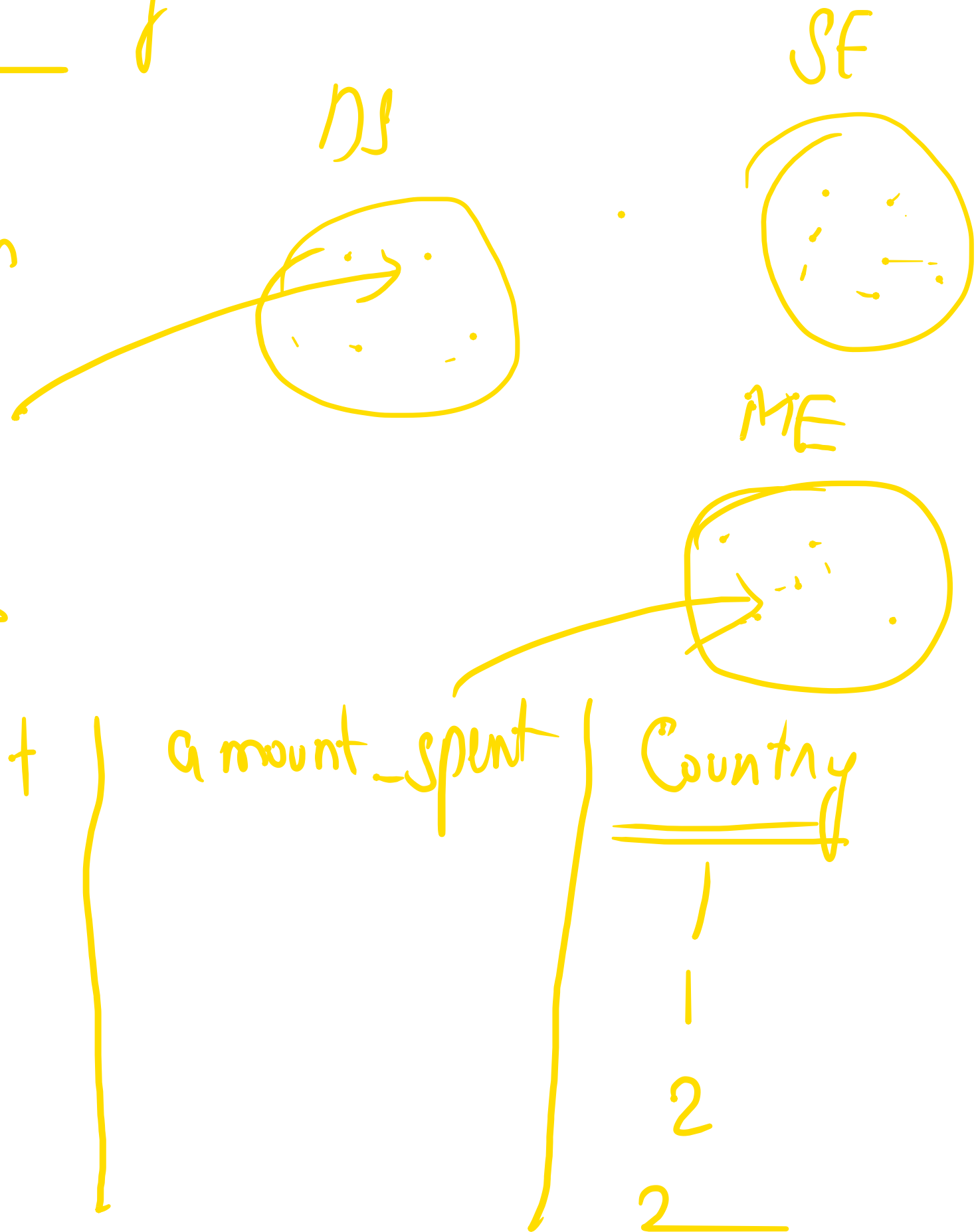


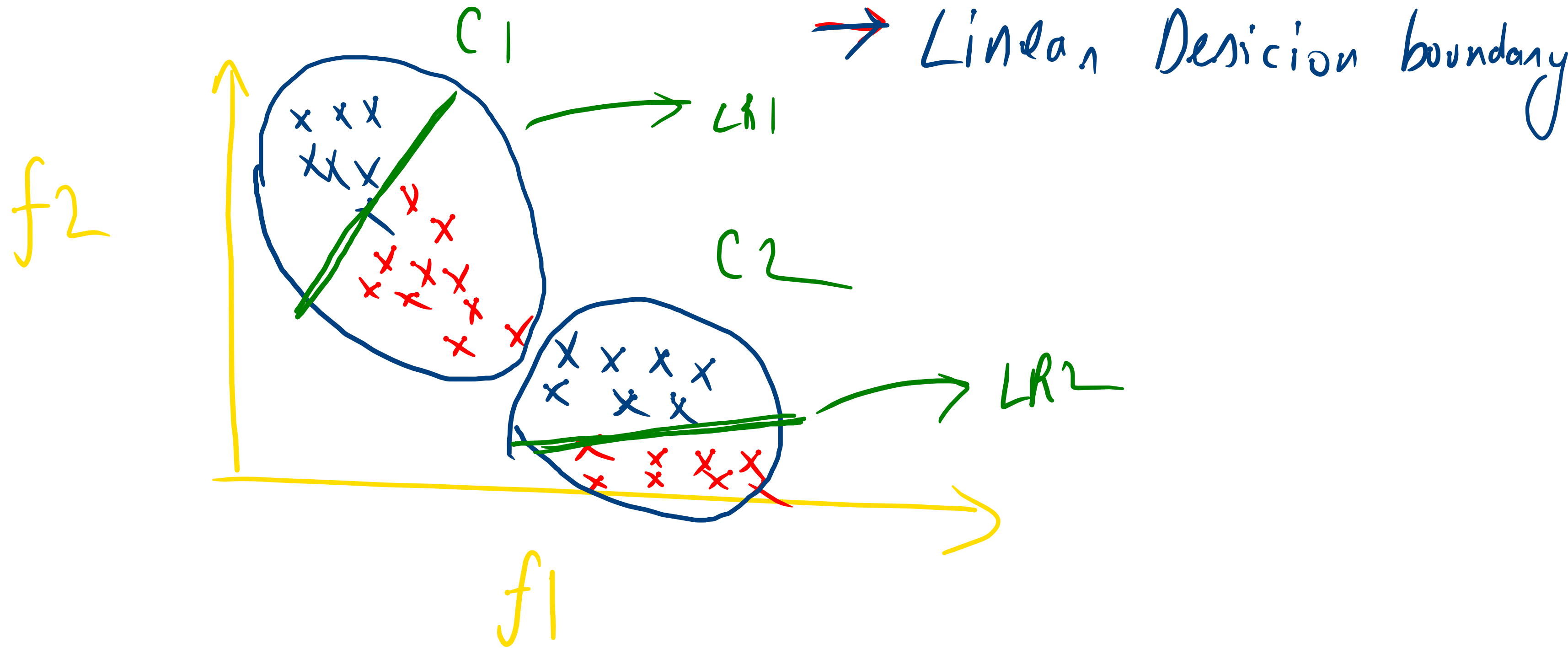
Application of Clustering

→ Search algorithm

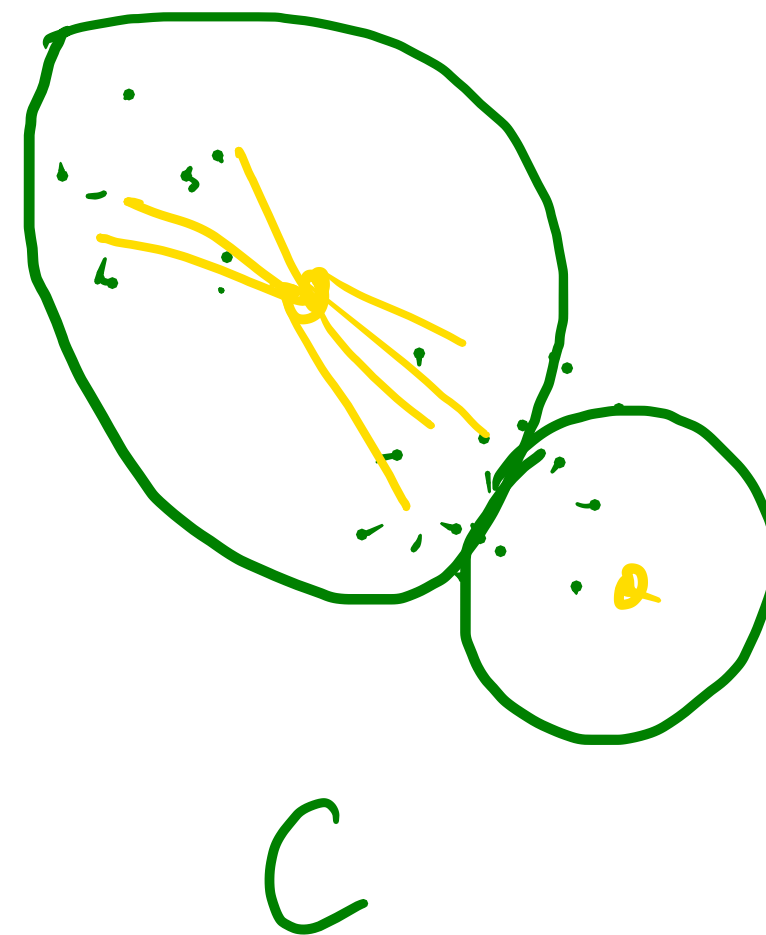
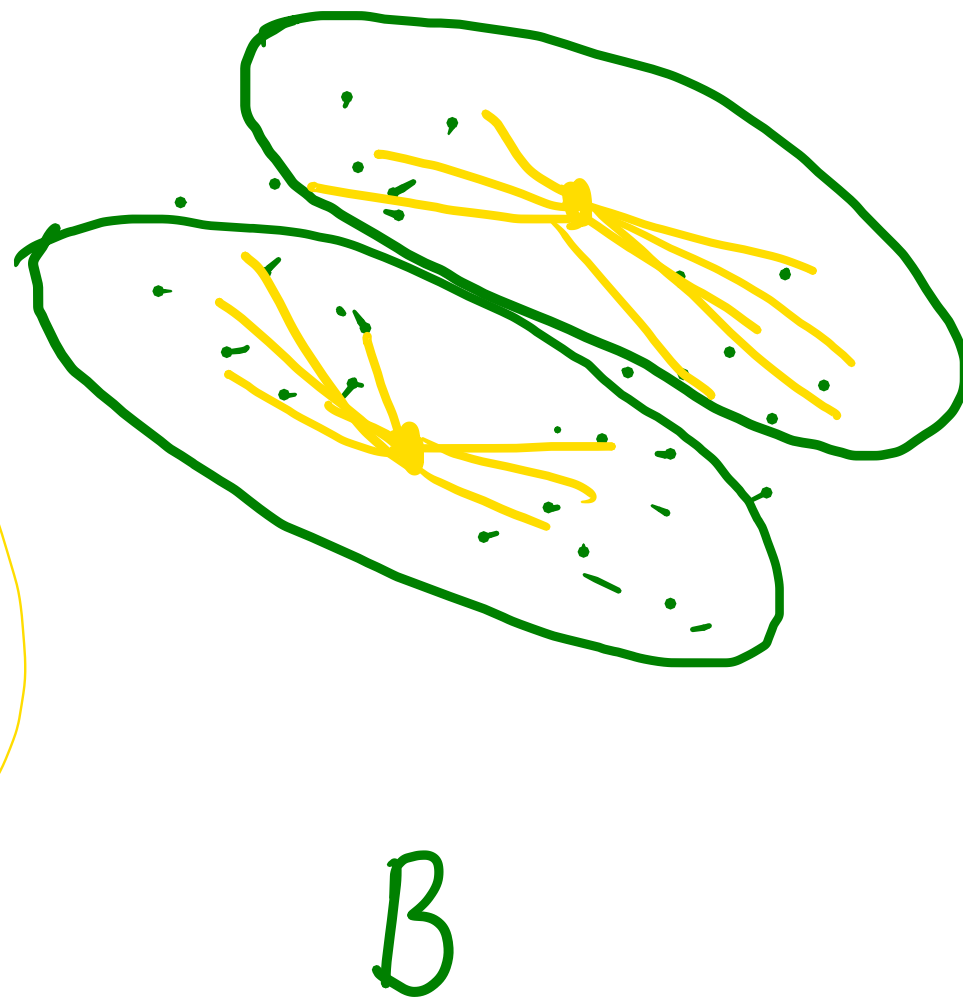
→ Feature c/lation

ID	n_clicks	n_visit	amount_spent	<u>Country</u>
				1
				1
				2
				2





Good Clustering



- Intra cluster dist minimize
- Inter cluster dist maximize

WCSS (Within Cluster Sum of Square)

$$WCSS = \sum_{P_i \text{ in } C_1} \text{dist}(P_i, C_1)^2 + \sum_{P_i \text{ in } C_2} \text{dist}(P_i, C_2)^2 + \sum_{P_i \text{ in } C_j} \text{dist}(P_i, C_j)^2$$

$$NCSS = \sum_{k=1}^K \sum_{i=1}^n 1(C_i = k) \|x_i - \underline{\underline{\mu_k}}\|^2$$

$$= \sum_{k=1}^K \sum_{i \in C_k} \|x_i - \underline{\underline{\mu_k}}\|^2$$

Gradient Descent

Update all parameters simultaneously

$$W \leftarrow W - \alpha \nabla_W L$$

↳
 w_1
 w_2
|
 w_n

Coordinate Descent

↳ Update only a subset
of parameters

↳ Fixing μ & we find
best c exactly

↳ Fixing c & find μ

K-means Clustering (Lloyd's Algorithm)

Steps:

→ Randomly initialize K centers

→ assign points to nearest center
to get your clusters ← C_i

→ find the centroids of these
clusters → because this will reduce WSS ← M_K

→ Re-assign points

→ Repeat until new centers = prev centers

Assignment
Step

Update

Elbow Method \rightarrow WCSS

\hookrightarrow Within Cluster sum of

Square

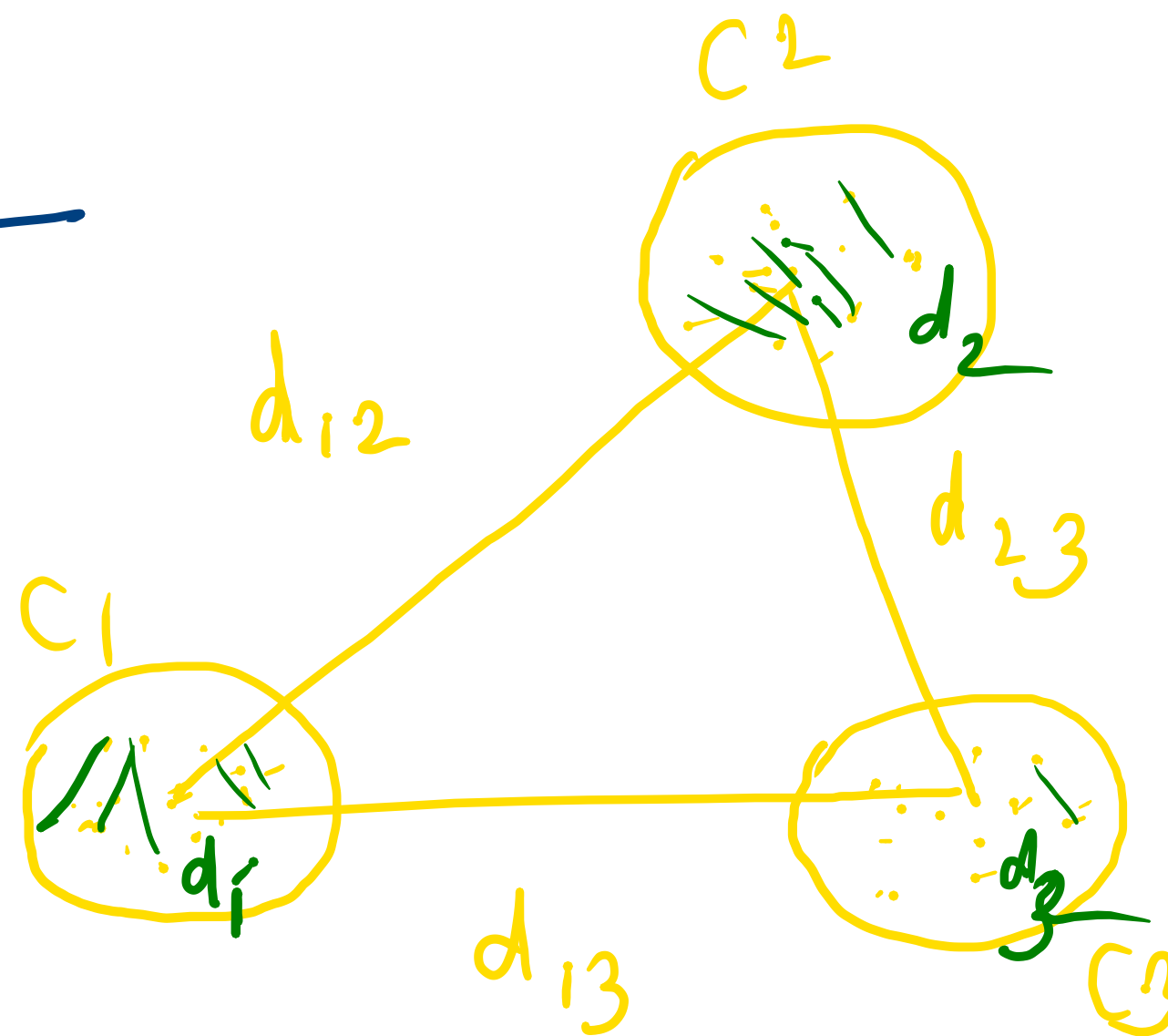
\rightarrow Intra cluster
distance

{ Dunn Index
Silhouette Score

Dunn Index

$$D = \frac{\text{Min Intercluster dist}}{\text{Max Intracluster dist}}$$

High \rightarrow good clustering



$$\text{Min Inter cluster dist} = \min(d_{12}, d_{13}, d_{23})$$

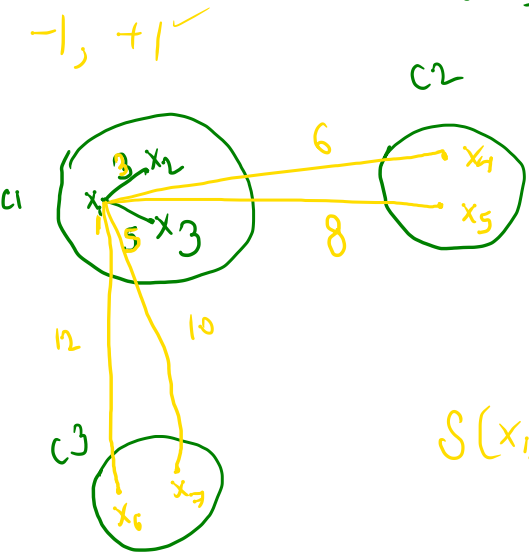
$$\text{Max Intra cluster dist} = \max(d_1, d_2, d_3)$$

Silhouette score

$$S(x_i) = \frac{b(x_i) - a(x_i)}{\max(a, b)}$$

$a(x_i)$ = Mean dist b/w i & all neighbouring point in a cluster

$b(x_i)$ = Min mean dist b/w i & all point in other cluster



$$S(x_1) = \frac{b(x_1) - a(x_1)}{\max(a(x_1), b(x_1))}$$

$$a(x_1) = \frac{3+5}{2} = 4$$

$$b(x_1) = \min\left(\frac{6+8}{2}, \frac{12+10}{2}\right)$$

$$= \min(7, 11) = 7$$

$$S(x_1) = \frac{7-4}{\max(7, 4)} = \frac{3}{7}$$

$$S(x_i) \rightarrow \begin{matrix} -1 & +1 \\ \text{~~~~~} & \text{~~~~~} \end{matrix}$$

+ve \rightarrow Good

$a > b \rightarrow S(x_i) = -ve \rightarrow$ Bad clustering

$= -1$

0 \rightarrow

Elbow Method Determine the best value of K (No. of clusters)
for K in range(10):

- K mean (K)
- WCSS



Red cluster

(1, 2)

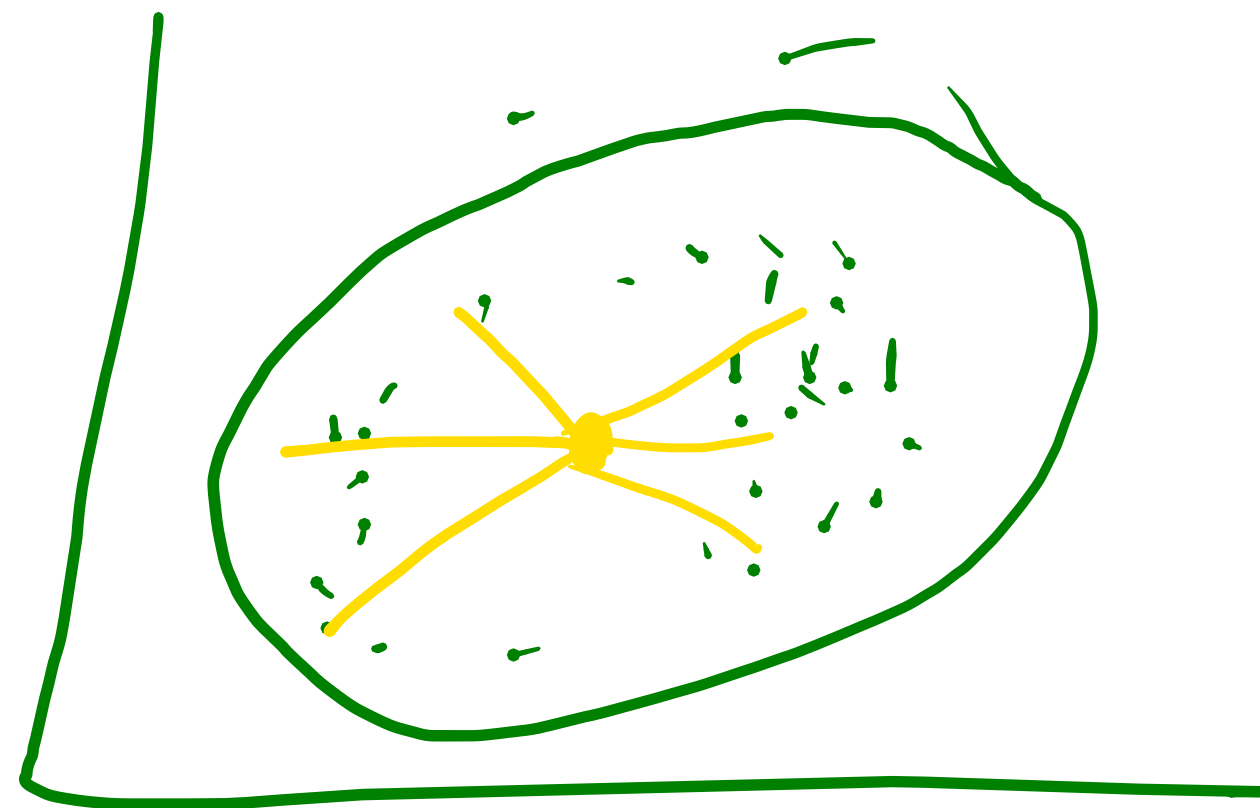
(2, 3)

(3, 4)

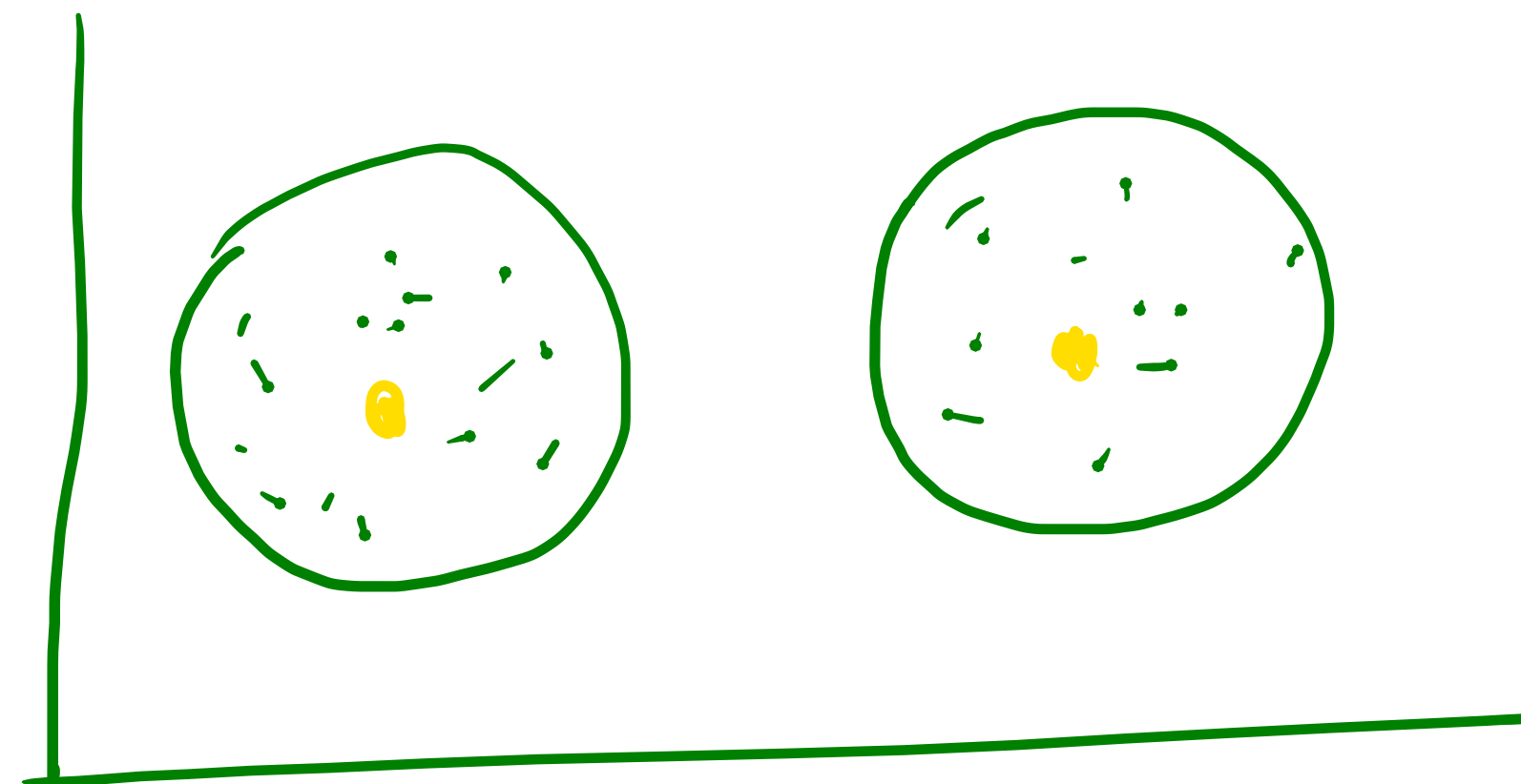
$$\mu_K = \frac{1}{n_K} \sum_{x_i \in C_K} x_i$$

$$\left(\frac{6}{3}, \frac{9}{3} \right) = (2, 3) \leftarrow \text{Centroid}$$

$$\mu_{\text{new}} - \mu_{\text{old}} < 10^{-4}$$



Cluster = 1



Cluster = 2



→ Break until 22:33

Business

- No of Ads
- Amount of Discount