→ Class start 9:05 PM

# ML 1.2    Unsupervised ML

→ Clustering (K Means, K means++, Hieanchial, DBSCAN, GMM)

→ Anomaly Detection like Isolation forest etc.

→ High Dimension Visualization (PCA, t-SNE, UMAP,

# Unsupervised ML

→ Supervised ML → $\begin{cases} \text{Features} \\ \text{target or labels, or ground} \\ \qquad\qquad\qquad\qquad \text{truth} \end{cases}$

→ Unsupervised ML → Features

→ No target or label

Classification

$$\begin{cases} \underline{\text{Binary}} \\ (x_i, y_i)_{i=1}^{m}, \ x_i \in \mathbb{R}^d, \ y_i \in (0,1) \end{cases}$$

$\underline{\text{Multi-class}}$ $y_i \in S$ → Set of Class

Regression

$y_i \in \mathbb{R}$

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

## Examples of Unsupervised ML

$\rightarrow$ Anamaly Detection / Fraud Detection

$\rightarrow$ Clustering problem

$\rightarrow$ Dimensionality Reduction like PCA

$\rightarrow$ Recom. System like MF

$\rightarrow$ Word - 2 -vec (NLP)
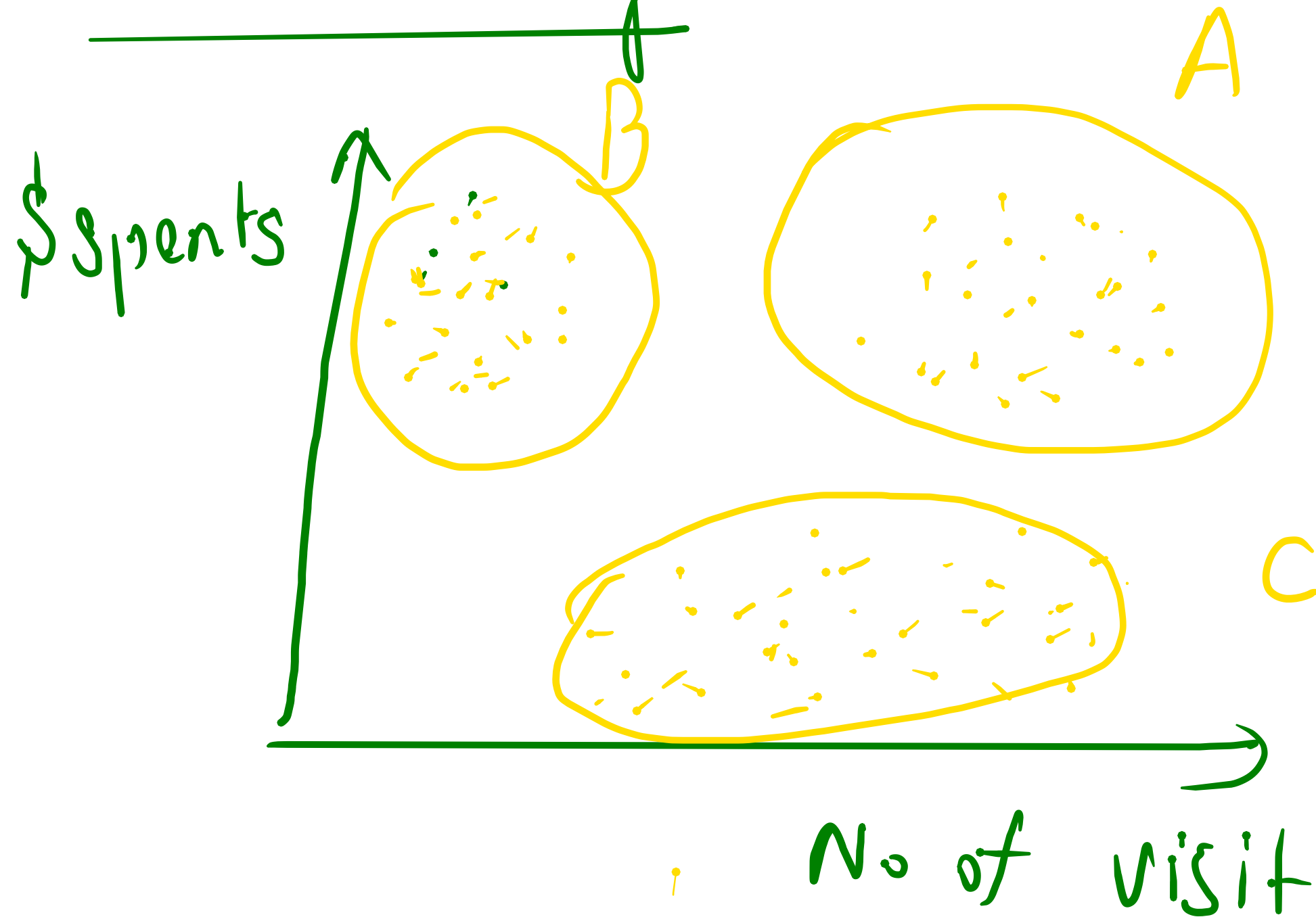
$\rightarrow$ Autoencoder (CV)

# Clustering

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

Ey
→ Customer Segmentation / Product Segmentation
→ Detecting similar stock
→ Google Photos groups similar in galaxy

# Clustering



$ spents (vertical axis)

No of visit (horizontal axis)

A
B
C

A: Heavy shopper, visit 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
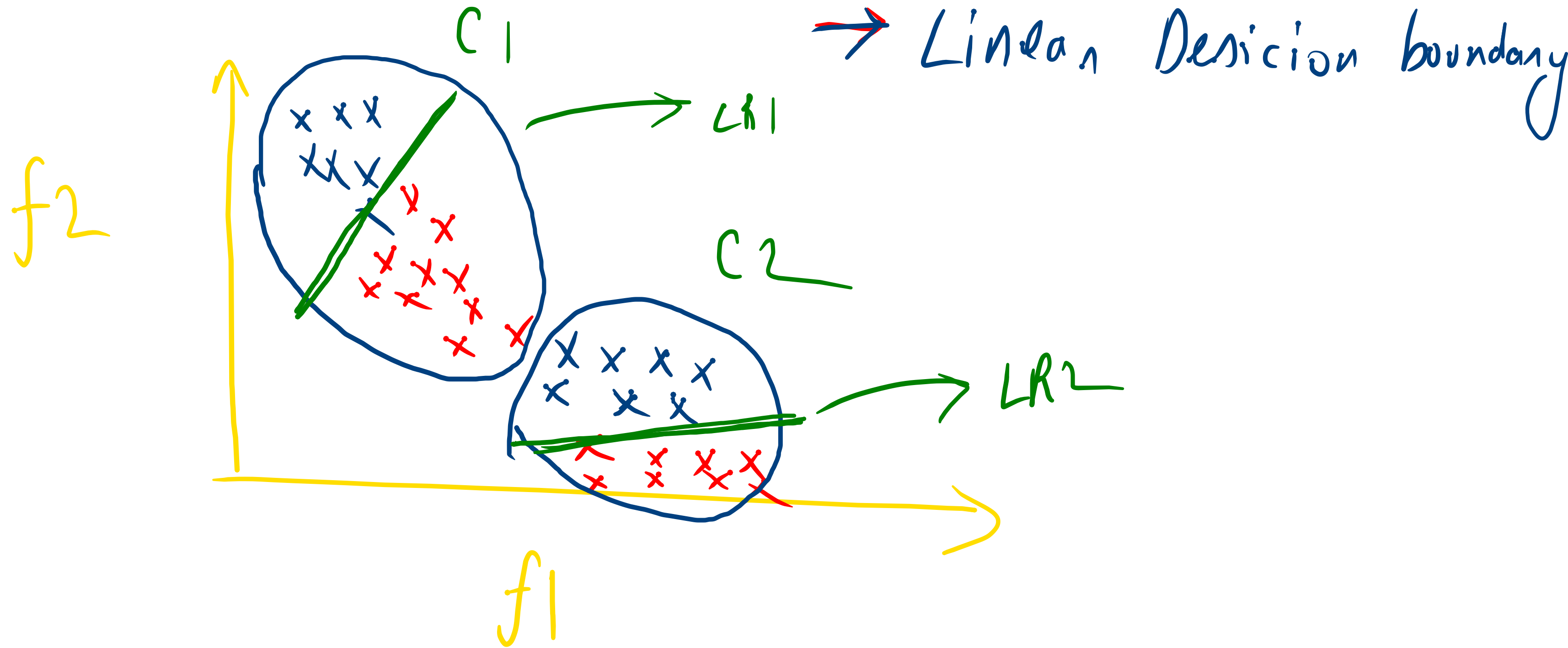


$ Spents

B

A

C

+More Ads

No of visit

# Application of Clustering

→ Search algorithm

DS

SE

→ Feature creation

ME

| ID | n_clicks | n_visit | amount_spent | Country |
|----|----------|---------|--------------|---------|
|    |          |         |              | 1       |
|    |          |         |              | 2       |
|    |          |         |              | 2       |

C1

LR1

C2

$\rightarrow$ Linear Desicion boundary

LR2

f2

f1

# Good Clustering



A

B

C

→ Intra cluster dist    minimize

→ Inter cluster dist    maximize

## WCSS (Within Cluster Sum of Square)

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

$C_1$

$C_2$

$$NCSS = \sum_{k=1}^{K} \sum_{i=1}^{n} \mathbb{1}(c_i = k) \, \|x_i - \mu_k\|^2 = L$$

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

$C_1 \quad C_2$

| Gradient Descent | Coordinate Descent |
|---|---|

Update all parameters simultaneously

$$W \leftarrow W - \alpha \nabla_w L$$

$\hookrightarrow$ $W_1$
$W_2$
$\vdots$
$W_n$

$\rightarrow$ Update only a subset of parameters

$\hookrightarrow$ Fixing $\mu$ & we find best $c$ exactly

$\hookrightarrow$ Fixing $c$ & find $\mu$

# K-means Clustering (Lloyd's Algorithm)

**Steps:**

→ Randomly initialize K centers

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

**Assignment Step**

**update**

→ find the centroids of these
clusters → because this will reduce WCSS

$M_K$

→ Re-assign points

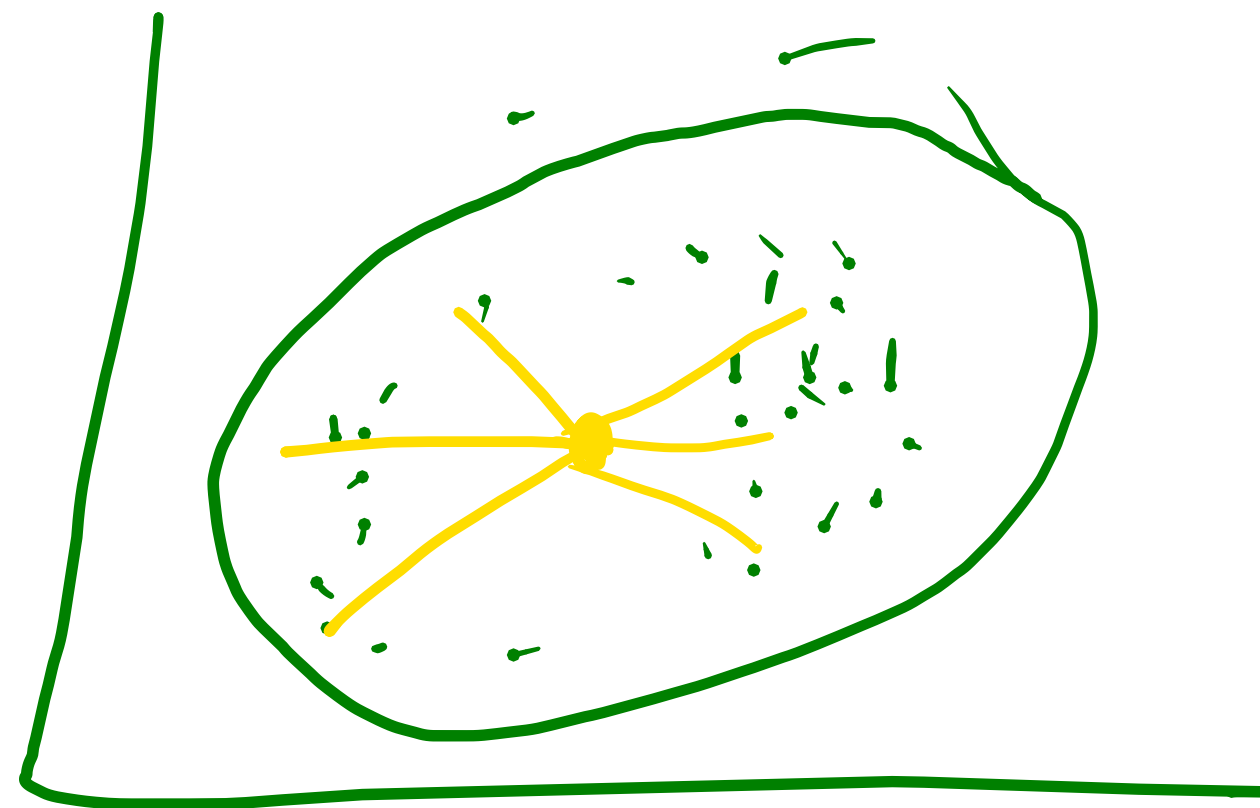→ Repeat until new centers = prev centers

Red cluster

(1, 2)

(2, 3)

(3, 4)

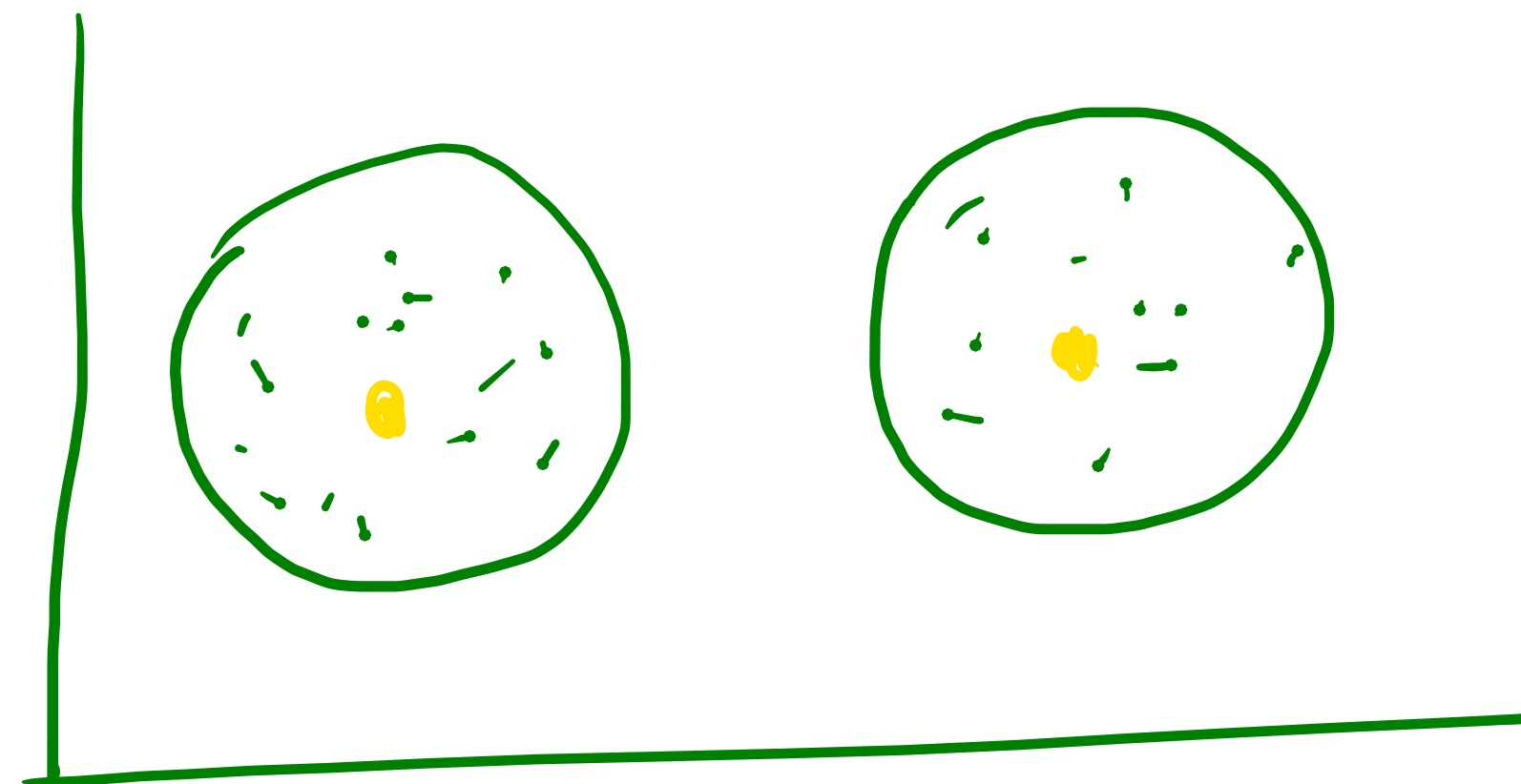$$M_K = \frac{1}{n_K} \sum x_i$$

$x_i \in C_K$

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

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

Cluster = 1

Cluster = 2

WCSS

$\rightarrow$ Break until 22:33

# Business

→ No of Ads

→ Amount of Discount