

# Introduction to Clustering

## KMeans Algorithm

[Unsupervised Learning]

- Module Intro
- Clustering intro
- Kmeans

## CASE STUDY

You are a data scientist at Amazon, and you need to decide different marketing strategies for different types of customers.

How will you identify the different customer groups and assign a group to every customer?

DATA:

ID	n_clicks	n_visits	amount_spent	amount_discount	days_since_registration	group
1476	130	65	213.905831	31.600751	233	1
1535	543	46	639.223004	5.689175	228	2
1807	520	102	1157.402763	844.321606	247	3
1727	702	83	1195.903634	850.041757	148	4
1324	221	84	180.754616	64.283300	243	5

→ Quiz

Q: What is the difference b/w this and a classification problem?

→ Here we have no target.

In binary classification, we need:

$$D = \left\{ (\vec{x}_i, y_i) \right\}_{i=1}^n ; \vec{x}_i \in \mathbb{R}^d; y_i \in \{0, 1\} \}$$

multi-classification  $\rightarrow \underline{y_i \in S}$   
finite set of classes

regression

$$\rightarrow \underline{y_i \in R}$$

## Unsupervised Learning

However, we do not have a 'target' in our problem.

$$D = \{ \vec{x}_i ; x_i \in \mathbb{R}^d \}$$

Such problems are called unsupervised learning

Q: Is this a valid use case then?

→ Of course.

Example:

→ customer / product segmentation

→ grouping similar images in gallery

→ detecting similar stocks automatically

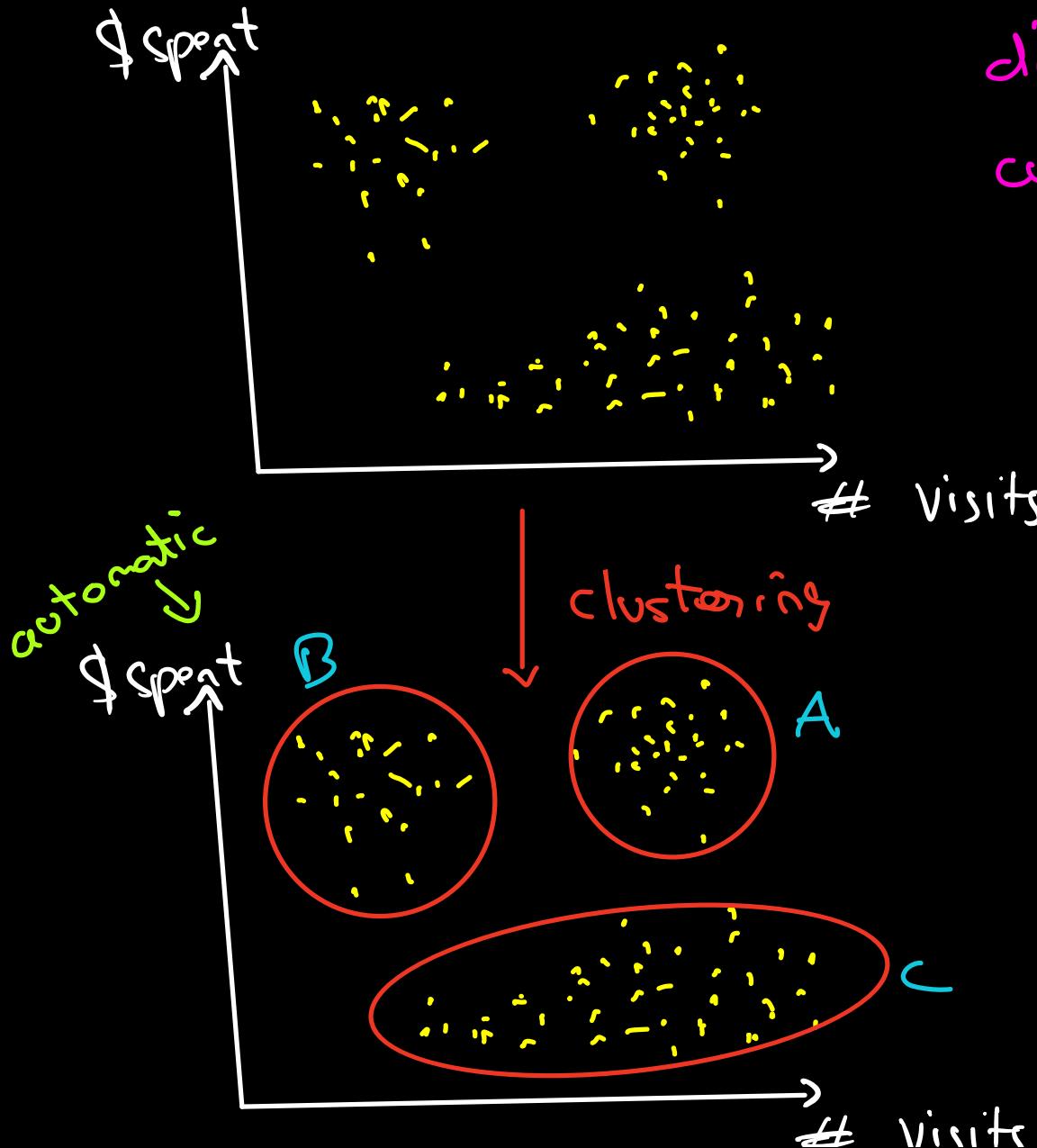
The process of grouping any kind of data based on similarity in their features, automatically, without human expertise, is called **Clustering**

This is one type of unsupervised - learning

Other examples of unsupervised learning:

- Association rules / Recommenders
- Dimensionality reduction (e.g.: PCA)
- Anomaly detection
- Encoding (word-to-vec)

## Clustering



Q: Can you spot different types of customers?

Intuitions

A → Heavy shoppers  
Visit a lot /  
Spend a lot

B → Rich people  
spend whenever they visit

C → Window shoppers

Q: What would you recommend for each group??

- No ads needed for 'A'
- Only ads needed for 'B'
- Ads + discount for 'C'

Intuitively, clustering is dividing a population into groups such that the pts in one group are similar to each other.

Each group is called a cluster

So the task is

- group similar points
- define "similarity"



es. distance on a scatter plot

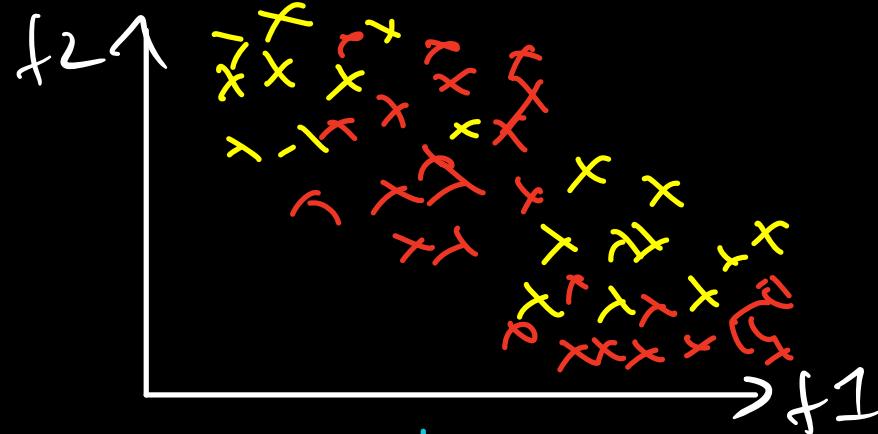
Q: How do you decide if a cluster is good or bad??

- There is no ground truth data!  
Hence nothing to compare with
- It should simply make business sense

## Advanced applications of clustering

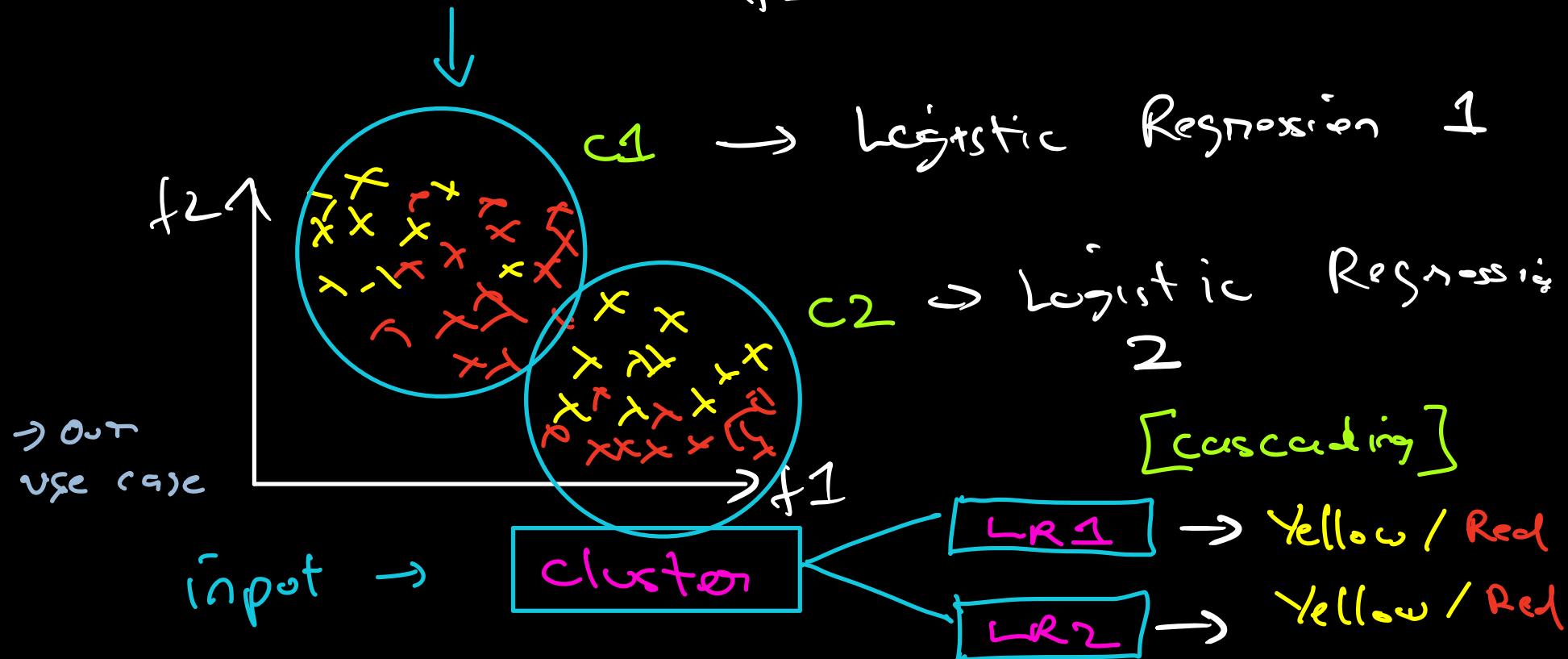
- Validation of domain knowledge / human opinion
- Search algorithms and recommender sys.
- Feature creation
- Clustering with supervised learning to improve model accuracy!!

Example:



Build a classifier  
for this problem.

Q: Any guesses??



## How to perform clustering

Q: Any ideas ??

So far we have seen a pattern

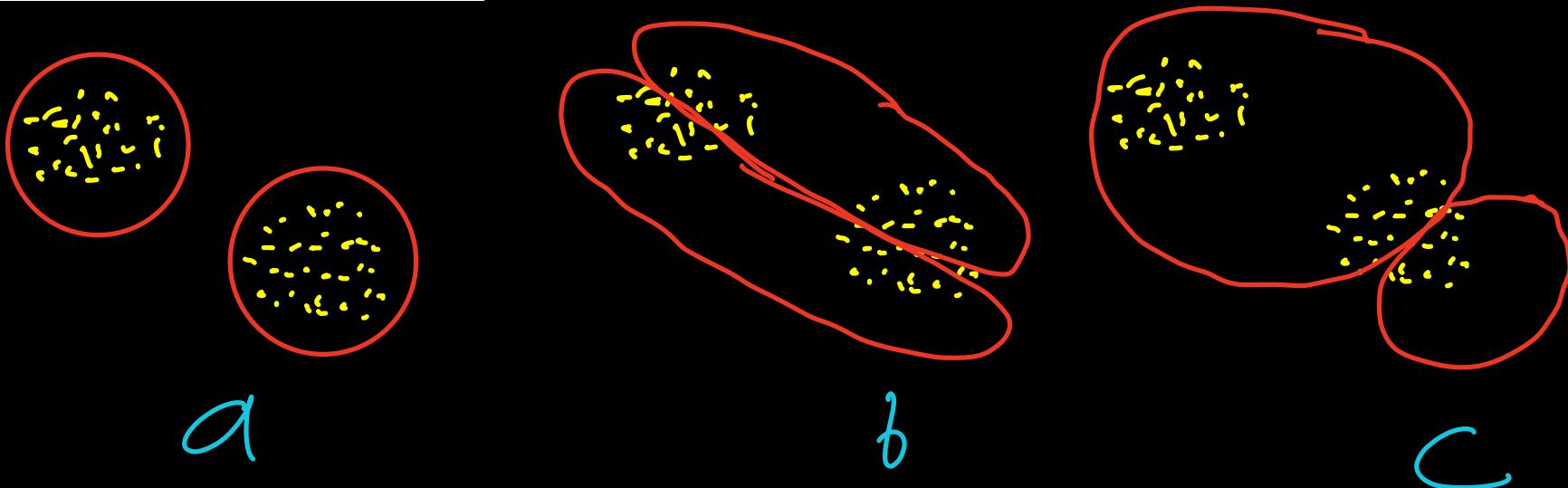
Geometrical idea

↳ Convert to math

↓  
Design an optimisation  
function

↓  
optimise

## Good clustering:



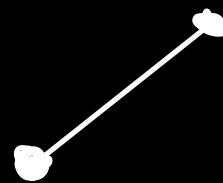
How to express this mathematically ??

- Within a cluster variance / distance should be as low as possible
- Between cluster distance (var of centroids) should be as high as possible

Q: What is distance mathematically ??

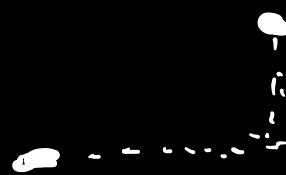
→ Euclidean distance

→ Low dimensions



→ Manhattan distance

→ Medium dimensions



→ Cosine distance

→ high dimensions



→ etc . .

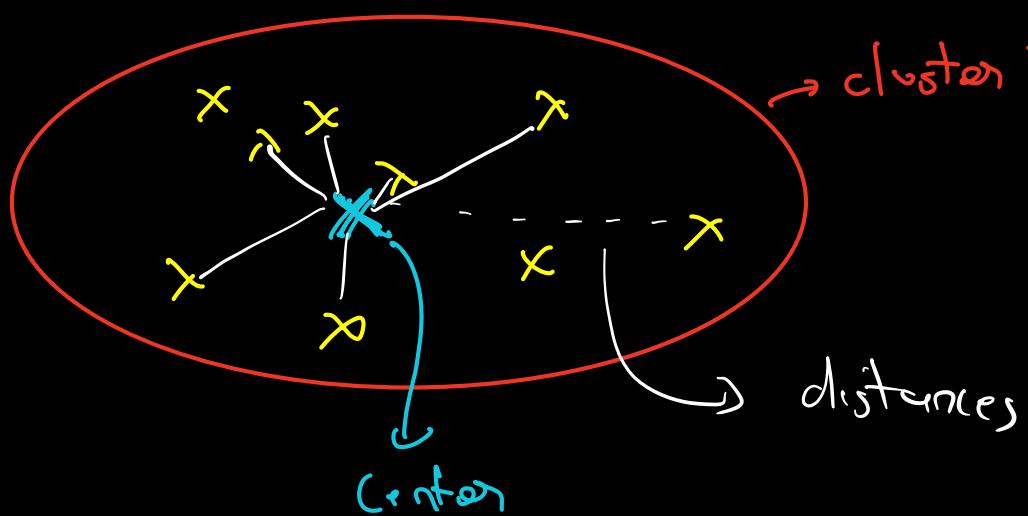
So the choice of the distance metric  
is a hyperparameter for many clustering  
algorithms. → pre-read.

So, we want to come up with a metric which can tell us which clustering is best!

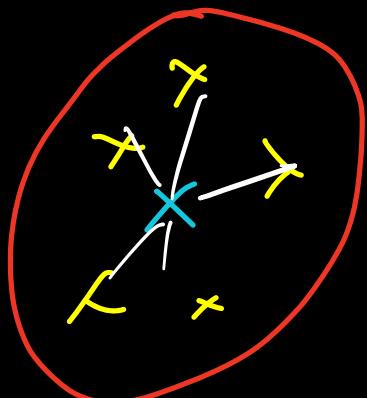
Idea :

→ Minimise within cluster distance.

## Within Cluster Sum of Squares (inertia)



$$\sum \text{dist}_{C_1} + \sum \text{dist}_{C_2}$$



$$WCSS = \sum_{j=1}^k \sum_{i \in C_k} (x_i - \bar{x}_j)^2$$

point  
 ↓  
 all clusters    all pts in j<sup>th</sup> cluster    center of i<sup>th</sup> cluster

idea: **Minimize** the **within cluster variance** (intra)

→ Here we are simplifying by completely ignoring between (inter) cluster distance.

Variation :

$$WCSS = \sum_{i=1}^k \sum_{j=1}^{m_k} \text{distance}(x_{ij}, c_i)^2$$

↓  
Any distance  
func<sup>n</sup> of  
your choice      ↓  
center  
of i<sup>th</sup>  
cluster

## Lloyd's Algorithm (K-Means Clustering)

Steps:

- Randomly initialize  $K$  centers
- assign points to nearest center  
to get your clusters
- find the centroids of these  
clusters → because this will reduce WSS
- Re-assign points
- Repeat until new centers =  $\underline{\text{prev centers}}$
- animation

→ animation

→ random initialise ' $K$ ' clusters  
↳ manually decide

while True:

$y = \text{closest}(x)$

for center in y.unique()

new\_center =  $x[y == c].mean()$

loss = metric( $x, c$ )

if prev\_loss - loss  $< 1e^{-4}$ :

break

→ code

Complexity:  $\uparrow n \nwarrow c \uparrow i \uparrow d$   
 $\downarrow \text{pts} \quad \uparrow K \quad \downarrow \text{iter} \quad \downarrow \text{dim}$

## Determining the best 'K'

→ Visually

→ Domain knowledge

└ Slack groups

→ Avg reply time, # users, # msgs

→ support channels (sales / ops)

→ discussion forums (students)

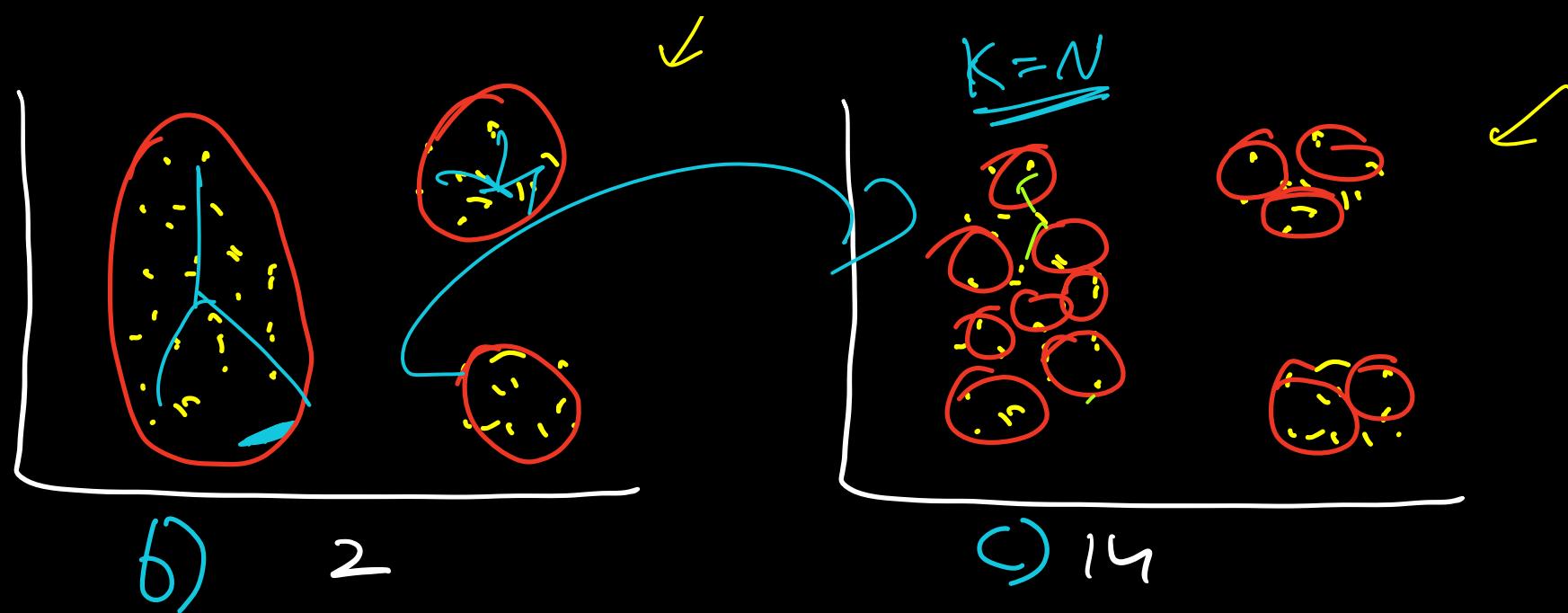
→ internal teams (employees)

→ Elbow method

Q: Which is best  
K ??

a) 3





→ We can visually see that its a) 3  
 ↳ How do you tell that numerically?

⇒ Metrics!!

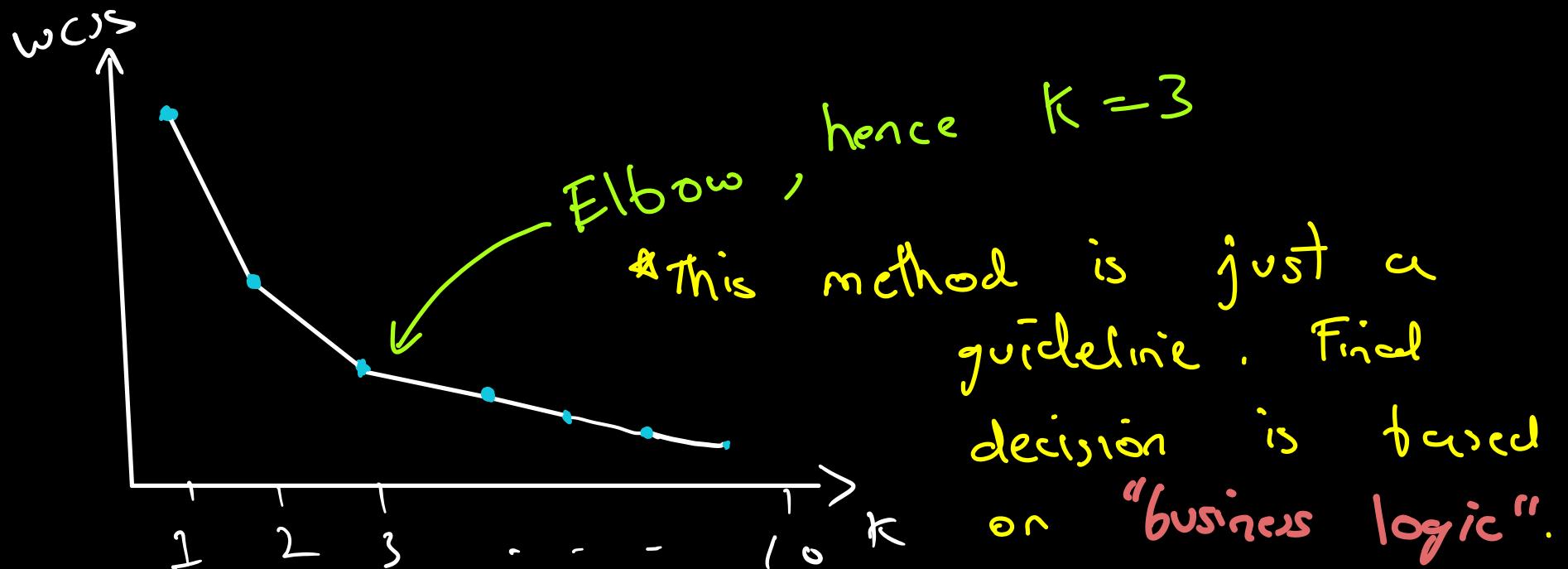
① WCSS

2) Silhouette Score

3) Dunn Index

--- Others

## Elbow Method



for  $K$  in range( $P$ ):

    res.append(metric(clustering( $K$ )))

plt.plot(res)