

Introduction to Clustering

KMeans Algorithm

[Unsupervised Learning]

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	profile_information
0	1476	130	65	213.905831	31.600751	233	235
1	1535	543	46	639.223004	5.689175	228	170
2	1807	520	102	1157.402763	844.321606	247	409
3	1727	702	83	1195.903634	850.041757	148	200
4	1324	221	84	180.754616	64.283300	243	259

→ Quiz

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

→ Here we have no target.

In binary classification, we had:

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

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

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

Unsupervised Learning

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

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

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

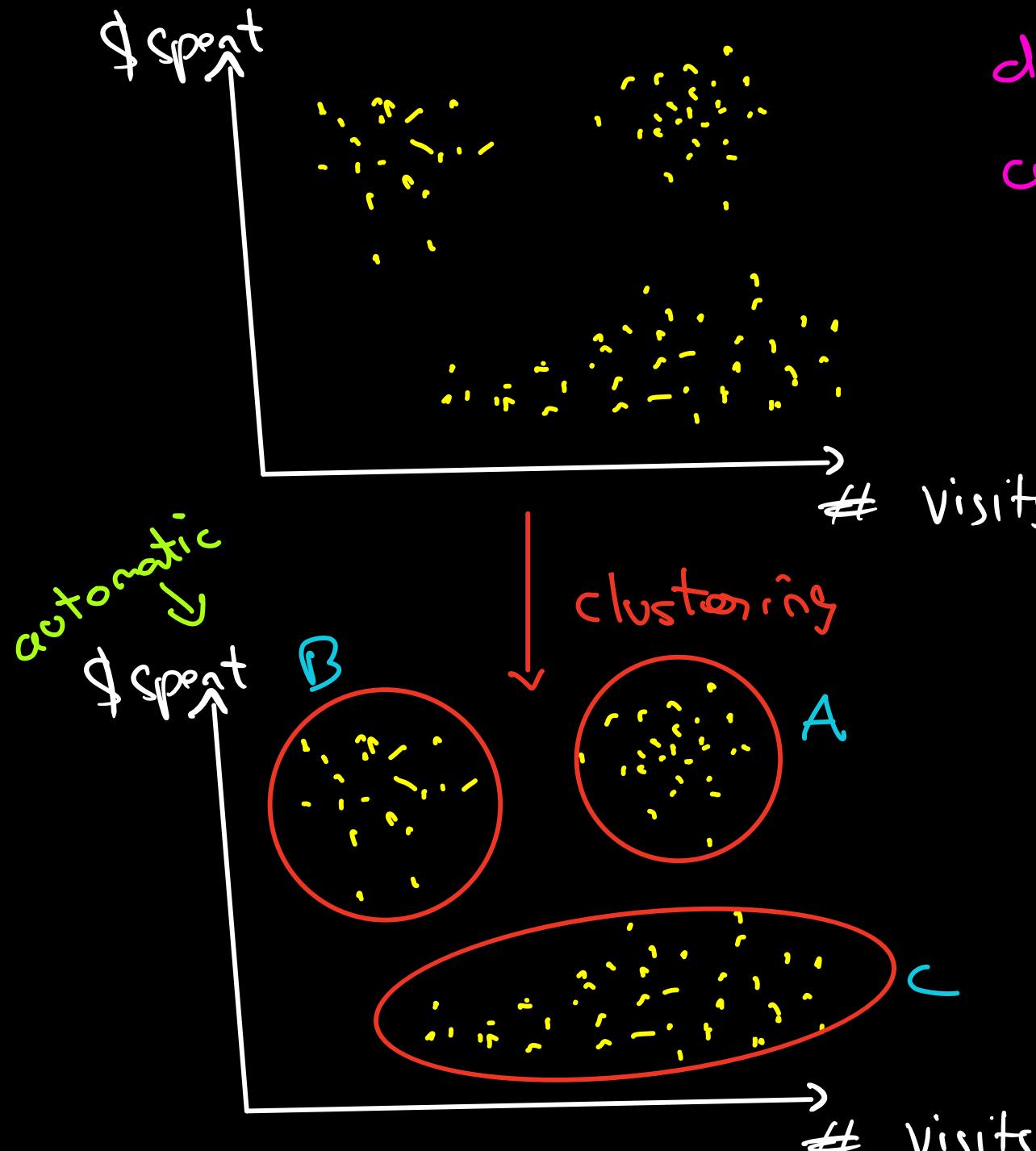
The process of grouping any kind of data based on similarity in their features, automatically, without human expertise, is called **Clustering**

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?

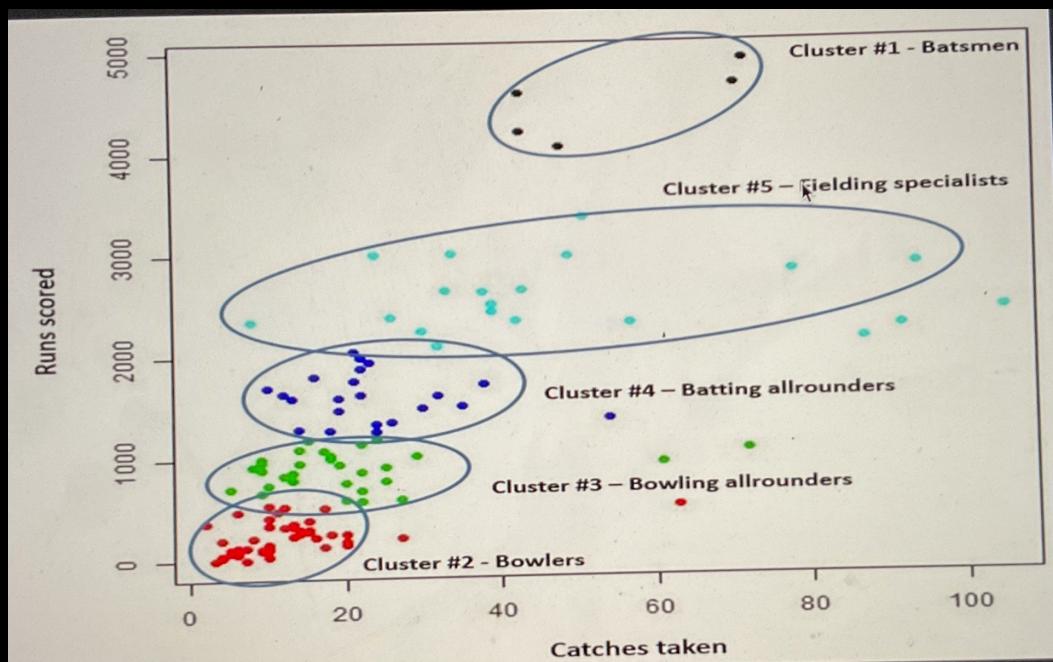
Intuitions

- 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'

Another example



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

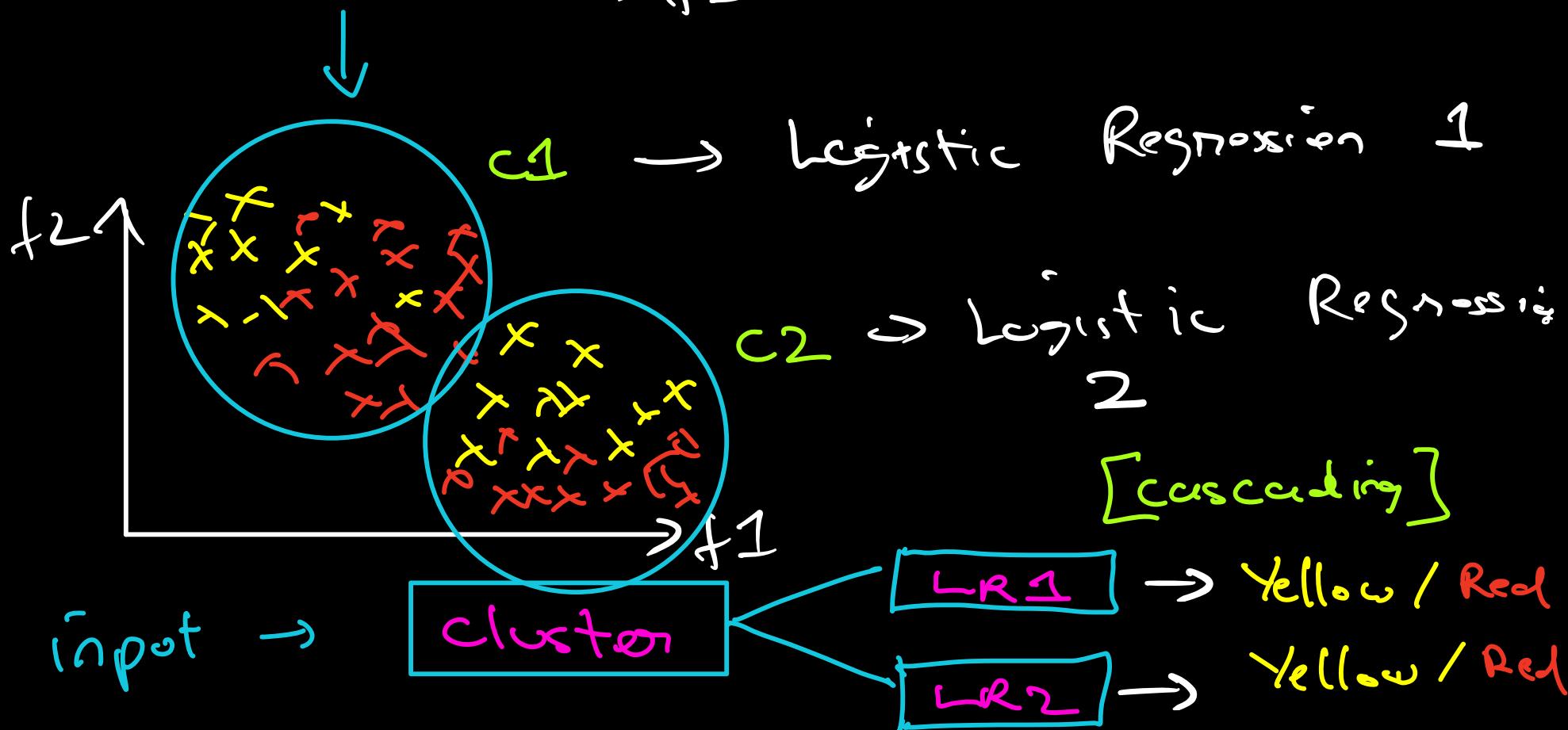
- Validation of domain knowledge / human opinion
- Search algorithms and recommender sys.
- Feature creation
- Clustering with supervised learning to improve model accuracy !!

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

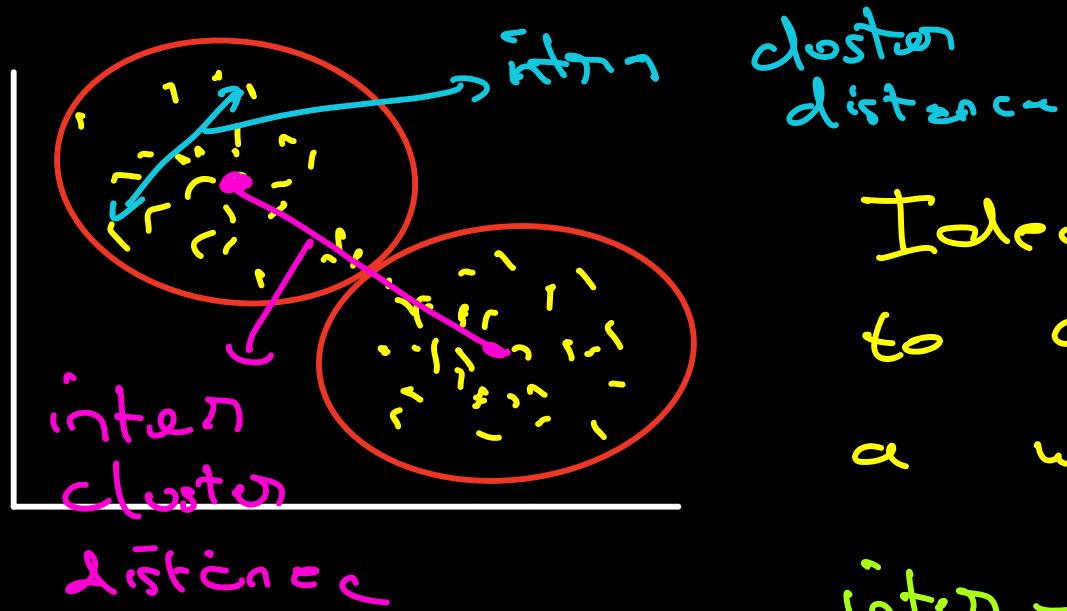
→ Convert to math

↓
Design an optimisation
function

↓
optimize

Q: Can we construct a loss (gain)
function for clustering ??
[Hint: ANOVA / f statistic]

→ We assuming that data comes from different populations here.



Ideally, we need to assign clusters in a way that:

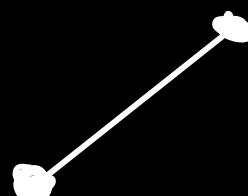
inter-cluster \uparrow (max)

intra-cluster \downarrow (min)

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.

Now,

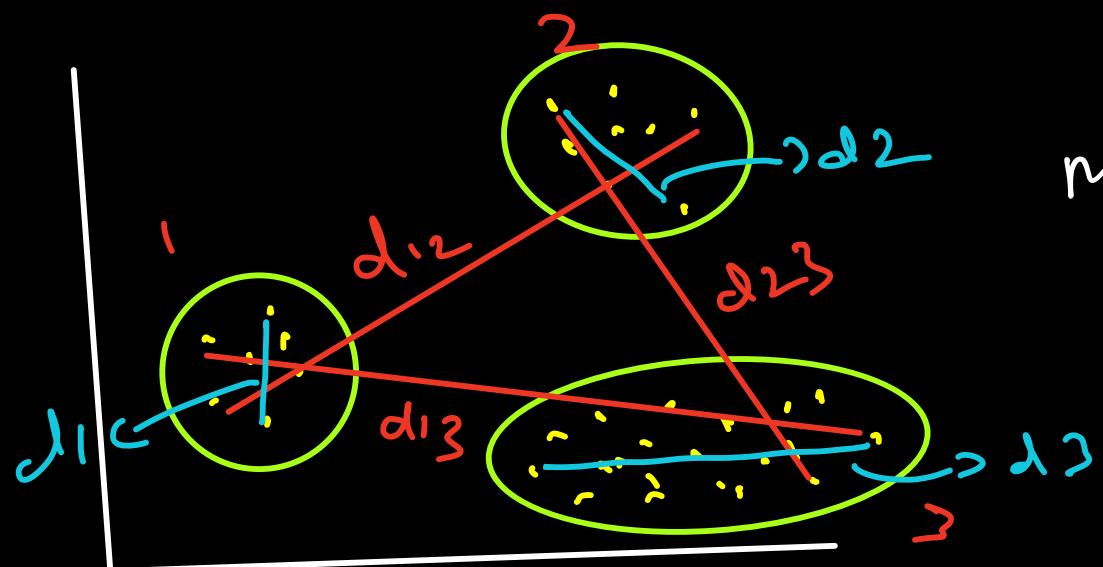
using this we can write some sort
of metric which we may optimise

Dunn Index

↳ slightly simpler metric than an f-statistic.

$$D = \frac{\min \text{ inter cluster dist}}{\max \text{ intra cluster dist}}$$

worst case scenarios are taken



$$\min(d_{12}, d_{13}, d_{23})$$

$$= d_{23}$$

$$\max(d_1, d_2, d_3)$$

$$= d_3$$

$$\therefore \text{Dunn Index} = \frac{d_{23}}{d_3}$$

If dunn index is high, clusters are well separated and spread out, else they are not.

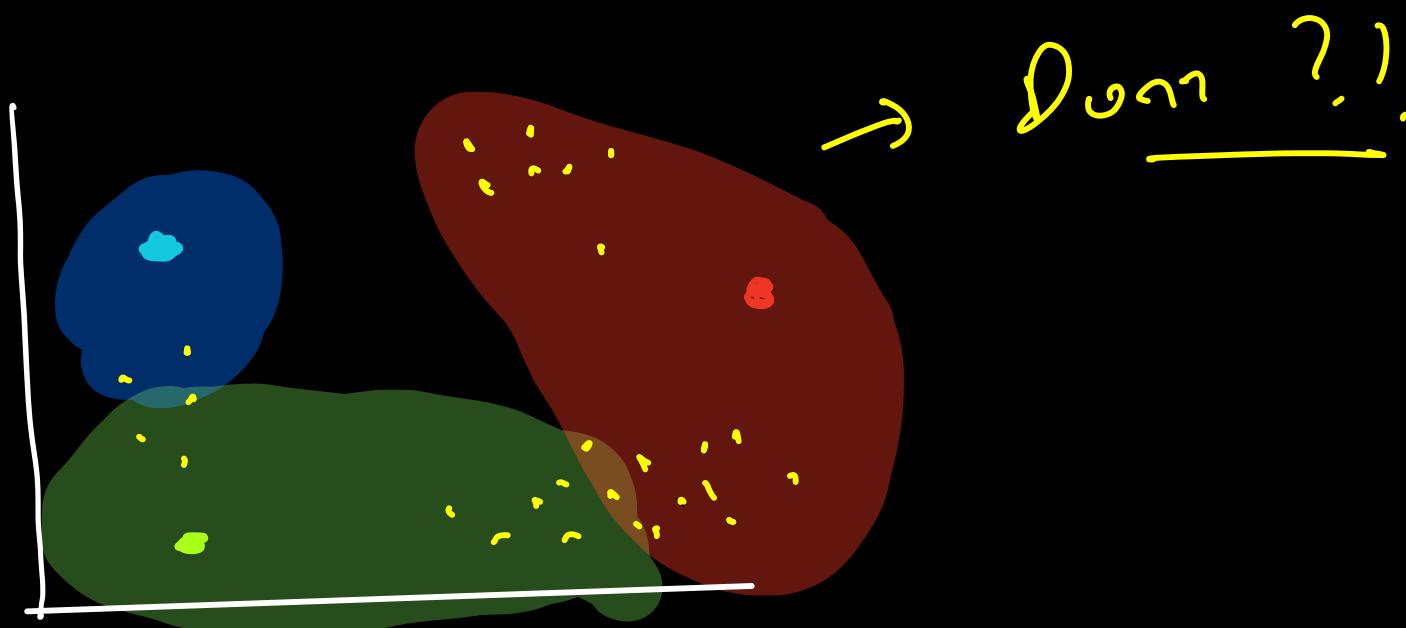
→ Note: dunn index is unbound, so it can only be interpreted in a relative sense. → Quiz: High or low?

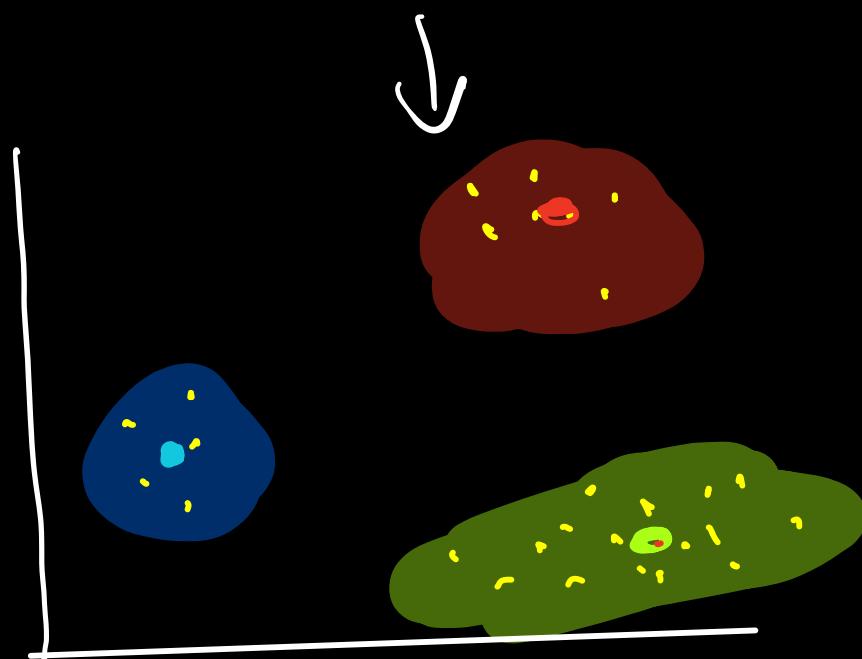
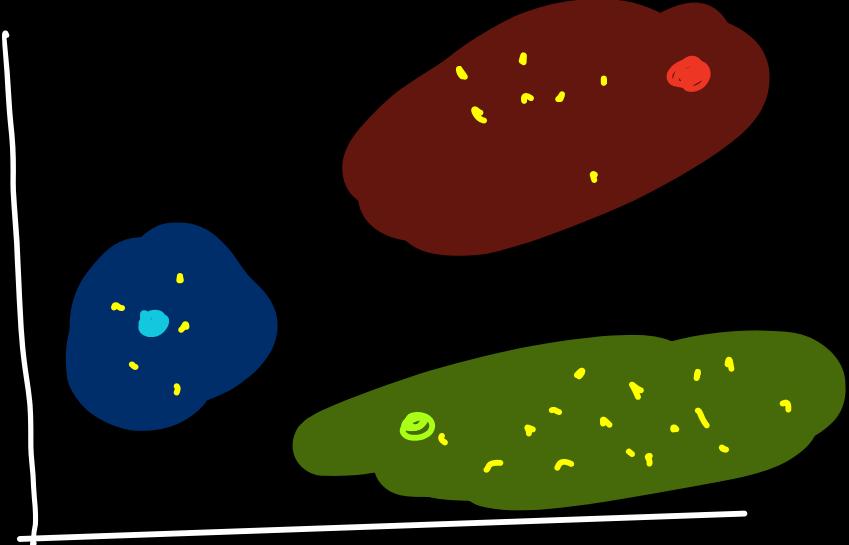
Optimisation

Now that we have a way to compare different cluster possibilities, we need to find a way to find the best clusters

Steps: → Decide # clusters (manual)

- Initialise cluster centers randomly
- To each point assign the cluster whose center is closest
- calculate sum index
- update the centroid location
- Repeat until best centroids are found.





↳ stop !!

Optimisation

max
[?]

$f(?)$

??

What is the parameter & the func' ?!

→ Params → centroids

→ $f = ??$ Dunn index (clusters)

Cluster centroids = $C = \begin{bmatrix} \vec{c}_1 \\ \vec{c}_2 \\ \vdots \\ \vec{c}_k \end{bmatrix}$

$g(C, X) =$ assignment
function based

on nearest centroid

$\arg\min(\text{distance}(C, X)) \mid d(x_i, c_i) = \text{dunn index}$

$$L_i: \max_{C_i} d(X, \underbrace{\text{argmin}}_{\substack{\downarrow \\ \text{not}}} (\underbrace{\text{distance}(C_i X)}_{\substack{\text{some distance} \\ \text{matrix are not} \\ \text{differentiable}})))$$

Hence we cannot solve this using calculus or gradient descent.

Clustering by nature is a discrete problem and hence can't be optimised easily.

[There are ways, but they are still hard]

Lloyd's Algorithm (K-Means Clustering)

Steps:

- Randomly initialize K centers
- assign points to this center
to get your cluster.
- find the centroids of these clusters
- Re-assign points
- Repeat until new centers = prev centers
- animation / code

Determining the best 'K'

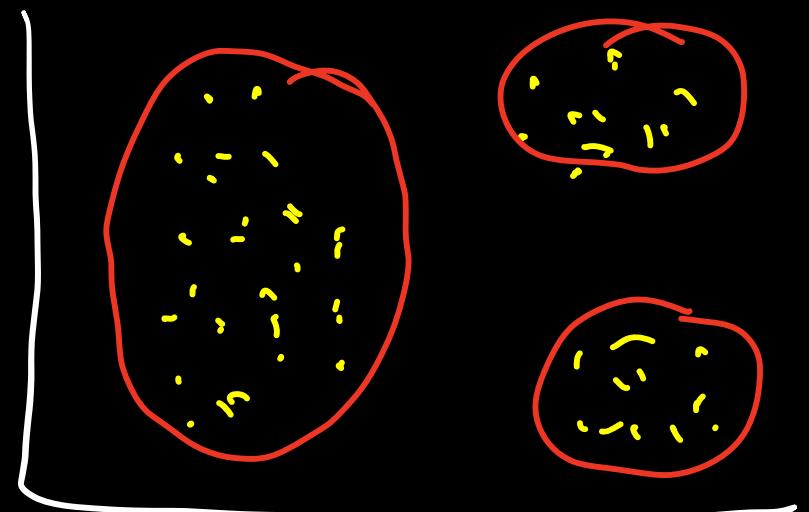
- Visually
- Domain knowledge
 - ↳ Slack groups
 - Avg reply time, # users, # msgs
 - support channels (sales / ops)
 - discussion forums (students)
 - internal teams (employees)

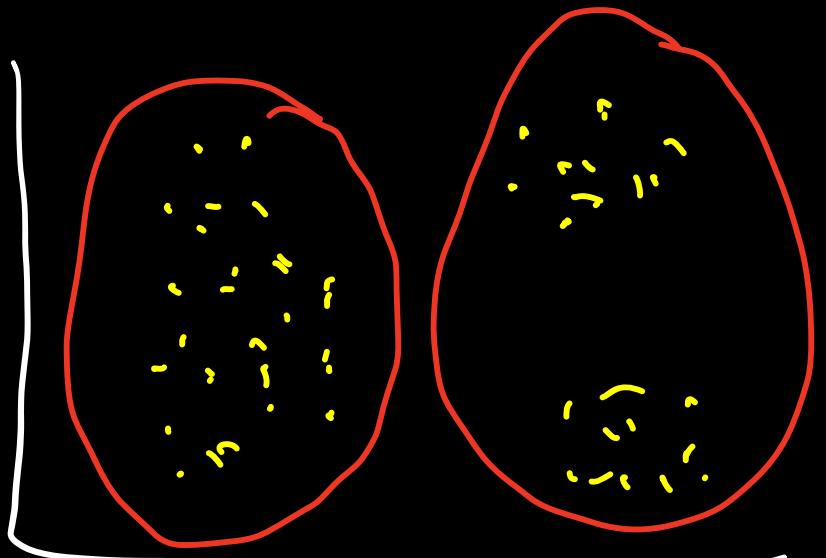
→ Elbow method

Q: Which is best

K ??

a) 3





b) 2



c) 14

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

⇒ Metrics !!

1) WCSS

2) Silhouette Score

WCSS (Inertia)

within - cluster sum of square

$$WCSS = \sum_{i=1}^K \sum_{j=1}^{m_i} \text{distance}(x_j, c_i)$$

distance of pts from their respective
cluster centroids (should be low)
[unbounded]

Ques: When will inertia go to zero?

Silhouette Score

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

a = avg distance of x_i with points in
its own cluster

b = avg distance of all point of nearest
cluster with x_i

How close is the point to the points
neighbouring cluster compared to own
cluster.

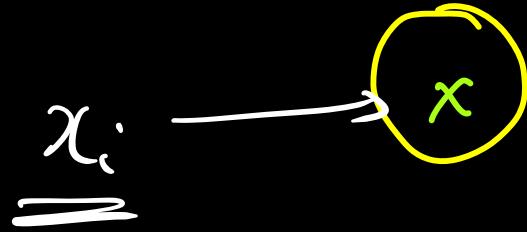
$$\max \sum_i s(x_i)$$

[Bounded]

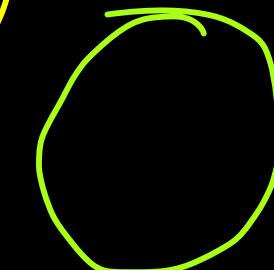
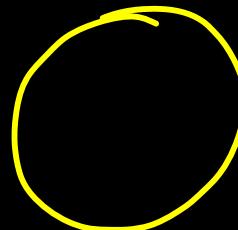
$$s_i(x_i) = 1, \quad b > a \quad \& \quad a = 0$$



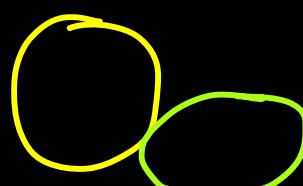
$s_i(x_i) = -1$, $a > b$ ~~&~~ $b = 0$
(incorrectly allotted clusters)



$$0 > s_i > 1$$



$$s_i \sim 0$$



$$-1 < s_i < 0$$

