

# Kmeans ++, Hierarchical [Clustering]

- Problem recap
- Kmeans algorithm
  - WCSS
- Code from scratch
- Kmeans ++
- Hierarchical clustering

## Recap

↳ Amazon customer segmentation

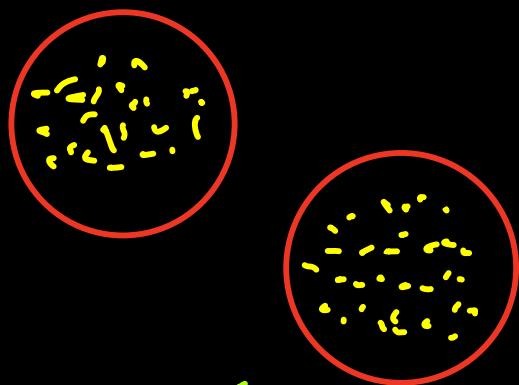
	ID	n_clicks	n_visits	amount_spent	amount_discount	days_since_registration	profile_information
0	1476	130	65	213.905831	31.600751	233	235
1	1535	543	46	639.223004	5.689175	228	170
2	1807	520	102	1157.402763	844.321606	247	409
3	1727	702	83	1195.903634	850.041757	148	200
4	1324	221	84	180.754616	64.283300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

→ Clustering :

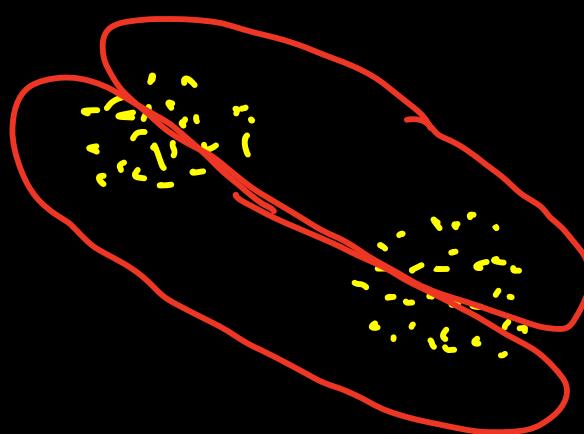
↳ Points belong to different populations

→ A good clustering algorithm  
will be able to detect these.

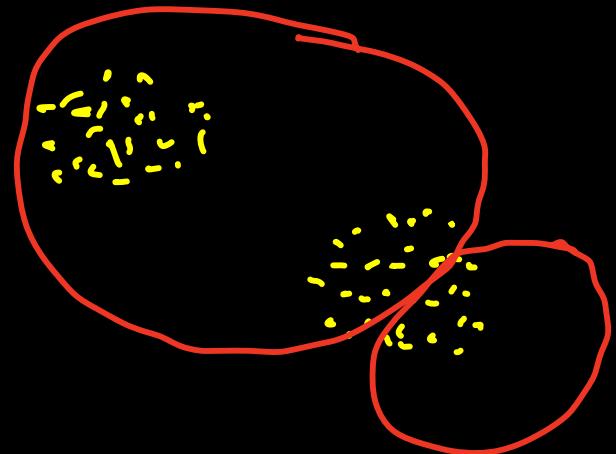
## Good clustering:



✓ Good



✗ Bad



✗ Bad

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

Hence, we can come up with metrics  
that can measure a clustering obj.

M →	inter cluster distance	(between)
	intra cluster distance	(within)

Possible metrics :

F-statistic / Dunn index / WCSS / Silhouette,  
etc,

Kind of an  
avg metric

$$\frac{\sum (\bar{x} - \bar{\bar{x}})^2 \cdot n/d_f}{\sum \sum (x - \bar{x})^2 / d_f}$$

↓  
extreme  
metric

$$\frac{\min(\text{inter})}{\max(\text{intra})}$$

will study

- ↓  
x not used much
- degrees of freedom
  - complex looking
  - tied to one distance metric

So we can use any good measure of clustering and proceed.

### Optimisation

params: cluster centers, metric: F, dunn, WCSS, etc

$$\max_C \text{metric}(\text{assignment}(X, C))$$

maximise over centroids, some metric measured on some assignment of points to closest centroid.

### Problem:

- ↳ "assignment of nearest cluster center" is not differentiable
- in some cases the loss metric may not be differentiable. eg: dunn / silhouette

Hence we cannot solve this using gradient descent.

# Lloyd's Algorithm (K Means)

→ animation

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

while True:

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

for center in  $y.\text{unique}()$

new\_center =  $x.\text{loc}[y == c].\text{mean}$ )

loss = metric( $x, c$ )

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

break

→ code

Complexity:  $\uparrow^{\text{pts}} \uparrow^{\text{K}} \downarrow^{\text{iter}} \downarrow^{\text{dim}}$

WCSS (inertia)  $\rightarrow$  Most popular

Within cluster sum of squares.

$$\text{WCSS} = \sum_{i=1}^k \sum_{j=1}^{n_k} (x_{ij} - \bar{x}_i)^2$$

↑  
point  
↓  
all clusters    all pts in  $i^{\text{th}}$  cluster    center of  $i^{\text{th}}$  cluster

Idea: Minimize the within cluster variance (intra)

distance -

$\rightarrow$  Here we are simplifying by completely ignoring between (inter) cluster distance.

## Variation :

$$WCSS = \sum_{i=1}^K \sum_{j=1}^{m_K} \text{distance}(x_{ij}, c_i)$$

↓  
 Any distance  
 func<sup>n</sup> of  
 your choice

↓  
 center  
 of *i*th  
 closure

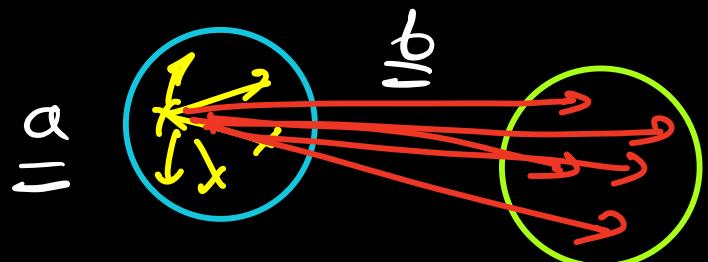
## Silhouette Score

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

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

*b* = avg distance of all point of nearest

cluster with  $x_i$



$\frac{b-a}{\max(a, b)}$  calculated for  
a single point.



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

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

[Bounded]

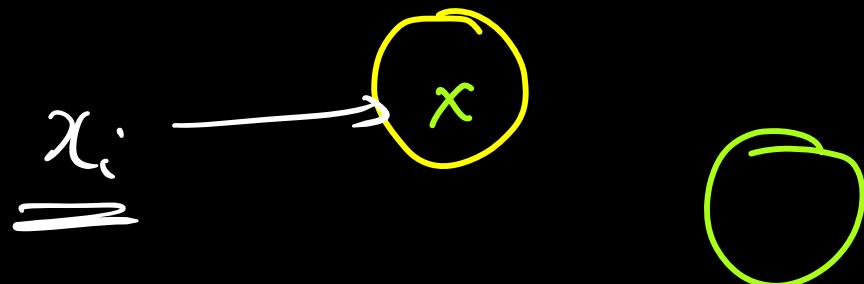
$i = 1$  to  $q$  &  $a = 0$

$$s_i(x_i) = 1, \quad c_1$$

$$s_i(x_i) = -1, \quad c_2$$

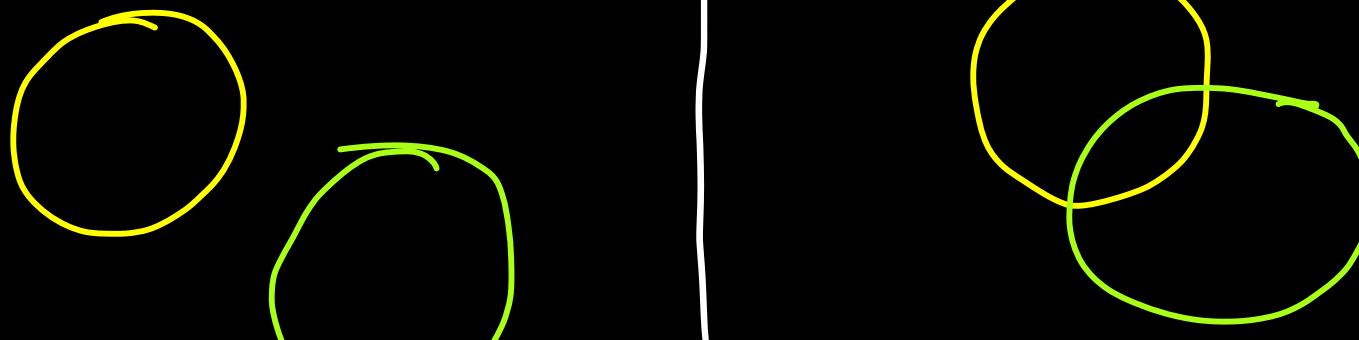
$$s_i(x_i) = -1, \quad a > b \quad \text{and} \quad b = 0$$

(incorrectly allotted clusters)

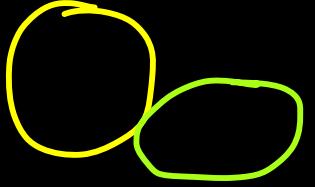


$$0 > s_i > 1$$

$$-1 < s_i < 0$$



$S_i \sim 0$



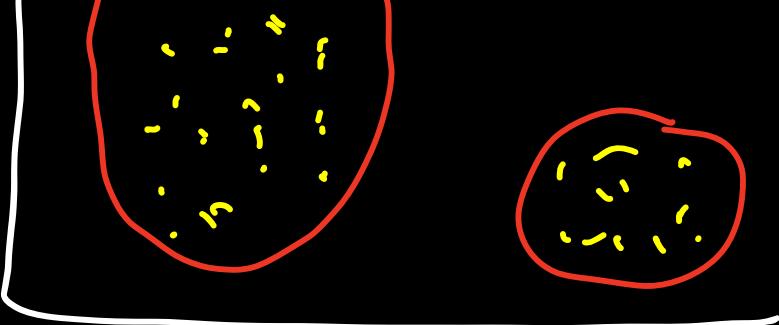
## Determining the best 'k'

- 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 it's a) 3  
↳ How do you tell that numerically?

⇒ Metrics!!

① WCSS

③ Dunn, etc.

2) Silhouette Score

## Elbow method

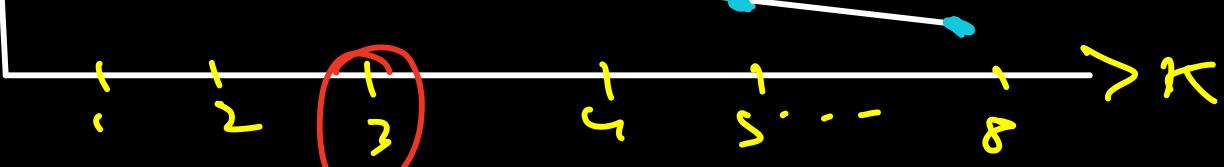
for  $K$  in range( $P$ ):

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

plt.plot(mes)

A metric  $\rightarrow$  e.g: WCSS  $\rightarrow$  after a while improvement  $\Delta$  becomes too small to increase ' $K$ ' further

Elbow



→ code

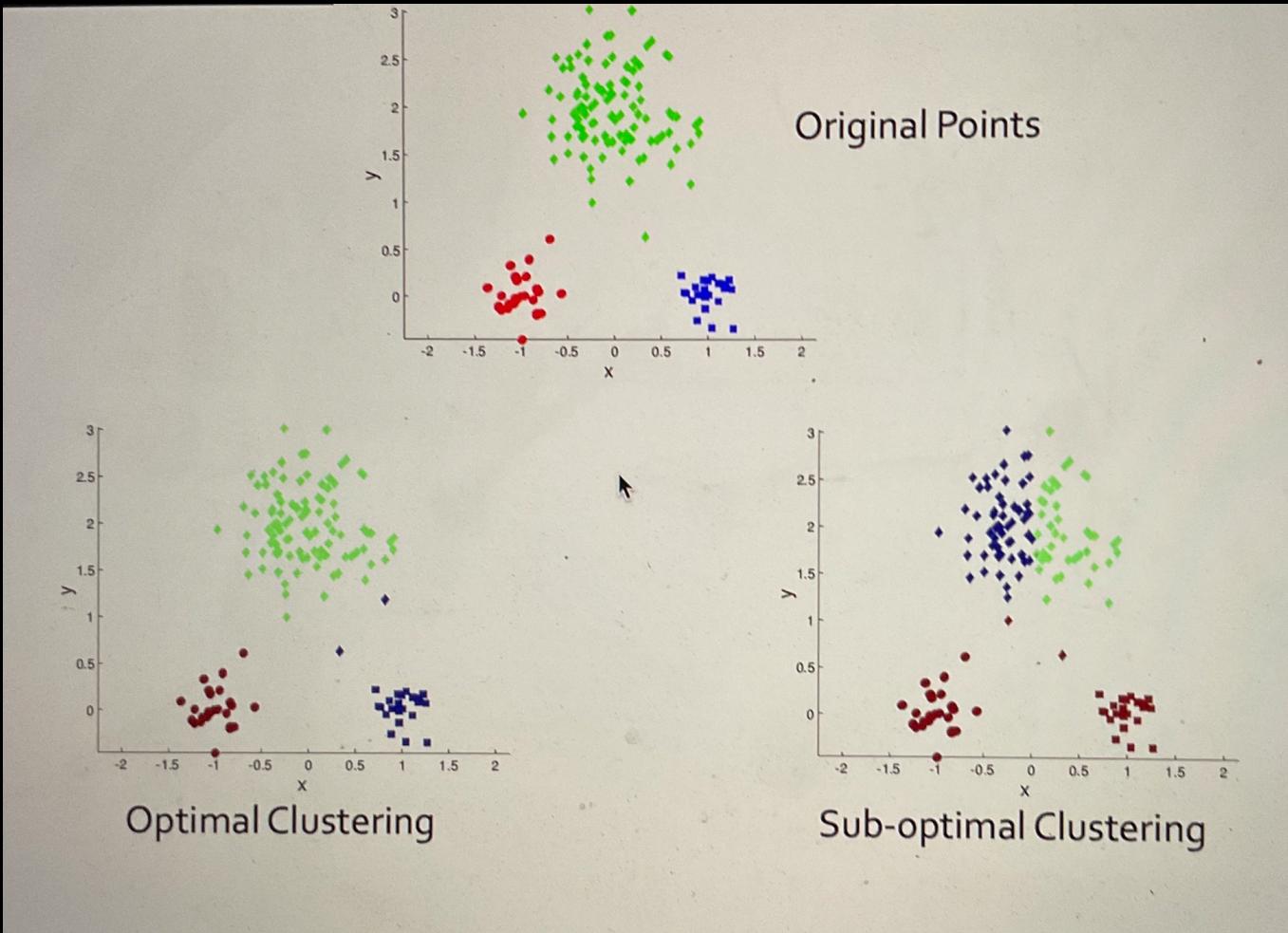
→ optimum K

## Kmeans ++

There is an issue with the Lloyd's algo.

→ Because, centroids are randomly initialized we may get different results some times.  
(bad)

→ animation



## Solution

→ choose centers smartly.

idea :

- 1st centroid initialised randomly
- then choose 2<sup>nd</sup> center as the furthest point from first
- 3<sup>rd</sup> to be far from 1<sup>st</sup> and 2<sup>nd</sup>

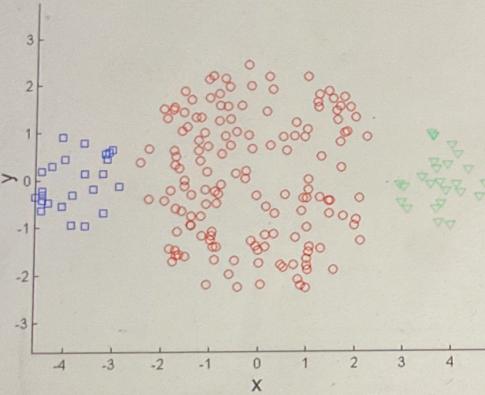
→ animation

What about outliers?

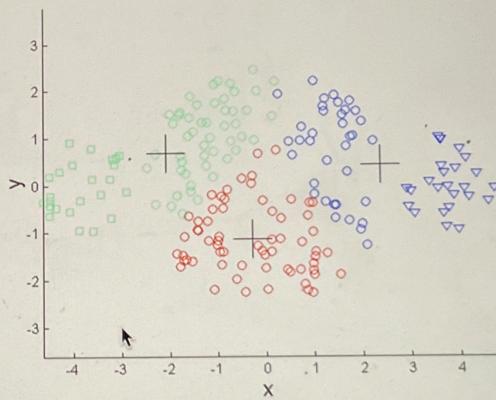
- Instead of choosing outliers we could choose furthest point as centroid with prob proportional to distance.
  - Remove outliers for learning centers.
- screen.

## Limitations of K means

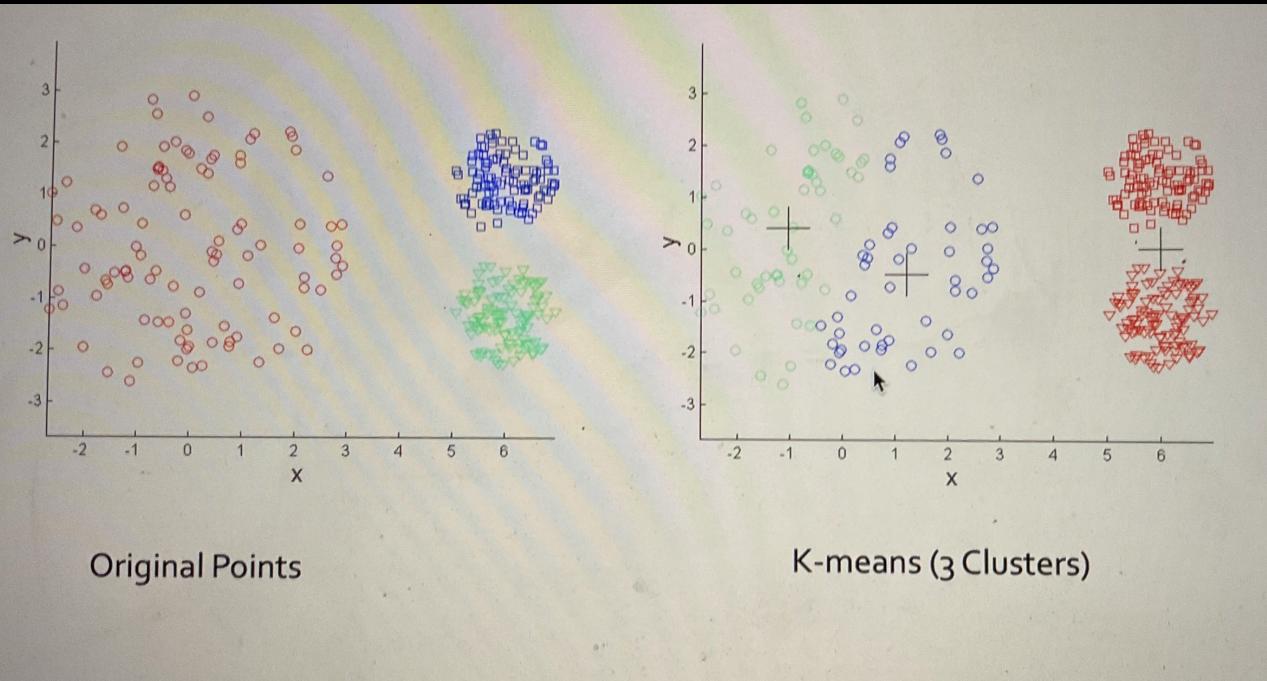
- Need to decide "K"
- May not be best for    clusters.
  - Diff size
  - diff density
  - non globular (hyper-spheres)



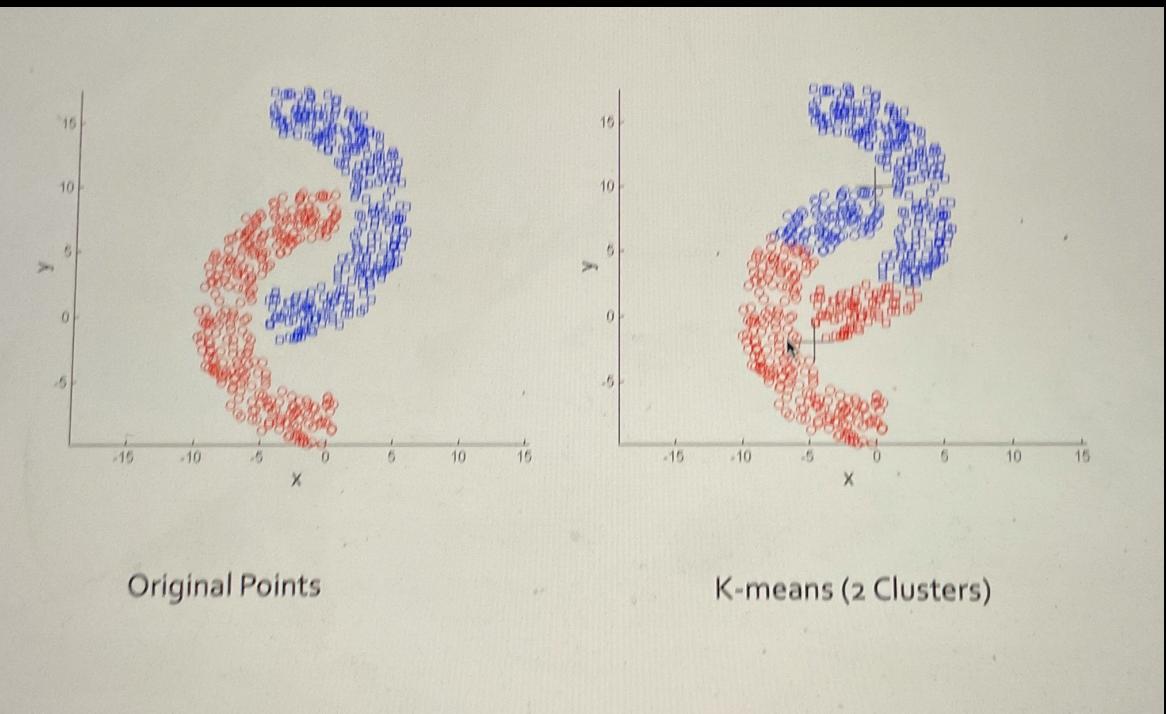
Original Points



K-means (3 Clusters)



Original Points



Original Points

K-means (2 Clusters)