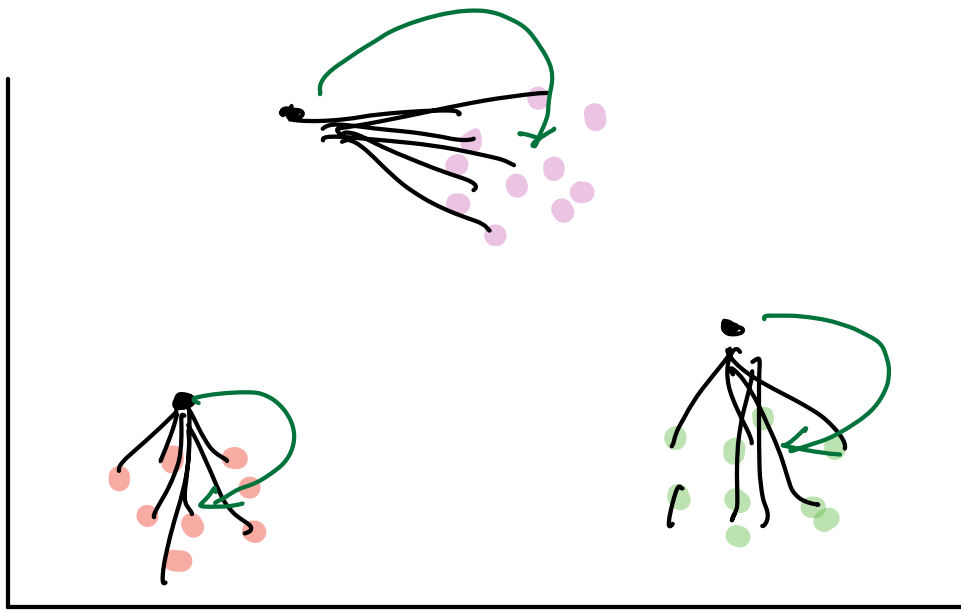


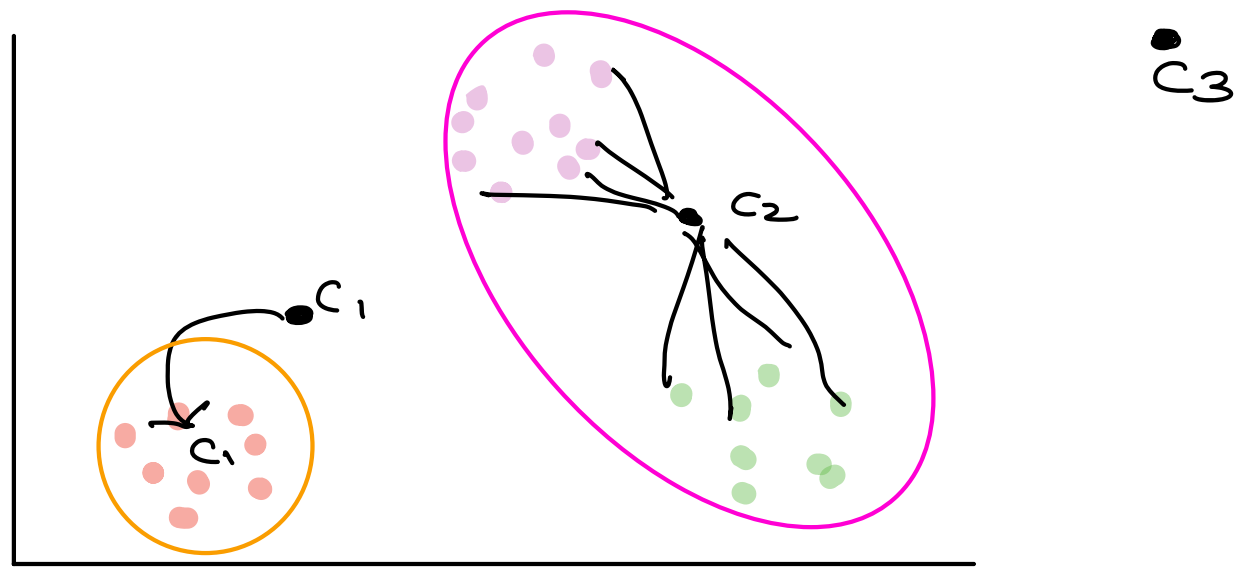
- Recap ✓
- WCSS and Elbow Method ✓
- Initialization Trap
- Kmeans ++
- Limitations of K-means
- K-median
- K-medoids

Limitations of K-means

Initialization Trap



Ideal Scenario



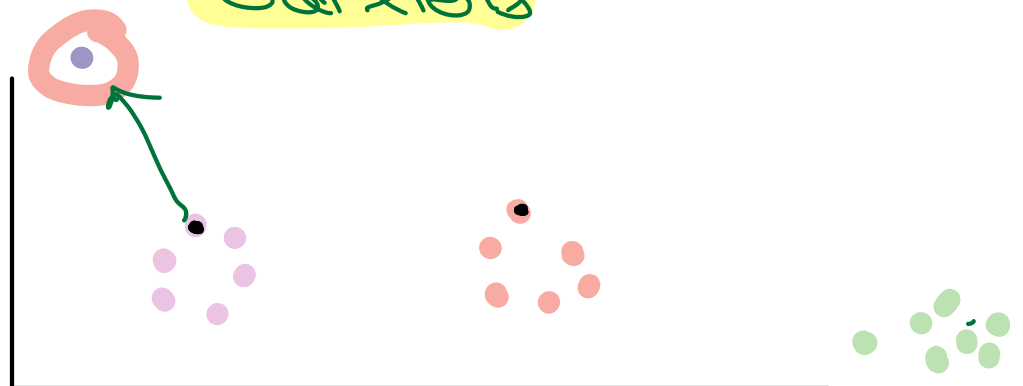
init - Trap

Kmean ++

Step - 1 : Consider any one data points as first Centroid randomly

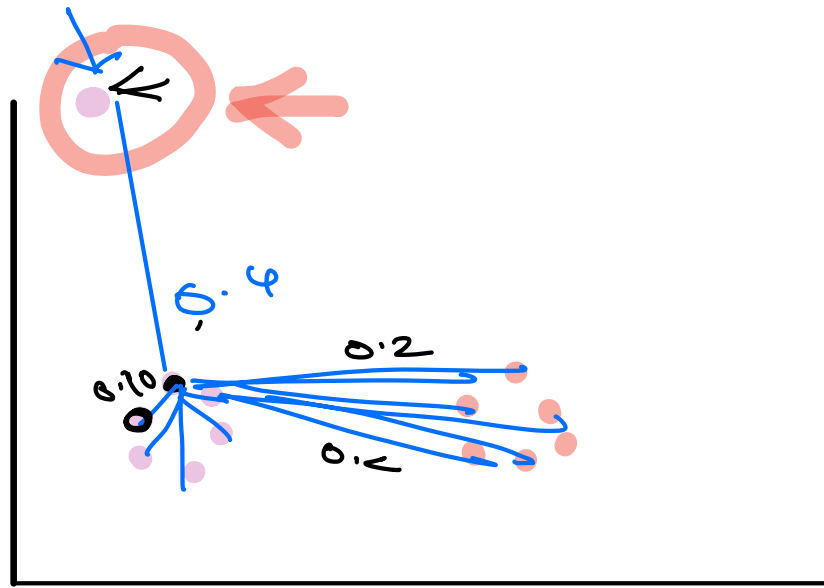
Step 2 : Pick the second centroid as far away as possible from first

Outliers



→ Handle outliers Bejeff hands

Pick centroids probabilistically

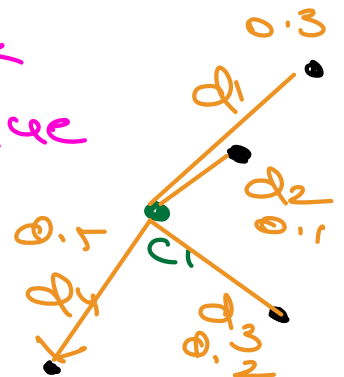


0.99
X

① Pick 1st Centroid Randomly

② Calculate $\text{dist}(C_1, x_i)$ and convert to probability

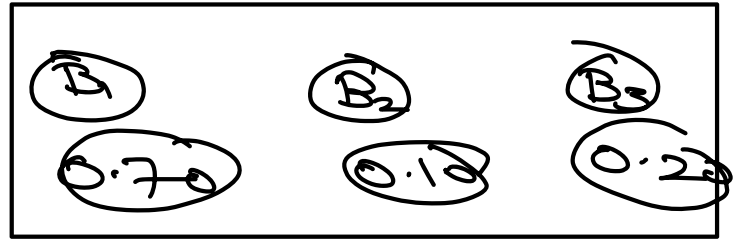
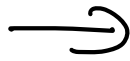
③ Based on P.O.D. pick the next Centroid Value



$$P_1 \Rightarrow \frac{d_1}{D} \quad P_2 \Rightarrow \frac{d_2}{D} \quad \dots$$

$$d_1 + d_2 + d_3 + d_4 \Rightarrow D$$

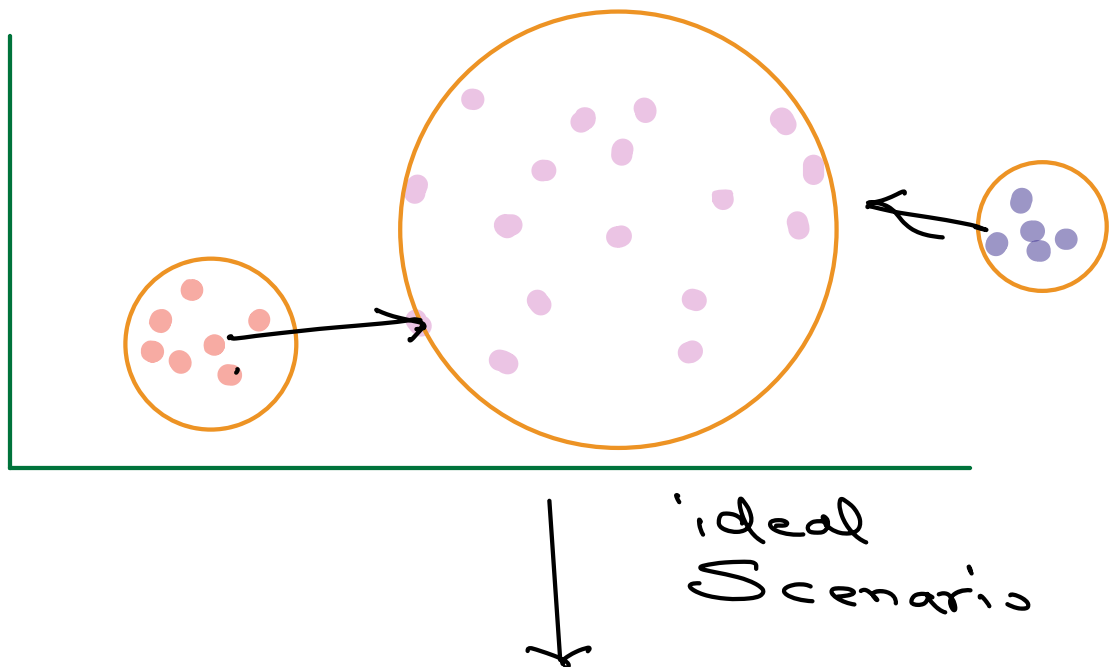
$B \Rightarrow 70\%$

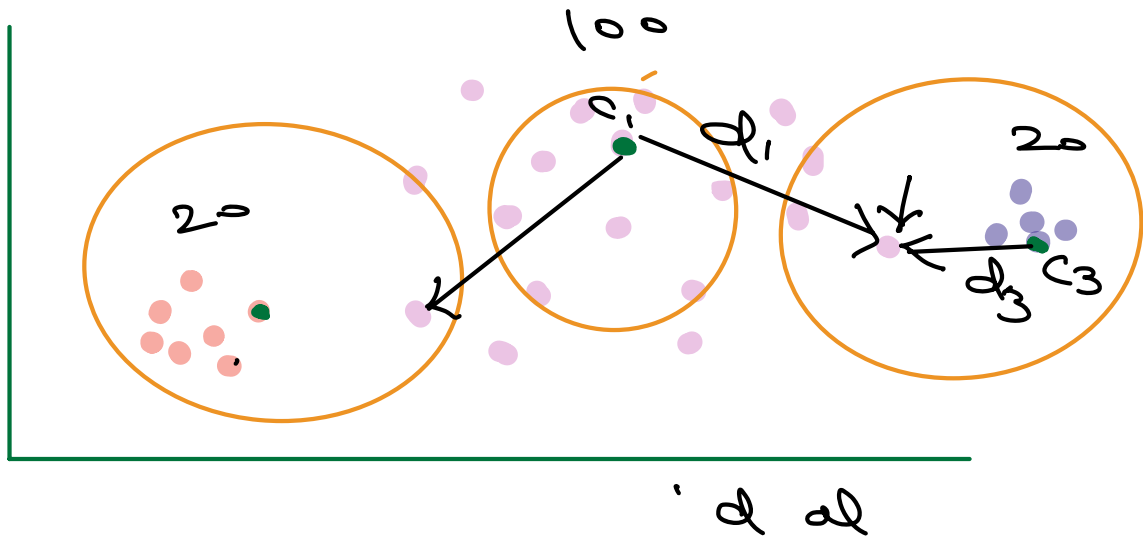


* The chances of a point getting picked will be proportional to its distance from previous Centroids

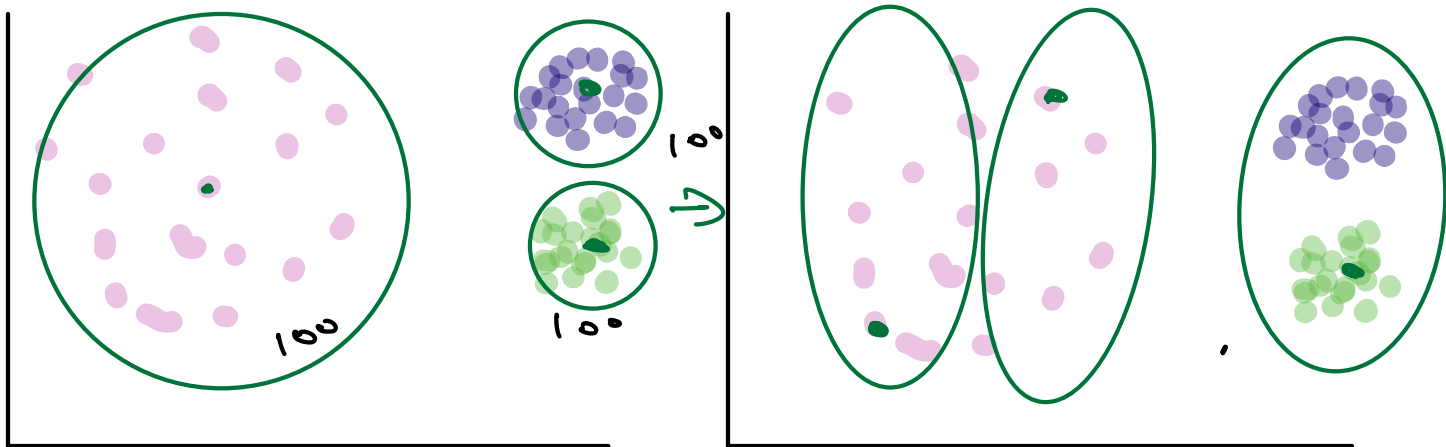
Limitations of kmeans/kmeant++

① Different sized Clusters

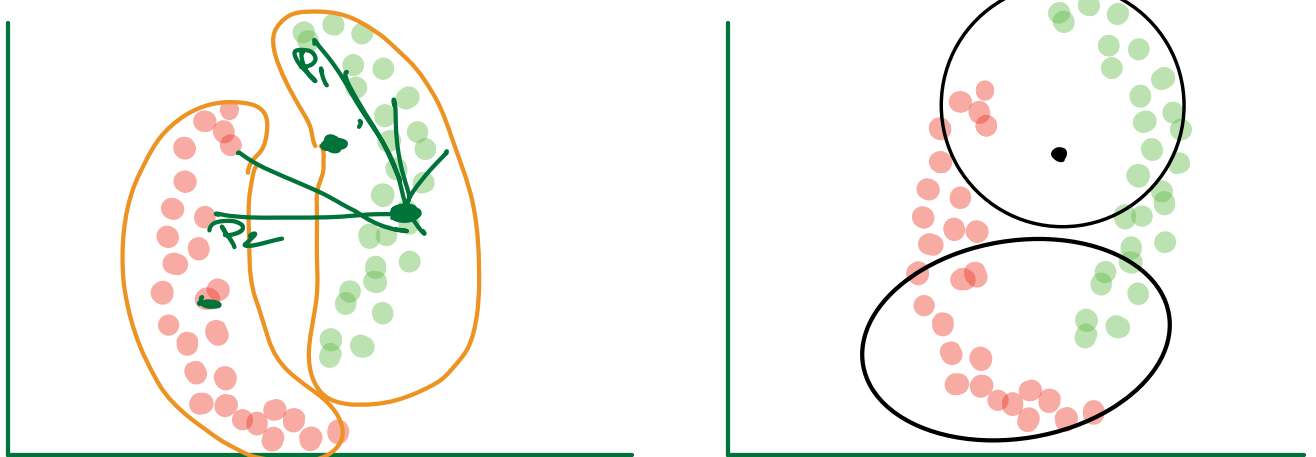




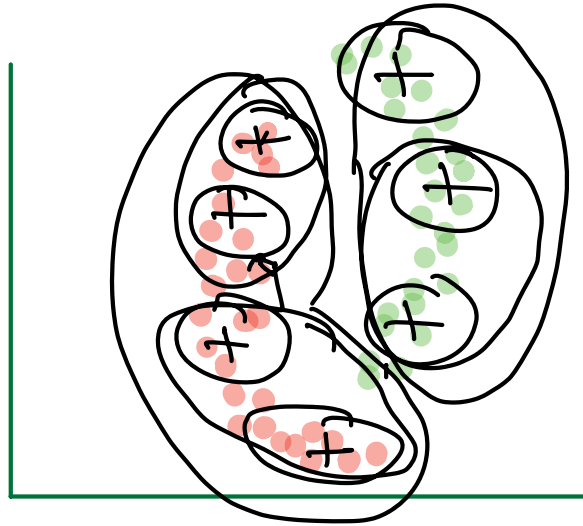
② Density of Clusters is Different



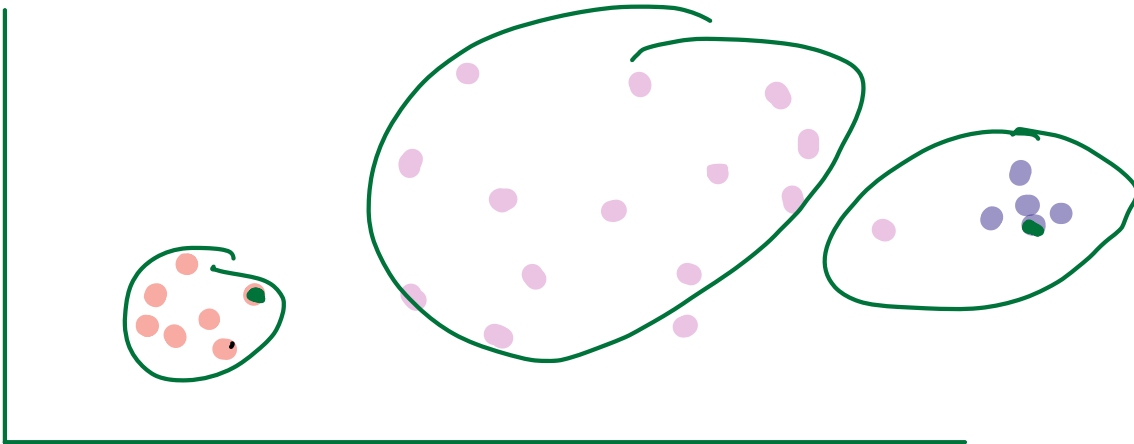
③ Clusters are Not Globular Shaped



Resolution



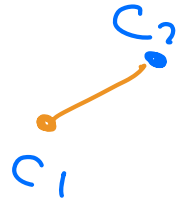
2



2 2

K - median

- ① Initialization
- ② Assignment
- ③ Update Step



In k-means, we were taking mean of data points assigned points

* In k-median, we use median of distance to update centroid.

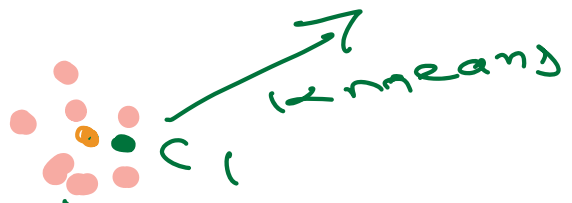
Sorted Distances

d_1 d_2 d_3

d_4 - - - - d_9

Which scenario can this help in?

Outliers



① What if you want the centroid to be one of existing point always

K-medoids

$n \text{ cluster} = 2 = k$



①

$$\text{Loss} \equiv \sum_{i=1}^k \sum_{x_i \in C_i} \text{distance}(x_i, m_i)$$

Iterate over all Data

② Find the value of m_i such that loss is minimum