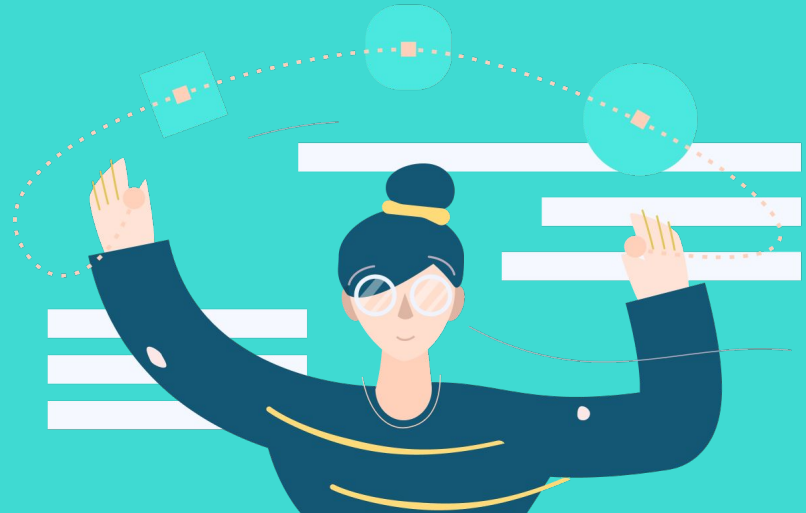




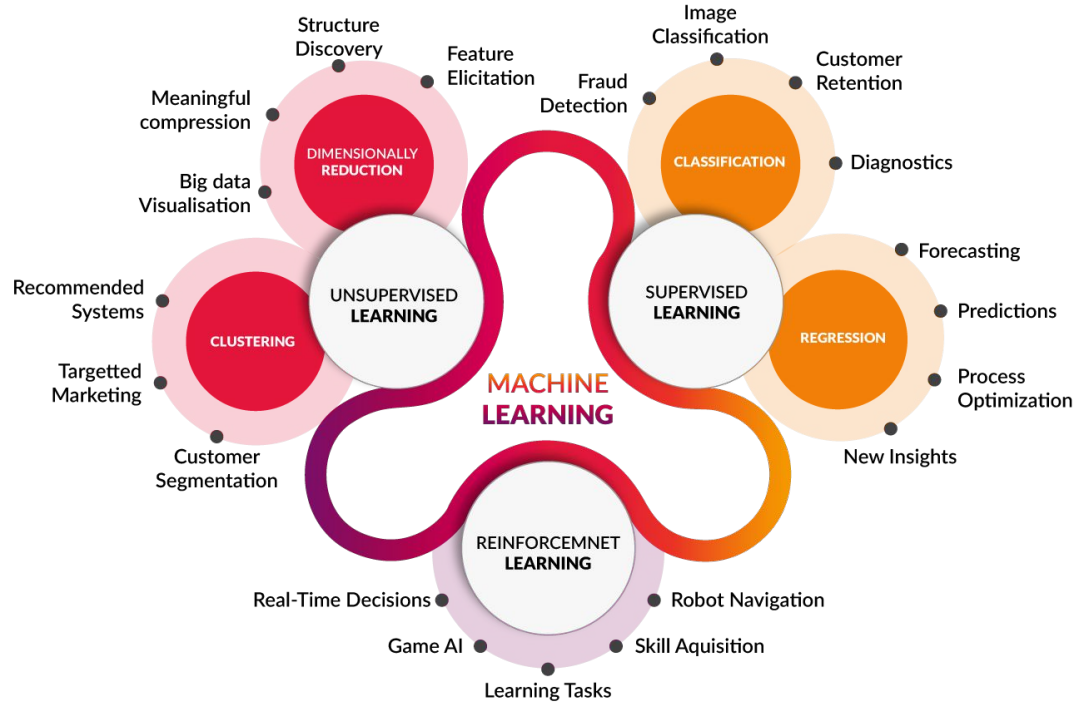
# K-Means

Clustering





# Machine Learning types






# Unsupervised Learning usecases

- **Dimensionality Reduction** ⇒ Compress datasets
- **Clustering** ⇒ Divide data into meaningful groups



# Supervised vs Unsupervised

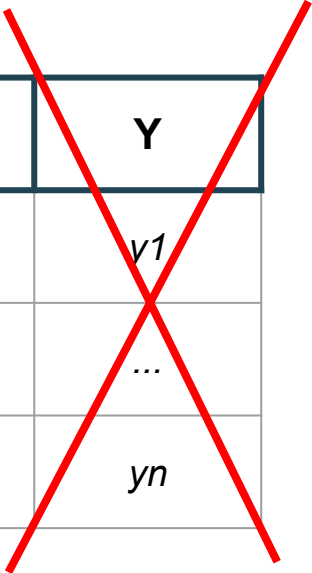
Predict



X1	X2	Y
$x_1$	$x_2$	$y_1$
...	...	...
$x_n$	$x_n$	$y_n$

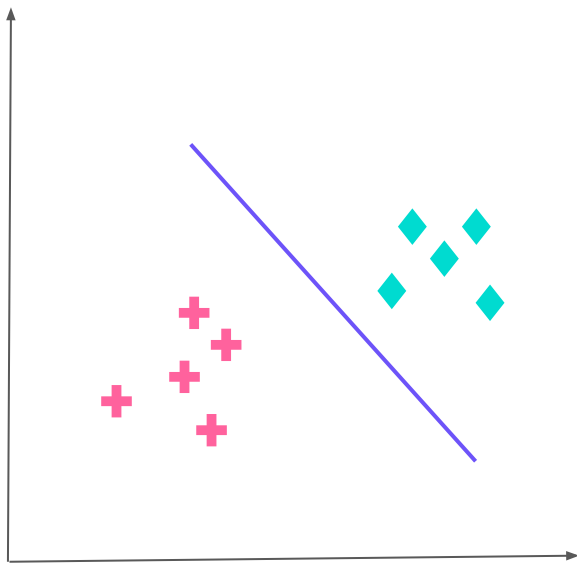
VS

X1	X2	Y
$x_1$	$x_2$	$y_1$
...	...	...
$x_n$	$x_n$	$y_n$

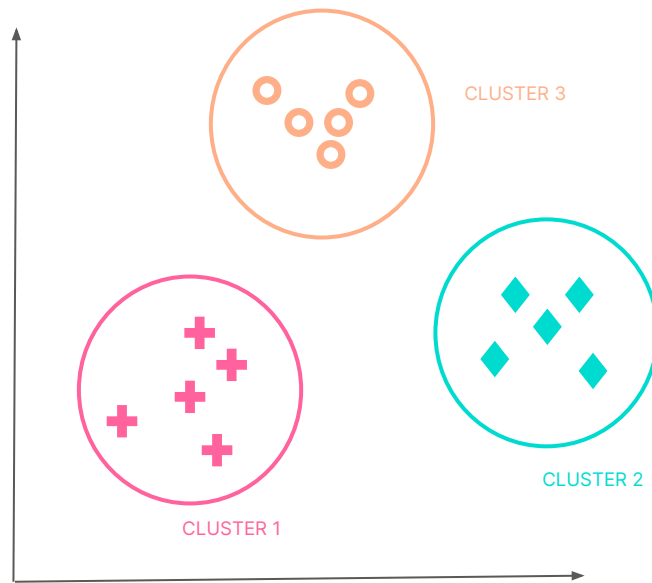




# Supervised vs Unsupervised



Supervised classification



Unsupervised Clustering



# K-Means



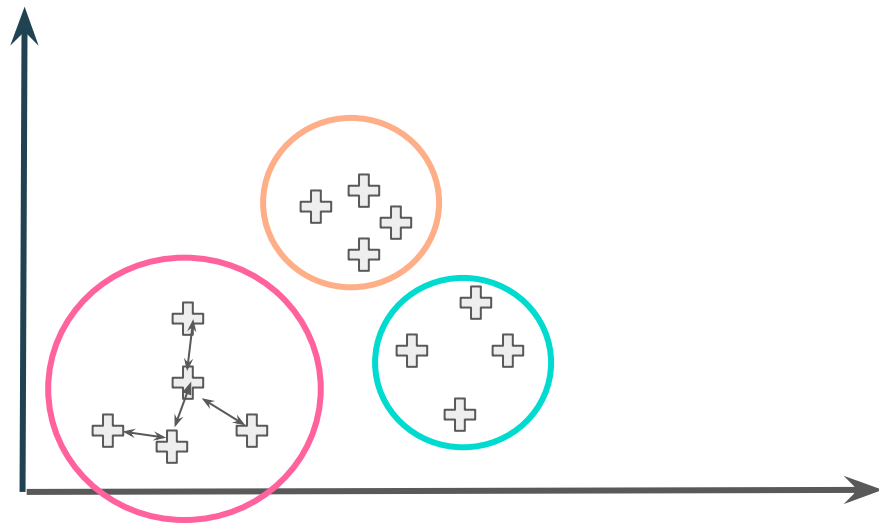
# K-Means Algorithm

The **KMeans** algorithm clusters data by trying to **separate samples in  $n$  groups** of equal variance, minimizing a criterion known as the **inertia or within-cluster sum-of-squares**

Source: [Sklearn](#)



# K-Means Algorithm



$$WCSS = \sum_{i=0}^n \min_{\mu_j \in C} (||x_i - \mu_j||^2)$$

$C$  = clusters

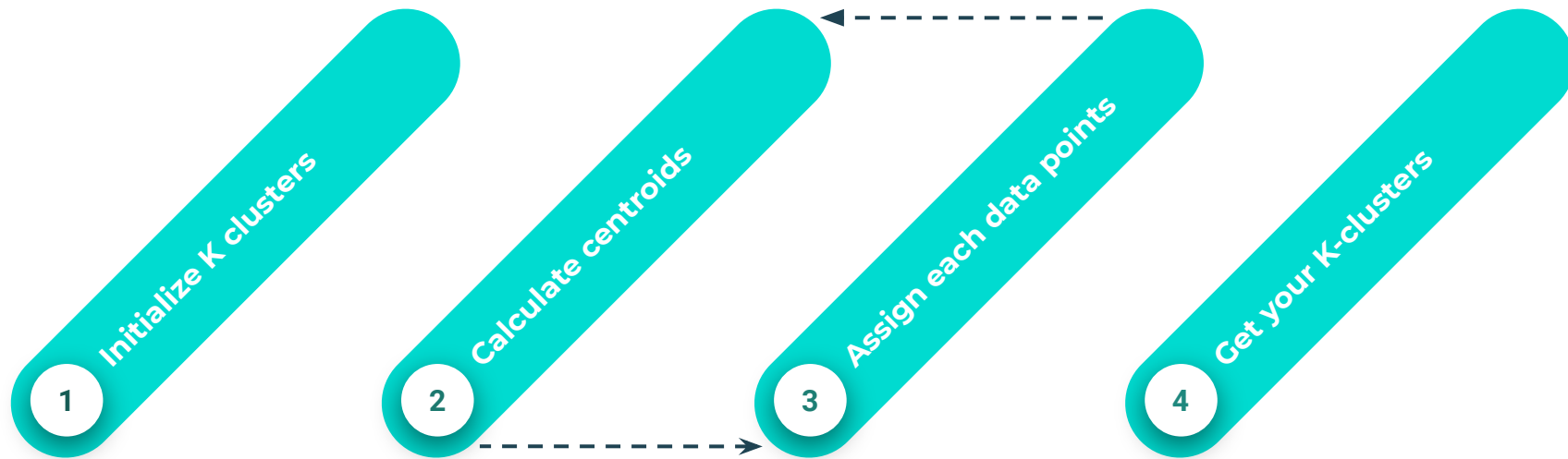
$\mu$  = mean within a given cluster

$x$  = data point





# Process



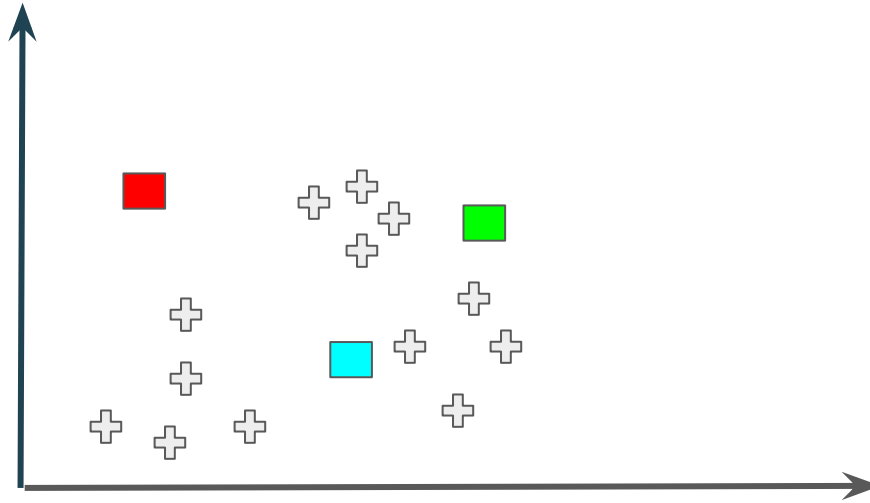


## Step 1 - Initialize K clusters

$$K = 3$$

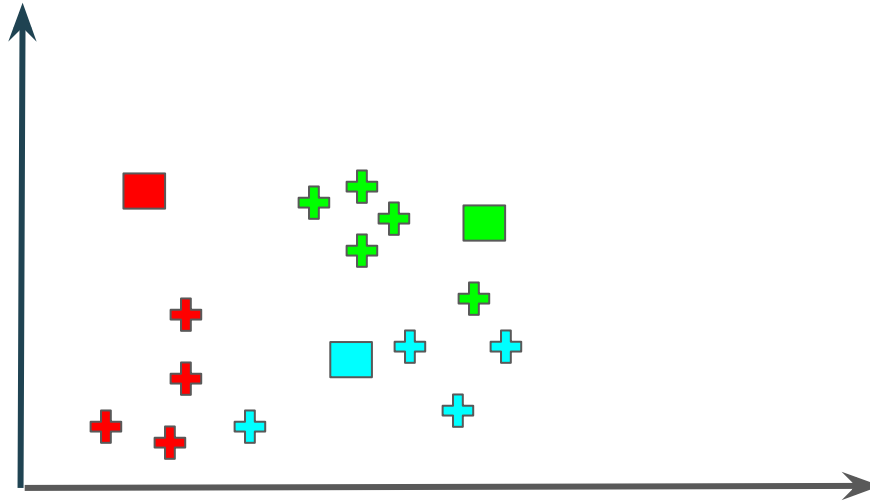


## Step 2 - Calculate centroids



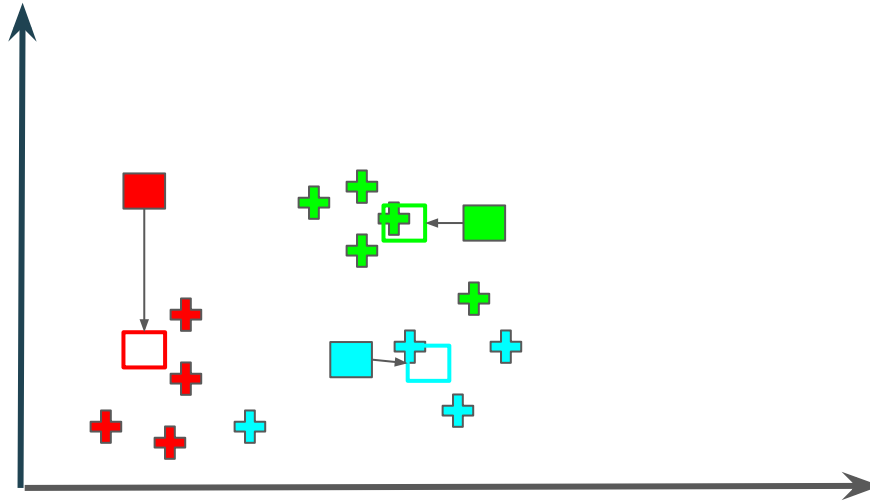


## Step 3 - Assign data points to the closest centroid



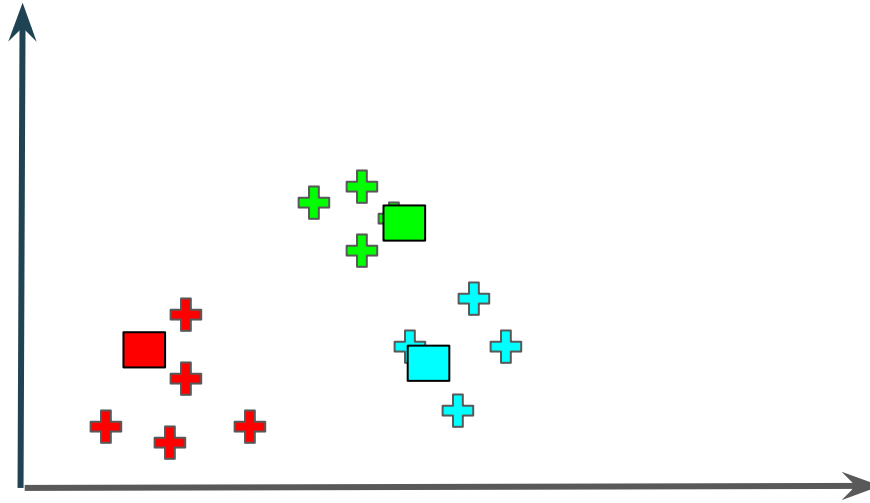


## Step 2 - Calculate centroids



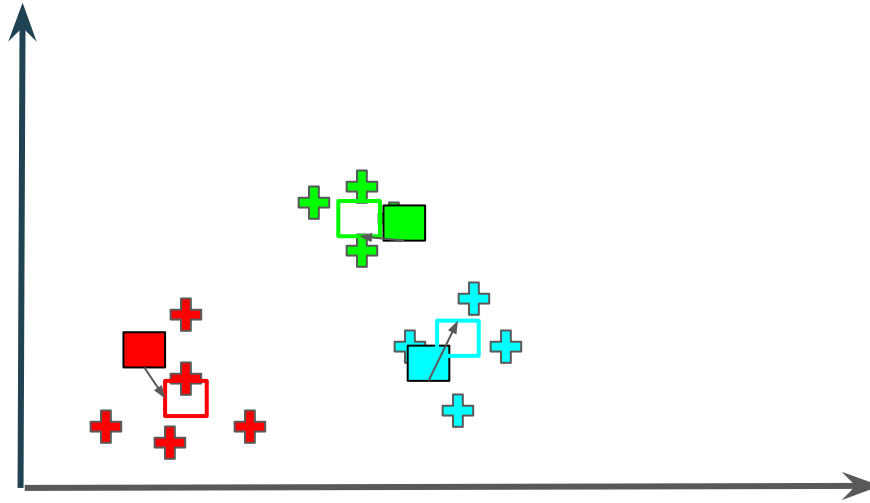


## Step 3 - Assign data points to the closest centroid



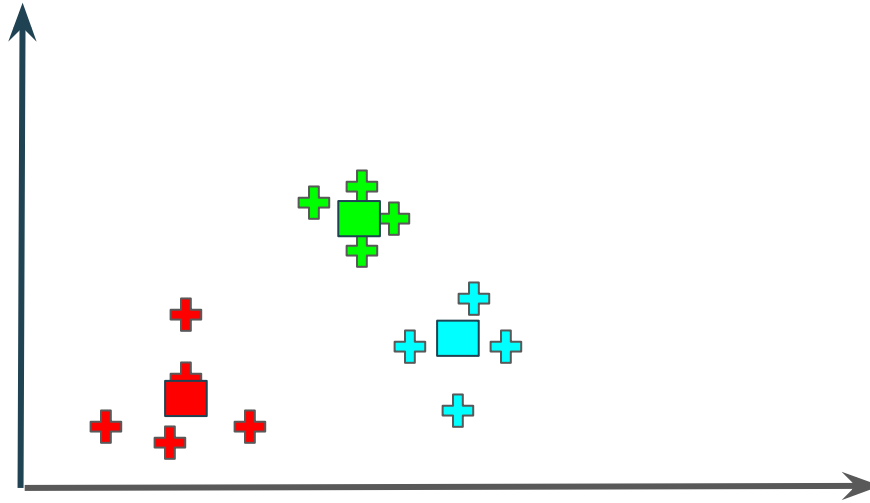


## Step 2 - Calculate centroids





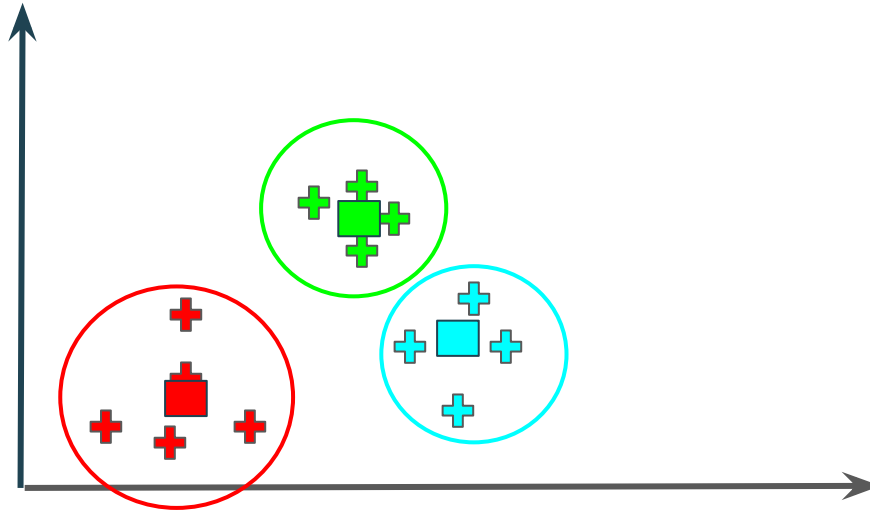
## Step 3 - Assign data points to the closest centroid







## Step 4 - Get our K clusters





**How to get the  
optimal number K?**

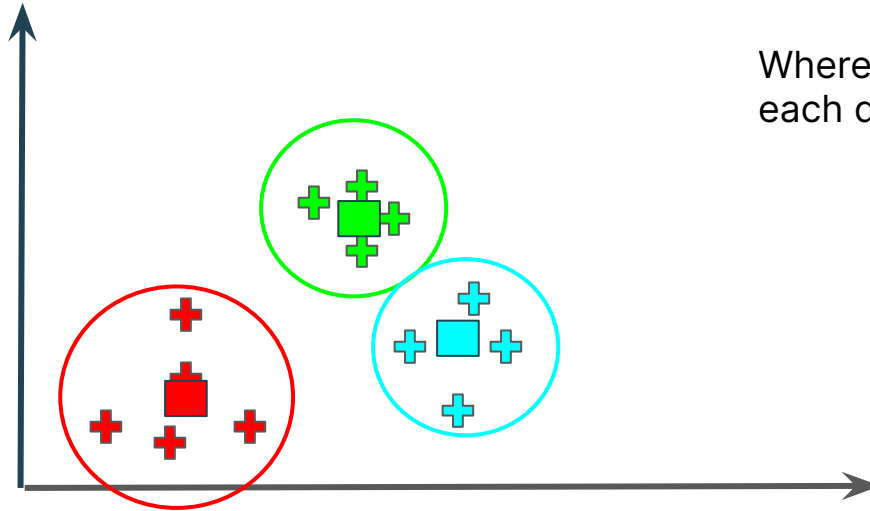


## Two methods

- **Elbow**  $\Rightarrow$  Check if data points within a cluster are close from their centroid
- **Silhouette**  $\Rightarrow$  Check if clusters are far from each other



# Elbow



$$WCSS = D(\text{+}, \blacksquare) + D(\text{+}, \blacksquare) + D(\text{+}, \blacksquare)$$

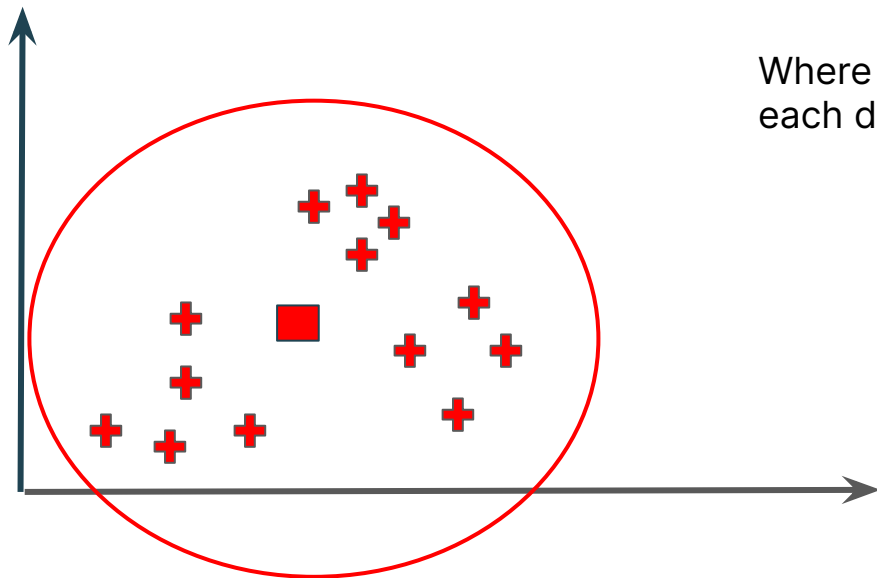
Where  $D()$  is the square distance from each data point to its centroid.



## Elbow - WCSS with 1 cluster

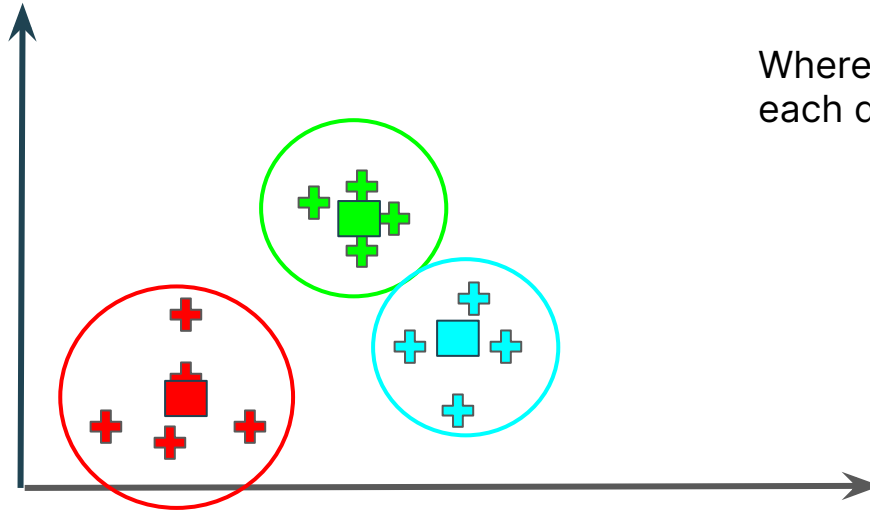
$$WCSS = D(+, \blacksquare)$$

Where  $D()$  is the square distance from each data point to its centroid.





## Elbow - WCSS with 3 clusters

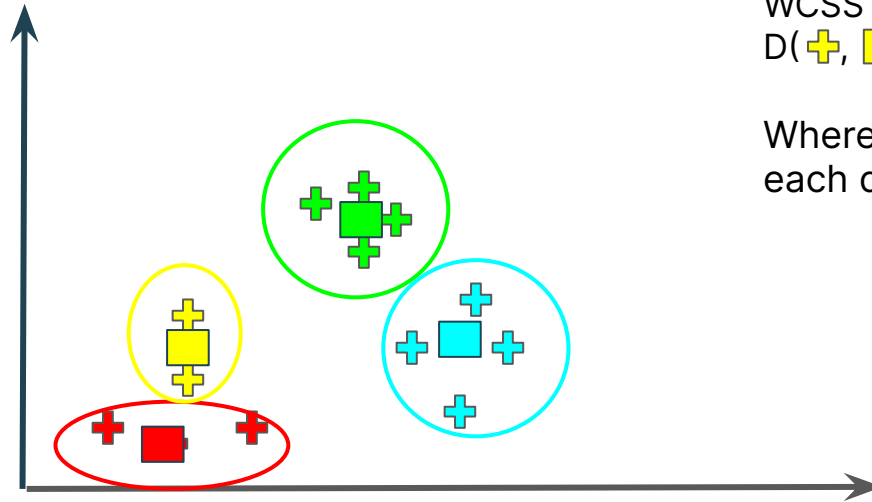


$$WCSS = D(\text{+}, \blacksquare) + D(\text{+}, \blacksquare) + D(\text{+}, \blacksquare)$$

Where  $D()$  is the square distance from each data point to its centroid.



## Elbow - WCSS with 4 clusters

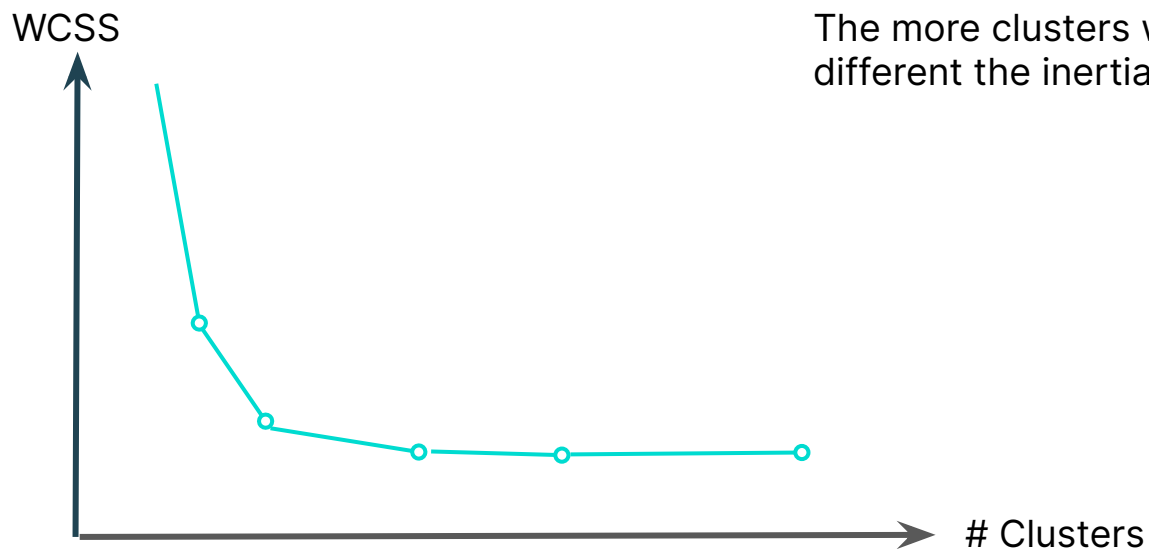


$$\text{WCSS} = D(\text{+}, \blacksquare) + D(\text{+}, \blacksquare) + D(\text{+}, \blacksquare) + D(\text{+}, \blacksquare)$$

Where  $D()$  is the square distance from each data point to its centroid.



## Elbow - WCSS with 4 clusters



The more clusters we get the less different the inertia (WCSS) is going to be





## Silhouette

$$\frac{b^i - a^i}{\max(a^i, b^i)}$$

### Where

a → average distance from all data points in the same cluster

b → average distance from all data points in the closest cluster



## Silhouette

- **Close to 0**  $\Rightarrow$  Clusters are close from each other
- **Close to 1**  $\Rightarrow$  Clusters are far from each other



# Thanks!

See you in the next course

