

# K-Means Clustering

# Centroid Based Clustering (**K-Means**)

- Select k number of classes and initialize their respective center points
- The center points are vectors of the same length as each data point vector
- Each data point is classified to be in the group whose center is closest to it
- Based on these classified points, recompute the group center
- Repeat these steps until the expected result
- Consider the following 2D points and find 3 cluster centroids for those

X	1	-3	0	-7	0	6	2	-1	0	10
Y	2	6	3	-4	7	9	3	-1	-100	10

# Centroid Based Clustering (**K-Means**)

<b>X</b>	<b>1</b>	<b>-3</b>	<b>0</b>	<b>-7</b>	<b>0</b>	<b>6</b>	<b>2</b>	<b>-1</b>	<b>0</b>	<b>10</b>
<b>Y</b>	<b>2</b>	<b>6</b>	<b>3</b>	<b>-4</b>	<b>7</b>	<b>9</b>	<b>3</b>	<b>-1</b>	<b>-100</b>	<b>10</b>
<b>Cluster</b>										

<b>Distance (C1)</b>										
<b>Distance (C2)</b>										
<b>Distance (C3)</b>										

<b>C1</b>	<b>(1,2)</b>
<b>C2</b>	<b>(-3,6)</b>
<b>C3</b>	<b>(0,3)</b>

**For simplicity, we will use squared distances**

# Centroid Based Clustering (**K-Means**)

<b>X</b>	<b>1</b>	<b>-3</b>	<b>0</b>	<b>-7</b>	<b>0</b>	<b>6</b>	<b>2</b>	<b>-1</b>	<b>0</b>	<b>10</b>
<b>Y</b>	<b>2</b>	<b>6</b>	<b>3</b>	<b>-4</b>	<b>7</b>	<b>9</b>	<b>3</b>	<b>-1</b>	<b>-100</b>	<b>10</b>
<b>Cluster</b>	<b>1</b>	<b>2</b>	<b>3</b>	<b>3</b>	<b>2</b>	<b>3</b>	<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>

<b>Distance (C1)</b>	<b>0</b>	<b>32</b>	<b>2</b>	<b>100</b>	<b>26</b>	<b>74</b>	<b>2</b>	<b>13</b>	<b>10405</b>	<b>145</b>
----------------------	----------	-----------	----------	------------	-----------	-----------	----------	-----------	--------------	------------

<b>Distance (C2)</b>	<b>32</b>	<b>0</b>	<b>18</b>	<b>116</b>	<b>10</b>	<b>90</b>	<b>34</b>	<b>53</b>	<b>11245</b>	<b>185</b>
----------------------	-----------	----------	-----------	------------	-----------	-----------	-----------	-----------	--------------	------------

<b>Distance (C3)</b>	<b>2</b>	<b>18</b>	<b>0</b>	<b>98</b>	<b>16</b>	<b>72</b>	<b>4</b>	<b>17</b>	<b>10609</b>	<b>149</b>
----------------------	----------	-----------	----------	-----------	-----------	-----------	----------	-----------	--------------	------------

<b>C1</b>	<b>(1,2)</b>
<b>C2</b>	<b>(-3,6)</b>
<b>C3</b>	<b>(0,3)</b>

**Update the cluster centroid**

<b>C1</b>	<b>(2.4,-17.2)</b>
<b>C2</b>	<b>(-1.5,6.5)</b>
<b>C3</b>	<b>(-0.3,2.7)</b>

# Centroid Based Clustering (**K-Means**)

X	1	-3	0	-7	0	6		-1	0	10
Y	2	6	3	-4	7			-100	10	
Cluster										

Distance (C1)										
---------------	--	--	--	--	--	--	--	--	--	--

Distance (C2)										
---------------	--	--	--	--	--	--	--	--	--	--

Distance (C3)										
---------------	--	--	--	--	--	--	--	--	--	--

C1	
C2	
C3	(-3, 2.7)

Update the cluster centroid

C1	
C2	
C3	

Do Yourself

# How to Find a Proper **k** value?

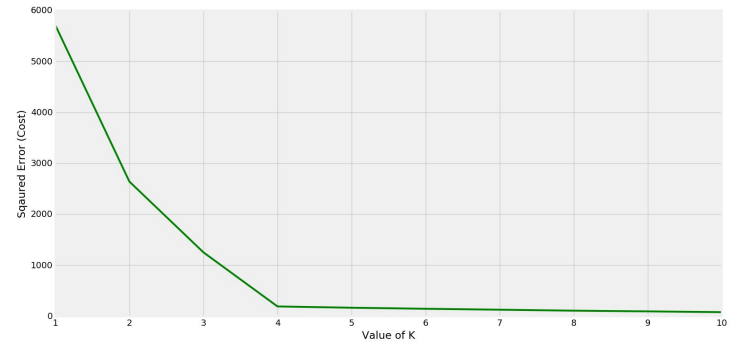
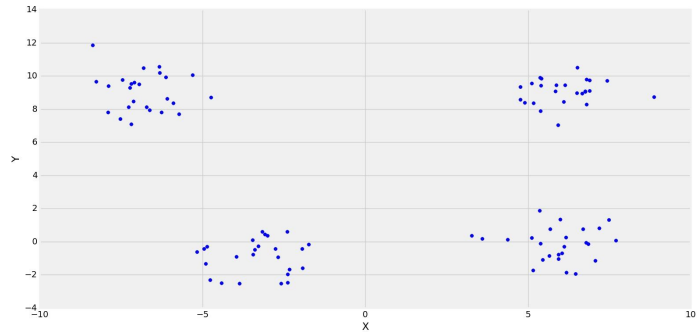
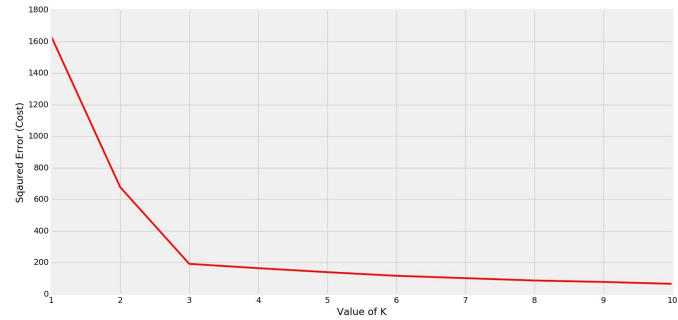
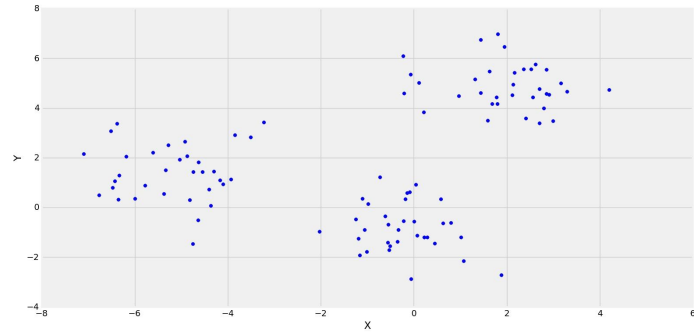
A number of analysis are used:

- Elbow Method
- Average Silhouette Method
- Gap Statistic Method

# Elbow Method

- Compute clustering algorithm for different values of  $k$
- For each  $k$ , calculate the sum of intra-cluster squared error (sse)
- Plot the curve of **sse vs  $k$**
- The location of a bend (knee) in the plot is generally considered as an indicator of the appropriate number ( $k$ ) of clusters

# Elbow Method





# Average Silhouette Method

- Compute clustering algorithm for different values of  $k$
- For each  $k$ , calculate the average silhouette of observations ( $\text{avg.sil}$ )
- Plot the curve of **avg.sil vs  $k$**
- The location of the maximum is considered as the appropriate number ( $k$ ) of clusters
- Silhouette Coefficient is calculated as:

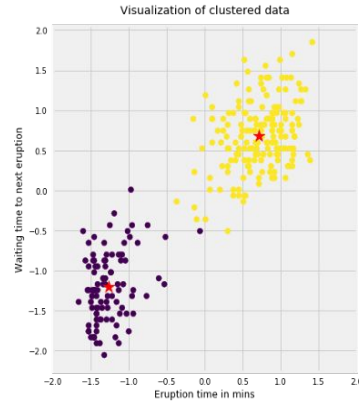
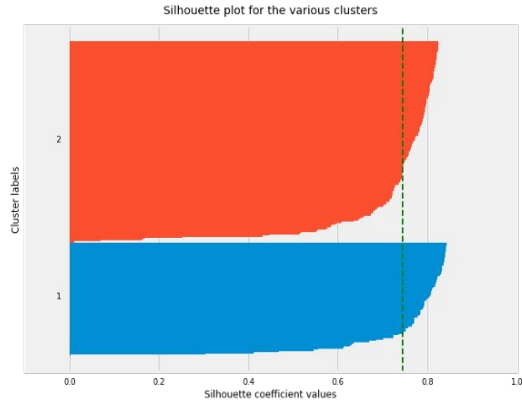
$$\text{silCof} = \frac{\beta - \alpha}{\max(\beta, \alpha)}$$

$\alpha = \text{average intra-cluster distance}$

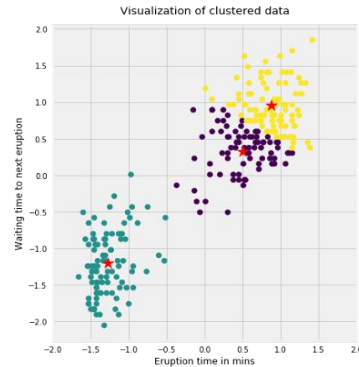
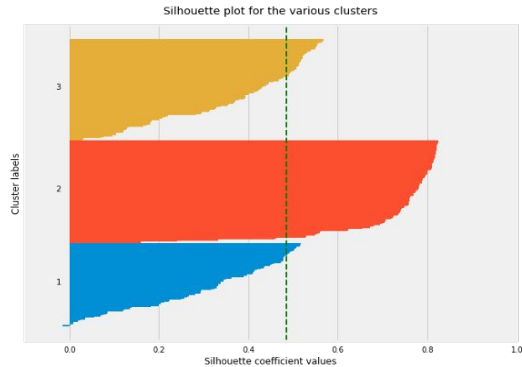
$\beta = \text{minimum average inter-cluster distance}$

# Average Silhouette Method

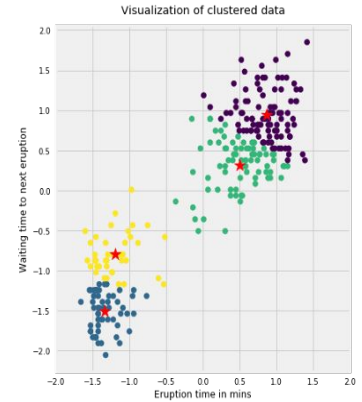
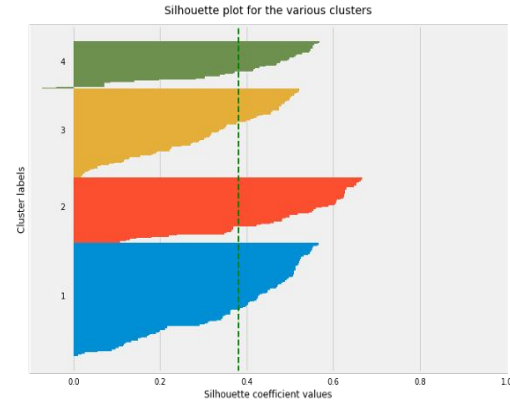
Silhouette analysis using  $k = 2$



Silhouette analysis using  $k = 3$



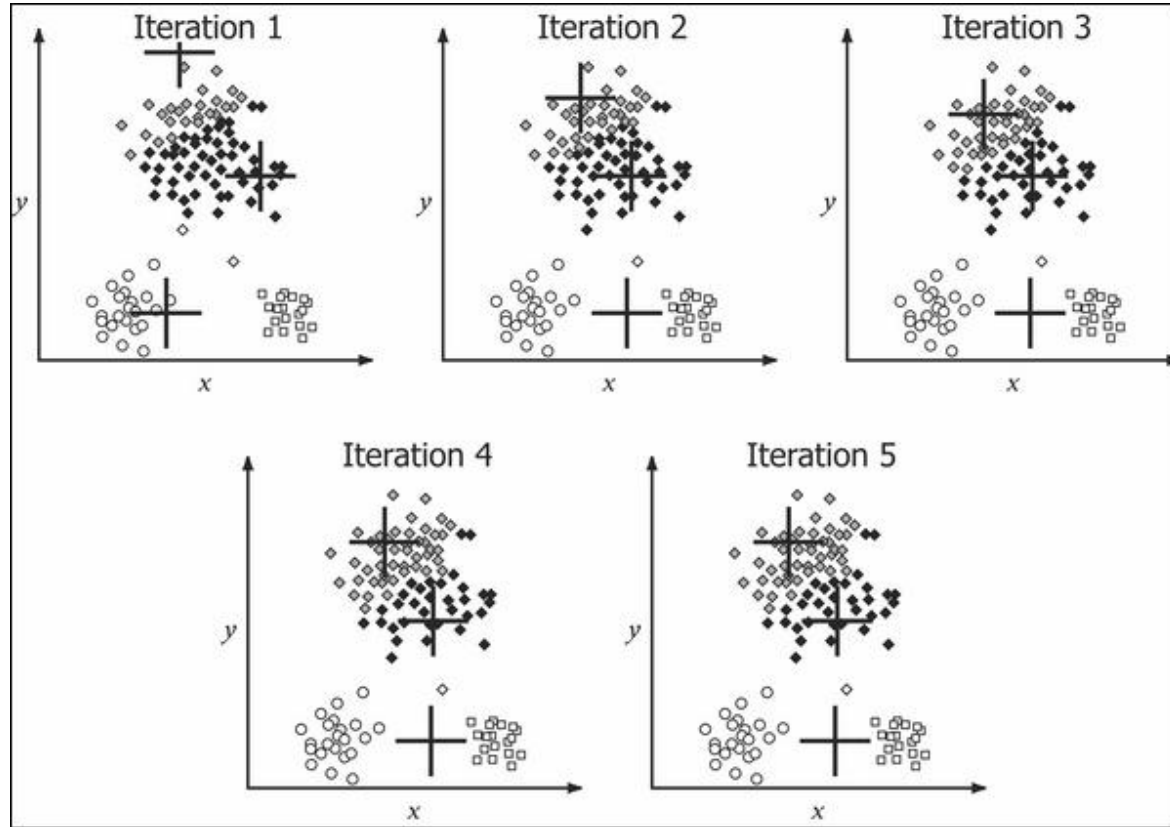
Silhouette analysis using  $k = 4$



# Drawbacks of K-Means

- Works poor with complex geometric shapes of clusters
- Can't handle outlier
- Can't guarantee to find the global optimum clusters
- Can't handle non-numerical data
- ...

# Drawbacks of K-Means



# Assignment

**Find the optimum number of cluster for the given dataset using Elbow Method.**

Dataset Preparation: If your student id is  
ABCD-E-FG-HIJ

X	-1	7	2	0	9	-3	5	8	-6	4
Y	A	B	C	D	E	F	G	H	I	J

*Thank You*

# References

- <https://towardsdatascience.com/understanding-k-means-clustering-in-machine-learning-6a6e67336aa1>
- <https://towardsdatascience.com/k-means-clustering-algorithm-applications-evaluation-methods-and-drawbacks-aa03e644b48a>
- <https://www.geeksforgeeks.org/ml-determine-the-optimal-value-of-k-in-k-means-clustering/#:~:text=There%20is%20a%20popular%20method,fewer%20elements%20in%20the%20cluster.>
- <http://user.ceng.metu.edu.tr/~akifakkus/courses/ceng574/k-means/>
-