

1000

① K-means clustering

- ↳ Clustering is the process of grouping observations in such way that members of same cluster are grouped together and members of different cluster are grouped together.
- ↳ Clustering is commonly used to explore the datasets & identify the patterns.
- ↳ K-means clustering algorithm is an iterative process of moving the centers to mean positions until there is no significant change.
- ↳ The cost function of K-means is determined by Euclidean distance between the observation.
- ↳ If there is one cluster ($K=1$), then the ~~distance b/w~~ distance b/w all obs. are compared with its single mean.
- ↳ If the cluster is increased to 2 ($K=2$), then, two means are calculated and few of them are assigned to cluster 1 and remaining to cluster 2.

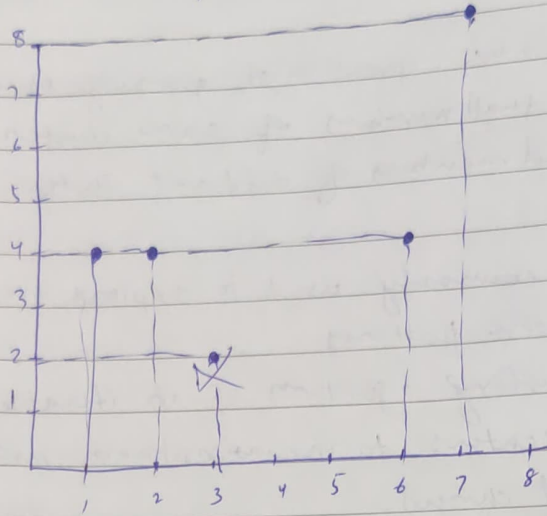
$$J = \sum_{k=1}^K \sum_{i \in C_k} \|x_i - \mu_k\|^2$$

example

Let's take 4 instances, with their x & y values

Instances	x	y
1	7	8
2	2	4
3	6	4
4	1	4

→ Plot the data points on 2D chart



→ Iteration 1: let's choose instance $(x=7, y=8)$ → centroid 2
 instance $(x=1, y=4)$ → centroid 1

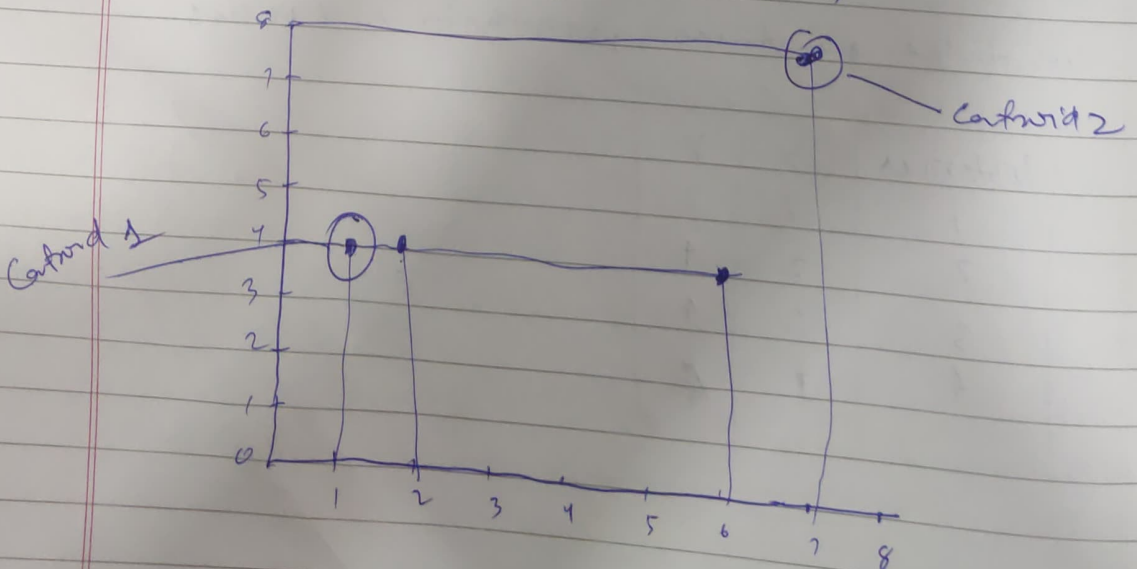
calculate euclidean distance

$$\frac{4+6}{10+36}$$

$$\frac{4+5}{10+36}$$

Instance	X	Y	Centroid 1 dist	Centroid 2 dist	Assign cluster
1	7	8	7.21 7.21	0	C2
2	2	4	2.81 2.81	6.40	C1
3	6	4	8.00 5.00	4.12	C2
4	1	4	0	7.21	C1
Centroid 1	1	4			
Centroid 2	7	8			

$$\text{Euclidean distance} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$



→ Iteration 2: calculate the new centroids

Instance	<u>X</u>	<u>Y</u>	Assigned cluster
1	7	8	C2
2	2	4	C1
3	6	4	C2
4	1	4	C1
Centroid 1	1.5	4.0	
Centroid 2	6.5	6	

Centroid 1 coordinates = Avg coordinates (instances 2, 4)

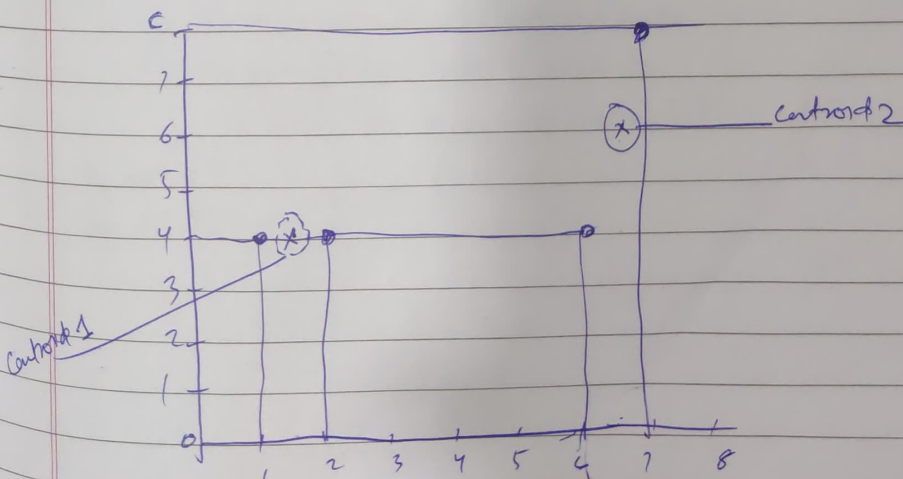
Centroid 2 coordinates = Avg coordinates (instances 1, 3)

$$\text{Centroid 1 } x = \frac{2+1}{2} = \frac{3}{2} = 1.5$$

$$\text{Centroid 1 } y = \frac{4+4}{2} = \frac{8}{2} = 4.0$$

$$\text{Centroid 2 } x = \frac{7+6}{2} = \frac{13}{2} = 6.5$$

$$\text{Centroid 2 } y = \frac{8+4}{2} = \frac{12}{2} = 6$$



30.25 +

Iteration 3: Check for the new cluster is new or same and calc. the centroids dist

$C1(1.5, 4.0)$ $C2(6.5, 6)$

<u>Instance</u>	<u>X</u>	<u>Y</u>	<u>C1 dist</u>	<u>C2 dist</u>	<u>old cluster</u>	<u>New cluster</u>	<u>Changed?</u>
1	7	8	6.80	2.06	C2	C2	No
2	2	4	0.5	4.92	C1	C1	No
3	4	4	4.5	1.5	C2	C2	No
4	1	4	0.5	5.85	C1	C1	No

Since there is no change
Hence soln is converged.

② PCA (Principal Component Analysis)

- ↳ It is dimensionality reduction method
- ↳ used to reduce the dimensions of large dataset by transforming the large dataset to smaller dataset
- ↳ While reduction, it may lead in reducing accuracy but using PCA it will not effect
- ↳ It maximizes the variance of the data
- ↳ It minimizes the mean squared distance

Need:

- Removes Inconsistent data
- Removes redundant data
- Makes data processing faster
- No loss of information

PCA steps:

- ① Standardization of the data
- ② Computing covariance matrix
- ③ calculating eigenvectors and eigenvalues
- ④ Computing principal components
- ⑤ Reducing the dimensions of dataset

Step 1: Standardization of the data

↳ It is about scaling the data

↳ Get some data & plot it.

↳ Take the $\text{mean}(x)$ & $\text{mean}(y)$

$$\text{mean}(x) = 3$$

$$\text{mean}(y) = 6.5$$

x	y
1.5	2.4
2.5	3.6
3.5	4.2
4.5	2.8

Date _____
Page _____

Step 1: computing covariance matrix

$$\begin{bmatrix} \text{cov}(a,a) & \text{cov}(a,b) \\ \text{cov}(b,a) & \text{cov}(b,b) \end{bmatrix}$$

→ The covariance values depends:

- (1) -ve covariance - indirectly proportional
- (2) +ve covariance - directly proportional

→ It defines the correlation b/w the different variable in dataset.

→ This is called DATA ADJUSTMENT

Step 3: calculating eigenvectors and eigenvalues

eigenvectors: These are the vectors when performed on them their direction does not changes

eigenvalues: denotes the scalar of eigenvectors.

They are the mathematical concepts that can be created using covariance matrix.

calculate covariant matrix

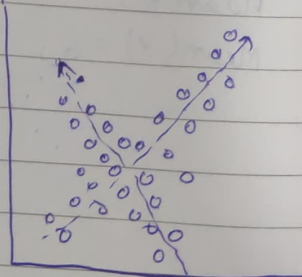
$$\text{cov}(X, X) = \frac{\sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})}{n-1}$$

	X	Y
X	$\text{cov}(X, X)$	$\text{cov}(X, Y)$
Y	$\text{cov}(X, Y)$	$\text{cov}(Y, Y)$

Step 4: Computing principal components.

→ PC1 is most significant and stores max possible info.

→ PC2 is second most significant and stores the remaining max info.



→ After calc. eigen vectors and values, we need to arrange them in descending order

Steps: Rearranging the dimensions:

↳ Rearrange the original data

↳ The original data will be redefined

↳ The eigen vectors with highest eigen value will be the PCA

PCA from first principles

① Take instances with (X and Y)

Instance	X	Y
1	0.1	0.6
2	0.2	0.7
3	0.3	0.8
4	0.4	0.9
5	0.5	0.1
Mean	0.3	0.62

$$\frac{3.1}{5 \times 10} = 0.2$$

② Remove the scale factors.

X	Y
$0.1 - 0.3 = -0.2$	$0.6 - 0.62 = -0.02$
$0.2 - 0.3 = -0.1$	$0.7 - 0.62 = 0.08$
$0.3 - 0.3 = 0$	$0.8 - 0.62 = 0.18$
$0.4 - 0.3 = 0.1$	$0.9 - 0.62 = 0.28$
$0.5 - 0.3 = 0.2$	$0.1 - 0.62 = -0.52$

$$\textcircled{3} \text{Cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

$$\text{Cov}(x, x) = \frac{(-0.2)^2 + (-0.1)^2 + (0)^2 + (0.1)^2 + (0.2)^2}{5-1} =$$

$$\text{Cov}(x, y) = \frac{(-0.2 \times -0.02) + (-0.1 \times 0.08) + \dots}{5-1} =$$

$$\text{Cov}(y, x) = \frac{(-0.02 \times -0.2) + (0.08 \times -0.1) + \dots}{5-1} =$$

$$\text{Cov}(y, y) = \frac{(-0.02)^2 + (0.08)^2 + \dots}{5-1} =$$

$$\text{Covariance matrix, } C = \begin{bmatrix} \text{Cov}(x, x) & \text{Cov}(x, y) \\ \text{Cov}(y, x) & \text{Cov}(y, y) \end{bmatrix}$$

Calculating the eigen vectors and values

$$\text{Eigen values} = \begin{bmatrix} \quad \quad \quad \end{bmatrix}$$

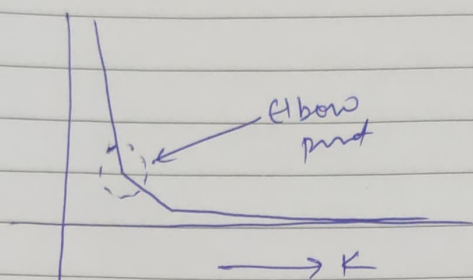
$$\text{Eigen vectors} = \begin{bmatrix} \quad \quad \quad \\ \quad \quad \quad \end{bmatrix}$$

The eigen vectors with highest eigen value will be the PCA.

am

① The elbow method

- ↳ It is the method used to find the optimal number of clusters
- ↳ It plots the values of cost function
- ↳ For value K , if $K \uparrow$, the avg. distortion will \downarrow
- ↳ The value of V where it is improved is called Elbow
- ↳ At this point we should stop dividing the data



② Evaluation of clusters.

- ↳ It can be done with ~~silhouette~~ silhouette coefficient.
- ↳ It is the measure of compactness and separation of clusters.
- ↳ Higher value sep. better quality of cluster
- ↳ values lies from -1 to $+1$
- ↳ Higher value will be better

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

a is mean dist in cluster
 b is mean dist of next cluster

③ PCA

Pros: Dimensionality Reduction
 : Data visualization
 : Improve data quality
 : easy to implement

Cons: Less accuracy
 : difficult to get result
 : slow