

There are two parts of Unsupervised Learning

- ✓ 1. KMeans Clustering ✓
2. Principal Component Analysis ✓

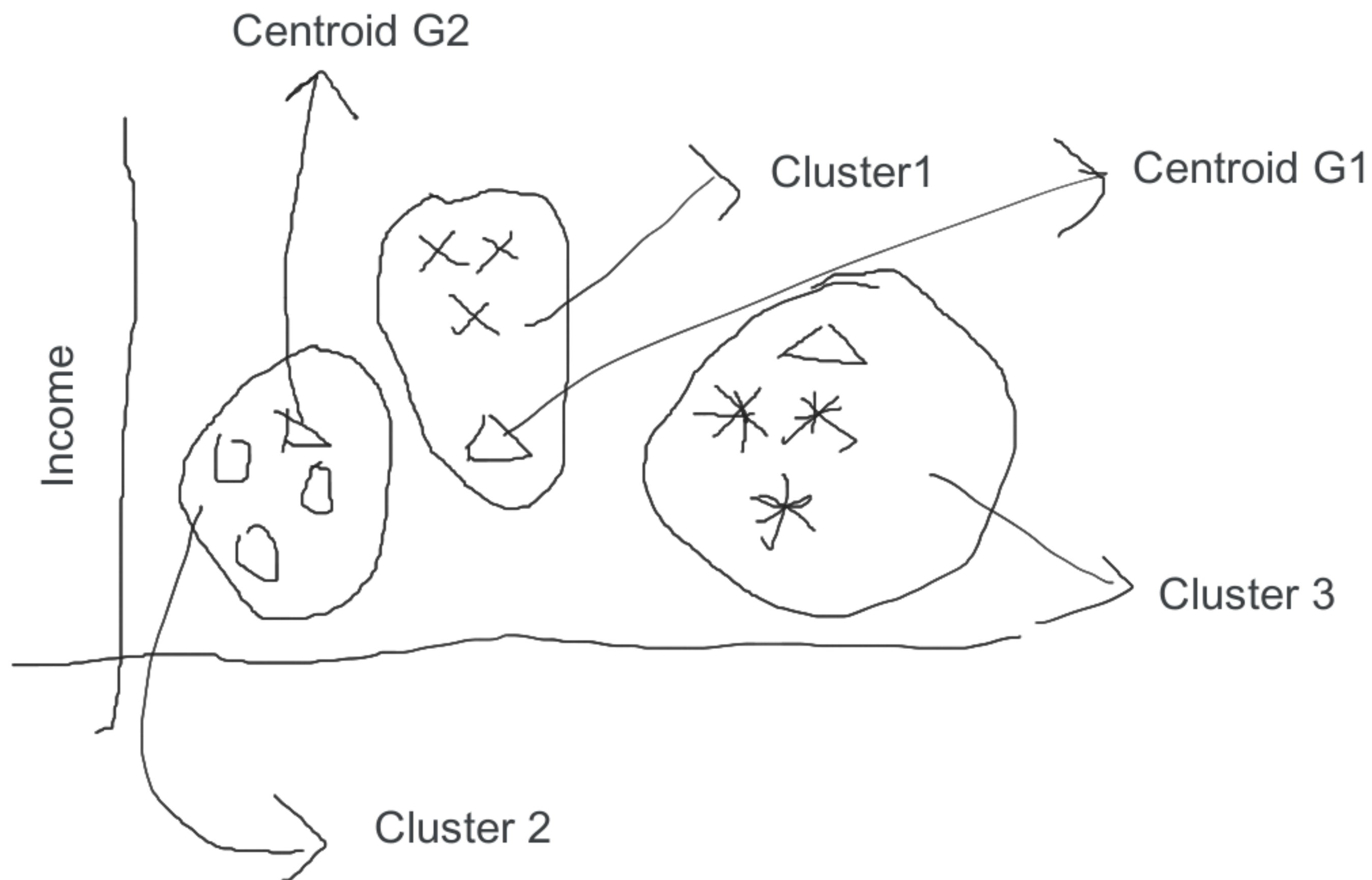
✓ K Means Clustering ✓

Important Case Studies

Walmart

Credit Card Segmentation

Working



Euclidean Distance Formula

$$d(x,y)(a,b) = \sqrt{(x-a)^2 + (y-b)^2}$$

(a, b)
 $(-2, 2)$

$\rightarrow 5$

✓

$$\sqrt{16 + 9}$$

$$= 5$$

$(2, -1)$
 (x, y)

✓ n_clusters=3 \Rightarrow K
V.Imp

How to decide how many clusters? \rightarrow how to choose the value of K

✓ 1. Business Objective - According to the client you will categorize/
clusterify/create groups. 100% of the time, follow the business objective.

✓ 2. Elbow Method

3. Silhouette Method

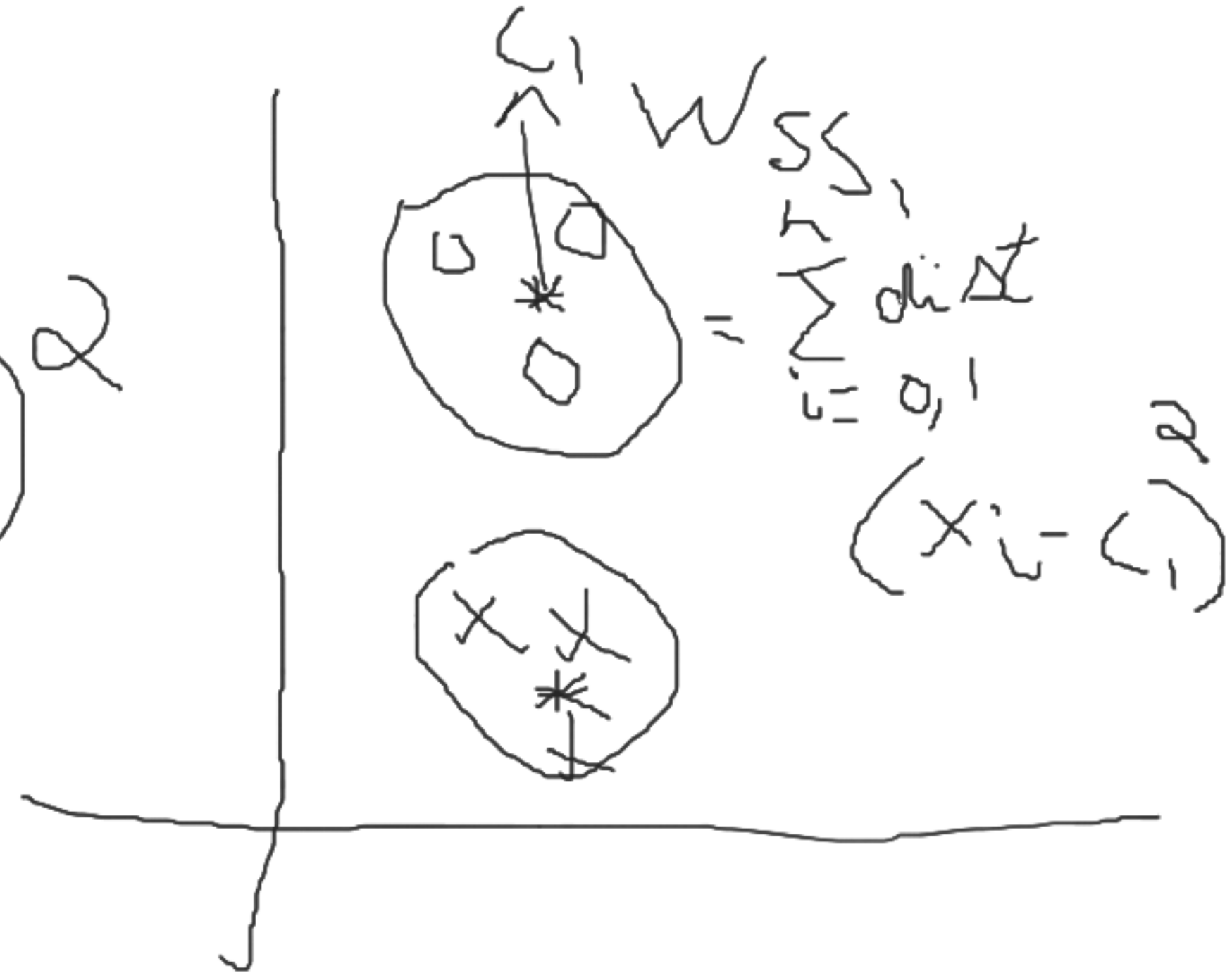
Elbow Method ✓

V.Imp

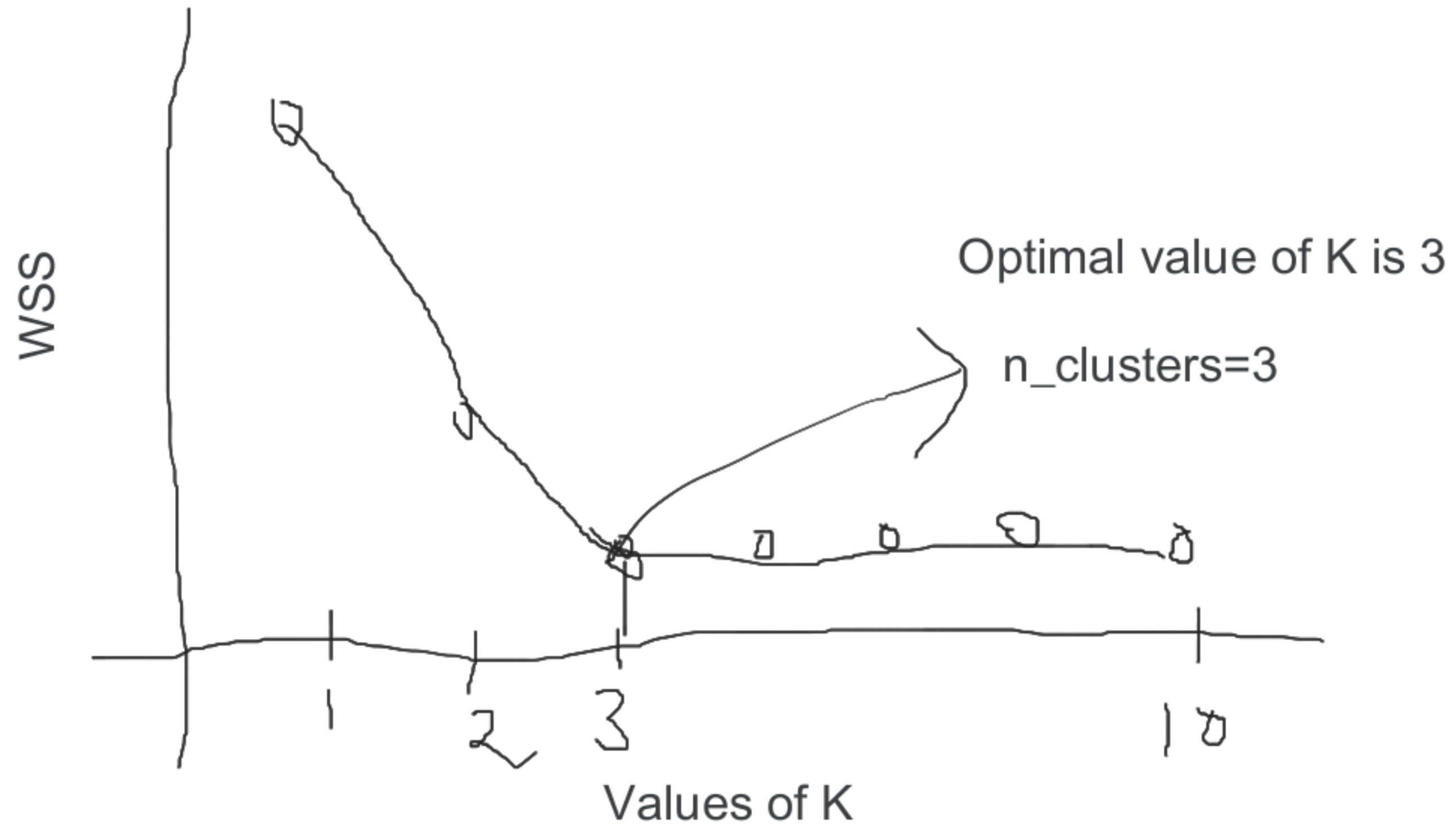
Within Sum of Squares

$$\underline{WSS} = WSS1 + WSS2 + \dots + WSSn$$

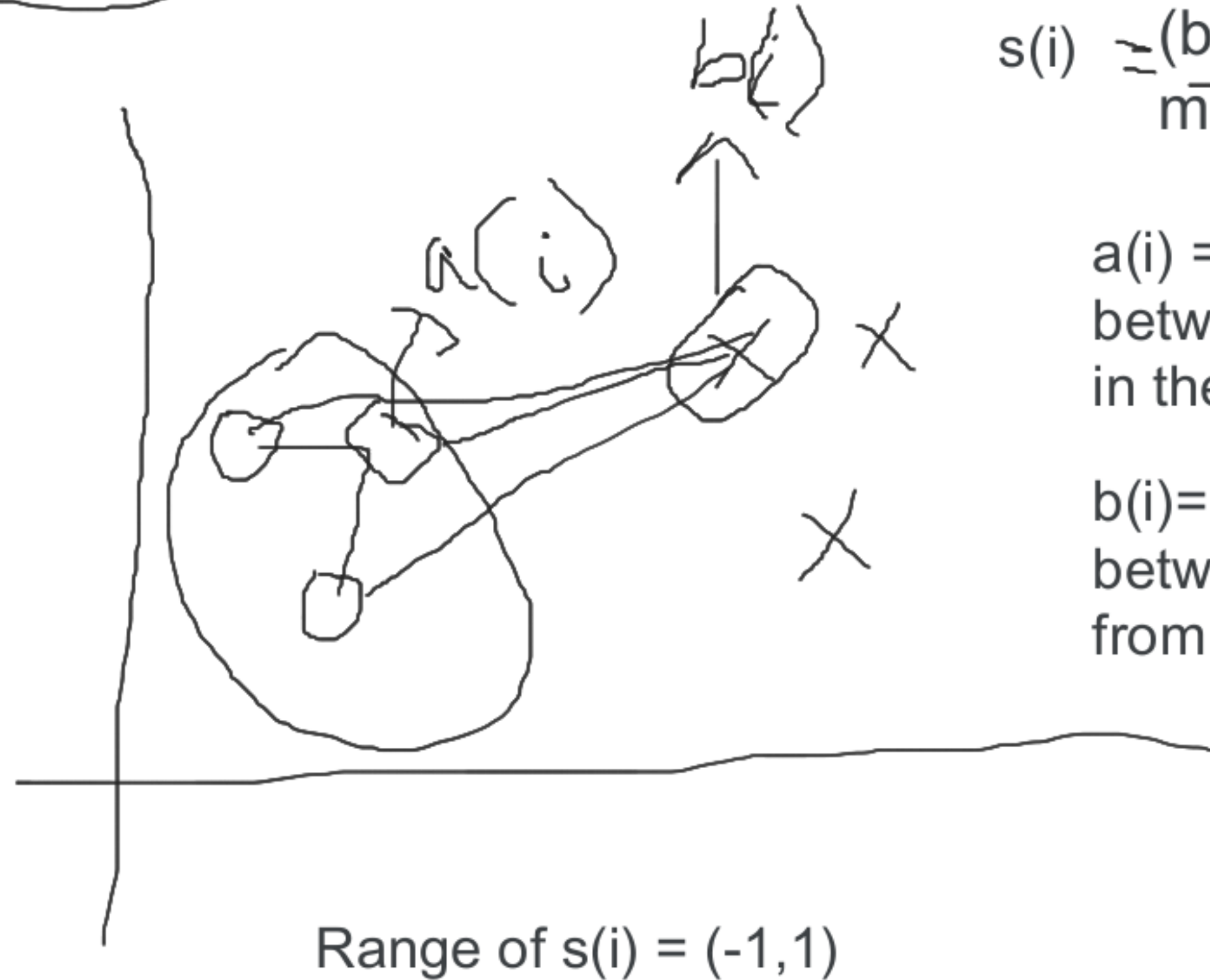
$$WSS1 = \sum_{i=0,1}^n \text{dist} (x_i - c_1)^2$$



Plotting of WSS values against values of K



Sillhoutte Method



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

$a(i)$ = Average distance
between every data point
in the same cluster.

$b(i)$ = Average distance
between every data point
from a different cluster

$s(i)$ = Silhouette Coefficient

If value is closer to -1 that means the data point is kind of dissimilar to the cluster or it does not belong to that cluster.

If value is closer to +1 that means the point belongs to the cluster.

Scaling of Data

MinMaxScaler → Range from 0 to +1

StandardScaler →

✓ Standard Scaler calculates Z Scores ✓

$$z = \frac{X - \bar{X}}{SD}$$

→ set of discrete values of x

→ Mean of X

→ $\frac{\sum (X - \bar{X})^2}{N}$

Need Of Scaling

X	Income(USD)	Cluster	Income(Rs.)	$X - \bar{X}$	$(X - \bar{X})^2$	$X - \bar{X}$
	6000	0	420000			
	7000	0	490000			
	4000	2	280000			
	5000	1	350000			
	5500	1	375000			