

K-means Clustering

Son Nguyen

What is clustering?

Clustering is grouping data points into groups where data points in one group are similar to each other.

What is clustering?

Machine Learning: Clustering



By color



By shape



By size



etc...

Methods of Clustering

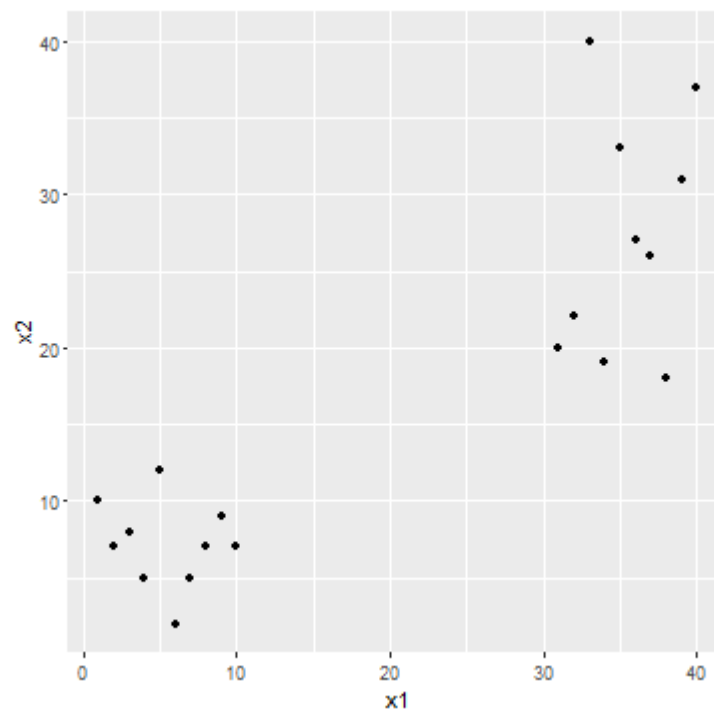
We will cover two clustering methods:

- K-means clustering and
- Hierarchical clustering

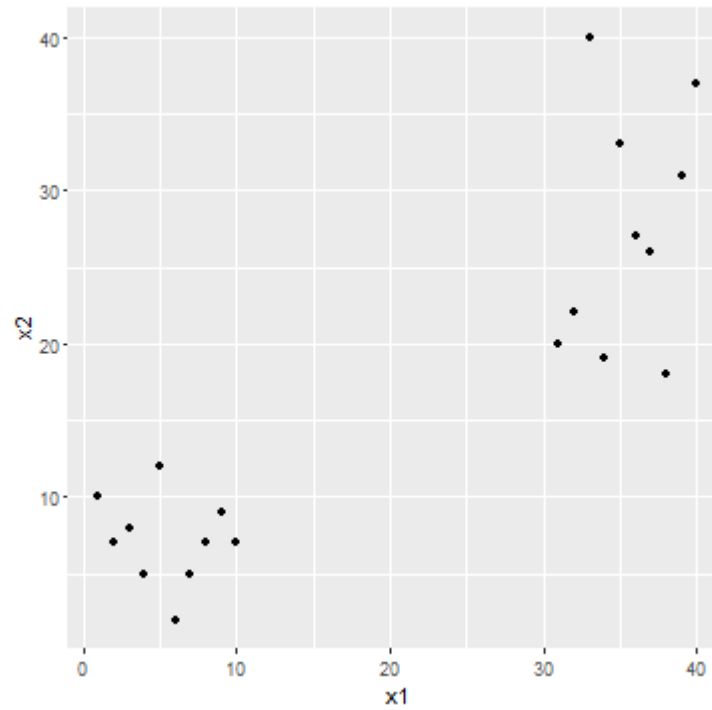
K-means Clustering

- Data
- Visualize Data
- Result of K-means clustering

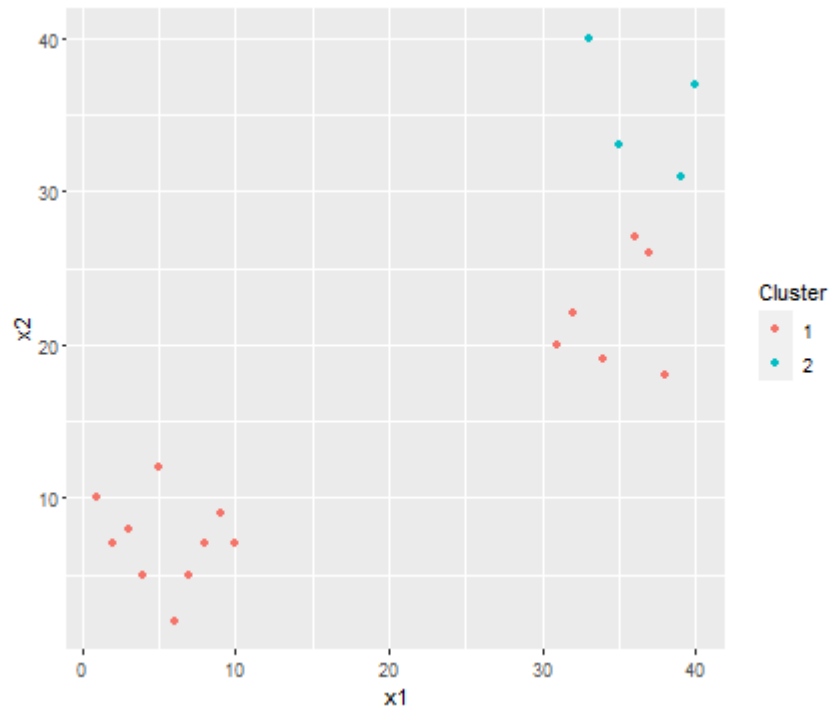
Step 1



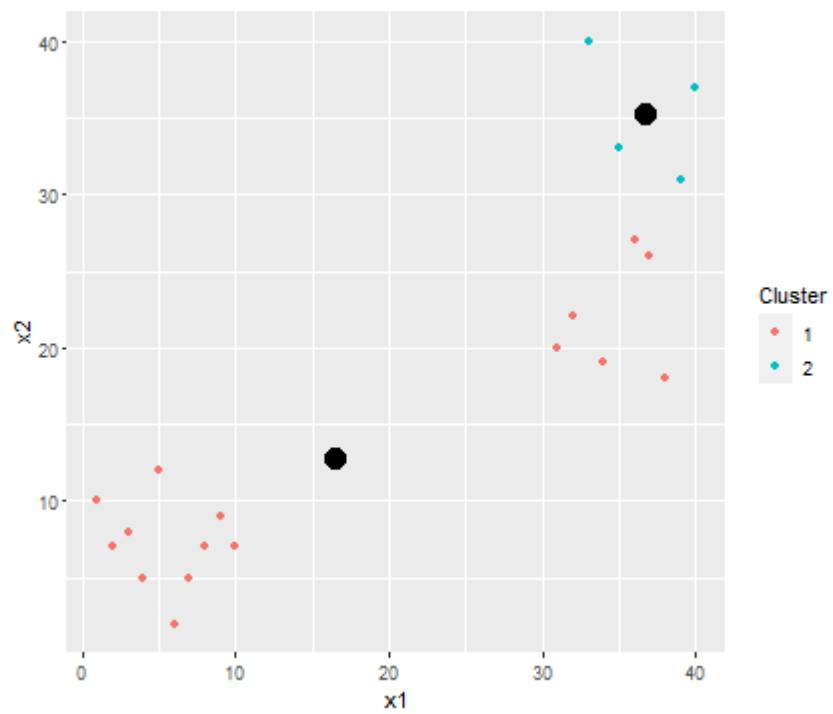
Step 1: Randomly select centroids



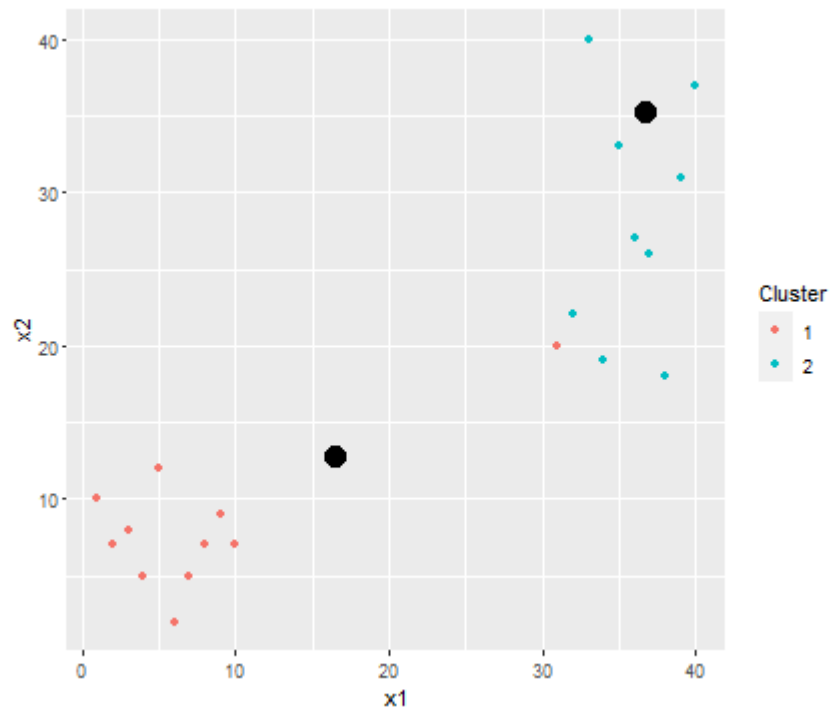
Step 1: Collect points for each clusters



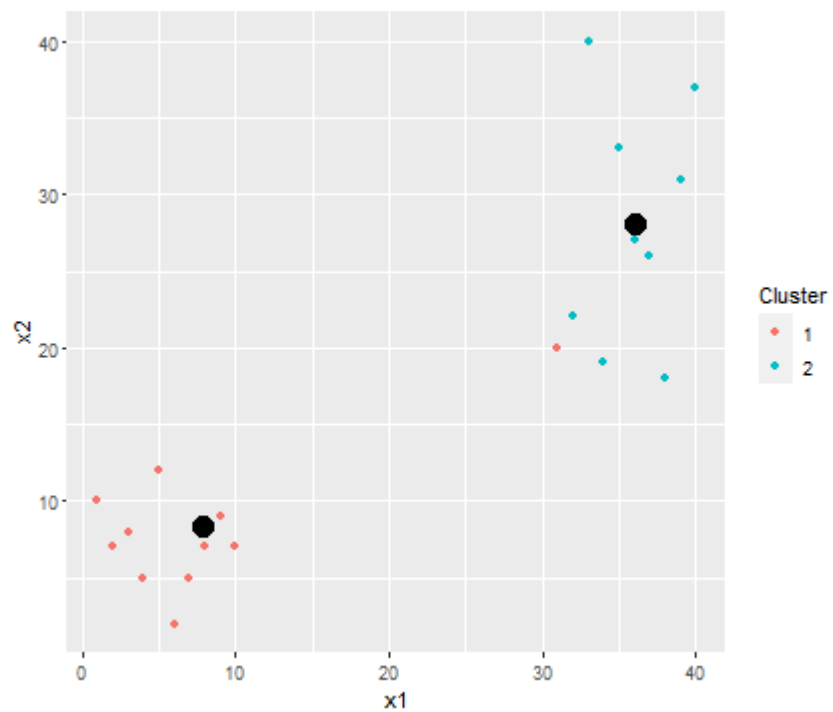
Locate centroids



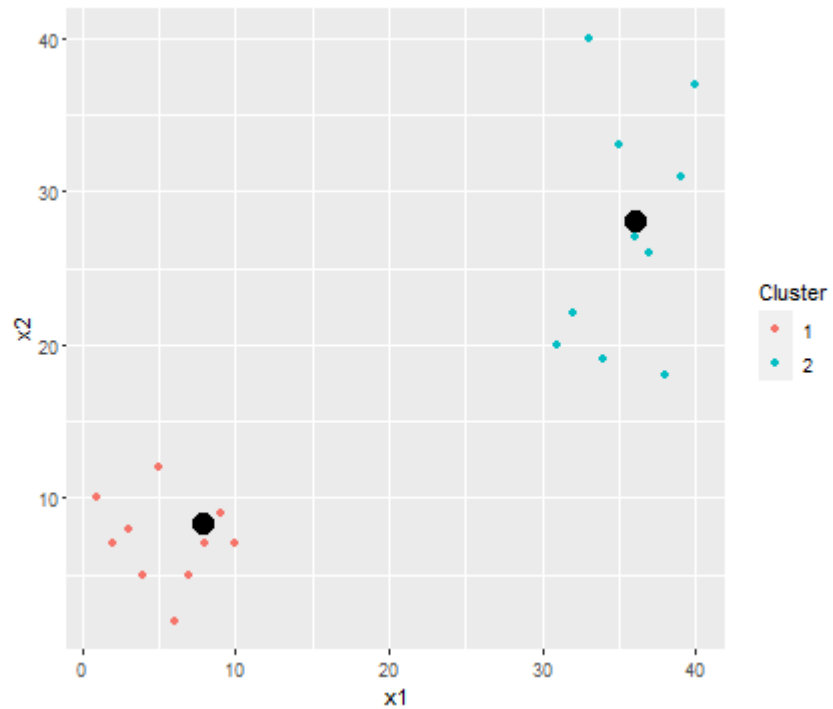
Collect points for each clusters



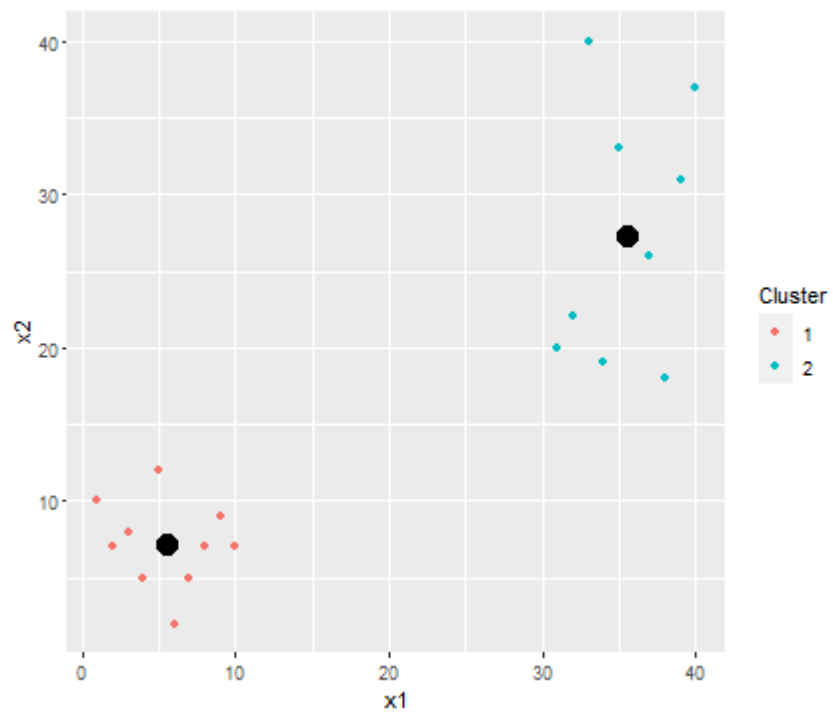
Relocate centroids



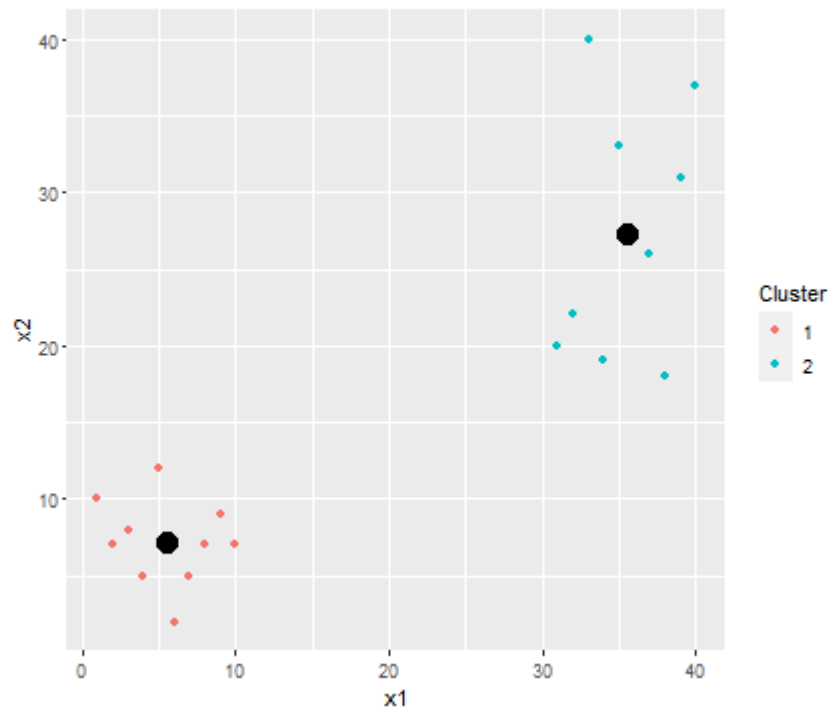
Collect points for each clusters



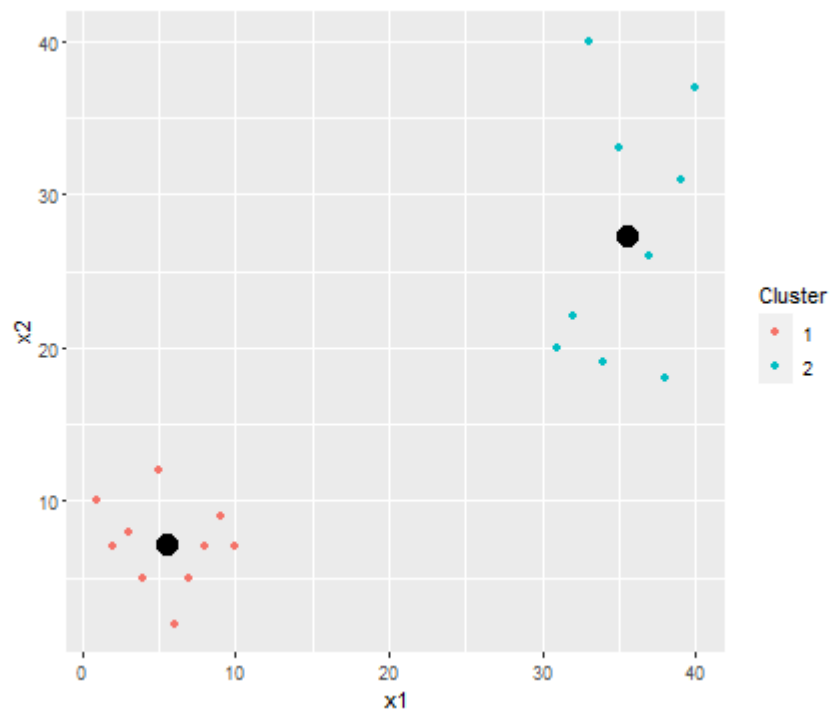
Relocate centroids



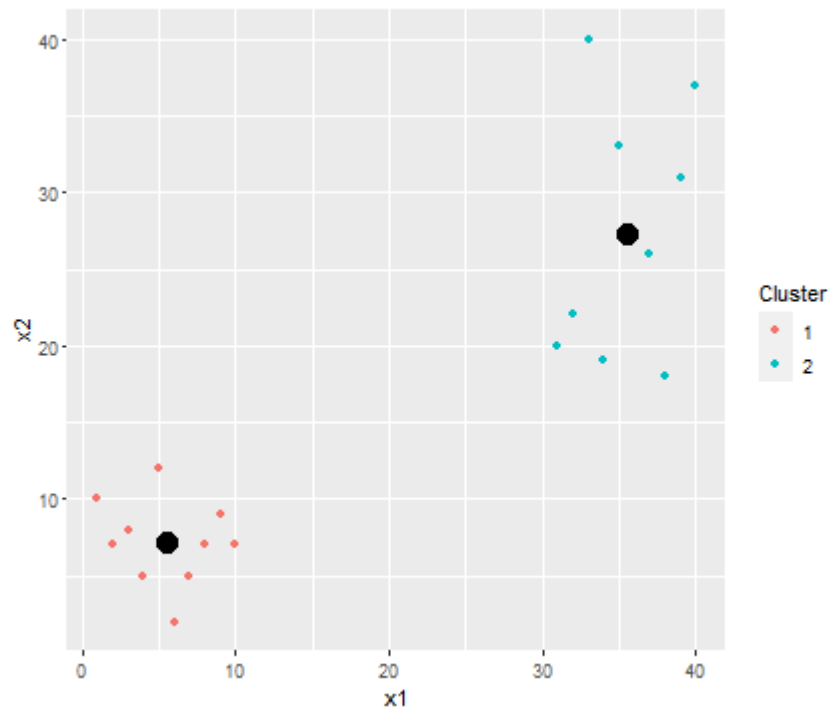
Collect points for each clusters



Relocate centroids

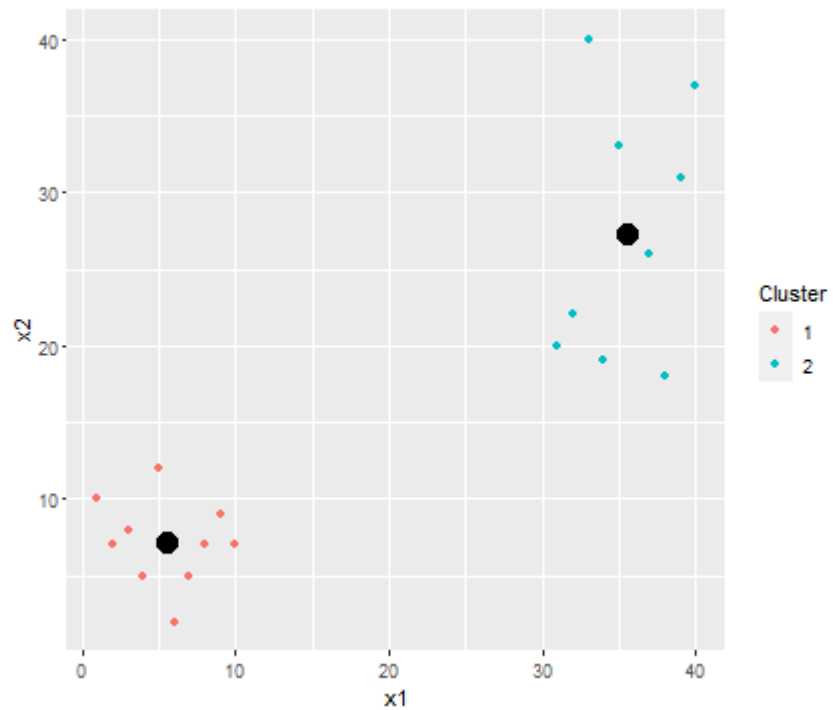


Step 2: Collect points for each clusters

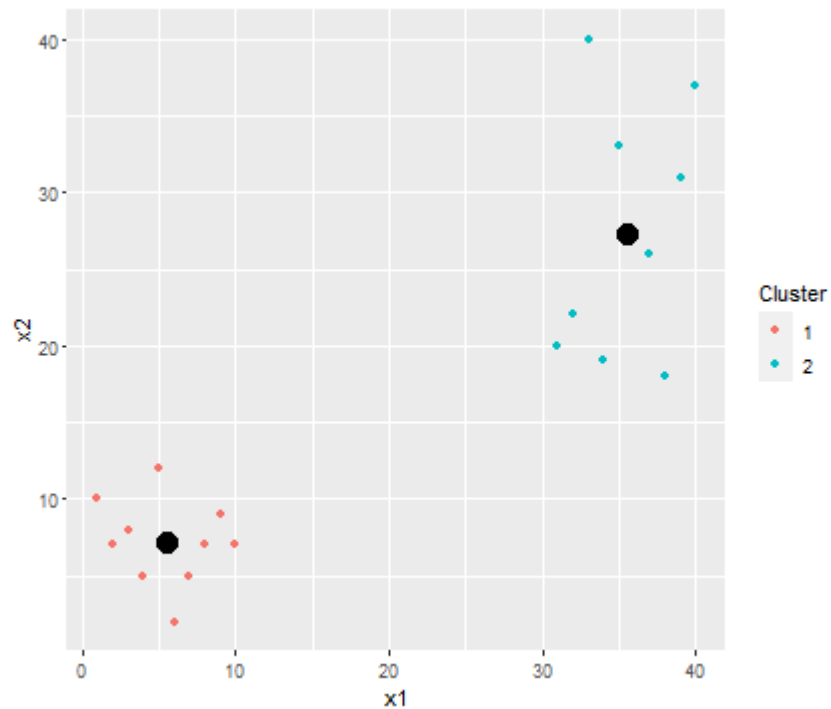


GPA	Screen Time
3.9	90
4.0	150
3.0	140
2.0	139

Step 2: Relocate centroids



Step 2: Collect points for each clusters



Centroids

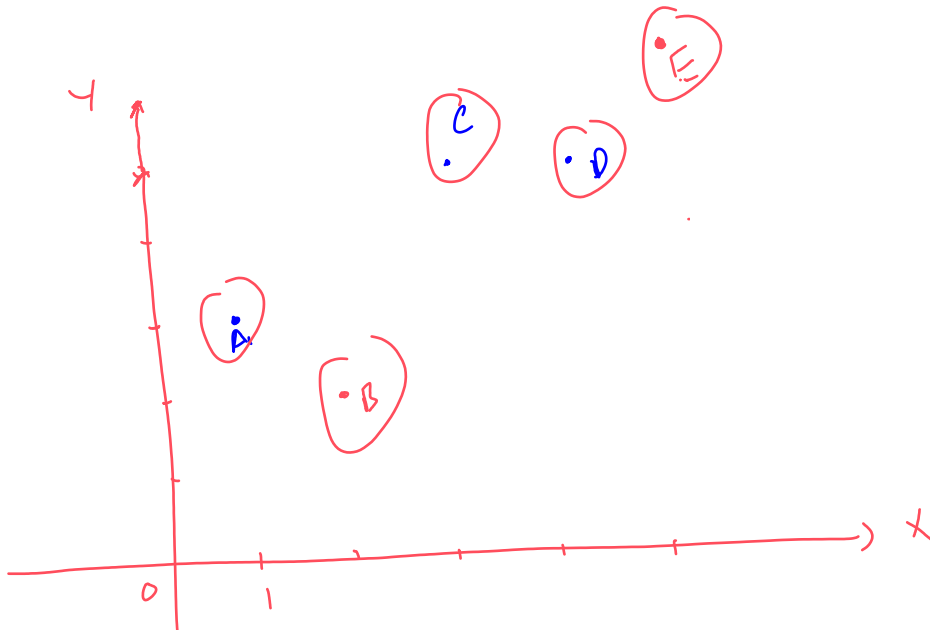
Cluster	x1	x2
1	5.5	7.2
2	35.5	27.3

K-means Algorithm

- 1. Randomly assign a number, from 1 to K, to each of the observations. These serve as initial cluster assignments for the observations.
- 1. Iterate until the cluster assignments stop changing:
 - (a) For each of the K clusters, compute the cluster centroid. The kth cluster centroid is the vector of the p feature means for the observations in the kth cluster.
 - (b) Assign each observation to the cluster whose centroid is closest (where closest is defined using Euclidean distance).

Dataset

Point	x	y
A	1	3
B	2	2
C	3	5
D	4	5
E	5	6



Randomly Assign Cluster to Points

Cluster	Point	x	y
1	A	1	3
2	B	2	2
1	C	3	5
1	D	4	5
2	E	5	6

Determine Centroids

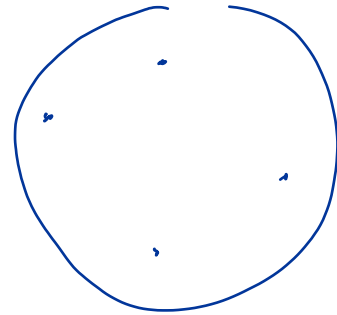
Cluster	Point	x	y	M_1x	M_1y	M_2x	M_2y
1	A	1	3	2.67	4.33	3.5	4
2	B	2	2	2.67	4.33	3.5	4
1	C	3	5	2.67	4.33	3.5	4
1	D	4	5	2.67	4.33	3.5	4
2	E	5	6	2.67	4.33	3.5	4

- Centroid 1: $M_1 = \frac{A+C+D}{3} = \frac{(1,3)+(3,5)+(4,5)}{3} = \frac{(8,13)}{3} = (2.67, 4.33)$
- Centroid 2: $M_2 = \frac{B+E}{2} = \frac{(2,2)+(5,6)}{2} = \frac{(7,8)}{2} = (3.5, 4)$

Distance to Centroids

Cluster	Point	x	y	M_1x	M_1y	M_2x	M_2y	dc1	dc2
1	A	1	3	2.67	4.33	3.5	4	2.13	2.69
2	B	2	2	2.67	4.33	3.5	4	2.42	2.50
1	C	3	5	2.67	4.33	3.5	4	0.75	1.12
1	D	4	5	2.67	4.33	3.5	4	1.49	1.12
2	E	5	6	2.67	4.33	3.5	4	2.87	2.50

- $AM_1 = \sqrt{(1 - 2.67)^2 + (3 - 4.33)^2} = 2.13$ and so on
- Variance within Cluster 1 = $V_1 = AM_1^2 + CM_1^2 + DM_1^2 = 7.3195$
- Variance within Cluster 2 = $V_2 = BM_2^2 + EM_2^2 = 12.5$
- Total Variance = $V_1 + V_2 = 19.83$



Compare Distances to Centroids

Cluster	Point	x	y	dc1	dc2	min_distance
1	A	1	3	2.13	2.69	2.13
2	B	2	2	2.42	2.50	2.42
1	C	3	5	0.75	1.12	0.75
1	D	4	5	1.49	1.12	1.12
2	E	5	6	2.87	2.50	2.50

Reassign Clusters

Cluster	Point	x	y	dc1	dc2	min_distance	New_Cluster
1	A	1	3	2.13	2.69	2.13	1
2	B	2	2	2.42	2.50	2.42	1
1	C	3	5	0.75	1.12	0.75	1
1	D	4	5	1.49	1.12	1.12	2
2	E	5	6	2.87	2.50	2.50	2

- New cluster 1 = {A, B, C}
- New cluster 2 = {D, E}
- Total Variance = $\underbrace{AN_1^2 + BN_1^2 + CN_1^2}_{\text{Cluster 1}} + \underbrace{DN_2^2 + EN_2^2}_{\text{Cluster 2}} = 7.67$

Reassign Clusters

Cluster	Point	x	y	dc1	dc2	min_distance	New_Cluster
1	A	1	3	2.13	2.69	2.13	1
2	B	2	2	2.42	2.50	2.42	1
1	C	3	5	0.75	1.12	0.75	1
1	D	4	5	1.49	1.12	1.12	2
2	E	5	6	2.87	2.50	2.50	2

- New cluster 1 = {A, B, C}
- New cluster 2 = {D, E}
- Total Variance = $AN_1^2 + BN_1^2 + CN_1^2 + DN_2^2 + EN_2^2 = 7.67$
- The process continues until there is no change in the total variance
- The total variance will be reduced to its minimum.

Step 1: Total Variance within

Cluster	Point	x	y
1	A	1	3
2	B	2	2
1	C	3	5
1	D	4	5
2	E	5	6

- Total Variance within 19.83

Step 2: Total Variance within

New_Cluster	Point	x	y
1	A	1	3
1	B	2	2
1	C	3	5
2	D	4	5
2	E	5	6

- Total Variance within 7.67

Step 2: Total Variance within

New_Cluster	Point	x	y
1	A	1	3
1	B	2	2
1	C	3	5
2	D	4	5
2	E	5	6

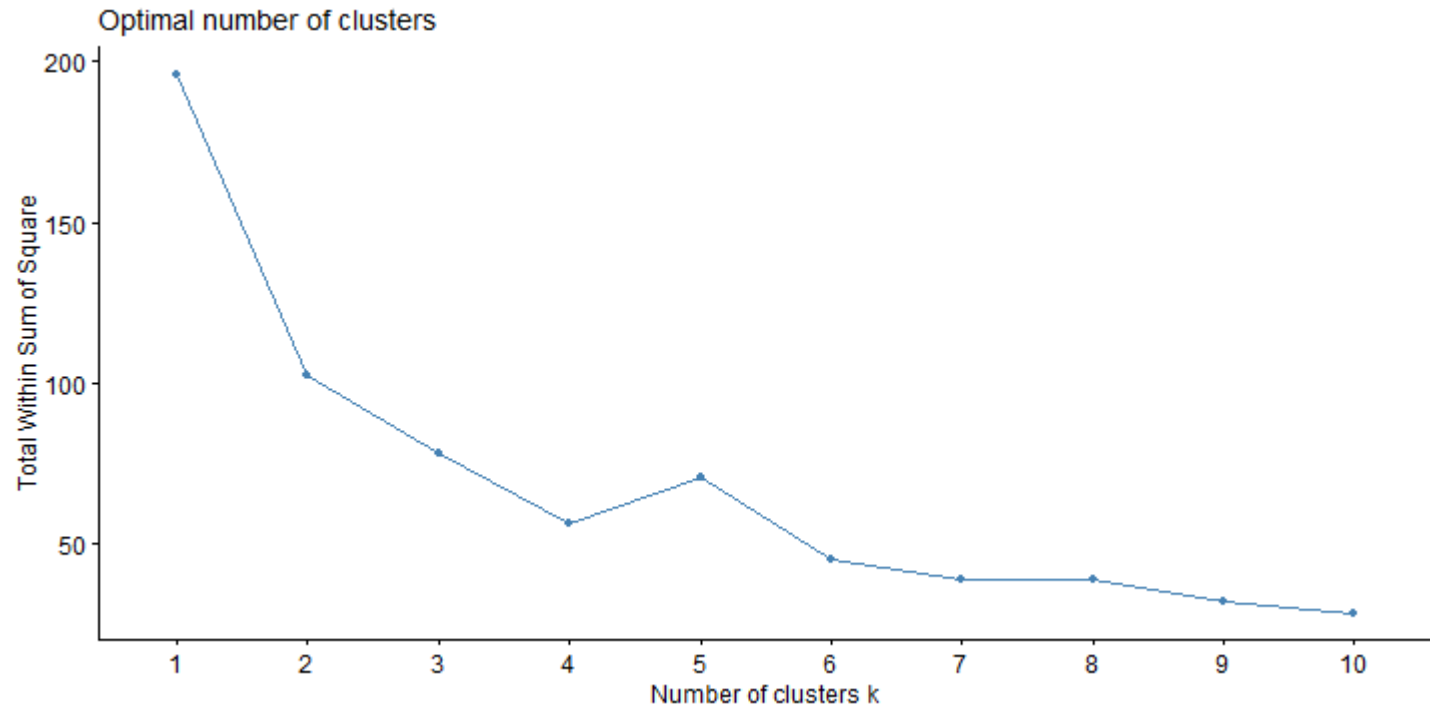
- Total Variance within 7.67
- The process continues until there is no change in the total variance
- The total variance will be reduced to its minimum.

Terminology

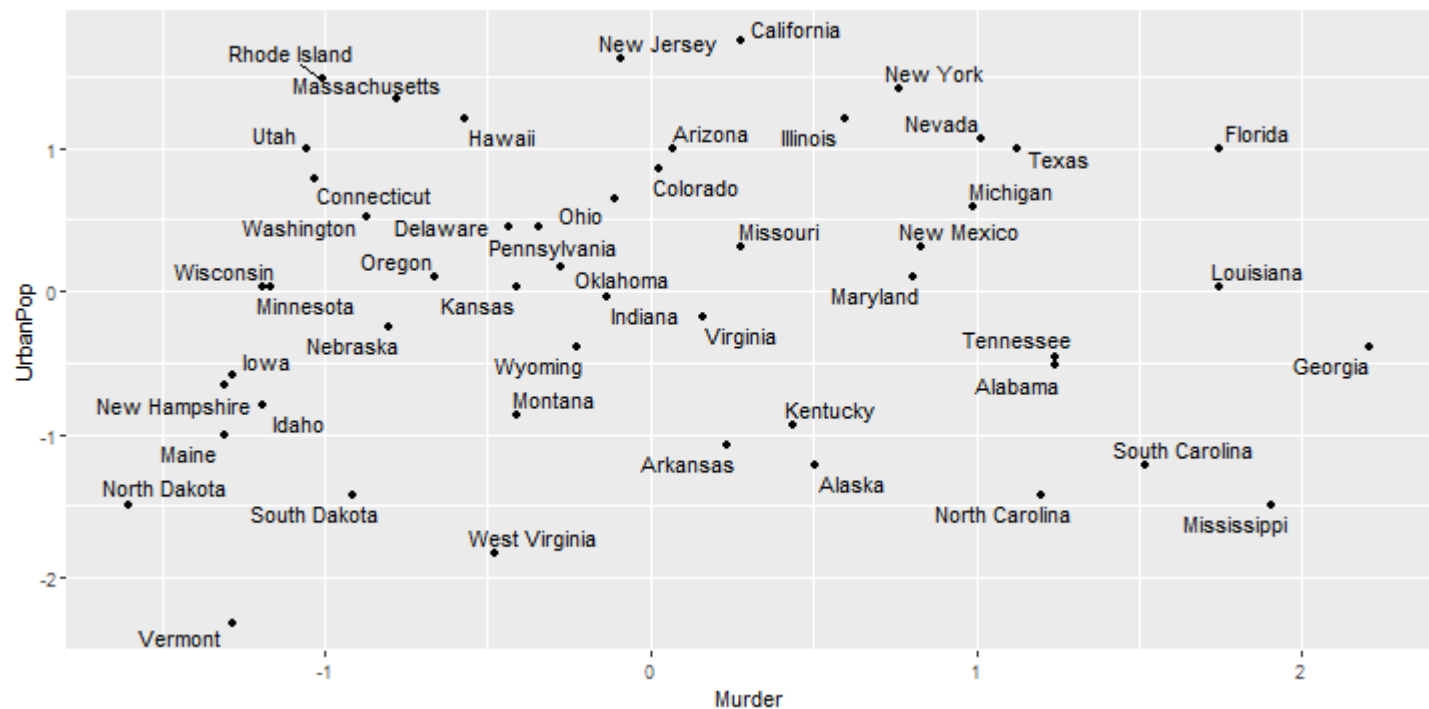
- Total Variance also called the total within sum square or the within-cluster sum of squares (WCSS) or WSS

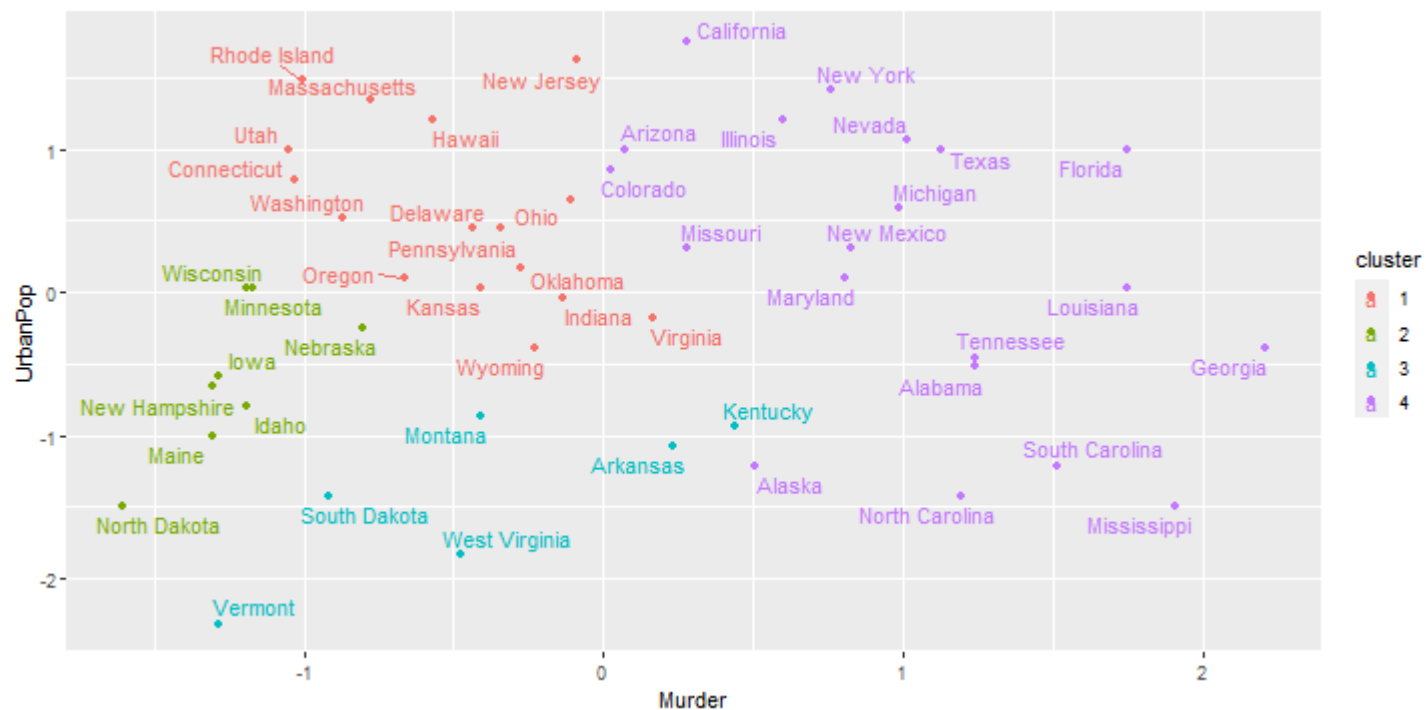
Number of clusters? Elbow method!

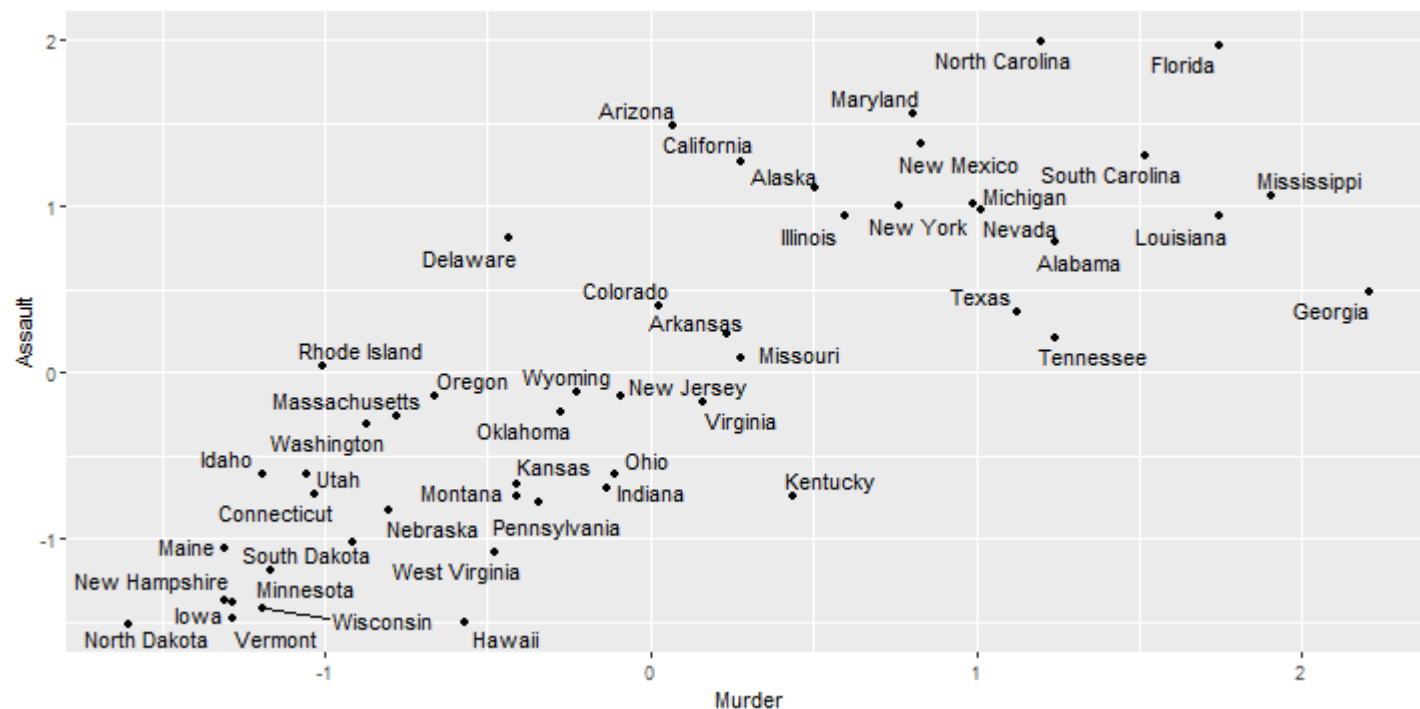
- Plot the WSS
- Decide the **elbow** of the graph

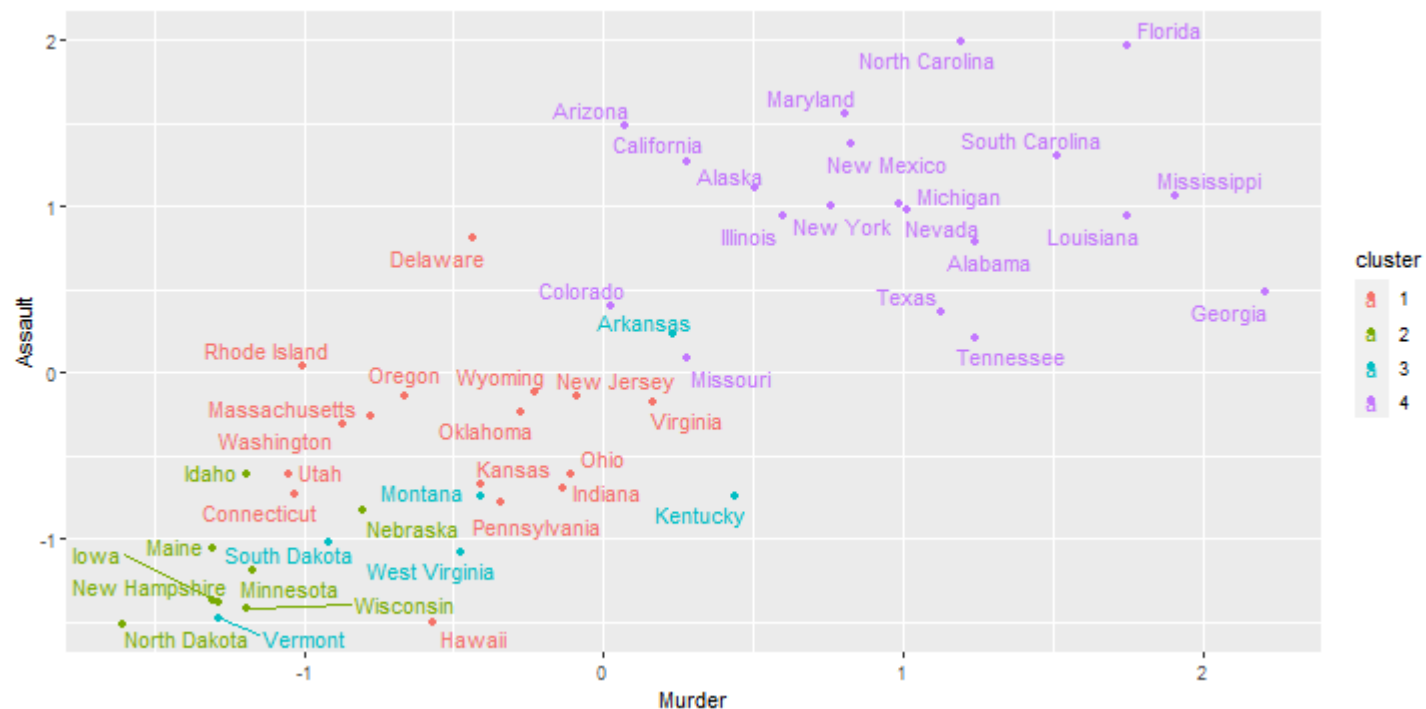


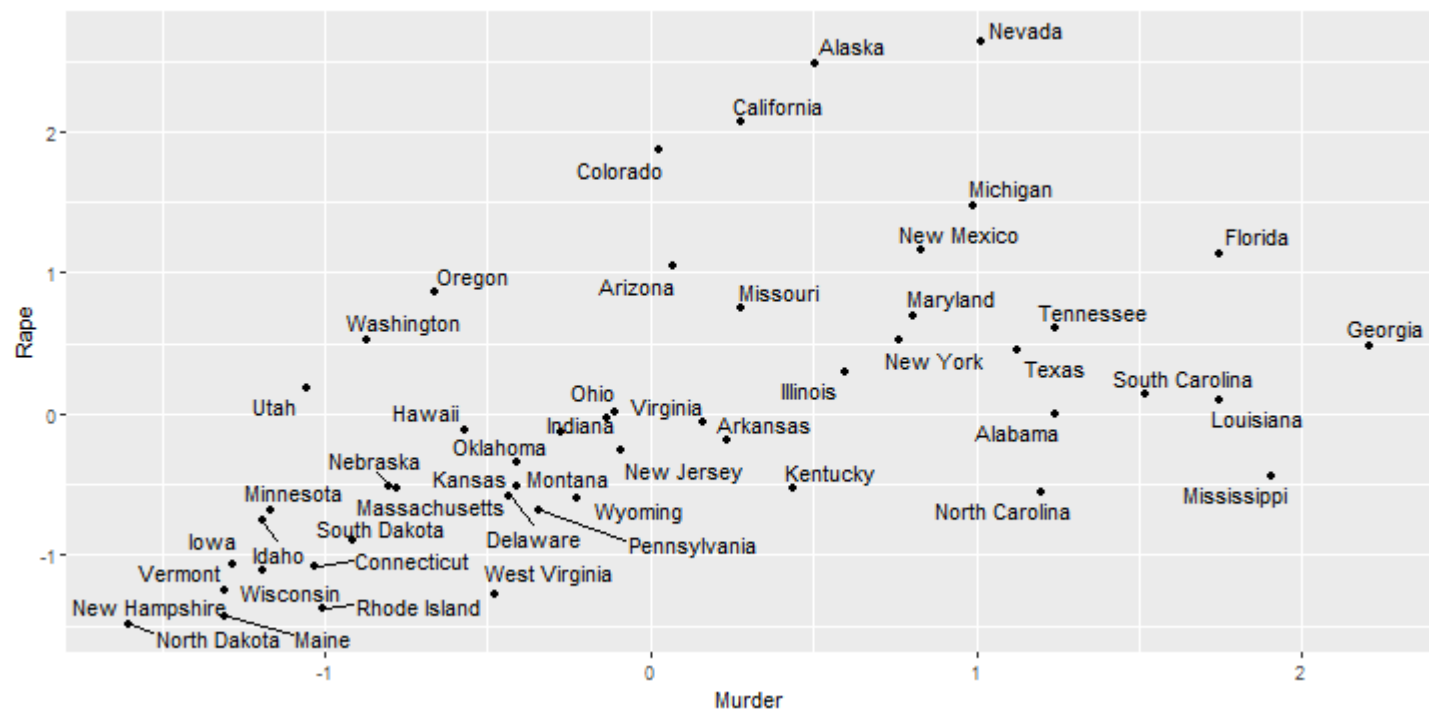
##		Murder	Assault	UrbanPop	Rape	cluster
##	Alabama	1.24256408	0.7828393	-0.5209066	-0.003416473	4
##	Alaska	0.50786248	1.1068225	-1.2117642	2.484202941	4
##	Arizona	0.07163341	1.4788032	0.9989801	1.042878388	4
##	Arkansas	0.23234938	0.2308680	-1.0735927	-0.184916602	3
##	California	0.27826823	1.2628144	1.7589234	2.067820292	4
##	Colorado	0.02571456	0.3988593	0.8608085	1.864967207	4

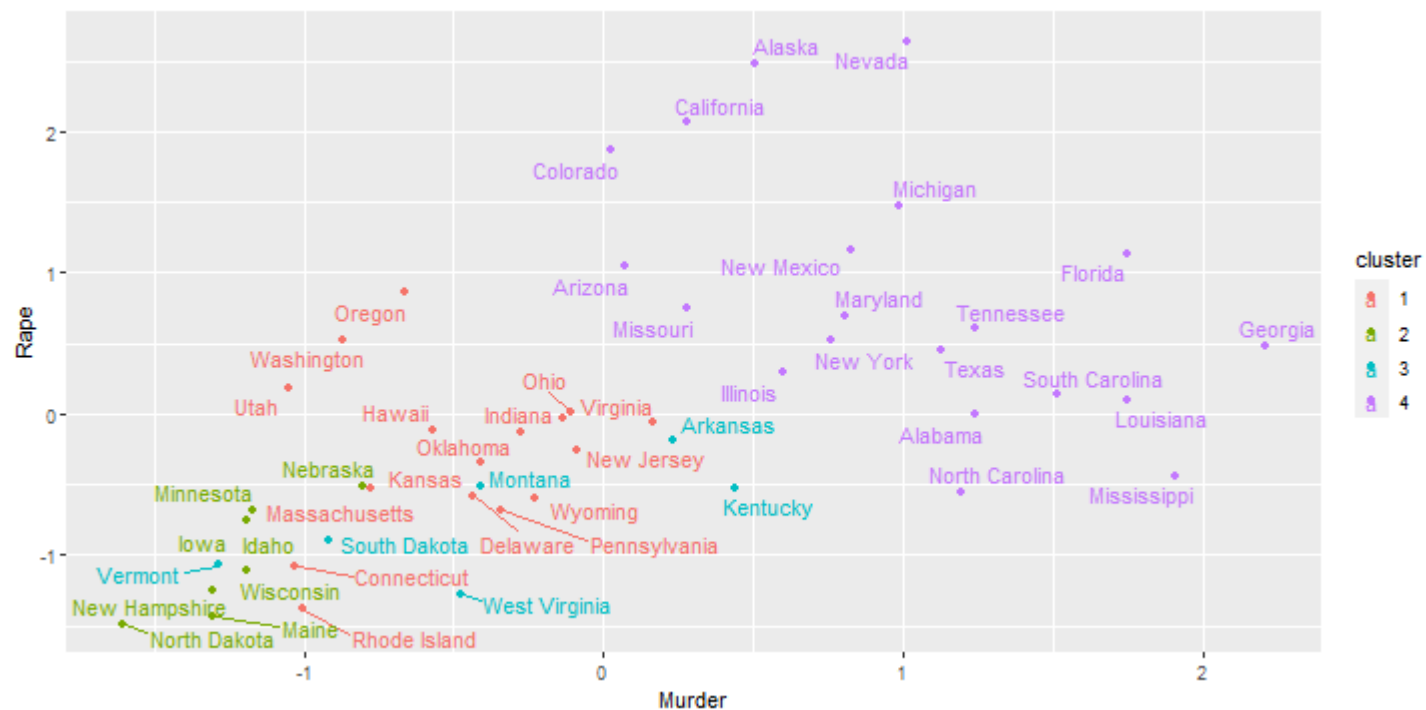


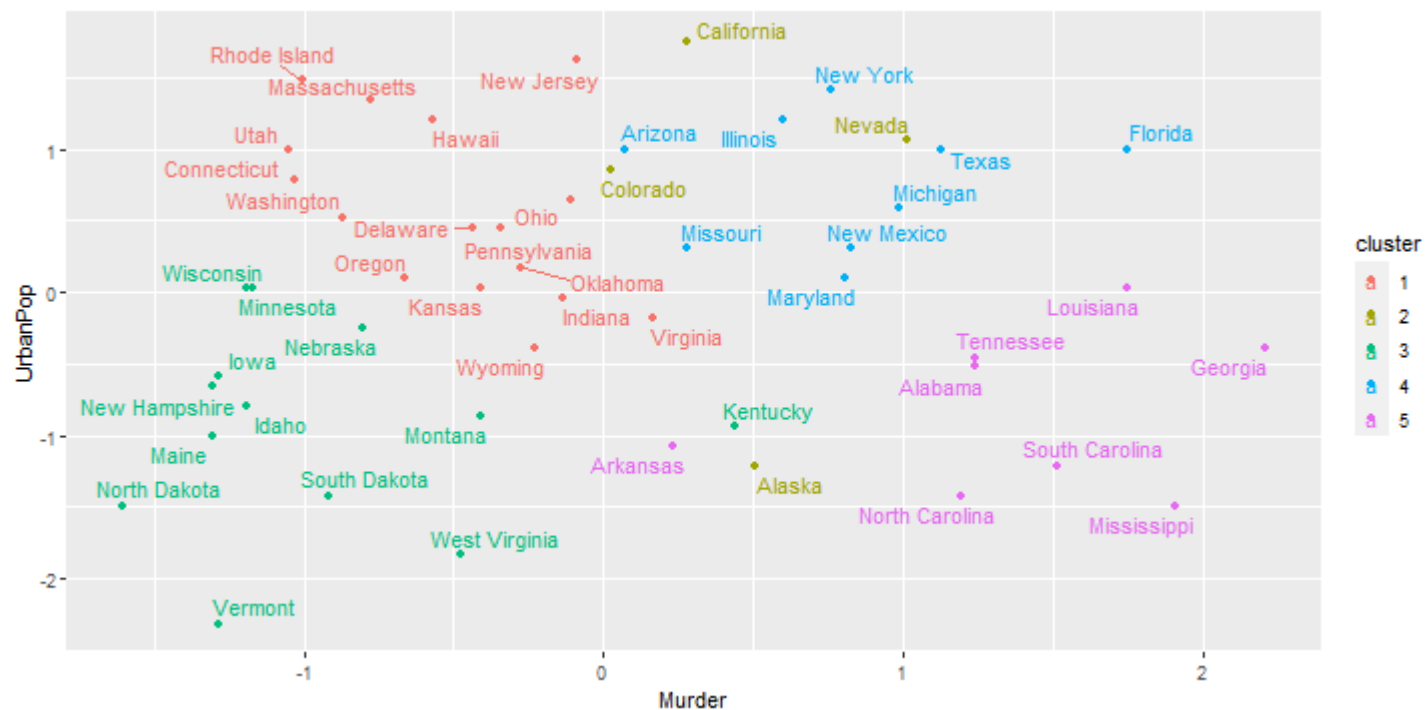


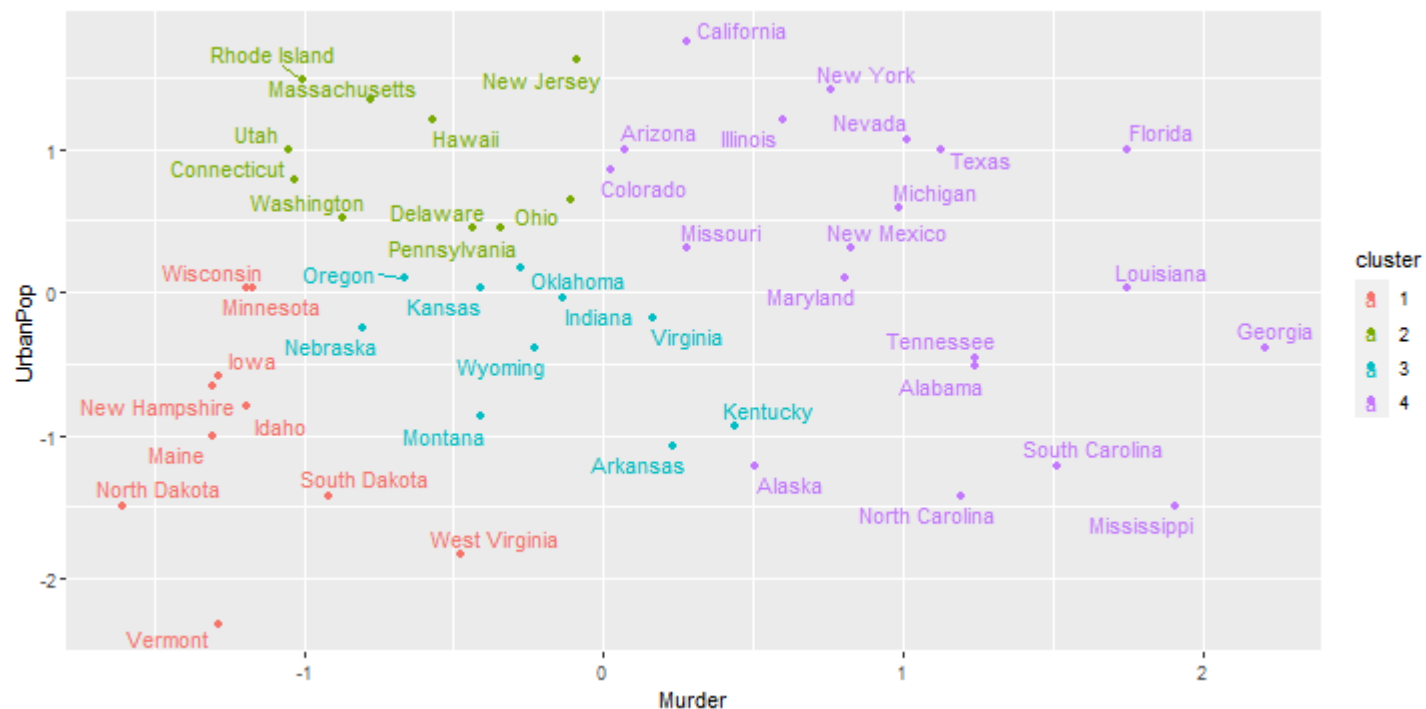


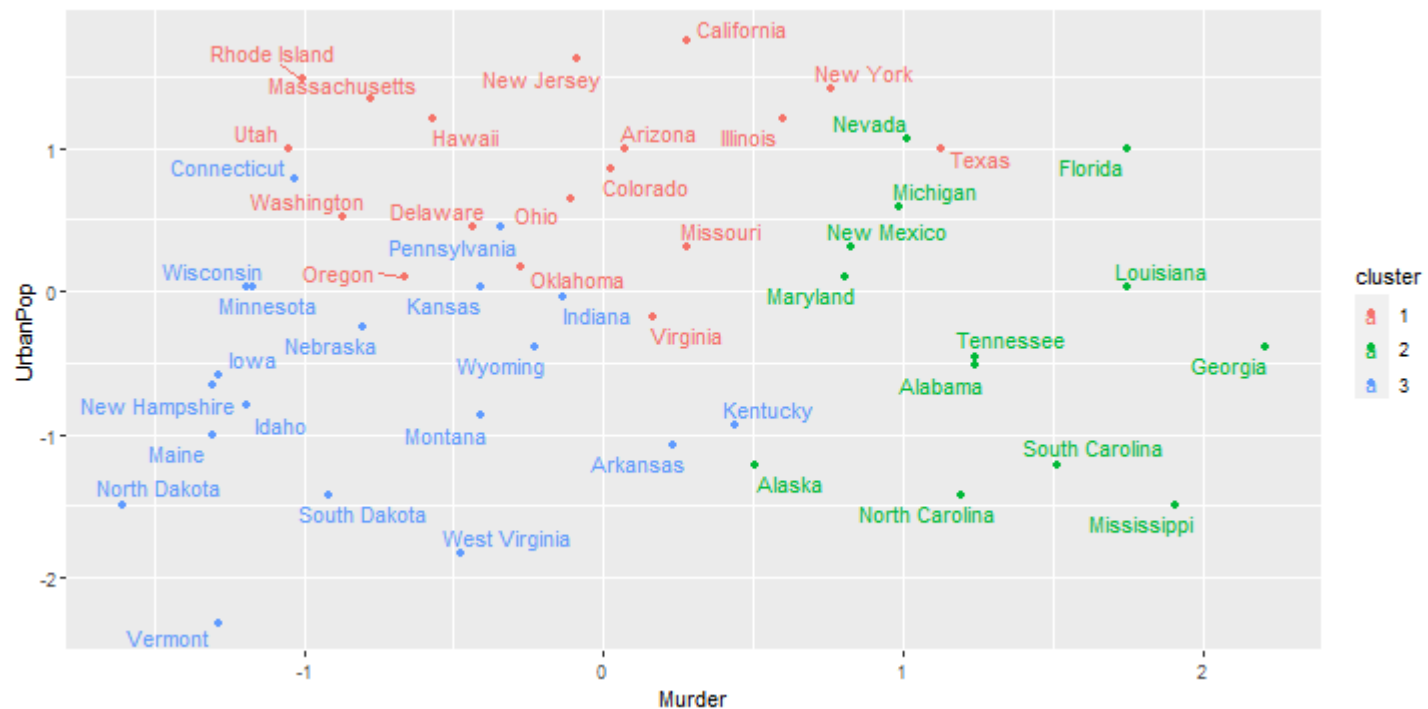












Principal Components

- Since we have four variables, we could have 6 different plots to visualize the clustering
- The more variables we have, the more plot we can have
- It is easier to contain all the variables in a few variables, then make plots.
- One way to do this is to use principal components analysis. The first two principals may contain most of the information of in the dataset.

