

## Table of Contents

K-means Clustering.....	pg. 2
Dataset .....	pg. 2
Step-1: Package Load .....	pg. 3
Step-2 : Import & Prepare Dataset .....	pg. 3-4
Step-3 : Find Optimal No. of Clusters .....	pg. 4
Step-3.1 : WSS .....	pg. 5
Step-3.2 : Gap Statics .....	pg. 6
Step-4 : Perform K-means .....	pg. 7
Visualizing the Clusters .....	pg. 8

## Applying K-means Clustering

### K-means Clustering

Using the K-means clustering technique, we assign each observation in a dataset to one of K clusters.

The ultimate goal is to have K clusters in which the observations are quite diverse from one another across clusters, but fairly similar to one another inside each cluster.

### Dataset

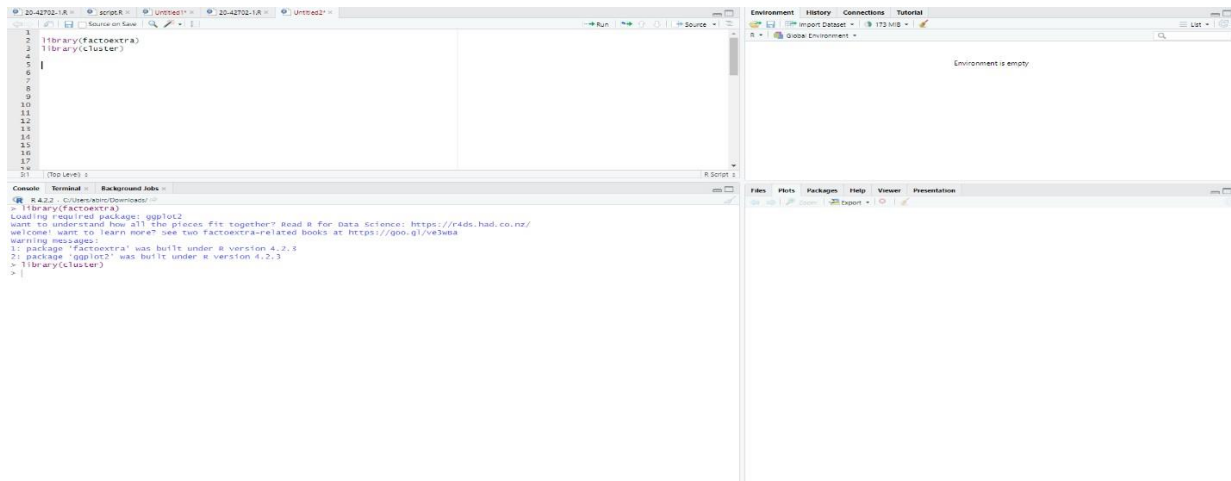
For this project I have chosen the kidney-stone dataset. Age, gender, weight, height, family history, and numerous blood and urine parameters are just a few of the 12 characteristics in this kidney stone dataset that are linked to the likelihood of having kidney stones. The target variable, which is a continuous value expressing the chance of acquiring kidney stones, is part of a dataset of 180 entries. The dataset is helpful for developing predictive models that can help identify people who have a high chance of developing kidney stones, which can help with the creation of preventative and treatment plans. The dataset includes details on 90 patients who had kidney stones, including age, gender, and the type of stone each patient had as well as the outcomes of laboratory tests. The link for the dataset:

<https://www.kaggle.com/datasets/harshghadiya/kidneystone?resource=download>.

The following shows all the necessary steps to perform K-means clustering on the chosen dataset in R.

## Step 1: Loading the Necessary Packages

I have first loaded two packages that have a number of helpful R functions for k-means clustering.

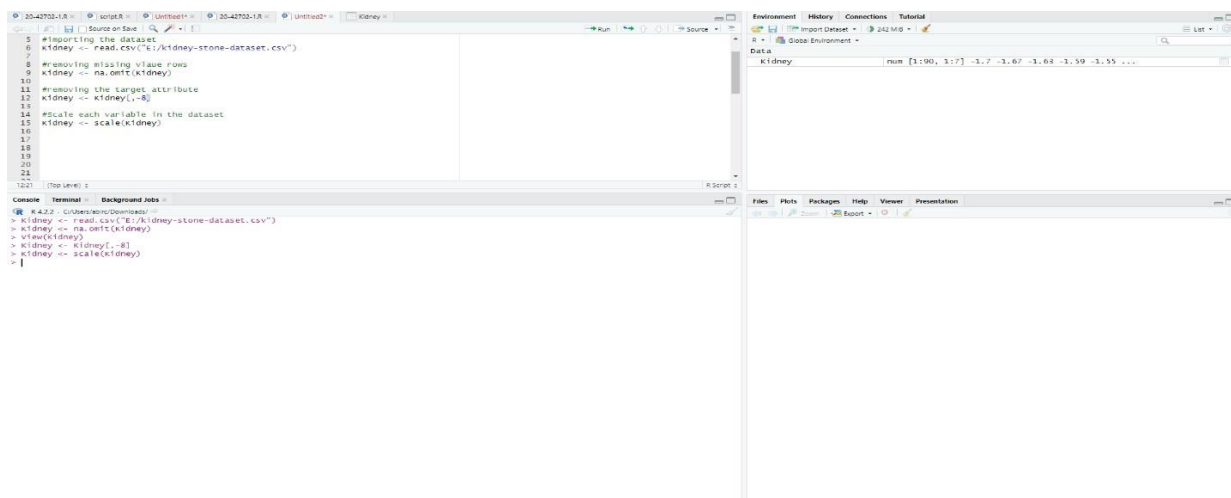


## Step 2: Importing and Preparing the Data

Here the dataset used is kidney-stone dataset. This dataset contains 12 features related to the risk of developing kidney stones, including age, gender, weight, height, family history, and various blood and urine measurements.

The code below demonstrates the following:

- Import the kidney-stone dataset
- Remove any rows with missing values
- Remove the target attribute
- Scale each variable in the dataset to have a mean of 0 and a standard deviation of 1



After step 2, the dataset:

	gravity	gR	mass	cond	urea	cAK
1	0.44958984	-1.562877024	0.51441234	-0.268377021	1.36530318	-0.51977682
2	-0.14038049	-0.41678918	-0.10623753	-0.06121907	0.27921174	0.15693486
3	-1.467738789	1.634274008	-1.17879352	-0.747496328	-1.161166229	-0.54961466
4	-1.025265916	-0.739884913	-0.81499216	-1.047977210	-0.25262015	-0.61623701
5	-1.910191663	0.476000882	-1.76173359	-1.741656470	-1.23020183	-0.94765685
6	0.302072705	-1.077027777	0.27837835	0.611188791	-0.04870663	-0.22471038
7	-0.877801623	-0.580347350	-0.59289246	-0.422890847	-0.46683021	-0.36788825
8	1.629471325	-0.519100033	0.11355647	0.99040409	0.15529643	-1.47897947
9	-0.435548751	-0.880273406	-0.24861901	0.167032618	-0.65148406	-0.94745665
10	0.449556996	0.131144987	0.74088560	0.663445958	0.81455336	-0.58934491
11	-1.025265916	0.215458203	-1.07914783	-1.191844501	-0.78445203	-0.68217471
12	1.039484160	-0.711787174	1.277164370	1.014162087	1.40136739	-0.91108780
13	-1.762707371	1.821983033	-1.51108425	-1.217813100	-1.43446878	-0.99056220
14	-1.615223080	-0.864668322	-1.33914952	-1.450713269	-0.82133480	-0.84480581
15	-1.025265916	-1.161030889	-0.63802168	-0.355656531	-0.71797329	-0.22478271
16	0.007134122	-1.588868849	0.36437581	0.718783185	0.18097039	0.93547588
17	-1.615223080	0.833084855	-1.48499589	-1.596677777	-0.92478656	-0.98382850
18	1.039484160	1.086481028	-1.44539685	1.564862433	1.01548039	-0.65921271
19	-1.467738789	-1.64810168	-0.88466689	0.715703185	-1.20548811	1.21415148
20	-0.582833042	0.145193857	-0.18656028	0.388095705	-0.32646869	-0.85131150
21	0.860208649	0.369875766	1.13823558	1.212145556	0.89966228	0.57868531
22	0.154858471	-0.798020290	0.68710760	1.721654228	-0.49728400	-1.06499290
23	-0.582833042	-0.788236586	-0.10623753	1.238275154	-1.26437804	-0.89441108
24	0.302072705	-0.107682762	0.12021593	-1.230873999	1.20891754	-0.61815801
25	0.744523578	-0.501841784	0.61103523	-0.98487680	0.14182184	-0.82815406
26	-0.14038049	-1.012223756	-0.67783595	1.544420560	0.08716134	-1.07870530
27	-0.14038049	0.210377851	-0.51891627	0.678510288	-1.35321649	-0.61280631
28	-1.172770207	0.805480716	-1.58237709	-1.413777568	-1.37537835	-1.27867628
29	-1.467738789	-0.33412816	-0.5117863	-0.03778963	-0.78176416	-1.08868320
30	0.302072705	-0.688226986	0.74823276	0.084567663	0.67069910	-0.33417087
31	-0.14038049	2.848880201	0.32870128	0.611188791	0.17879986	-0.88061015
32	0.154858471	-0.078588053	-0.39180517	-0.688112832	0.38898929	-0.02675475
33	-0.14038049	0.735044889	-0.18172175	-0.628919833	0.40433937	0.45162100
34	-1.467738789	-0.135783350	-1.45227615	-1.638870675	-0.94495622	-0.16171873

### Step 3: Finding the Optimal Number of Clusters

The built-in `kmeans()` function in R can be used to do k-means clustering and has the following syntax:

**kmeans(data, centers, iter.max, nstart)**

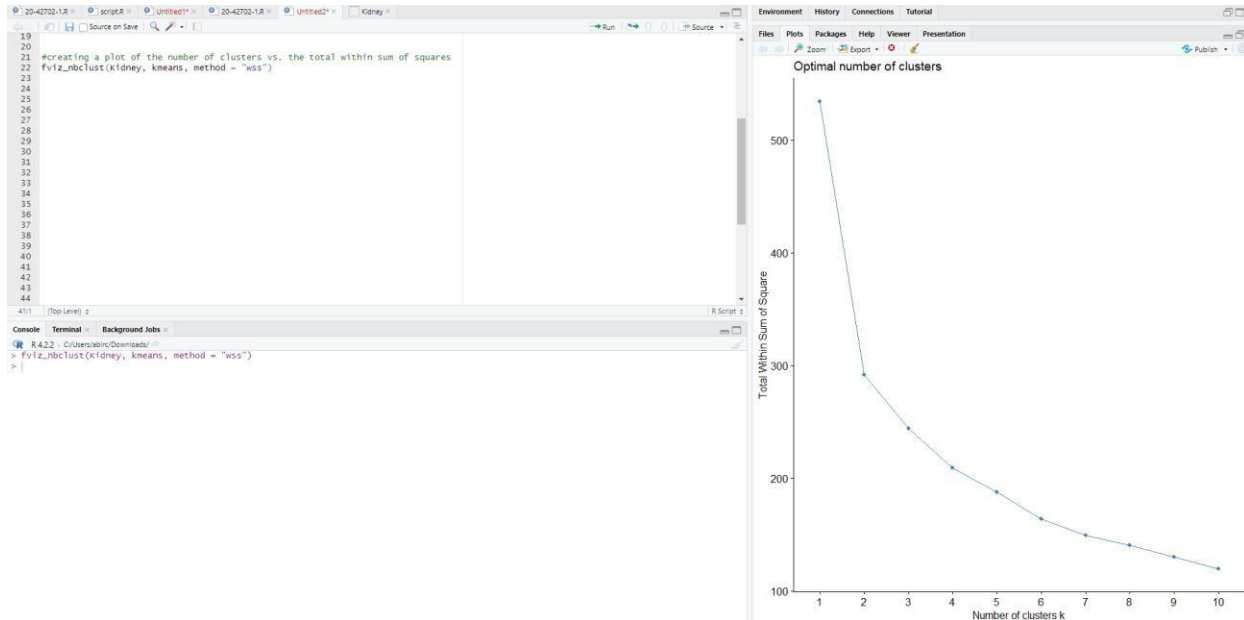
Here:

- **data:** Name of the dataset.
- **centers:** The number of centers denoted by k.
- **iter.max:** The maximum number of iterations allowed. The default value is 10.
- **nstart:** The number of random starting partitions.

Since how many clusters will be optimal is unknown, two different plots were chosen here for the help of deciding:

### Step 3.1: Total Within Sum of Squares

At first the **fviz\_nbclust()** function was used to plot the number of clusters vs the total within sum of squares.



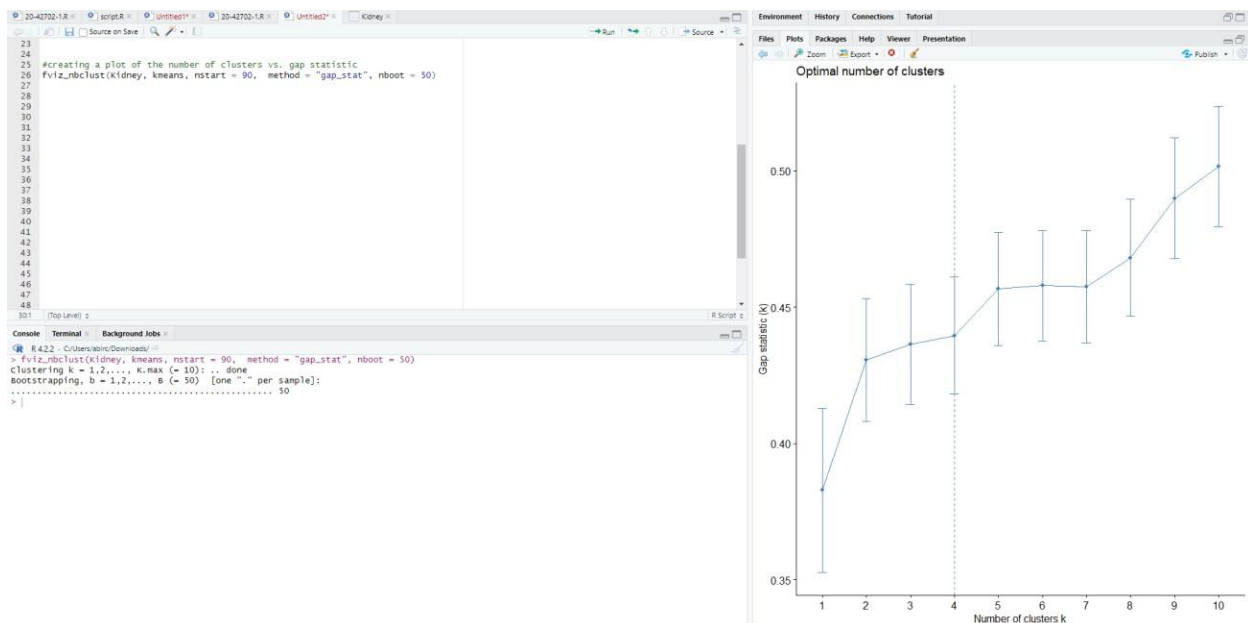
Usually, when making this kind of plot, we search for an "elbow" shape where the total number of squares starts to bend or level off. The ideal number of clusters is usually this.

There appears to be a slight elbow or bend in this figure at the k = 4 clusters.

## Step 3.2: Gap Statistics

Utilizing the gap statistic, which contrasts the total intra-cluster variation for various values of  $k$  with their anticipated values for a distribution with no clustering, is another method for figuring out the ideal number of clusters.

It can be calculated using the **fviz\_nbclust()** function.



From the plot we can see that gap statistic suggests  $k = 4$  clusters, which matches the Within Sum of Squares method.

## Step 4: Performing K-Means Clustering with Optimal Clusters

Lastly, k-means clustering was performed on the dataset using the optimal value for  $k$  of 4:

```

28
29
30 #performing k-means clustering on the dataset using the optimal value for k of 6
31 km <- kmeans(Kidney, centers = 4, iter.max = 90, nstart = 90)
32
33 #displaying the result
34 km
35
36 |
37
38
39
40
41
42
43
44

```

```

R 4.2.2 - C:/Users/abirc/Downloads/
> km <- kmeans(Kidney, centers = 4, iter.max = 90, nstart = 90)
> km
K-means clustering with 4 clusters of sizes 28, 20, 27, 15

cluster means:
  gravity      ph      osmo      cond      urea      calc
1 -0.8149479 -0.1486525 -1.0414054 -1.0736606 -0.8770478 -0.77188896
2 -0.5060479  1.1617310 -0.3691747  0.0693985 -0.5894470  0.19131089
3  0.4768689 -0.3459962  0.6453452  0.7379609  0.5757682 -0.07809824
4  1.3376027 -0.6486967  1.2745684  0.5833056  1.3867024  1.32635503

clustering vector:
 1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44
 3  3  2  1  1  3  1  4  1  3  1  3  1  1  1  3  1  3  2  2  3  3  2  3  3  1  2  1  1  3  2  3  2  1  3  3  1  3  1  3  3  1  2  1
45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88
 3  3  4  3  3  4  2  1  2  3  1  4  2  3  4  4  4  4  1  1  2  1  3  3  1  4  3  4  2  1  4  2  2  1  4  2  1  2  2  4  2  4  2  4
89 90
 1  1

within cluster sum of squares by cluster:
[1] 59.67405 59.67220 56.94724 32.72046
(between_SS / total_SS = 60.9 %)

Available components:

[1] "cluster"      "centers"      "totss"        "withinss"     "tot.withinss" "betweenss"    "size"         "iter"
[9] "ifault"
> |

```

The result showed the following:

- 28 instances in cluster 1
- 20 instances in cluster 2
- 27 instances in cluster 3
- 15 instances in cluster 4

## Visualizing the Clusters

The **fviz\_cluster()** function can be used to see the clusters on a scatterplot with the first two principal components shown on the axes:

