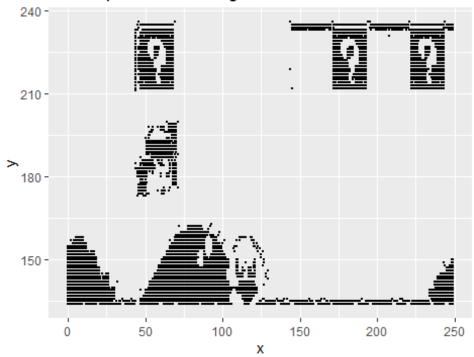
```
# Assignment: ASSIGNMENT 11.2.2
# Name: Anjale, Jiteshwar
# Date: 2021-05-29
# Clustering
## Load the package
library(tidyverse) # data manipulation
## Warning: package 'tidyverse' was built under R version 4.0.5
## -- Attaching packages ----- tidyverse 1.
3.1 --
## v ggplot2 3.3.3 v purrr 0.3.4
## v tibble 3.1.0 v dplyr 1.0.5
## v tidyr 1.1.3 v stringr 1.4.0
## v readr 1.4.0 v forcats 0.5.1
## Warning: package 'tidyr' was built under R version 4.0.5
## Warning: package 'readr' was built under R version 4.0.5
## Warning: package 'purrr' was built under R version 4.0.5
## Warning: package 'dplyr' was built under R version 4.0.5
## Warning: package 'stringr' was built under R version 4.0.5
## Warning: package 'forcats' was built under R version 4.0.5
## -- Conflicts ----- tidyverse conflict
s() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag() masks stats::lag()
library(cluster) # clustering algorithms
## Warning: package 'cluster' was built under R version 4.0.5
library(factoextra) # clustering algorithms & visualization
## Warning: package 'factoextra' was built under R version 4.0.5
## Welcome! Want to learn more? See two factoextra-related books at https://g
oo.gl/ve3WBa
setwd('C:/Users/anjal/OneDrive/Desktop/MS/DSC520/dsc520')
# Load the `data/clustering-data.csv` to
cluster_df <- read.csv("C:/Users/anjal/OneDrive/Desktop/MS/DSC520/dsc520/data</pre>
/clustering-data.csv")
```

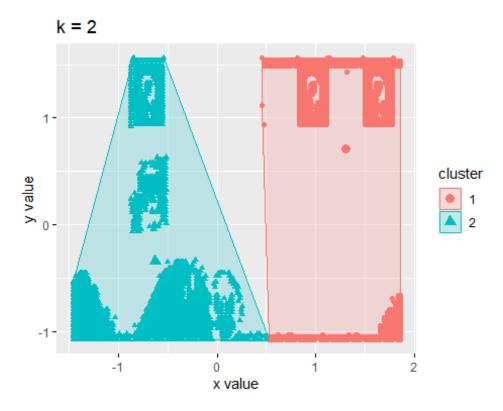
```
# Examine the structure of `clustering-data.csv` using `str()`
str(cluster_df)
                    4022 obs. of 2 variables:
## 'data.frame':
## $ x: int 46 69 144 171 194 195 221 244 45 47 ...
## $ y: int 236 236 236 236 236 236 236 235 235 ...
# Show the top rows of clustering-data.csv
head(cluster_df)
##
       Х
## 1 46 236
## 2 69 236
## 3 144 236
## 4 171 236
## 5 194 236
## 6 195 236
# i.Plot the dataset using a scatter plot.
# scatter plot of data
library(ggplot2)
ggplot(data = cluster_df, aes(x=x, y=y)) +
  geom_point(size = 0.4) +
  ggtitle("Scatterplot of clustering data")
```

## Scatterplot of clustering data

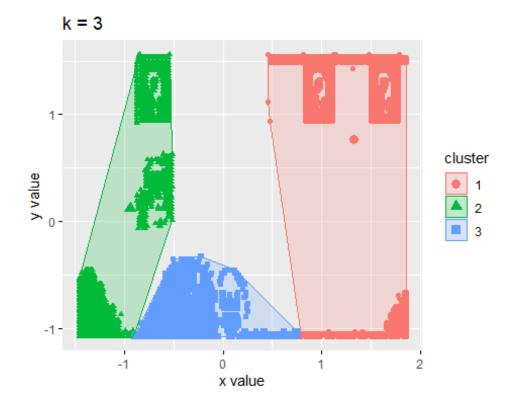


# ii.Fit the dataset using the k-means algorithm from k=2 to k=12. Create a s catter plot of the resultant clusters for each value of k.

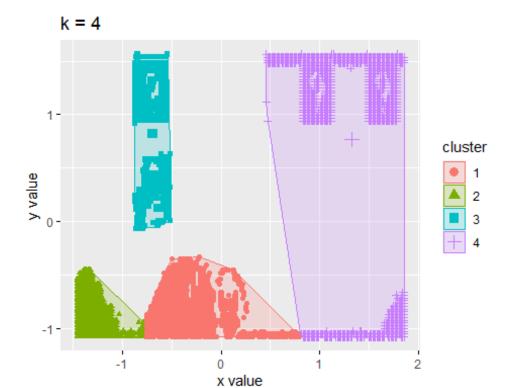
```
#Kmeans for k=2
set.seed(123)
kmeans_2 <- kmeans(cluster_df, 2, iter.max = 300, nstart = 10)</pre>
#Kmeans for k=3
set.seed(123)
kmeans_3 <- kmeans(cluster_df, 3, iter.max = 300, nstart = 10)</pre>
#Kmeans for k=4
set.seed(123)
kmeans 4 <- kmeans(cluster df, 4, iter.max = 300, nstart = 10)</pre>
#Kmeans for k=5
set.seed(123)
kmeans_5 <- kmeans(cluster_df, 5, iter.max = 300, nstart = 10)</pre>
#Kmeans for k=6
set.seed(123)
kmeans_6 <- kmeans(cluster_df, 6, iter.max = 300, nstart = 10)</pre>
#Kmeans for k=7
set.seed(123)
kmeans_7 <- kmeans(cluster_df, 7, iter.max = 300, nstart = 10)</pre>
#Kmeans for k=8
set.seed(123)
kmeans_8 <- kmeans(cluster_df, 8, iter.max = 300, nstart = 10)</pre>
#Kmeans for k=9
set.seed(123)
kmeans_9 <- kmeans(cluster_df, 9, iter.max = 300, nstart = 10)</pre>
#Kmeans for k=10
set.seed(123)
kmeans_10 <- kmeans(cluster_df, 10, iter.max = 300, nstart = 10)</pre>
#Kmeans for k=11
set.seed(123)
kmeans 11 <- kmeans(cluster df, 11, iter.max = 300, nstart = 10)</pre>
#Kmeans for k=12
set.seed(123)
kmeans_12 <- kmeans(cluster_df, 12, iter.max = 300, nstart = 10)</pre>
# plots to compare
#Scatter plot for k=2
fviz_cluster(kmeans_2, geom = "point", data = cluster_df) + ggtitle("k = 2")
```



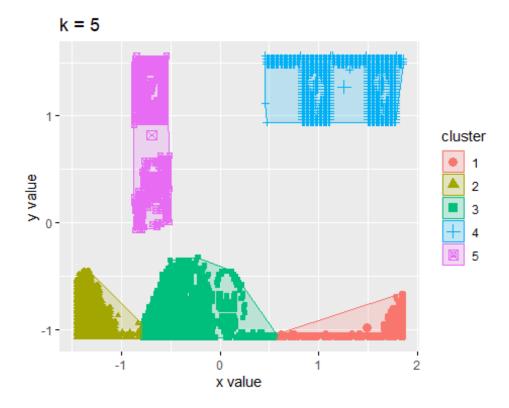
#Scatter plot for k=3
fviz\_cluster(kmeans\_3, geom = "point", data = cluster\_df) + ggtitle("k = 3")



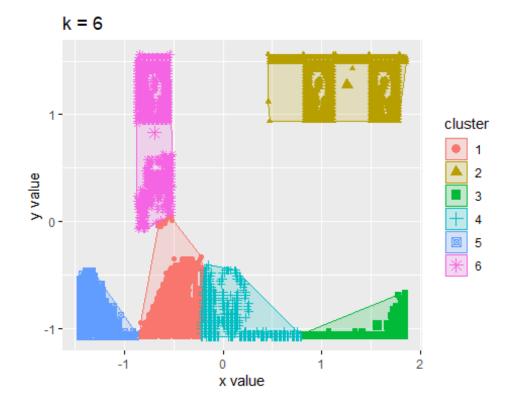
## #Scatter plot for k=4 fviz\_cluster(kmeans\_4, geom = "point", data = cluster\_df) + ggtitle("k = 4")



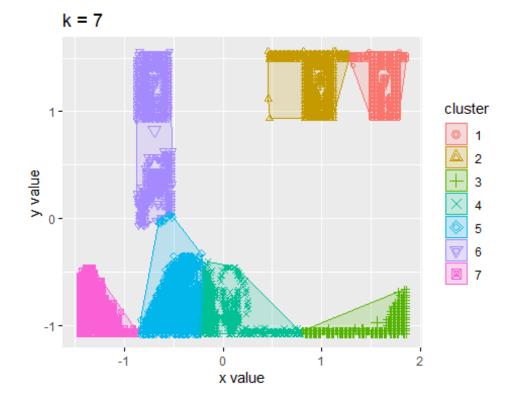
#Scatter plot for k=5
fviz\_cluster(kmeans\_5, geom = "point", data = cluster\_df) + ggtitle("k = 5")



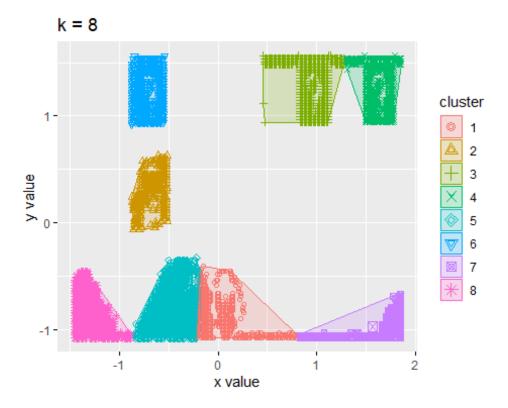
#Scatter plot for k=6
fviz\_cluster(kmeans\_6, geom = "point", data = cluster\_df) + ggtitle("k = 6")



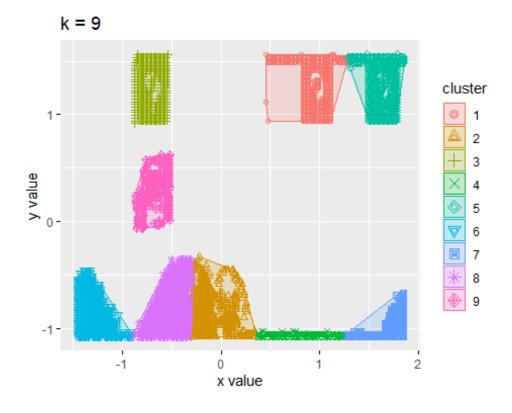
## #Scatter plot for k=7 fviz\_cluster(kmeans\_7, geom = "point", data = cluster\_df) + ggtitle("k = 7")



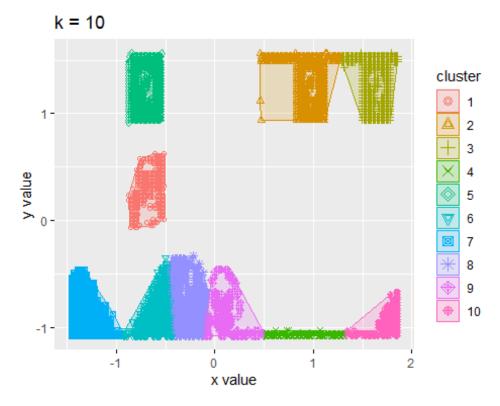
#Scatter plot for k=8
fviz\_cluster(kmeans\_8, geom = "point", data = cluster\_df) + ggtitle("k = 8")



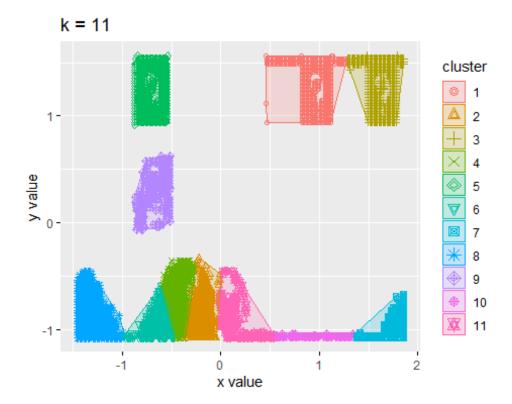
#Scatter plot for k=9
fviz\_cluster(kmeans\_9, geom = "point", data = cluster\_df) + ggtitle("k = 9")



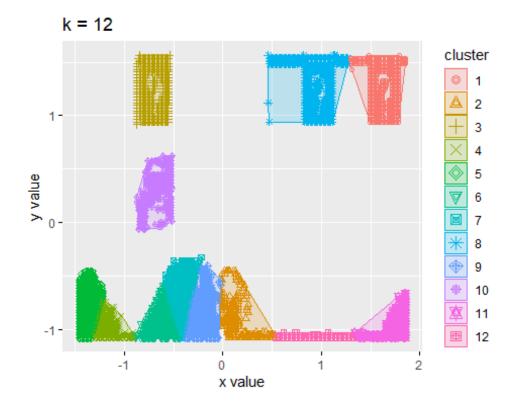
```
#Scatter plot for k=10
fviz_cluster(kmeans_10, geom = "point", data = cluster_df) + ggtitle("k = 10"
)
```

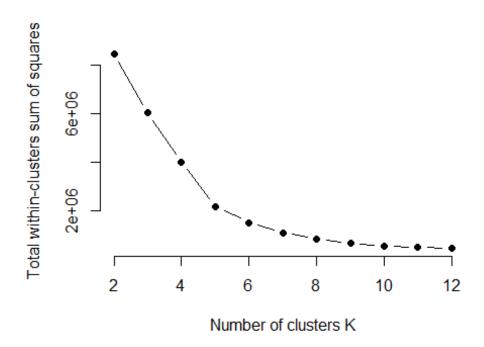


```
#Scatter plot for k=11
fviz_cluster(kmeans_11, geom = "point", data = cluster_df) + ggtitle("k = 11
")
```



#Scatter plot for k=12
fviz\_cluster(kmeans\_12, geom = "point", data = cluster\_df) + ggtitle("k = 12")





#The results suggest that 6 is the optimal number of clusters as it appears to be the bend in elbow.