

BDA Prac 1 K-Means algorithm

Code ▼

Hide

```
install.packages("plyr")
```

WARNING: Rtools is required to build R packages but is not currently installed. Please download and install the appropriate version of Rtools before proceeding:

```
https://cran.rstudio.com/bin/windows/Rtools/  
Installing package into 'C:/Users/User/AppData/Local/R/win-library/4.2'  
(as 'lib' is unspecified)  
trying URL 'https://cran.rstudio.com/bin/windows/contrib/4.2/plyr_1.8.8.zip'  
Content type 'application/zip' length 1153074 bytes (1.1 MB)  
downloaded 1.1 MB
```

```
package 'plyr' successfully unpacked and MD5 sums checked  
Warning in install.packages :  
  cannot remove prior installation of package 'plyr'  
Warning in install.packages :  
  problem copying C:\Users\User\AppData\Local\R\win-library\4.2\00LOCK\plyr\libs\x64\plyr.dll  
to C:\Users\User\AppData\Local\R\win-library\4.2\plyr\libs\x64\plyr.dll: Permission denied  
Warning in install.packages :  
  restored 'plyr'
```

```
The downloaded binary packages are in  
  C:\Users\User\AppData\Local\Temp\RtmpiwDCp8\downloaded_packages
```



Hide

```
install.packages("ggplot2")
```

WARNING: Rtools is required to build R packages but is not currently installed. Please download and install the appropriate version of Rtools before proceeding:

```
https://cran.rstudio.com/bin/windows/Rtools/  
Installing package into 'C:/Users/User/AppData/Local/R/win-library/4.2'  
(as 'lib' is unspecified)  
trying URL 'https://cran.rstudio.com/bin/windows/contrib/4.2/ggplot2_3.4.1.zip'  
Content type 'application/zip' length 4226768 bytes (4.0 MB)  
downloaded 4.0 MB
```

```
package 'ggplot2' successfully unpacked and MD5 sums checked
```

```
The downloaded binary packages are in  
  C:\Users\User\AppData\Local\Temp\RtmpiwDCp8\downloaded_packages
```

Hide

```
install.packages("cluster")
```

WARNING: Rtools is required to build R packages but is not currently installed. Please download and install the appropriate version of Rtools before proceeding:

```
https://cran.rstudio.com/bin/windows/Rtools/  
Installing package into 'C:/Users/User/AppData/Local/R/win-library/4.2'  
(as 'lib' is unspecified)  
trying URL 'https://cran.rstudio.com/bin/windows/contrib/4.2/cluster_2.1.4.zip'  
Content type 'application/zip' length 585941 bytes (572 KB)  
downloaded 572 KB
```

```
package 'cluster' successfully unpacked and MD5 sums checked  
Warning in install.packages :  
  cannot remove prior installation of package 'cluster'  
Warning in install.packages :  
  problem copying C:\Users\User\AppData\Local\R\win-library\4.2\00LOCK\cluster\libs\x64\cluster.dll to C:\Users\User\AppData\Local\R\win-library\4.2\cluster\libs\x64\cluster.dll: Permission denied  
Warning in install.packages :  
  restored 'cluster'  
  
The downloaded binary packages are in  
  C:\Users\User\AppData\Local\Temp\RtmpiWDCp8\downloaded_packages
```

Hide

```
install.packages("lattice")
```

WARNING: Rtools is required to build R packages but is not currently installed. Please download and install the appropriate version of Rtools before proceeding:

```
https://cran.rstudio.com/bin/windows/Rtools/  
Installing package into 'C:/Users/User/AppData/Local/R/win-library/4.2'  
(as 'lib' is unspecified)  
trying URL 'https://cran.rstudio.com/bin/windows/contrib/4.2/lattice_0.20-45.zip'  
Content type 'application/zip' length 1171761 bytes (1.1 MB)  
downloaded 1.1 MB
```

```
package 'lattice' successfully unpacked and MD5 sums checked  
Warning in install.packages :  
  cannot remove prior installation of package 'lattice'  
Warning in install.packages :  
  problem copying C:\Users\User\AppData\Local\R\win-library\4.2\00LOCK\lattice\libs\x64\lattice.dll to C:\Users\User\AppData\Local\R\win-library\4.2\lattice\libs\x64\lattice.dll: Permission denied  
Warning in install.packages :  
  restored 'lattice'  
  
The downloaded binary packages are in  
  C:\Users\User\AppData\Local\Temp\RtmpiWDCp8\downloaded_packages
```

[Hide](#)

```
install.packages("grid")
```

WARNING: Rtools is required to build R packages but is not currently installed. Please download and install the appropriate version of Rtools before proceeding:

```
https://cran.rstudio.com/bin/windows/Rtools/  
Installing package into 'C:/Users/User/AppData/Local/R/win-library/4.2'  
(as 'lib' is unspecified)
```

Warning in install.packages :
package 'grid' is a base package, and should not be updated

[Hide](#)

```
install.packages("gridExtra")
```

WARNING: Rtools is required to build R packages but is not currently installed. Please download and install the appropriate version of Rtools before proceeding:

```
https://cran.rstudio.com/bin/windows/Rtools/  
Installing package into 'C:/Users/User/AppData/Local/R/win-library/4.2'  
(as 'lib' is unspecified)  
trying URL 'https://cran.rstudio.com/bin/windows/contrib/4.2/gridExtra_2.3.zip'  
Content type 'application/zip' length 1109591 bytes (1.1 MB)  
downloaded 1.1 MB
```

package 'gridExtra' successfully unpacked and MD5 sums checked

The downloaded binary packages are in
C:\Users\User\AppData\Local\Temp\RtmpiwDCp8\downloaded_packages

[Hide](#)

```
library(plyr)  
library(ggplot2)  
library(cluster)  
library(lattice)  
library(grid)  
library(gridExtra)
```

[Hide](#)

```
grade_input=as.data.frame(read.csv("F:/GitHub/Practical_BscIT_MscIT_Ninad/MscIT/Semester 2/Bi  
gDataAnalytics/Dataset/grades_km_input.csv"))  
kmdata_orig=as.matrix(grade_input[, c ("Student","English","Math","Science")])  
kmdata=kmdata_orig[,2:4]  
kmdata[1:10,]
```

	English	Math	Science
[1,]	99	96	97
[2,]	99	96	97
[3,]	98	97	97
[4,]	95	100	95
[5,]	95	96	96
[6,]	96	97	96
[7,]	100	96	97
[8,]	95	98	98
[9,]	98	96	96
[10,]	99	99	95

Hide

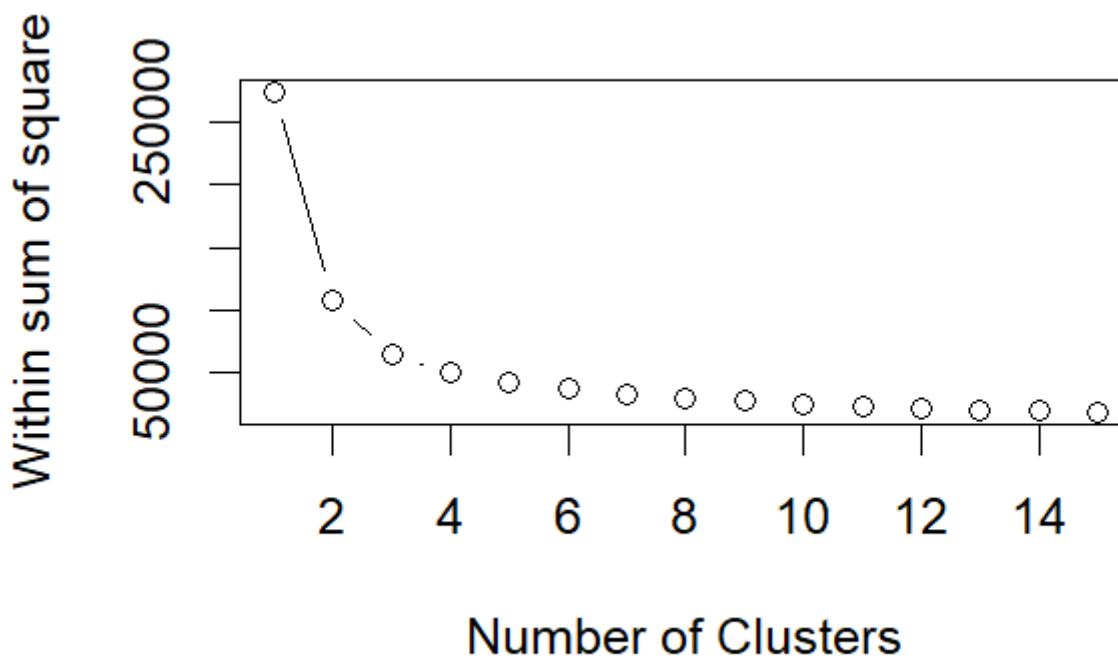
the k-means algorithm is used to identify clusters for $k = 1, 2, \dots, 15$. For each value of k , the WSS is calculated.

```
wss=numeric(15)
```

Hide

the option `n start=25` specifies that the k-means algorithm will be repeated 25 times, each starting with k random initial centroids

```
for(k in 1:15)wss[k]=sum(kmeans(kmdata,centers=k,nstart=25)$withinss)
plot(1:15,wss,type="b",xlab="Number of Clusters",ylab="Within sum of square")
```



Hide

#As can be seen, the WSS is greatly reduced when k increases from one to two. Another substantial reduction in WSS occurs at $k = 3$. However, the improvement in WSS is fairly linear for $k > 3$.

```
km = kmeans(kmdata,3,nstart=25)
```

```
km
```

Hide

```
g1=ggplot(data=df, aes(x=English, y=Math, color=cluster )) + geom_point() + theme(legend.position="right") + geom_point(data=centers,aes(x=English,y=Math, color=as.factor(c(1,2,3))),size=10, alpha=.3, show.legend =FALSE)
```

```
g2=ggplot(data=df, aes(x=English, y=Science, color=cluster )) + geom_point ( ) +geom_point(data=centers,aes(x=English,y=Science, color=as.factor(c(1,2,3))),size=10, alpha=.3, show.legend=FALSE)
```

```
g3 = ggplot(data=df, aes(x=Math, y=Science, color=cluster )) + geom_point ( ) + geom_point(data=centers,aes(x=Math,y=Science, color=as.factor(c(1,2,3))),size=10, alpha=.3, show.legend=FALSE)
```

```
tmp=ggplot_gtable(ggplot_build(g1))
```

Hide

```
grid.arrange(arrangeGrob(g1 + theme(legend.position="none"),g2 + theme(legend.position="none"),g3 + theme(legend.position="none"),top ="High School Student Cluster Analysis" ,ncol=1))
```

