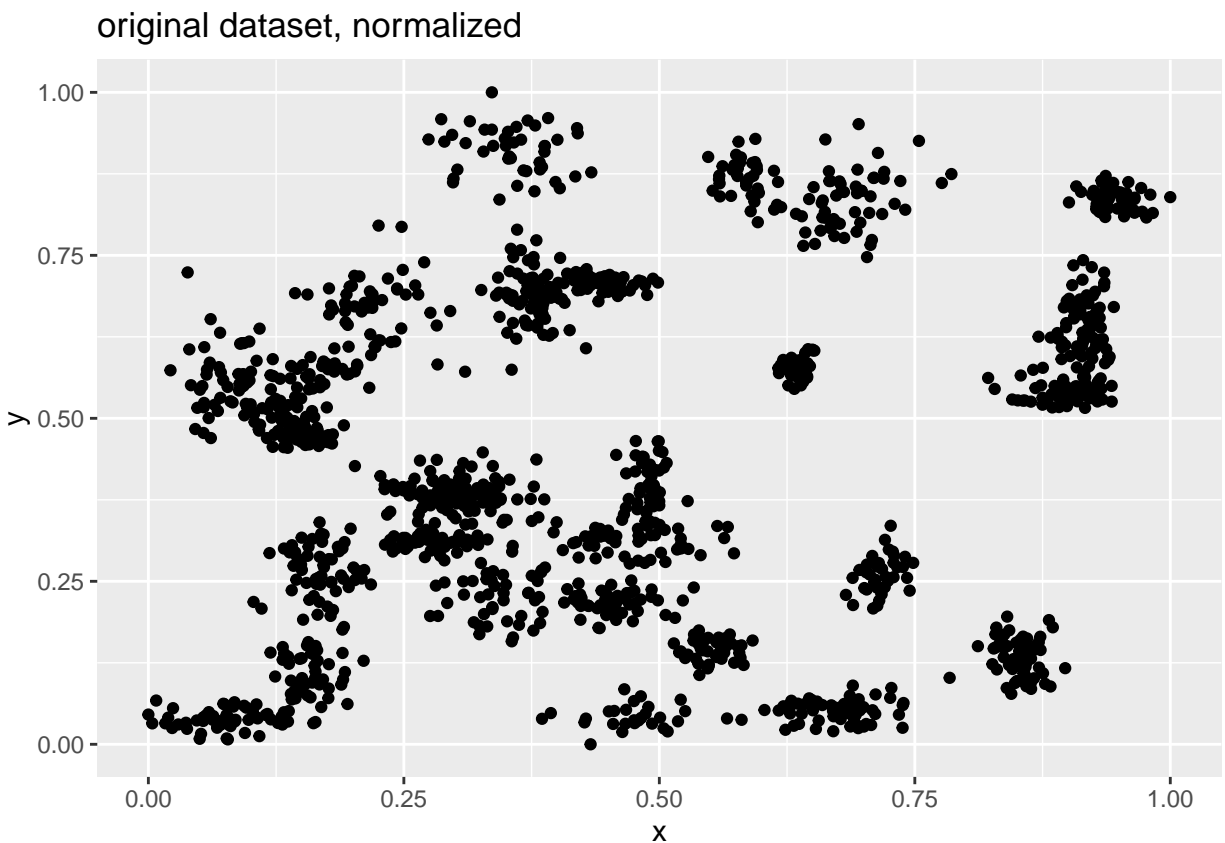title: "DSC20 week10 assignment" author: "xin tang" date: "2023-08-08" output: pdf_document: default editor_options: markdown: wrap: 72

## use the nearest neighbors algorithm to fit a model on two simplified datasets.

**Plot the data from each dataset using a scatter plot.**

```
##
## Attaching package: 'dplyr'

## The following objects are masked from 'package:stats':
##
##     filter, lag

## The following objects are masked from 'package:base':
##
##     intersect, setdiff, setequal, union

##
## Attaching package: 'scales'

## The following object is masked from 'package:readr':
##
##     col_factor
```

original dataset, normalized

from data, the dataset is too scattered to use a linear regression

now start to fit a nearest neighbor model using K =3, 5,15, 20 and 25.

```
## [1] "0.711 is the accuracy when k =3"
```

```
## [1] "0.708 is the accuracy when k =5"
```

```
## [1] "0.691 is the accuracy when k = 10"
```
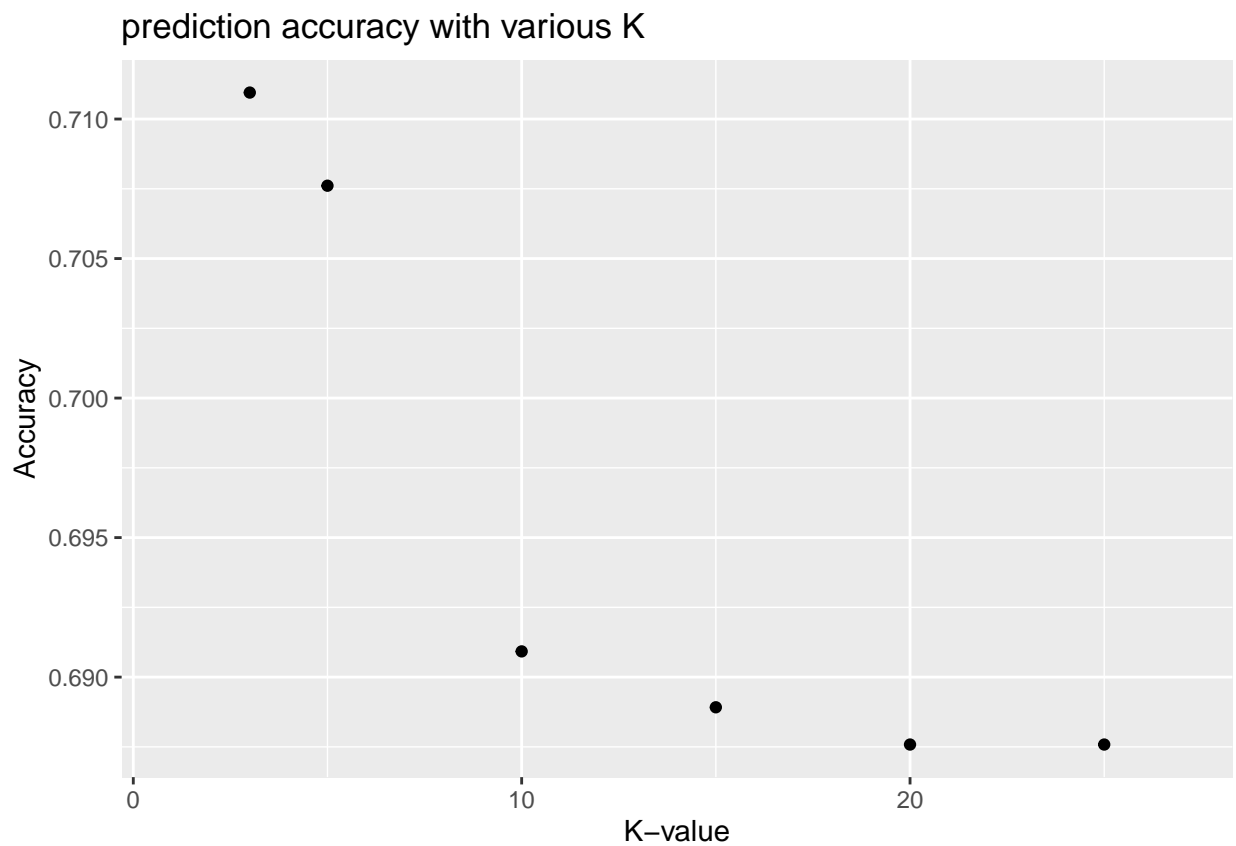
```
## [1] "0.689 is the accuracy when k = 15"
```

```
## [1] "0.688 is the accuracy when k = 20"
```

```
## [1] "0.688 is the accuracy when k = 25"
```
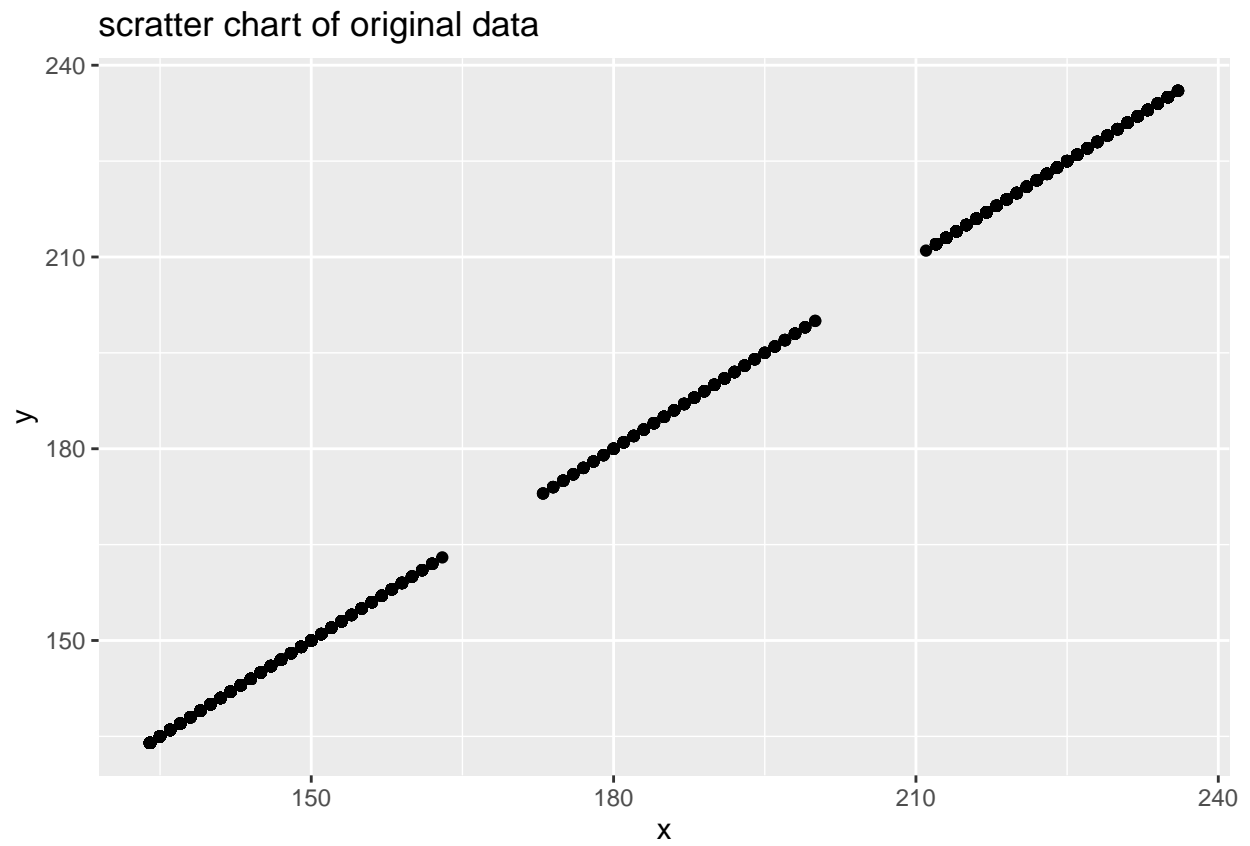
From the output, the accuracy greatly improved (last week is only 58%), which prove this is a better model.

Plot a comparision chart for different K value
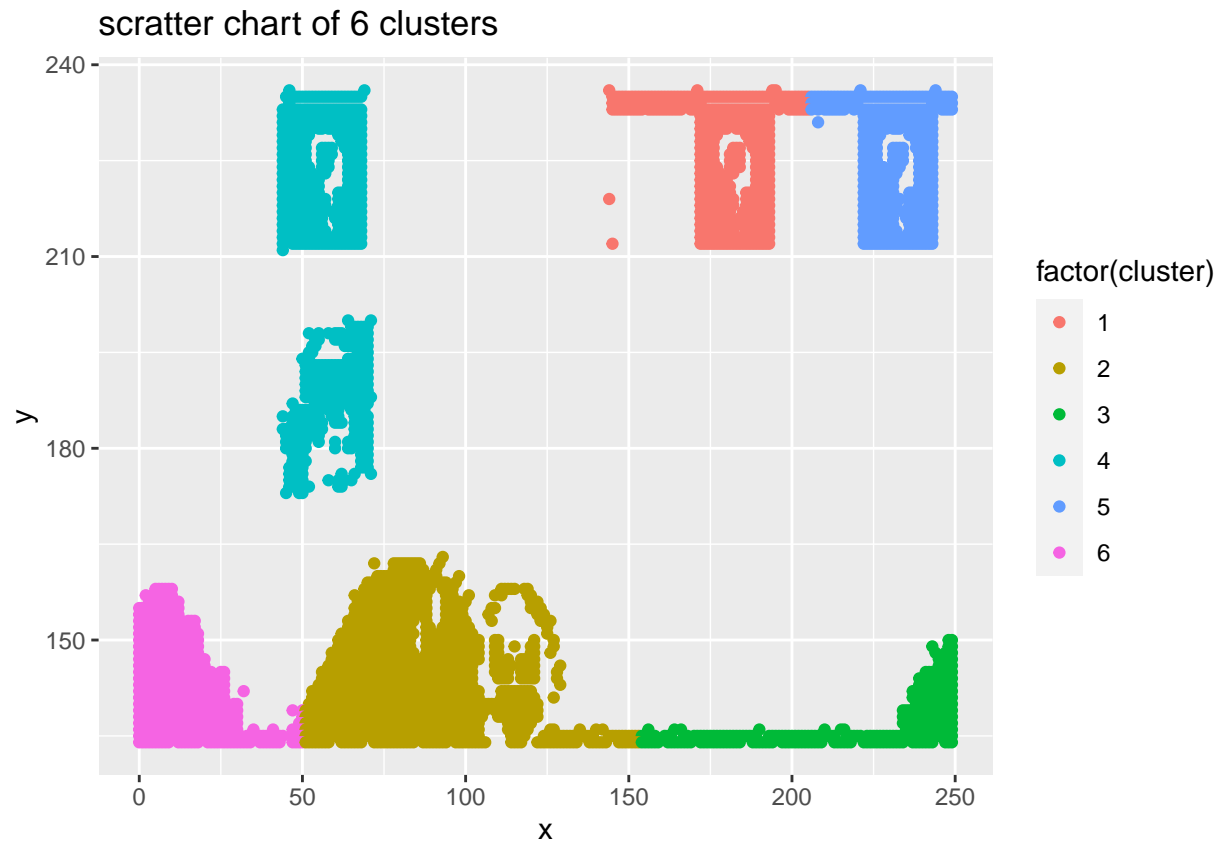


prediction accuracy with various K

# Begin of the cluster data analysis using K means clustering

##plot the original data in scatter plot

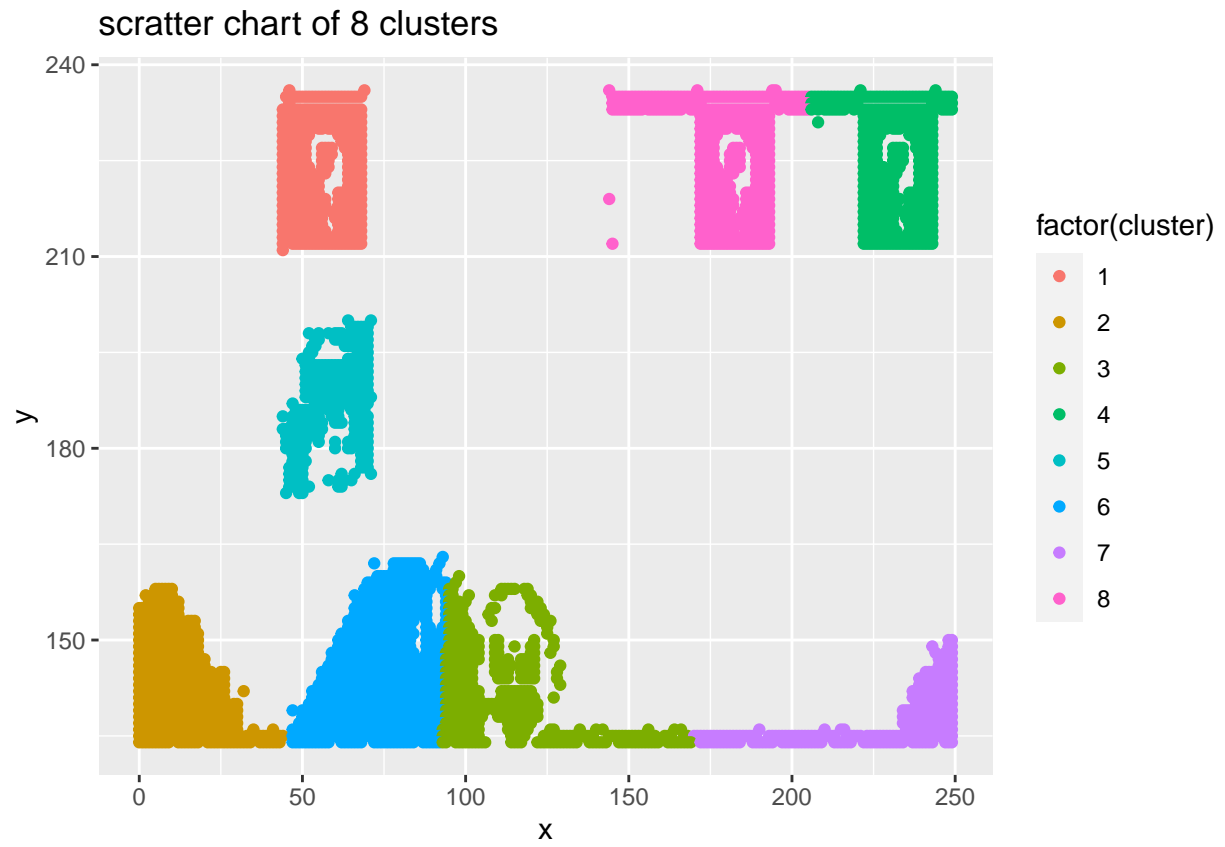scratter chart of original data

start to model using different K-value, visualize using scatter chart

## scratter chart of 2 clusters



## scratter chart of 4 clusters

scratter chart of 6 clusters

scratter chart of 8 clusters

scratter chart of 10 clusters

scratter chart of 12 clusters

**Finally do a comparison of average distance under different K**



Avg Distance for different K

**Finally, using Factoextra package to visualize the optimal K-value**

```
## Welcome! Want to learn more? See two factoextra-related books at https://goo.gl/ve3WBa
```

Optimal number of clusters

Optimal number of clusters

Cluster plot