

ETX2250/ETF5922: Data Visualization and Analytics

K means clustering

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

Lecturer: Lauren Kennedy

Department of Econometrics and Business Statistics

✉ lauren.kennedy1@monash.edu

CALENDAR Week 8



Partitional Clustering

1. Random assignment of points in the midst of cluster
2. Repeatedly re-assign points cluster
3. Minimise within-cluster distances
4. Maximise between-cluster distances

Assignment Project Exam Help

One example is k-means, which we focus on today

<https://tutorcs.com>

WeChat: cstutorcs

Comparing clustering algorithms

i

Hierarchical clustering

- Suitable for small data set (<500 observations)
- Easily examine solutions with increasing numbers of clusters
- Convenient method if you want to observe how clusters are nested

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

i

K-means clustering

- Suitable if you know how many clusters you want
- Large data set
- Good for summarising data with k “average” observations that describe the data with the minimum amount of error

Obtaining clusters: K means method

- Non-hierarchical clustering method
- Require number of clusters be specified before assigning observations
- No single objective method for determining the number of clusters

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

Obtaining clusters: K means method

Consider the following

- Want cohesion with, diversity between
- How much improvement does one more cluster provide?
- Single member or very small clusters are not usually useful
- All clusters should be significantly different across the set of clustering variables
- The utility of the clusters ultimately depends on external validation
 - Other variables not used in the analysis
 - External information, e.g. location, geopolitical information

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

K-means Clustering

- Given the value of K, the K-means algorithm **randomly** partitions the observations into K clusters
- After all observations have been assigned to a cluster, the resulting cluster centroids are calculated
- Using the updated cluster centroids, all observations are reassigned to the cluster with the closest centroid (usually Euclidean distance)
- Continue until the assignments no longer change or a predetermined maximum number of iterations is reached

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

K-means Clustering

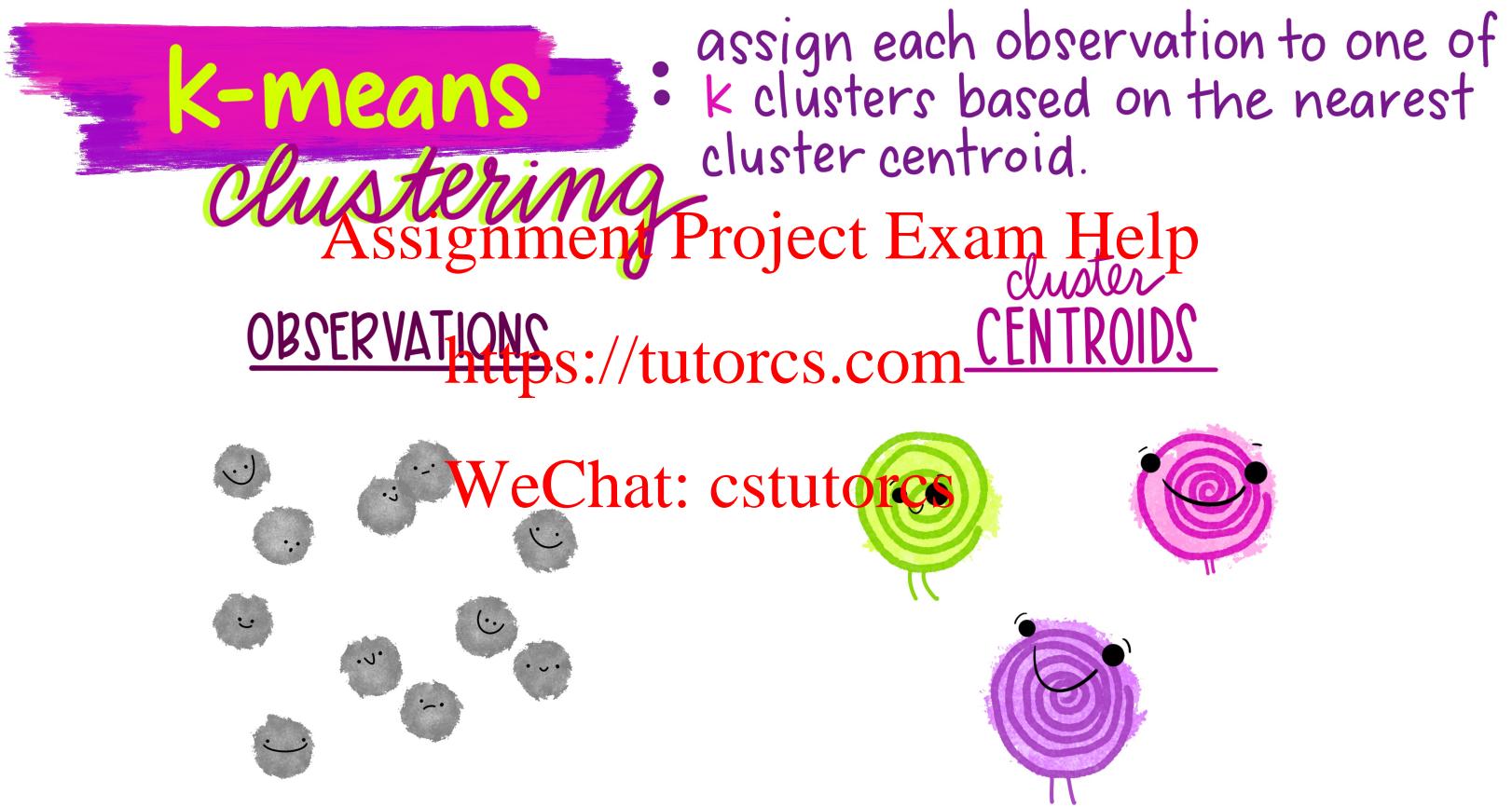
- Number of K clusters needs to be specified, usually advisable to try several values of K
- Initial K clusters are chosen randomly and the results do depend strongly on this particular starting point
 - procedure should be repeated a number of times (>10) with different starting points
- The “best” result is chosen

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

The algorithm in pictures



@allison_horst

The algorithm in pictures

①

Specify the number of clusters (in this example, $k=3$).

Then imagine k cluster centroids are created.



@allison_horst

The algorithm in pictures

②

Those k centroids get randomly placed
in your space.

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

observations,
not currently assigned
to any cluster.

@allison_horst

The algorithm in pictures

③

Each observation gets temporarily "assigned" to its closest centroid.
(e.g. by Euclidean distance)

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

Note: observations stay put, but their color updates to nearest centroid!



@allison_horst

The algorithm in pictures

④

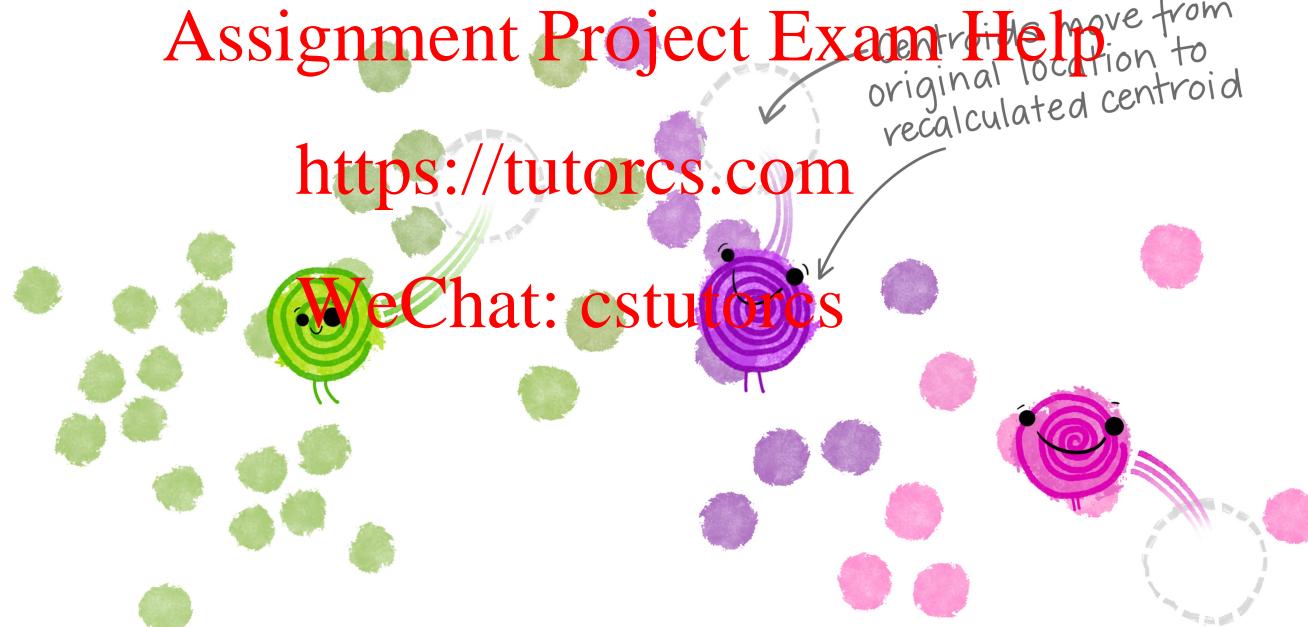
Then the centroid of each cluster is calculated based on all observations assigned to that cluster...

Assignment Project Exam Help

<https://tutorcs.com>

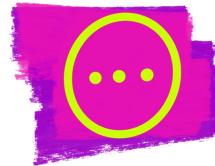
WeChat: cstutorcs

centroids move from original location to recalculated centroid

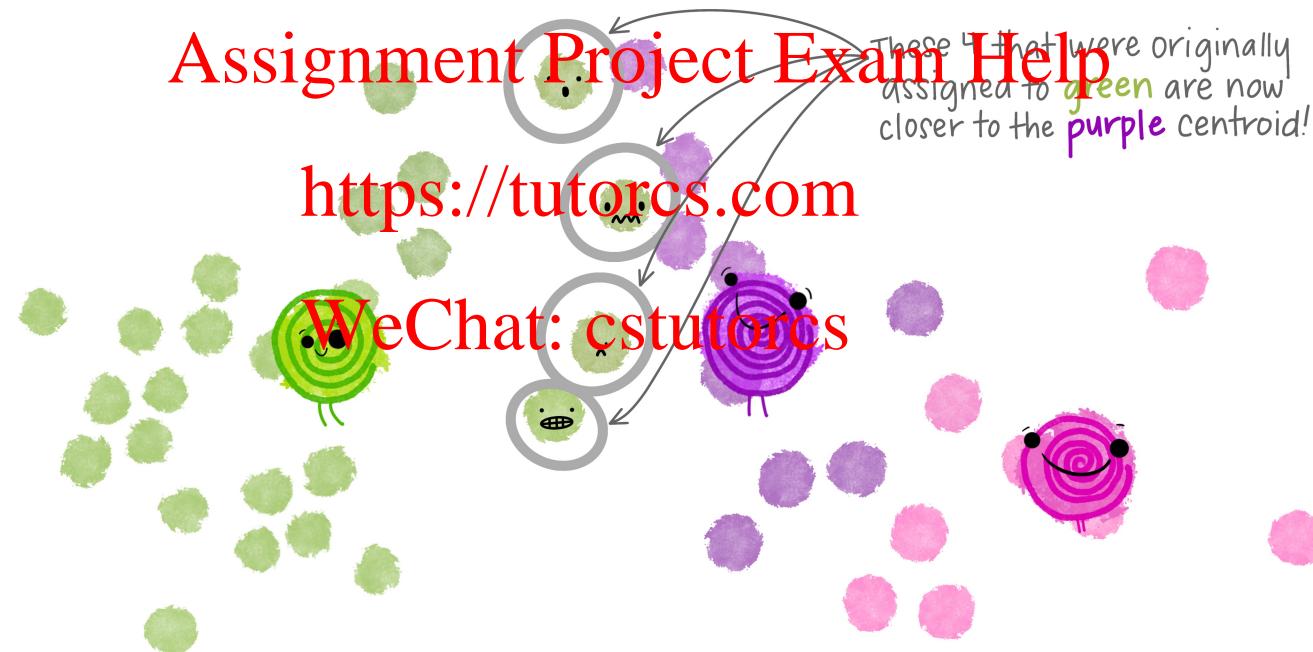


@allison_horst

The algorithm in pictures



UH OH. Now that the cluster centroids have moved, some of the observations are now closer to a different centroid!



@allison_horst

The algorithm in pictures

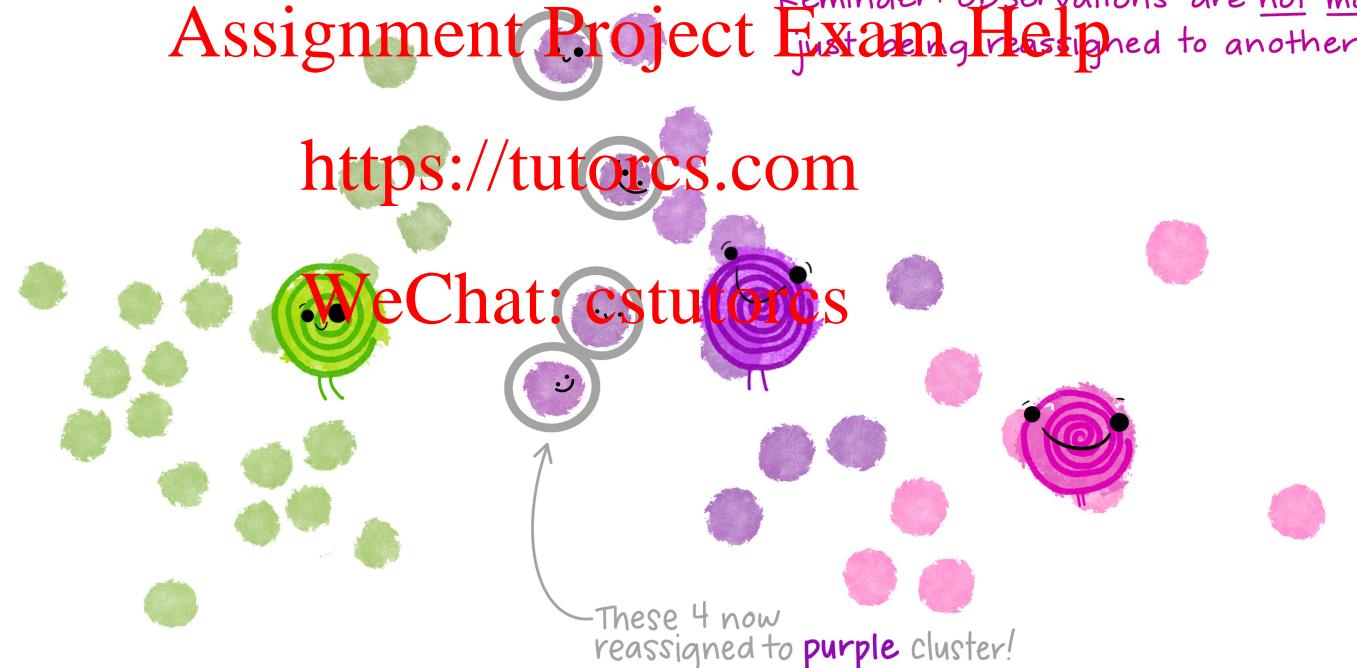


NO PROBLEM!

Observations get reassigned* to a different cluster based on the recalculated centroid.

*Reminder: observations are not moving, just being reassigned to another cluster.

Assignment Project Exam Help



@allison_horst

The algorithm in pictures



But now that observations have been reassigned,
the centroids need to move again [recalculate
centroids from updated clusters]

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

beep beep
RECALCULATING
CENTROIDS

@allison_horst

The algorithm in pictures

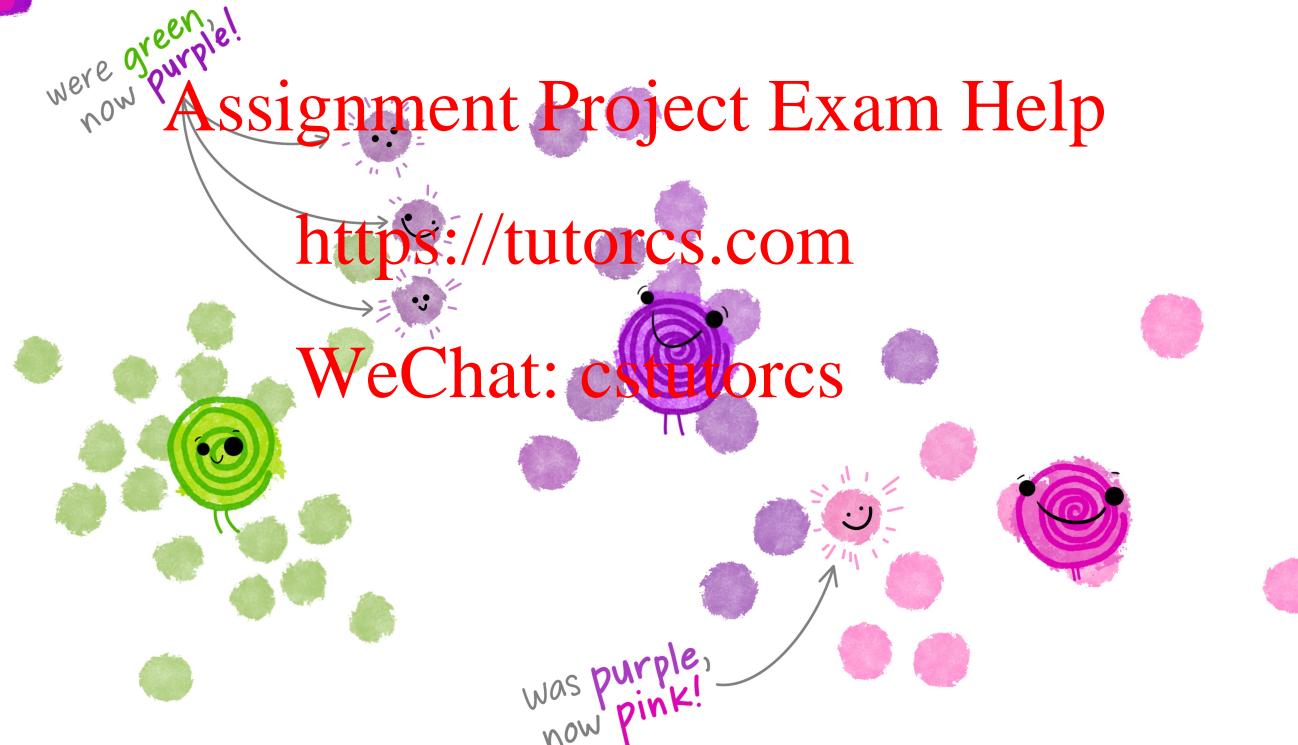


Again, now observations are reassigned as needed to the closest centroid.

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs



@allison_horst

The algorithm in pictures



Then the centroid for each cluster
is recalculated...

Assignment Project Exam Help

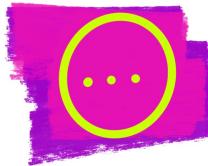
<https://tutorcs.com>

WeChat: cstutorcs

...which means observations will be reassigned...

@allison_horst

The algorithm in pictures



That iterative process of

Assignment Project Exam Help
Recalculate cluster centroids
→ Reassign observations to nearest centroid

https://tutortcs.com
↳ Recalculate cluster centroids

Reassign observations to nearest centroid

WeChat: cstutorcs
↳ Recalculate cluster centroids

Reassign observations to nearest centroid



Continues until nothing is moving
or being reassigned anymore!

@allison_horst

The algorithm in pictures



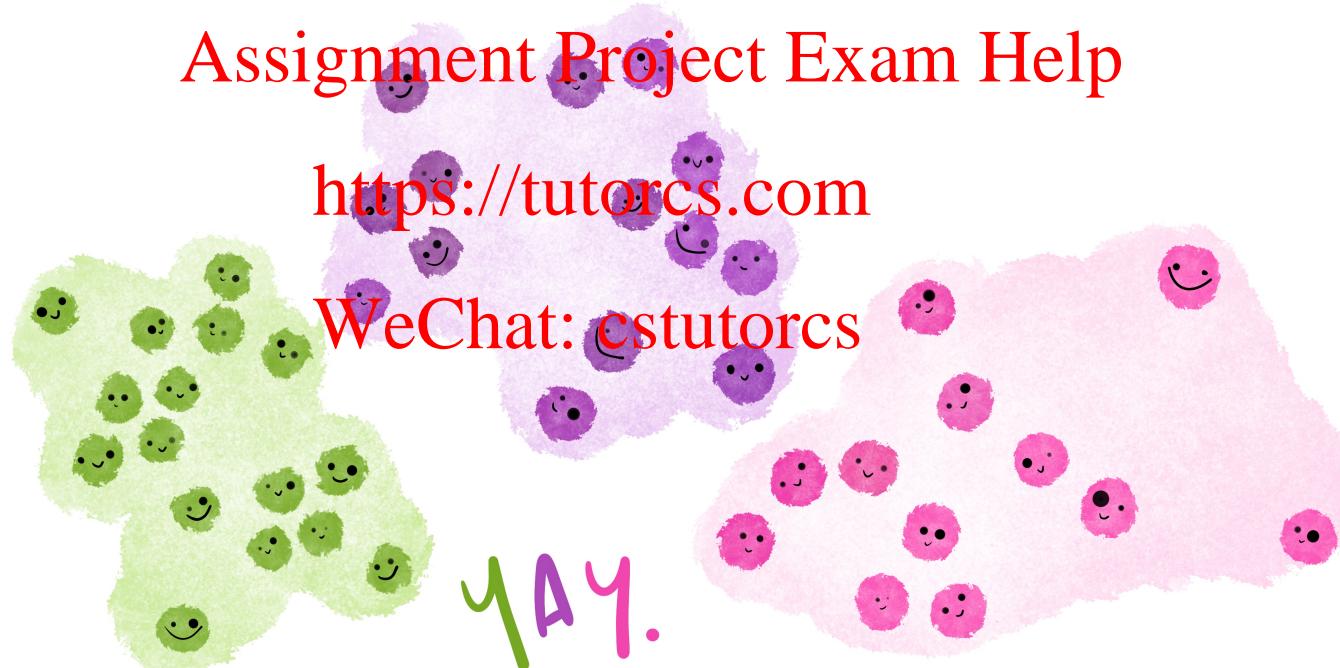
Which means the iteration is done and each observation is assigned to its final cluster.

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

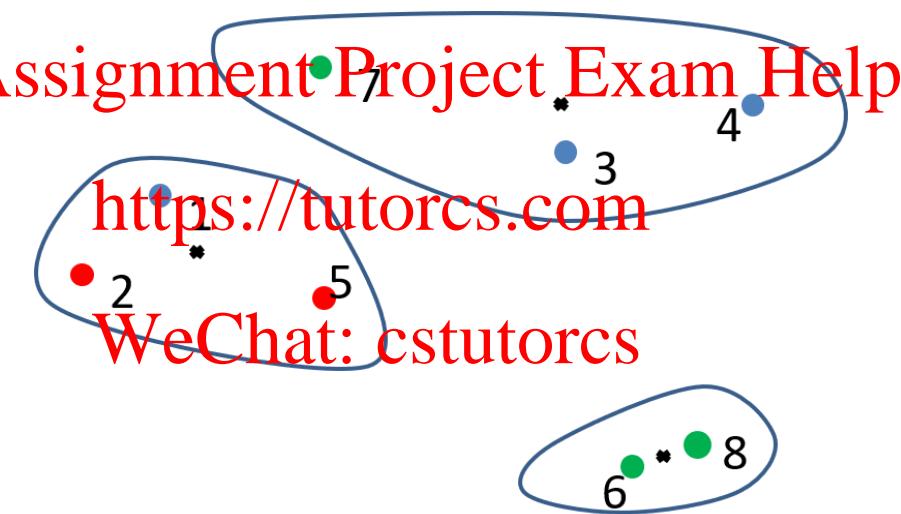
yay.



@allison_horst

Centroid

- A centroid is the average observation of a cluster
- The average value for each variable across the observations in the cluster
- Centroid does not have to be an actual observation



Recall KTC Example

- Example: continuous variables in KTC data
- These are: Age, Income, Children (columns 2, 4 and 6)

ID	Age	Female	Income	Married	Children	CarLoan	Mortgage
1	49	1	17546	0	1	0	0
2	40	0	30085	1	3	1	1
3	51	1	16375	1	0	1	0
4	23	1	20375	1	3	0	0

- Use Euclidean distance: the variables will need standardising so that each of the three variables is on a similar scale. This is easily done in R.

Within Sum of Squares

- Use $\| \cdot \|$ to represent the Euclidean distance in n dimensional space, where n is the number of variables
- Suppose we have k clusters, call them S_1, \dots, S_k
- We are trying to select the clusters in such a way as to minimise the “within sum of squares”

Assignment Project Exam Help

$$\sum_{j=1}^k \sum_{x \in S_j} \|x - \bar{m}_j\|^2$$

WeChat: cstutorcs

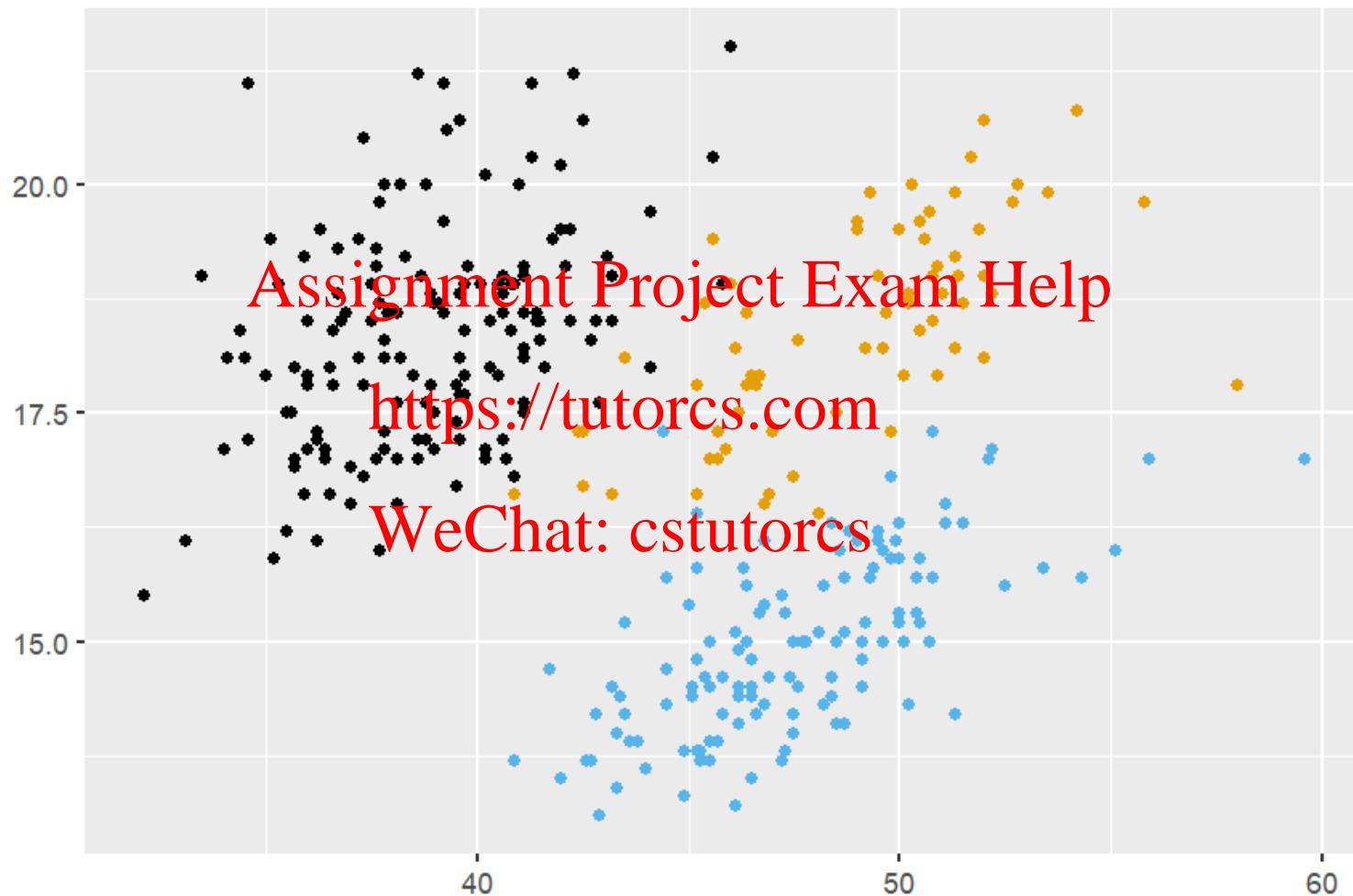
Total SS and between SS

Total SS = the sum of squared distances between each data point to the global sample mean

Obtain between SS:

1. For each cluster calculate the squared distance from the centroid of the cluster to the global sample mean
2. Multiply the number of points in the cluster
3. Add up the results for each cluster <https://tutorcs.com>
4. If there were no discernible pattern of clustering, the three means of the three groups would be close to the global mean, and between SS would be a very small fraction of total SS5

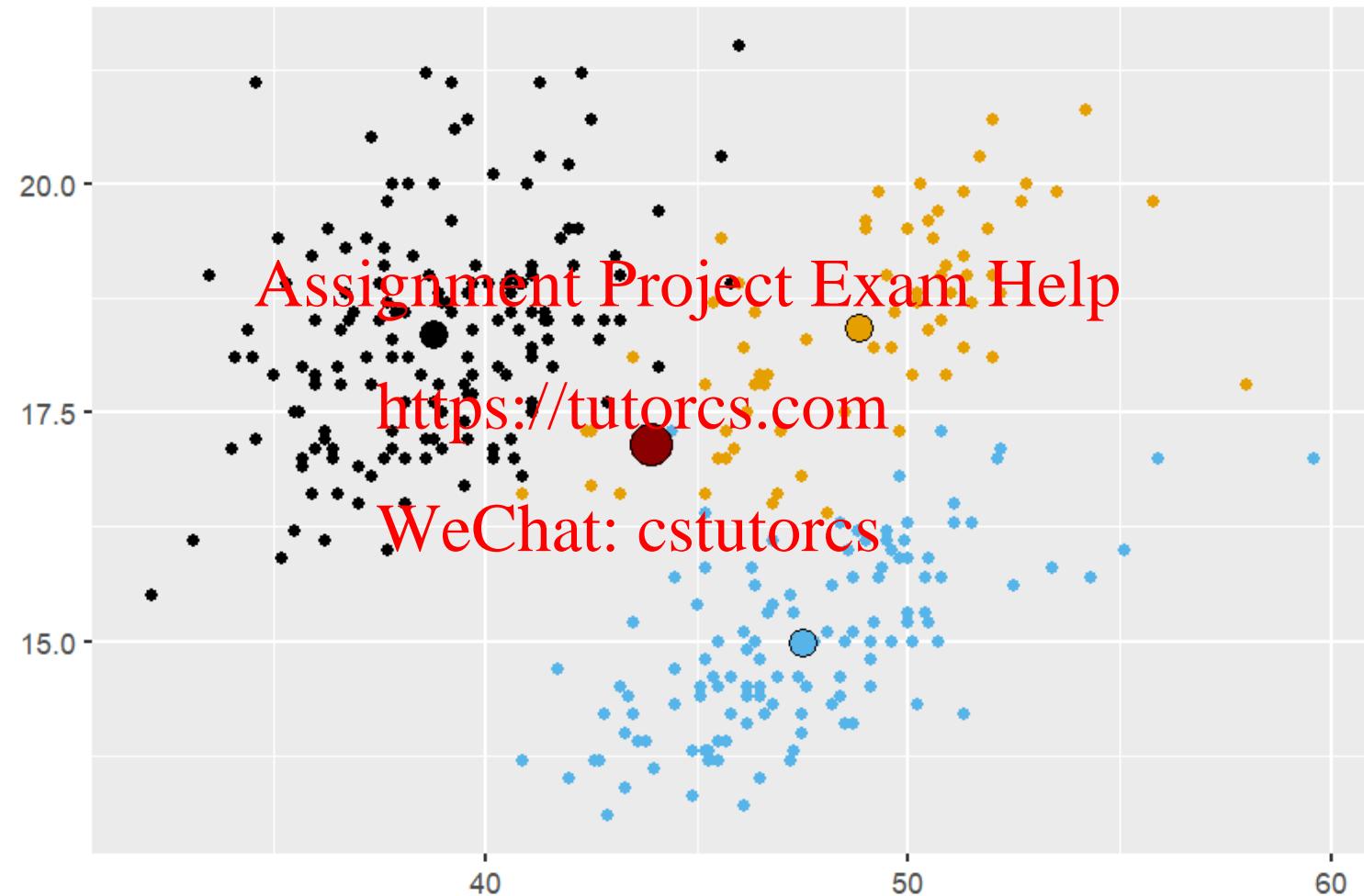
Visualise



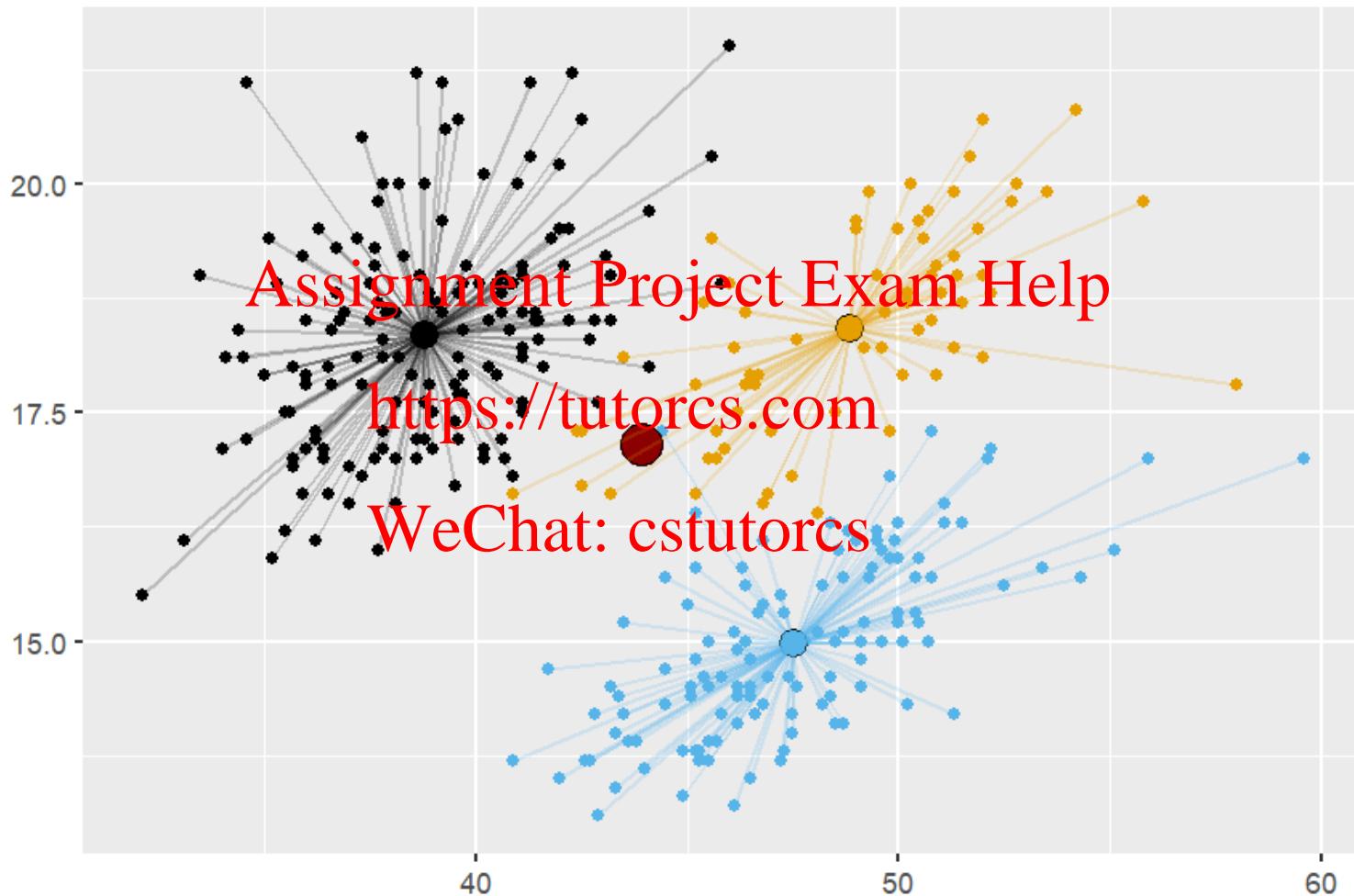
Group centroids



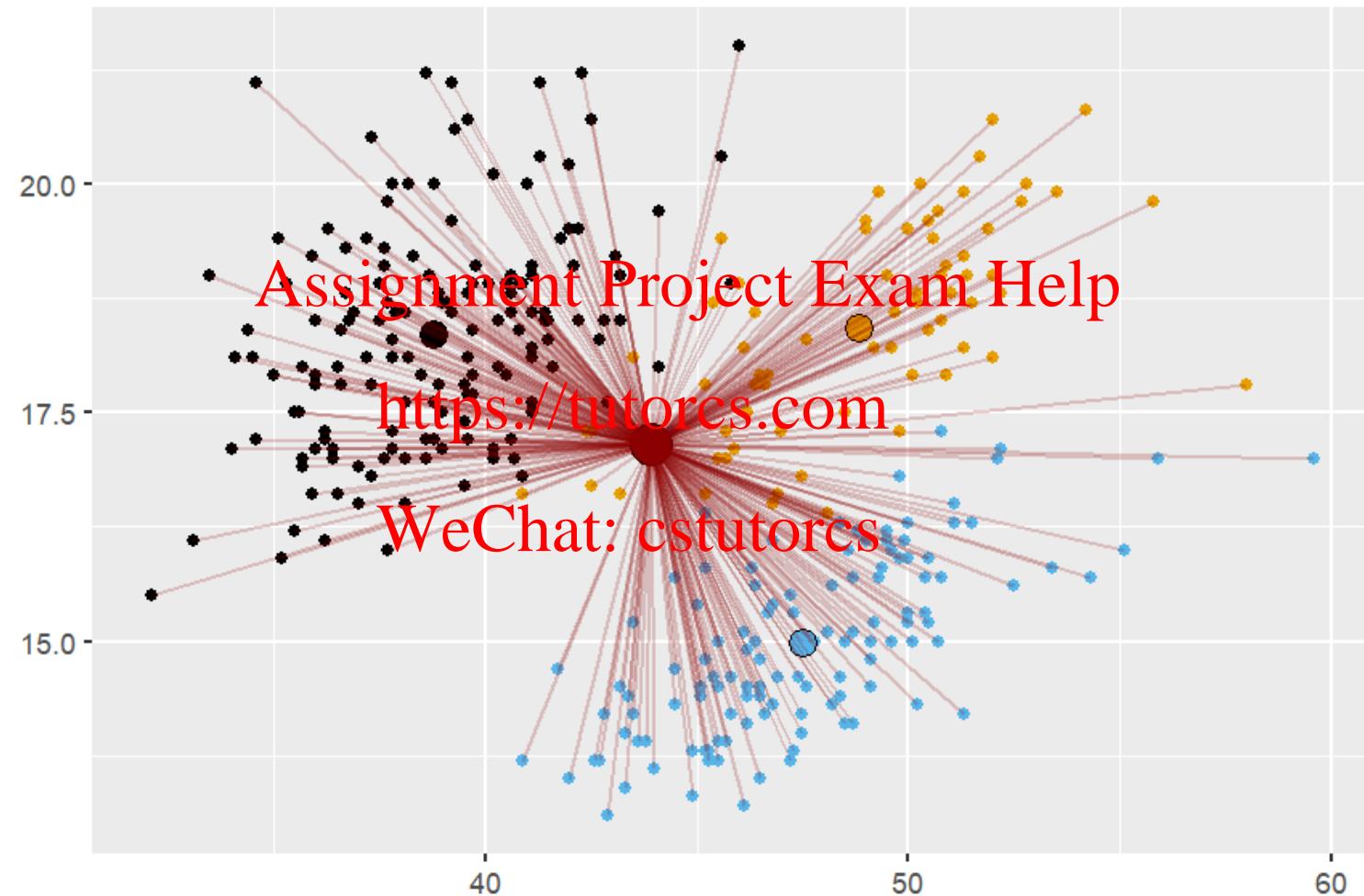
Overall centroid



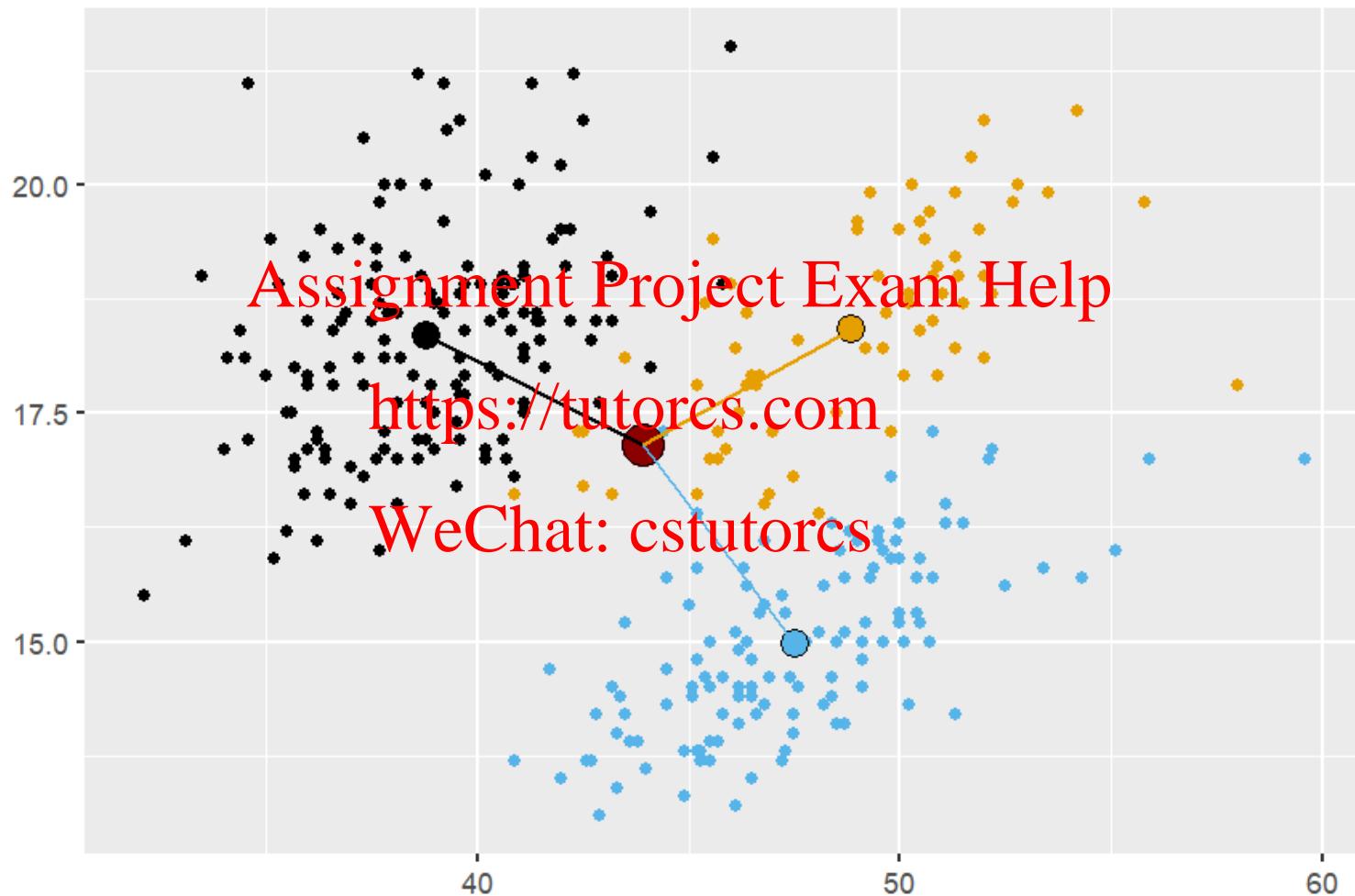
Within sum of squares



Total sum of squares



Between sum of squares



Outcome measure

- The algorithm repeats this process (calculate cluster centroid, assign observations to the cluster with the nearest centroid) until there is no change in the clusters or a specified maximum number of iterations has been reached
- Since we want little variation within clusters and a lot of variations between clusters, the ratio

Assignment Project Exam Help

$$\frac{SS_{BTWN}}{SS_{TOT}}$$

<https://tutorcs.com>

- Is a possible measure of how good the clustering that k-means has found
- One rule of thumb is that the ratio

WeChat: cstutorcs

$$\frac{SS_{BTWN}}{SS_{WTHN}}$$

should exceed one (i.e. there is more variation between clusters than within.)

Assignment Project Exam Help

How to in R

<https://tutorcs.com>

WeChat: cstutorcs

Know thy Customer

- Remember from your tutorials last week

```
ktc_df <- read_csv(here::here("data/KTC.csv"))
head(ktc_df)
```

```
## # A tibble: 6 x 8
##   ID    Age Female Income Married Children CarLoan Mortgage
##   <dbl> <dbl>  <dbl>  <dbl>   <dbl>   <dbl>   <dbl>
## 1 1     48      1     17546     0       1       0       0
## 2 2     40      0     30055     1       3       1       1
## 3 3     51      1     16575     1       0       1       0
## 4 4     23      1     20375     1       3       0       0
## 5 5     57      1     50576     1       0       0       0
## 6 6     57      1     37870     1       2       0       0
```

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

Focus only on the numeric predictors:

```
ktc_numeric <- ktc_df %>%
  select(ID, Age, Income, Children) %>%
  # We move the ID to rownames so we can keep track of the rows but ID
  # doesn't get included in the clustering (which it would if we left
  # it in the data frame)
  column_to_rownames(var = 'ID') %>%
  # Remember last week? k-means has the same sensitivity to scale,
  # so we need to scale all our numeric variables.
  scale
```

Assignment Project Exam Help

WeChat: cstutors

```
head(ktc_numeric)
```

```
##          Age      Income   Children
## 1  0.1559026 -0.7637480  0.06169096
## 2 -0.4574848  0.1512843  1.91241969
## 3  0.3859229 -0.8346067 -0.86367341
## 4 -1.7609332 -0.5573020  1.91241969
## 5  0.8459635  1.6466130  0.86367341
## 6  0.8459635  0.7193939  0.98705533
```

Assignment Project Exam Help
<https://tutorcs.com>

WeChat: cstutorcs

How do we do kmeans?

```
set.seed(145)  
km <- kmeans(ktc_numeric, centers = 3, nstart = 25)
```

- centers is the number of clusters. Note that for k-means (unlike hierarchical clustering), we need to specify this. You can also specify a set of cluster centers
- nstart is the number of times you run through the procedure. You only need to specify this if you specify the number of clusters (why? more on that later...)

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

What does km look like?

km

```
## K-means clustering with 3 clusters of sizes 8, 11, 11
```

```
##
```

```
## Cluster means:
```

```
##          Age      Income   Children
```

```
## 1 -0.6300000 -0.4368752  1.3340670
```

```
## 2  0.9644588  0.9939719 -0.3589292
```

```
## 3 -0.5062770 -0.6762445 -0.6113013
```

```
##
```

```
## 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
```

```
## 3  1  3  1  2  2  3  2  1  1  2  3  3  2  3  3  1  2  2  3  2  1  2  3  1  2  3  1  2
```

```
## 27 28 29 30
```

```
## 3  3  1  2
```

```
##
```

```
## Within cluster sum of squares by cluster:
```

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

What do the clusters look like?

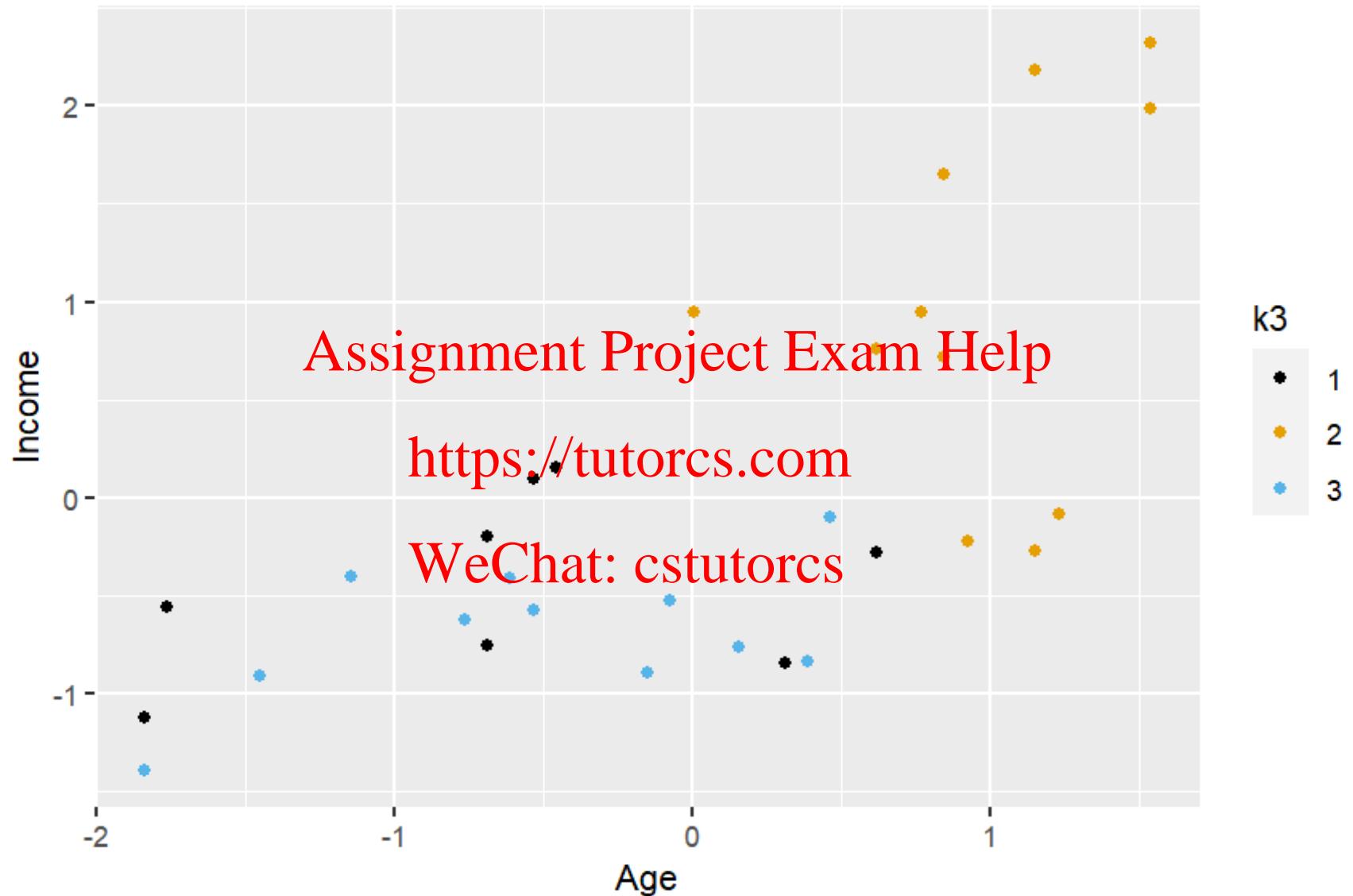
- Let's use our visualisation skills!

```
ktc_numeric %>%
  #This is technically not a dataframe, so let's make it so
  as.data.frame() %>%
  # Add in the cluster allocations
  mutate(k3 = factor(km$cluster))%>%
  ggplot(data = .,
         aes (x =Age,
              y= Income,
              colour = k3))+
  geom_point()+
  ggthemes::scale_color_colorblind()
```

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs



How can we summarise the clusters?

- Cluster 1 characterised by younger, low income customers
- Cluster 2 characterised by older, high income customers
- Cluster 3 is characterised by younger, low income customers
- It appears that the separation between 1 and 3 is not good because they appear very similar...

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

What do the clusters look like (pt 2)?

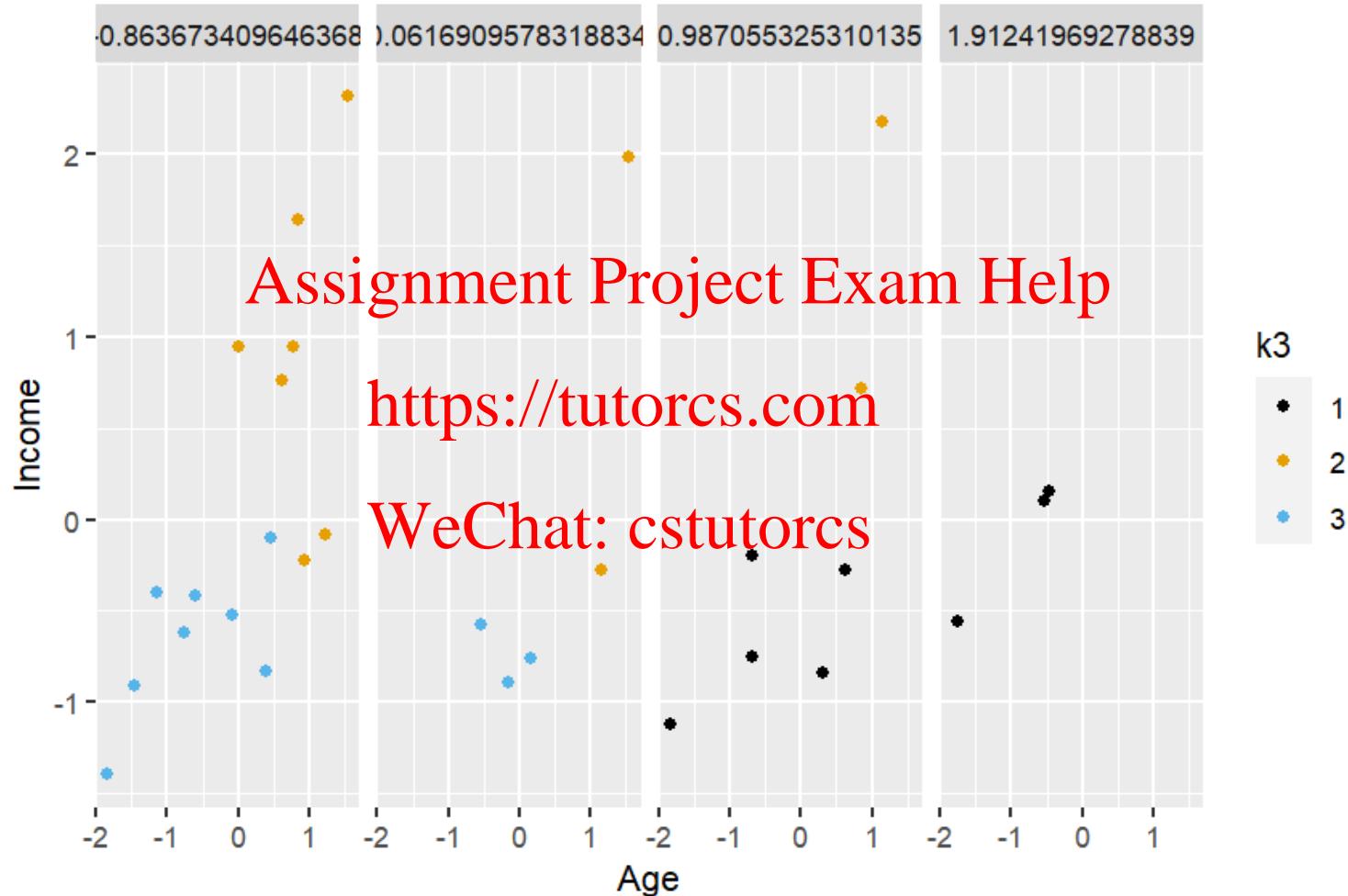
```
ktc_numeric %>%
  #This is technically not a dataframe, so let's make it so
  as.data.frame() %>%
  # Add in the cluster allocations
  mutate(k3 = factor(km$cluster))%>%
  ggplot(data = .,
         aes (x =Age,
              y= Income,
              colour = k3))+  
  geom_point()+
  facet_grid(.~Children)+
  ggthemes::scale_color_colorblind()
```

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

What do the clusters look like (pt 2)?



How can we summarise the clusters (pt 2)?

- Cluster 1 characterised by younger, low income customers with more children
- Cluster 2 characterised by older, high income customers with less children
- Cluster 3 is characterised by younger, low income customers with no children
- The separation appears much better!

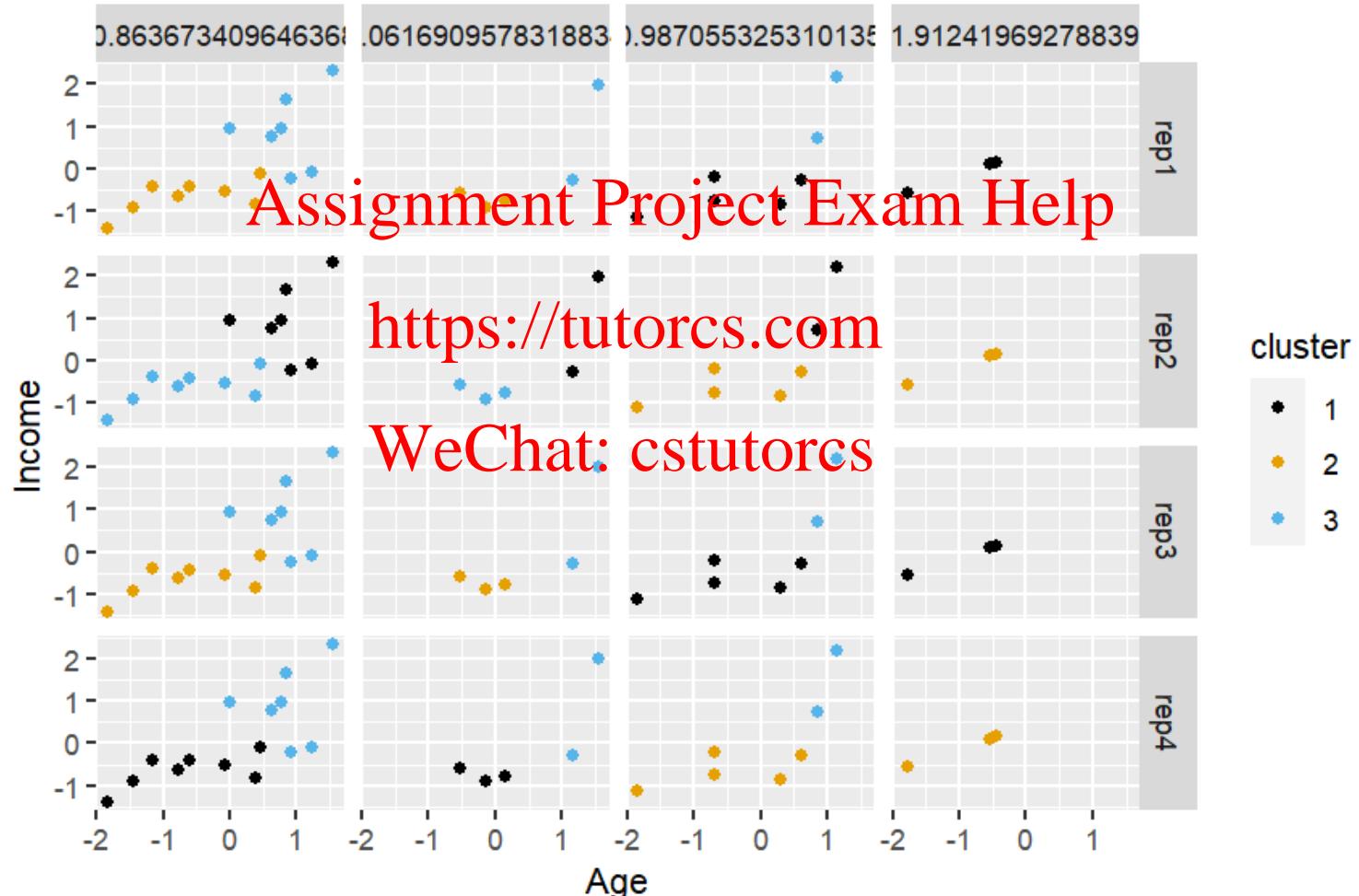
Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

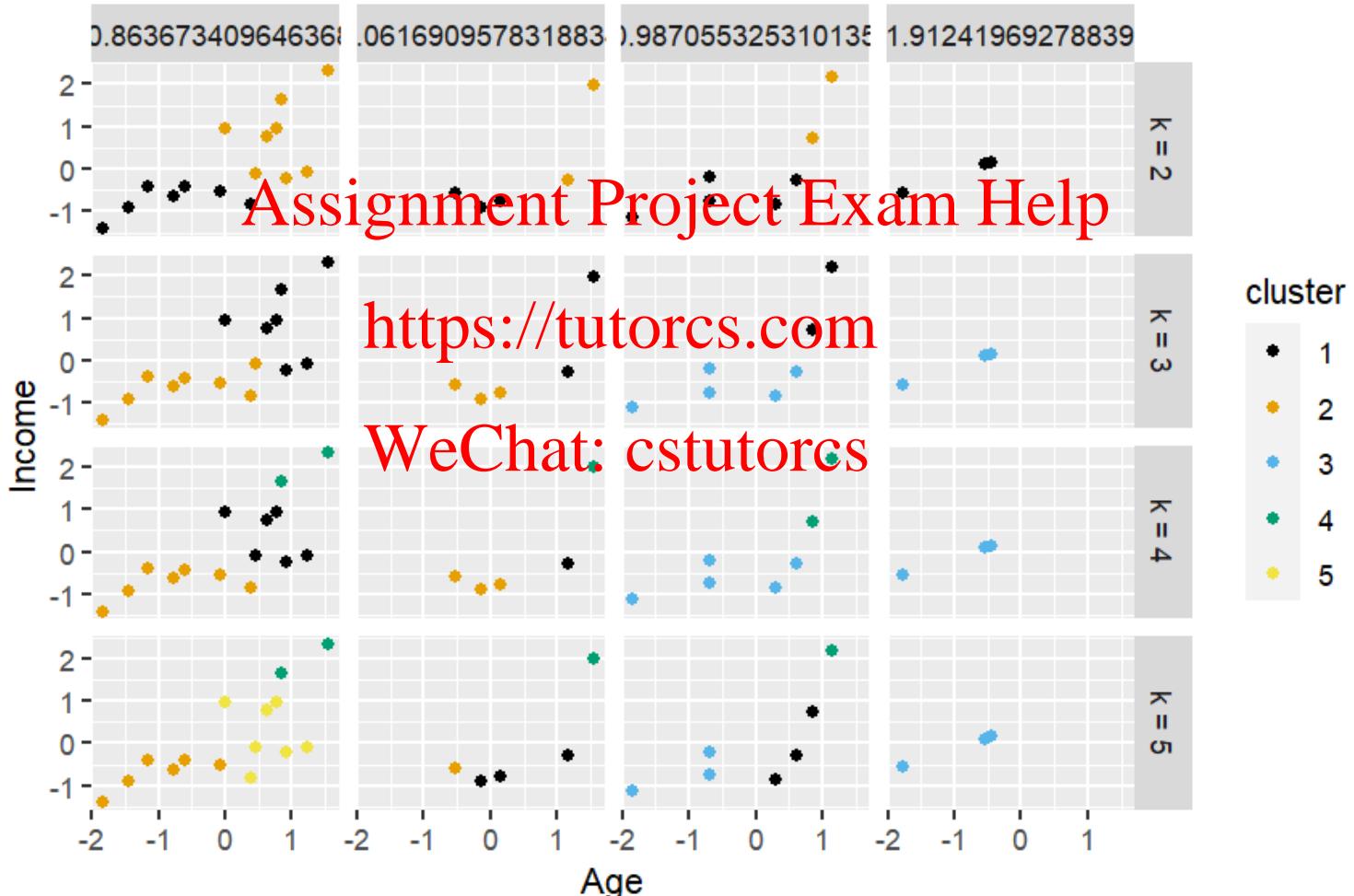
How stable is k-means?

- We can explore by running the same kmeans, with the same number of clusters a few times over:



But wait...why are we using $k = 3$?

- What if we tried different values of k



But wait...why are we using $k = 3$?

- It's difficult to tell from visual inspection which cluster is best (sometimes it is).
- Can we visualise a diagnostic of some kind instead?
- Consider:
 - The total within sum of squares. This measures the compactness (i.e., goodness) of the clustering.
Smaller is better

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

To do this let's run for a number of different k

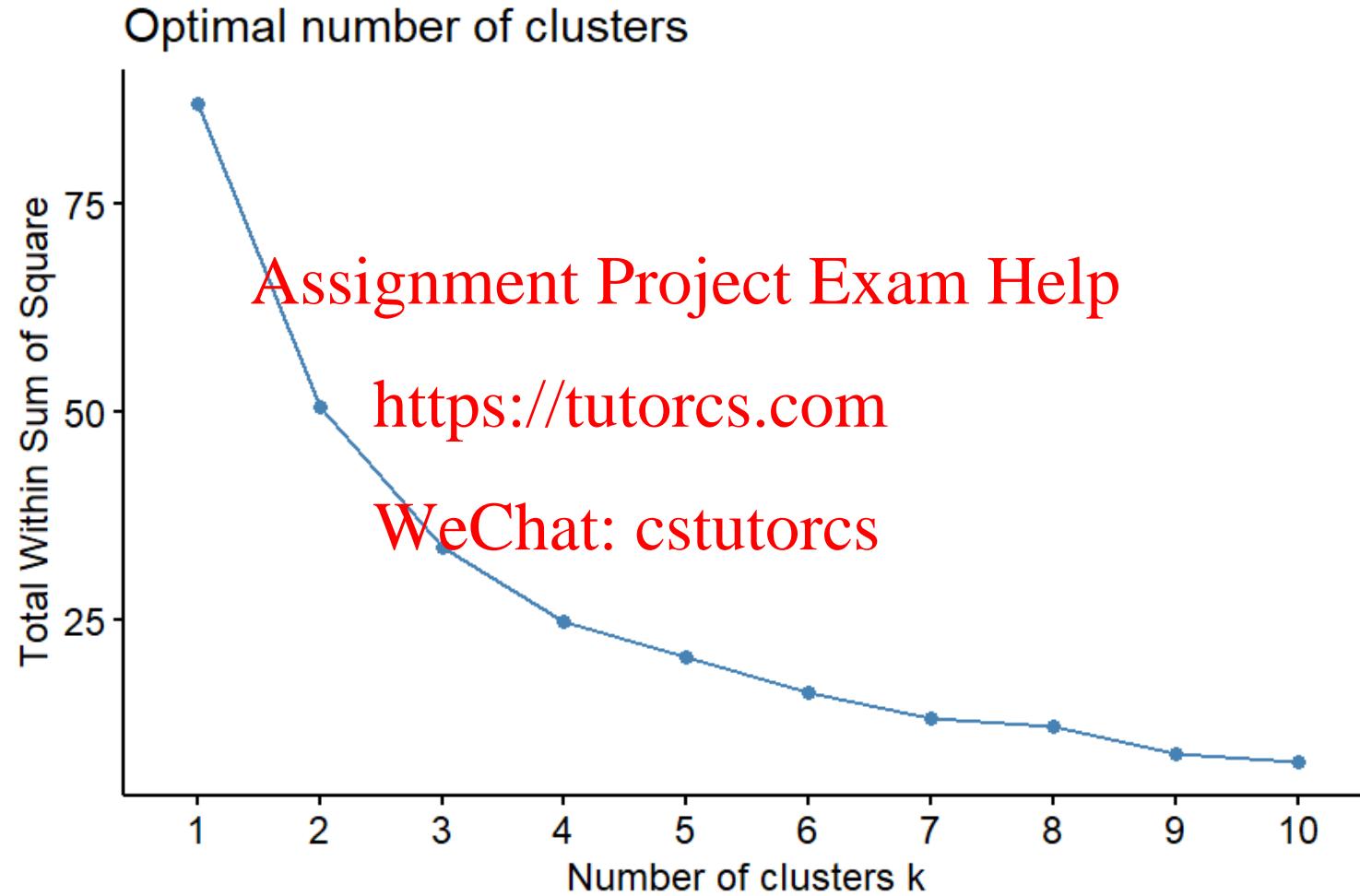
- We could do this by hand, but it is a bit of a pain
- Fortunately there's a function in the `factoextra1` package to help us!

```
library(factoextra)
#specify the dataframe, method (kmeans) and within sum of squares (wss)
fviz_nbclust(ktc_numeric, kmeans, method = "wss")
```

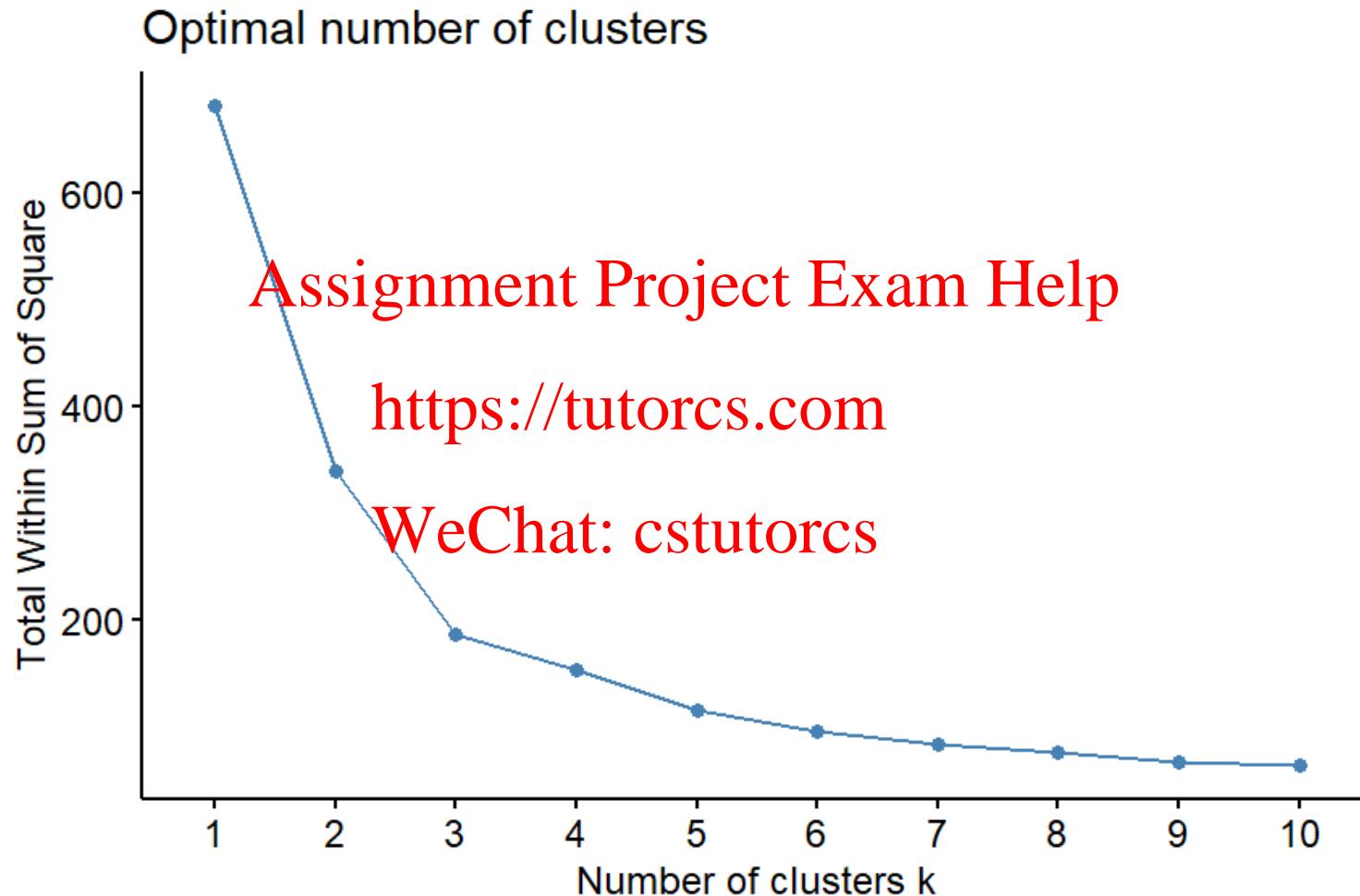
[Assignment Project Exam Help
https://tutorcs.com](https://tutorcs.com)

WeChat: cstutorcs

Elbow plot



A sharper elbow



Summary

- Unlike hierarchical clustering, kmeans requires us to specify the number of clusters directly
- kmeans is also iterative, which can mean we are less likely to get strange shaped clusters (see the snake like cluster from the previous lecture)
- Choosing k can be difficult, but we can use the Elbow plot to help to understand what number is best.

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs



This work is licensed under a [Creative Commons Attribution-ShareAlike 4.0 International License](#).

Assignment Project Exam Help

<https://tutorcs.com>

WeChat: cstutorcs

Lecturer: Lauren Kennedy

Department of Econometrics and Business Statistics

✉ lauren.kennedy1@monash.edu

CALENDAR Week 8

