

K-Means Problem

Jack Reddan

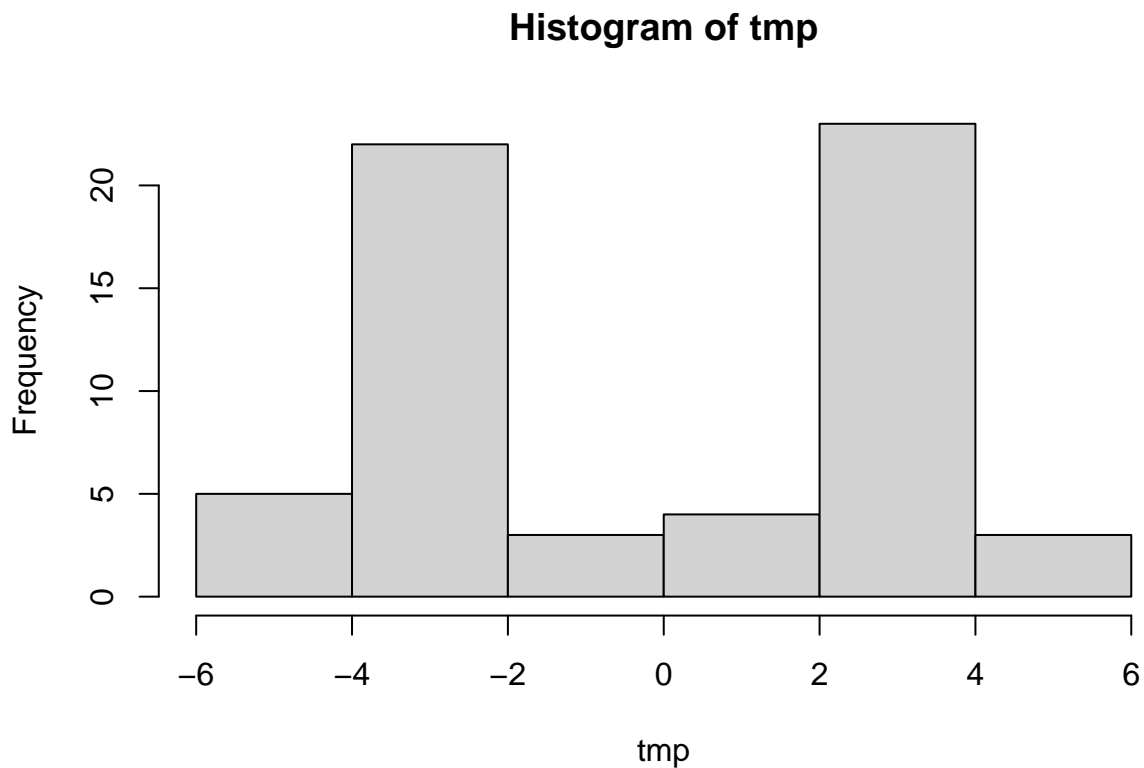
10/21/2021

Try K-Means clustering

Generate fake data and explore how the method works.

Generate example data

```
tmp <- c(rnorm(30,-3), rnorm(30,3))  
hist(tmp)
```



Generate multidimensional example data

```
x <- cbind(x = tmp, y = rev(tmp))  
plot(x)
```


[Q] What component of your results object details:

Cluster size

```
clusters$size
```

```
## [1] 30 30
```

Cluster assignment

```
clusters$cluster
```

[illegible]

Cluster center

```
clusters$centers
```

```
##           x           y
## 1 -3.110046  2.835387
## 2  2.835387 -3.110046
```

Plot x colored by the kmeans cluster centers as blue points

Load ggplot2

```
library(ggplot2)
```

Convert matrices to be used in ggplot to data frames.

```
df <- data.frame(x)
centroids <- data.frame(clusters$centers)
```

Plot the original data colored by kmeans clusters and add blue centroids. IBM's colorblind palette is used.

```
ggplot(data = df) +
  aes(x = x, y = y, color = factor(clusters$cluster)) +
  geom_point() +
  scale_color_manual(values = c("#785EF0", "#FE6100"), name = "Cluster") +
  geom_point(data = centroids, aes(x = x, y = y), color = "#648FFF", shape = 8)
```

