

Clustering (2020 May)

Hwang Seong-Yun

2022 9 15

SOM cluster

reference1 : <https://data-make.tistory.com/91> (<https://data-make.tistory.com/91>)

reference2 : <https://www.statmethods.net/advstats/cluster.html>
(<https://www.statmethods.net/advstats/cluster.html>)

```
water <- read.csv("C:/Users/HSY/Desktop/영산강 수질악화 관련 데이터 정리_결과 포함(220915)/월별 평균 자료/2020년 5월.csv",
sep=",", header=T)
water_name <- water[,1]
water <- water[,-1]
rownames(water) <- water_name
```

Distance matrix

```
water_scale <- scale(water)
d <- dist(water_scale, method="euclidean")
as.matrix(d)
```

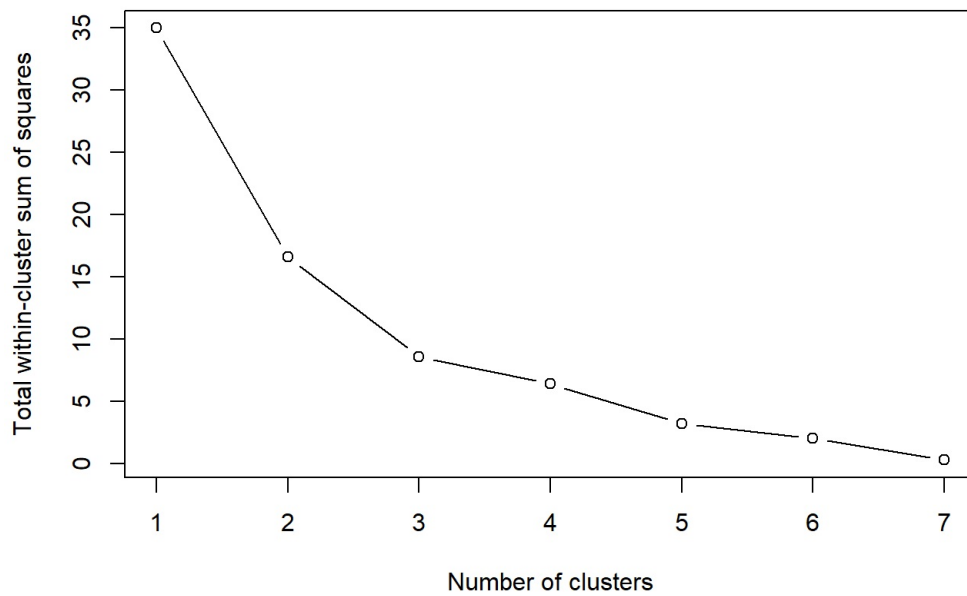
```
##           우치   광주1   방류수   광주천2   광주2   광주3   황룡강5
## 우치      0.0000000  1.872549  3.011595  2.133140  3.866692  4.468859  0.7396523
## 광주1     1.8725486  0.000000  3.905016  1.862021  2.771558  3.607066  1.6989452
## 방류수     3.0115947  3.905016  0.000000  3.099345  4.226119  4.198521  3.3760459
## 광주천2    2.1331402  1.862021  3.099345  0.000000  2.889395  3.027282  2.4352696
## 광주2     3.8666918  2.771558  4.226119  2.889395  0.000000  1.859893  3.7776854
## 광주3     4.4688594  3.607066  4.198521  3.027282  1.859893  0.000000  4.6689785
## 황룡강5    0.7396523  1.698945  3.376046  2.435270  3.777685  4.668979  0.0000000
## 광산      3.5215701  2.728719  4.113089  3.094849  1.540072  2.059765  3.5619629
##           광산
## 우치      3.521570
## 광주1     2.728719
## 방류수     4.113089
## 광주천2    3.094849
## 광주2     1.540072
## 광주3     2.059765
## 황룡강5    3.561963
## 광산      0.000000
```

Decide number of clusters

find the optimal number of clusters using Total within-cluster sum of squares

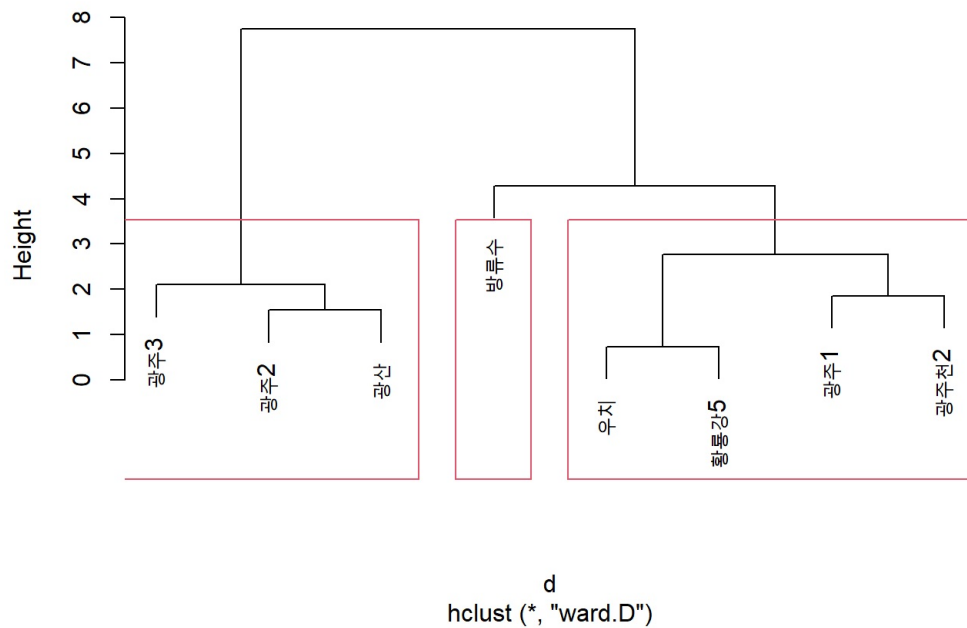
```
tot_withinss <- c()
for (i in 1:7){
  set.seed(1004) # for reproducibility
  kmeans_cluster <- kmeans(water_scale, centers = i, iter.max = 1000)
  tot_withinss[i] <- kmeans_cluster$tot.withinss}
plot(c(1:7), tot_withinss, type="b",
     main="Optimal number of clusters",
     xlab="Number of clusters",
     ylab="Total within-cluster sum of squares")
```

Optimal number of clusters



```
fit <- hclust(d, method="ward.D")
plot(fit)
rect.hclust(fit, k=3)
```

Cluster Dendrogram



SOM cluster

```
library(SOMbrero)
```

```
## Warning: 패키지 'SOMbrero'는 R 버전 4.1.3에서 작성되었습니다
```

```
## 필요한 패키지를 로딩중입니다: igraph
```

```
## Warning: 패키지 'igraph'는 R 버전 4.1.2에서 작성되었습니다
```

```
##
## 다음의 패키지를 부착합니다: 'igraph'
```

```
## The following objects are masked from 'package:stats':
##
##   decompose, spectrum
```

```
## The following object is masked from 'package:base':
##
##   union
```

```
## 필요한 패키지를 로딩중입니다: markdown
```

```
##
```

```
## *****
```

```
##
```

```
##   This is 'SOMbrero' package, v 1.4.1
```

```
##
```

```
## Citation details with citation('SOMbrero')
```

```
##
```

```
## Further information with help(SOMbrero)...
```

```
##
```

```
## Use sombreroGUI() to start the Graphical Interface.
```

```
##
```

```
## *****
```

```
library(kohonen)
```

```
## Warning: 패키지 'kohonen'는 R 버전 4.1.3에서 작성되었습니다
```

Normalization of data

```
water_scale <- data.frame(scale(water))
water_scale_matrix <- as.matrix(water_scale)
```

Training the SOM model

```
som_grid <- somgrid(xdim=1, ydim=3, topo="hexagonal")
som_model1 <- som(water_scale_matrix, grid=som_grid)
som_model2 <- trainSOM(x.data=water_scale, dimension=c(1,3),
                      nb.save=10, maxit=2000, scaling="none",
                      radius.type="letremy")
```

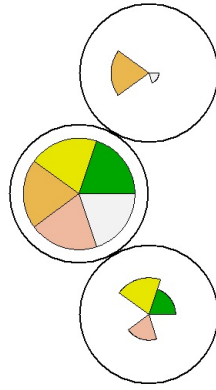
Visualization

```
table(som_model2$clustering)
```

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

```
plot(som_model1, main="feature distribution")
```

feature distribution



```
plot(som_model2, what="obs", type="names", print.title=T, scale=c(1,1))
```

```
## Warning in plot.somRes(som_model2, what = "obs", type = "names", print.title =  
## T, : 'print.title' will be deprecated, please use 'show.names' instead
```

Observations overview

repartition of row.names values

