

goophi - clustering

2022.08.24.

Table of contents

1	Introduction	1
2	Import sample data	1
3	K-means clustering	2
4	K-means clustering without hyperparameters	4
5	Visualize	5

1 Introduction

- 1) 본 문서는 goophi 패키지를 Shiny app에서 사용하는 것을 상정해 작성했습니다.
- 2) 본 문서의 케이스 스타일은 Camel case와 Snake case가 혼용되어 있습니다.
 - Camel case : goophi의 함수명 및 파라미터명
 - Snake case: 유저로부터 받는 입력, shiny app의 server에서 사용(될 것이라고 예상)하는 object명, snake case로 작성된 dependencies의 함수명 등

2 Import sample data

- 1) 전처리가 완료된 샘플데이터를 불러옵니다.
 - NA가 없어야 함
 - string value가 있는 열은 factor로 변환

- 한 열이 모두 같은 값으로 채워져 있을 경우 제외해야 함
- Date type column이 없어야 함
- Outcome 변수는 classification의 경우 factor, regression의 경우 numeric이어야 함 (clustering은 outcome변수를 사용하지 않음)

```
library(goophi)

cleaned_data <- read.csv(file = "~/git/goophi/data/boston_c.csv",
                        stringsAsFactors = TRUE
                        )
cleaned_data$Pcrime <- NULL
cleaned_data$chas <- as.numeric(cleaned_data$chas)

str(cleaned_data)
```

```
'data.frame':  506 obs. of  13 variables:
 $ zn      : num  18 0 0 0 0 0 12.5 12.5 12.5 12.5 ...
 $ indus   : num  2.31 7.07 7.07 2.18 2.18 2.18 7.87 7.87 7.87 7.87 ...
 $ chas    : num  1 1 1 1 1 1 1 1 1 1 ...
 $ nox     : num  0.538 0.469 0.469 0.458 0.458 0.458 0.524 0.524 0.524 0.524 ...
 $ rm      : num  6.58 6.42 7.18 7 7.15 ...
 $ age     : num  65.2 78.9 61.1 45.8 54.2 58.7 66.6 96.1 100 85.9 ...
 $ dis     : num  4.09 4.97 4.97 6.06 6.06 ...
 $ rad     : int  1 2 2 3 3 3 5 5 5 5 ...
 $ tax     : int  296 242 242 222 222 222 311 311 311 311 ...
 $ ptratio : num  15.3 17.8 17.8 18.7 18.7 18.7 15.2 15.2 15.2 15.2 ...
 $ black   : num  397 397 393 395 397 ...
 $ lstat   : num  4.98 9.14 4.03 2.94 5.33 ...
 $ medv    : num  24 21.6 34.7 33.4 36.2 28.7 22.9 27.1 16.5 18.9 ...
```

3 K-means clustering

```
# user input
max_k <- "15" # k = 2:max_k, <= number of columns
n_start <- "25" # attempts 25 initial configurations, <= 175
iter_max <- "10" # <= 5000
n_boot <- "100" # Used only for determining the number of clusters using gap statistic
algorithm = "Hartigan-Wong" ## "Hartigan-Wong", "Lloyd", "Forgy", "MacQueen"
select_optimal <- "silhouette" # "silhouette", "gap_stat" // there's no mathematical definit.
```



```
[1] 2873179 2868624
(between_SS / total_SS = 70.3 %)
```

Available components:

```
[1] "cluster"      "centers"      "totss"        "withinss"     "tot.withinss"
[6] "betweenss"    "size"         "iter"         "ifault"
```

4 K-means clustering without hyperparameters

```
# K-means clustering
km_model <- goophi::kMeansClustering(data = cleaned_data)

km_model$result
```

K-means clustering with 2 clusters of sizes 137, 369

Cluster means:

	zn	indus	chas	nox	rm	age	dis	rad
1	0.00000	18.451825	1.058394	0.6701022	6.006212	89.96788	2.054470	23.270073
2	15.58266	8.420894	1.073171	0.5118474	6.388005	60.63225	4.441272	4.455285

	tax	ptratio	black	lstat	medv
1	667.6423	20.19635	291.0391	18.67453	16.27226
2	311.9268	17.80921	381.0426	10.41745	24.85718

Clustering vector:

```
[1] 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
[38] 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
[75] 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
[112] 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
[149] 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
[186] 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
[223] 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
[260] 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
[297] 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
[334] 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 1 1 1 1 1 1 1 1 1 1
[371] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
[408] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
[445] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
```

```
[482] 1 1 1 1 1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 2 2 2 2 2 2 2
```

Within cluster sum of squares by cluster:

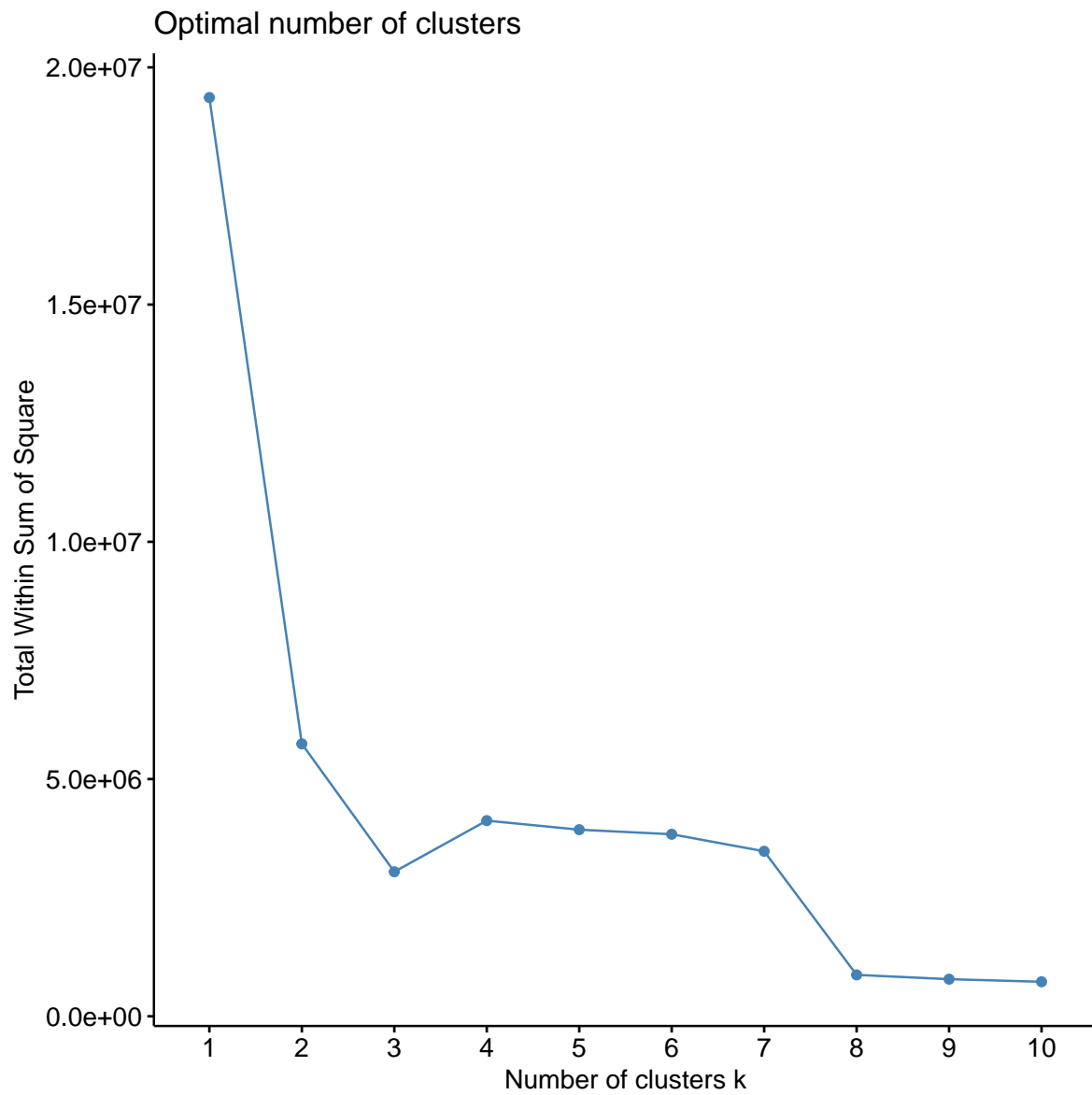
```
[1] 2873179 2868624  
(between_SS / total_SS = 70.3 %)
```

Available components:

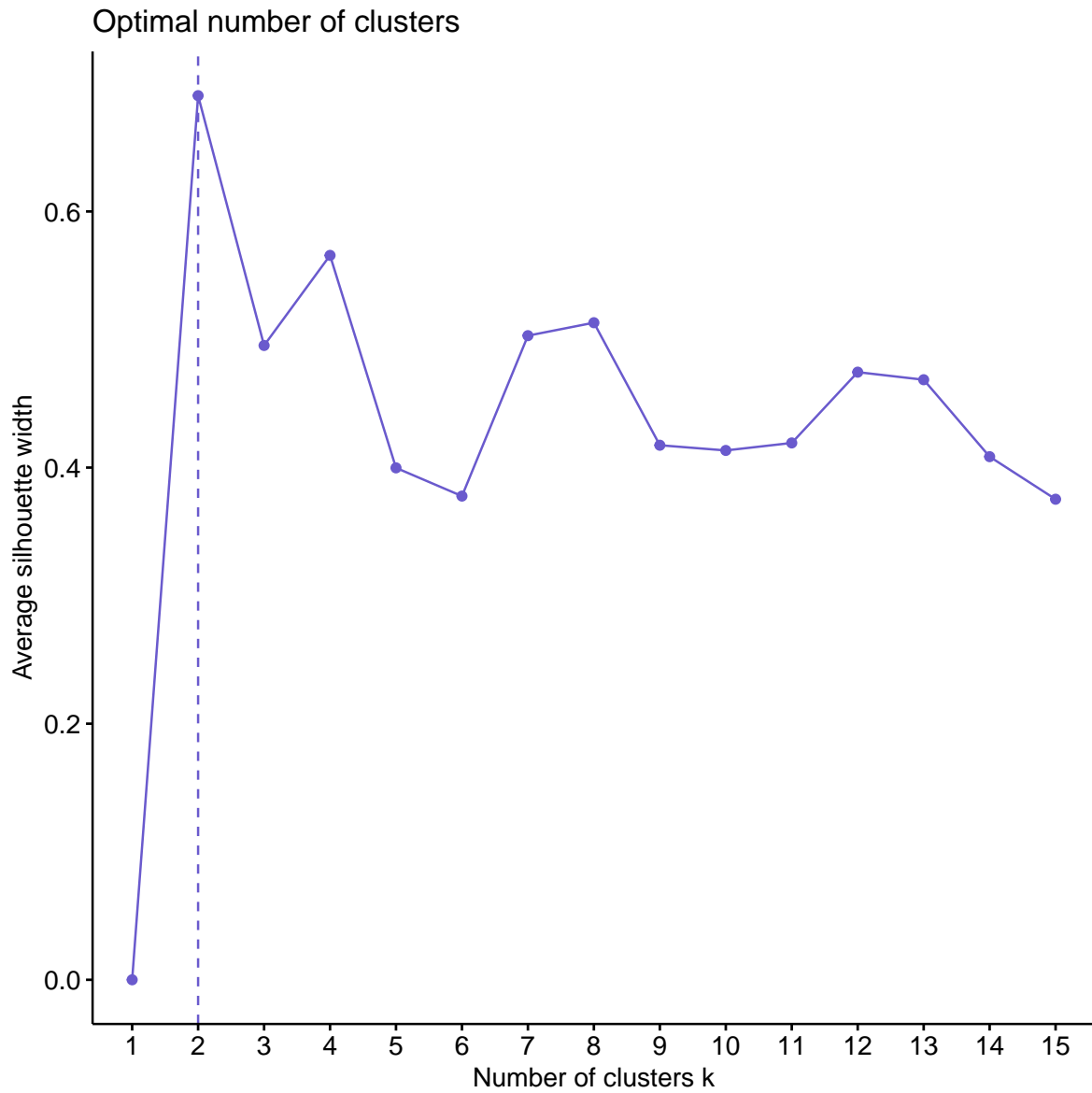
```
[1] "cluster"      "centers"      "totss"        "withinss"     "tot.withinss"  
[6] "betweenss"    "size"         "iter"         "ifault"
```

5 Visualize

```
km_model$elbowPlot
```



```
km_model$optimalK
```



```
km_model$clustVis
```

