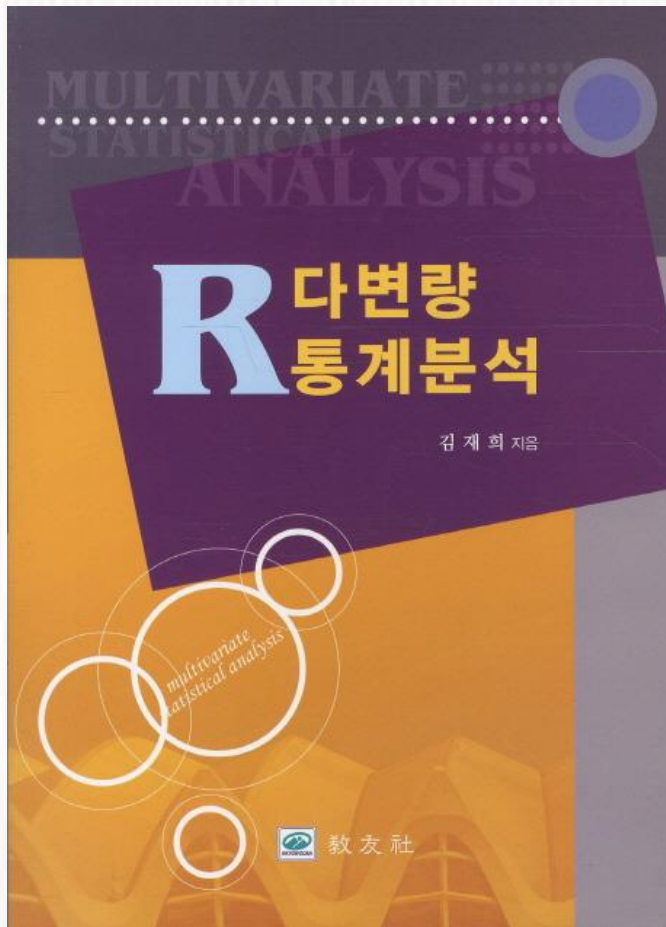




팀	이름	발표일	팀	이름	발표일
19B	김의영	월요일	18	권용재	목요일
19A	정다솜	월요일	18	김나연	
16B	김해수	월요일	16A	허정인	목요일
15A	최효인	월요일	15B	김동원	목요일
14	김사곤	월요일	13	이돈녕	목요일
14	김남현		13	박기완	
12	최민기	월요일	11	강두훈	목요일
12	이경주		11	정재정	
10	김다인	월요일	9	허정민	목요일
10	이다은		9	오제윤	
8	김예은	월요일	7	한상혁	목요일
8	김병관		7	홍관표	
5	최재승	월요일	6	윤다정	목요일
5	유시은		6	이제은	
4	남기수	월요일	3	장연지	목요일
4	박태현		3	오세규	
1	강지애	월요일	2	김재문	목요일
1	이혜린		2	이재우	

s)

군집분석: Clustering Analysis



1. 군집분석 소개
2. 위계적 군집분석 방법
3. 분할 군집분석 방법

1. 군집분석(Cluster Analysis) 소개

- 개체의 특징을 나타내는 데이터로부터 구조를 찾아내고, 통계적 특성이 서로 다른 군집으로 분리할 수 있는지 알아내는 방법
 - 개체 간의 유사성을 측정하여 **유사한 것들끼리** 분류하는 방법
 - 군집의 개수나 구조에 대한 가정없이 다변량 데이터로부터 거리 (유사성) 기준에 의해 **자발적인 군집화**를 유도
 - 개체가 속할 군집이 미리 정해져 있지 않기 때문에 비지도학습에 해당 (Unsupervised Learning) cf. Logistic regression
- 예제
 - 구매성향/인구통계학 속성에 따른 고객세분화
 - 심리실험결과에 의거한 집단분류
 - 사회경제활동지표를 근거로 한 계급분류

1. 군집분석(Cluster Analysis) 소개

- 유사성(similarity) 척도: 거리

Name	Formula
Euclidean metric	$d_E(\mathbf{x}_i, \mathbf{x}_j) = \{\sum_g w_g (x_{gi} - x_{gj})^2\}^{1/2}$
Unstandardized	$w_g = 1$
Standardized by s.d. (Karl Pearson distance)	$w_g = 1/s_g^2$
Standardized by range	$w_g = 1/R_g^2$
Mahalanobis metric	$d_{Ml}(\mathbf{x}_i, \mathbf{x}_j) = \{(\mathbf{x}_i - \mathbf{x}_j)S^{-1}(\mathbf{x}_i - \mathbf{x}_j)'\}^{1/2}$ $= \{\sum_g \sum_{g'} s_{gg'}^{-1} (x_{gi} - x_{gj})(x_{g'i} - x_{g'j})\}^{1/2}$ where $S = (s_{gg'})$ is any $G \times G$ positive definite matrix, usually the sample covariance matrix of the variables. When the matrix is the identity, this reduces to the unstandardized Euclidean distance.
Manhattan metric	$d_{Mn}(\mathbf{x}_i, \mathbf{x}_j) = \sum_g w_g x_{gi} - x_{gj} $
Minkowski metric	$d_{Mk}(\mathbf{x}_i, \mathbf{x}_j) = \{\sum_g w_g x_{gi} - x_{gj} ^\lambda\}^{1/\lambda}, \lambda \geq 1.$ $\lambda = 1$: Manhattan distance $\lambda = 2$: Euclidean distance
Canberra metric	$d_C(\mathbf{x}_i, \mathbf{x}_j) = \sum_g \frac{ x_{gi} - x_{gj} }{(x_{gi} + x_{gj})}$
One minus Pearson correlation	$d_{corr}(\mathbf{x}_i, \mathbf{x}_j) = 1 - \frac{\sum_g (x_{gi} - \bar{x}_{.i})(x_{gj} - \bar{x}_{.j})}{\{\sum_g (x_{gi} - \bar{x}_{.i})^2\}^{1/2} \{\sum_g (x_{gj} - \bar{x}_{.j})^2\}^{1/2}}$

The formulae refer to distances between observations (arrays).

1. 군집분석(Cluster Analysis) 소개

- 군집방법의 유형

- 위계적 군집(Hierarchical clustering)

- ✓ 응집형(Agglomerative, "**Bottom-up**"): 각 개체를 군집으로 정의하고 가까운 것들부터 묶어 나감
 - 단일(최단) 연결법(Single Linkage Method): **두 군집의** 개체 간 거리 중 가장 짧은 거리를 기준
 - 완전(최장) 연결법(Complete Linkage Method): 두 군집의 개체 간 거리 중 가장 긴 거리를 기준
 - 평균 연결법(Average Linkage Method): 두 군집의 모든 개체 간 거리의 평균을 기준
 - 중심 연결법(Centroid Linkage Method): 두 군집의 중심(Centroid) 간 거리를 기준
- ✓ 분리형(Divisive, "**Top-down**"): 하나의 군집에서 출발하여 세분화
 - DIANA 방법(DIANA Method): 모든 세분화 경우를 고려

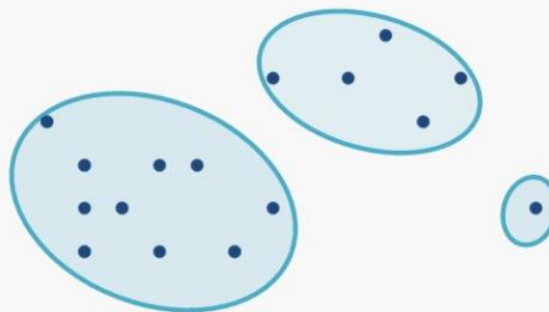
1. 군집분석(Cluster Analysis) 소개

- 군집방법의 유형
 - 비위계적 군집/분할 군집(Non-hierarchical/Partitional clustering)
 - ✓ K-중심 군집(K-Centroid Clustering)
 - K-평균 군집(K-means clustering)
 - K-중앙값 군집(K-median clustering)
 - K-메도이드 군집(K-medoid clustering)

위계적 군집

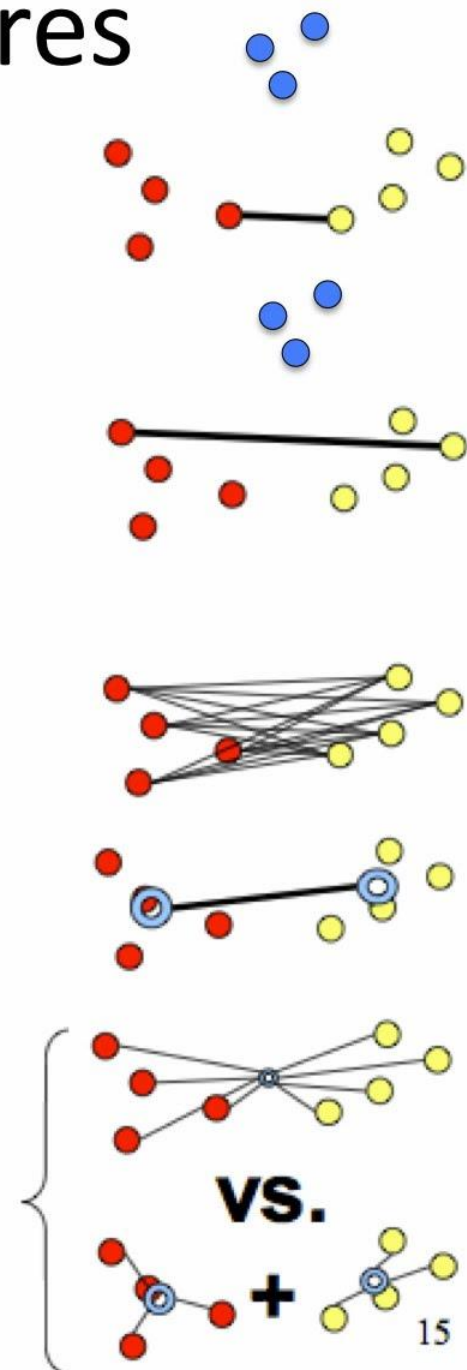


비위계적/분할 군집



Cluster distance measures

- **Single link:** $D(c_1, c_2) = \min_{x_1 \in c_1, x_2 \in c_2} D(x_1, x_2)$
 - distance between closest elements in clusters
 - produces long chains $a \rightarrow b \rightarrow c \rightarrow \dots \rightarrow z$
- **Complete link:** $D(c_1, c_2) = \max_{x_1 \in c_1, x_2 \in c_2} D(x_1, x_2)$
 - distance between farthest elements in clusters
 - forces "spherical" clusters with consistent "diameter"
- **Average link:** $D(c_1, c_2) = \frac{1}{|c_1|} \frac{1}{|c_2|} \sum_{x_1 \in c_1} \sum_{x_2 \in c_2} D(x_1, x_2)$
 - average of all pairwise distances
 - less affected by outliers
- **Centroids:** $D(c_1, c_2) = D\left(\left(\frac{1}{|c_1|} \sum_{x \in c_1} \vec{x}\right), \left(\frac{1}{|c_2|} \sum_{x \in c_2} \vec{x}\right)\right)$
 - distance between centroids (means) of two clusters
- **Ward's method:** $TD_{c_1 \cup c_2} = \sum_{x \in c_1 \cup c_2} D(x, \mu_{c_1 \cup c_2})^2$
 - consider joining two clusters, how does it change the total distance (TD) from centroids?





1) Single Linkage Method

```
# 1) single linkage method
set.seed(29) #to make the result always the same
x <- c(rnorm(7, mean = 0, sd = 1), rnorm(7, mean = 5, sd = 1))
y <- c(rnorm(7, mean = 2, sd = 1), rnorm(7, mean = 7, sd = 1))
plot(x, y, main = "Single Linkage Method")
distance <- max(dist(cbind(x, y)))
distance
for(i in 1:7){ # for each element of cluster i
  for(j in 8:14){ #for each element of cluster j
    shortest <- sqrt((x[i]-x[j])^2 + (y[i]-y[j])^2)
    if(distance > shortest){
      distance <- shortest
      shortest_i <- i #position of chosen element from cluster i
      shortest_j <- j #position of chosen element from cluster j
    }
  }
}
distance
shortest_i
shortest_j
lines(x[c(shortest_i, shortest_j)], y[c(shortest_i, shortest_j)],
col = "red")
```




2) Complete Linkage Method

```
# 2) complete linkage method
# set.seed(29) #to make the result always the same
# x <- c(rnorm(7, mean = 0, sd = 1), rnorm(7, mean = 5, sd = 1))
# y <- c(rnorm(7, mean = 2, sd = 1), rnorm(7, mean = 7, sd = 1))
plot(x, y, main = "Complete Linkage Method")
distance <- min(dist(cbind(x, y)))
distance
for(i in 1:7){ # for each element of cluster i
  for(j in 8:14){ #for each element of cluster j
    longest <- sqrt((x[i]-x[j])^2 + (y[i]-y[j])^2)
    if(distance < longest){
      distance <- longest
      longest_i <- i #position of chosen element from cluster i
      longest_j <- j #position of chosen element from cluster j
    }
  }
}
distance
longest_i
longest_j
lines(x[c(longest_i, longest_j)], y[c(longest_i, longest_j)], col
= "blue")
```



3) Average Linkage Method

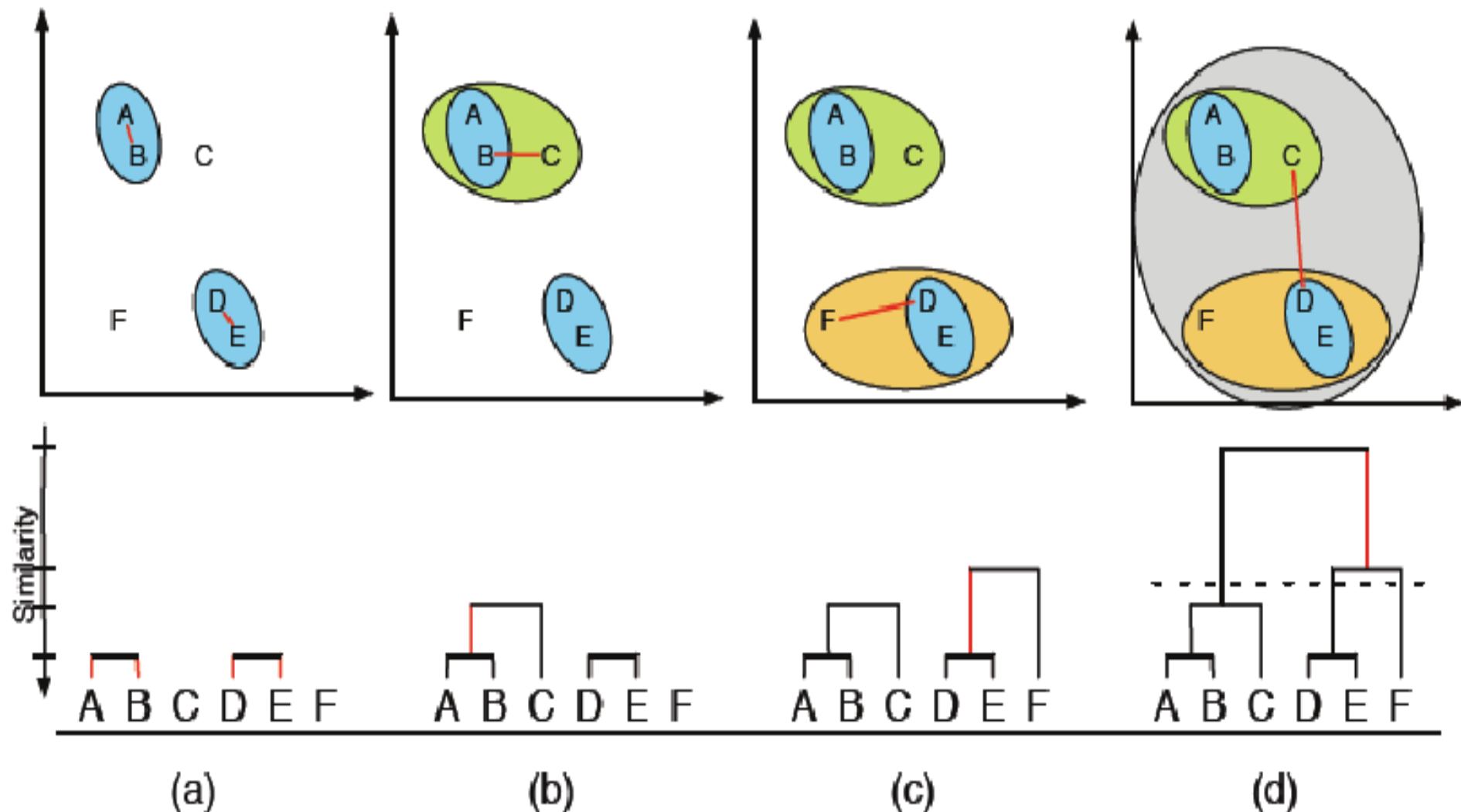
```
# 3) average linkage method
# set.seed(29) #to make the result always the same
# x <- c(rnorm(7, mean = 0, sd = 1), rnorm(7, mean = 5, sd = 1))
# y <- c(rnorm(7, mean = 2, sd = 1), rnorm(7, mean = 7, sd = 1))
plot(x, y, main = "Average Linkage Method")
for(i in 1:7){
  for(j in 8:14) lines(x[c(i, j)], y[c(i, j)], col = "gray")
}
```

2. 위계적 군집분석 방법

- 위계적 군집분석 절차

1. 알맞은 속성 선택
2. 데이터 표준화
3. 이상치 선별
4. 거리 계산
5. 군집 알고리즘 선택
6. 군집의 개수 결정
7. 최종결과 획득
8. 분석 결과의 시각화
9. 군집분석 결과의 해석

Example: Hierarchical Agglomerative Clustering



2. 위계적 군집분석 방법

1. 미국 주 별 폭력 범죄율 (USArrests)

- 미국 50개 주에 관한 4개 변수 값을 포함
 - ✓ Murder numeric Murder arrests (per 100,000) in 1973
 - ✓ Assault numeric Assault arrests (per 100,000) in 1973
 - ✓ UrbanPop numeric Percentage of urban population in 1973
 - ✓ Rape numeric Rape arrests (per 100,000) in 1973

```
#자료 준비
USArrests
#자료 표준화
USArrests.scaled <- scale(USArrests) #scale = T (default)
#거리 계산
USArrestsDist <- dist(USArrests.scaled, method = "euclidean")
str(USArrestsDist)
```



2. 위계적 군집분석 방법

#위계적 군집분석 시행하고 Dendrogram으로 확인 (5가지 방법)

```
hc1 <- hclust(USArrestsDist, method="single")  
plot(hc1, hang=-1, main="Single Linkage Method")
```

```
hc2 <- hclust(USArrestsDist, method="complete")  
plot(hc2, hang=-1, main="Complete Linkage Method")
```

```
hc3 <- hclust(USArrestsDist, method="average")  
plot(hc3, hang=-1, main="Average Linkage Method")
```

```
hc4 <- hclust(USArrestsDist, method="centroid")  
plot(hc4, hang=-1, main="Centroid Linkage Method")
```

```
hc5 <- hclust(USArrestsDist, method = "ward.D2")  
plot(hc5, hang = -1, main = "Ward Method")
```

#군집 개수 결정 (참고자료)

```
#install.packages("NbClust")
```

```
library(NbClust)
```

```
devAskNewPage(ask = T) #새화면 나오기 전에 사용자 대기
```

```
nc <- NbClust(USArrests.scaled, distance = "euclidean", min.nc = 2,  
max.nc = 15, method = "average") #각종 추천 알고리즘 결과
```

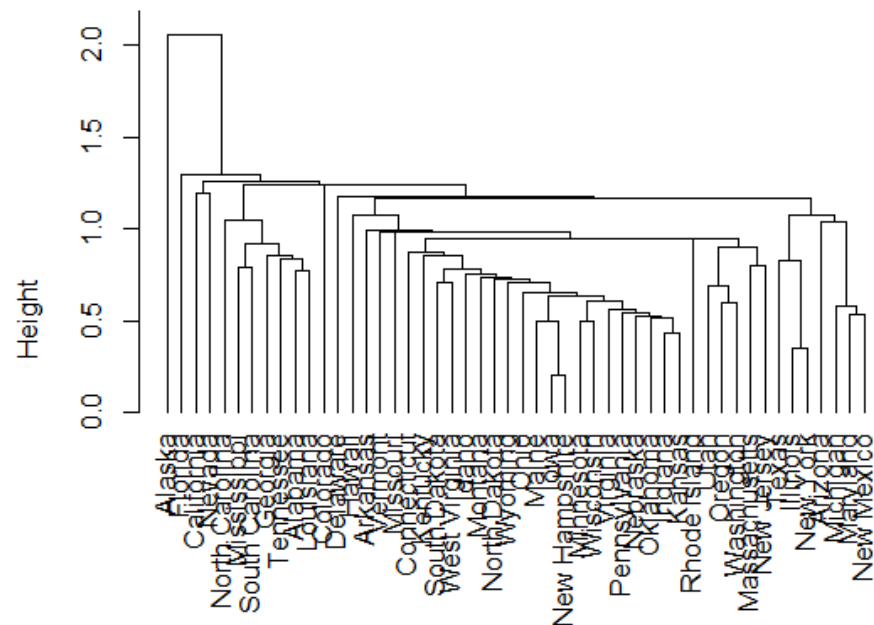
```
devAskNewPage(ask = F)
```

```
par(mfrow = c(1,1)) # 쪼개진 plot 화면 하나로 복원
```

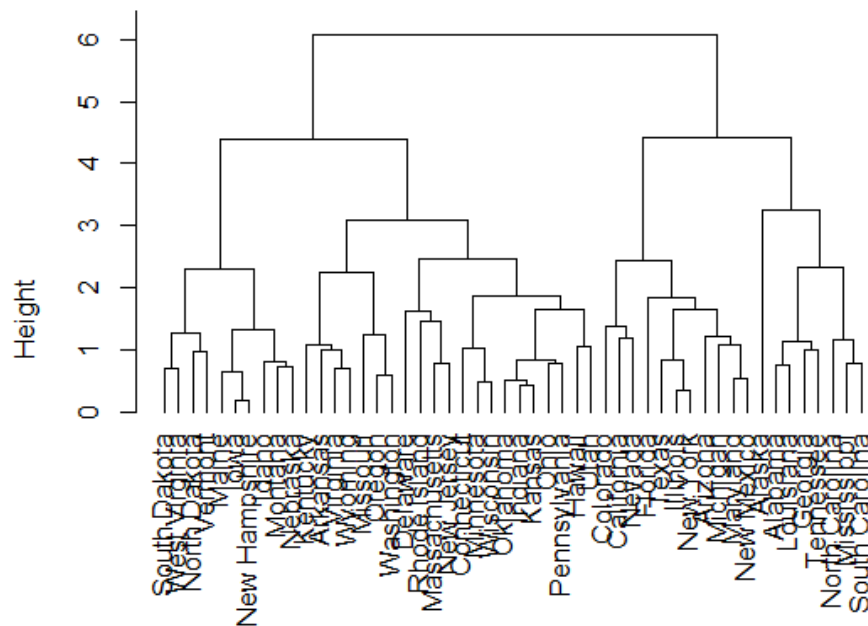
```
table(nc$Best.nc[1,]) #각종 추천 알고리즘 결과 요약
```

2. 위계적 군집분석 방법

Single Linkage Method

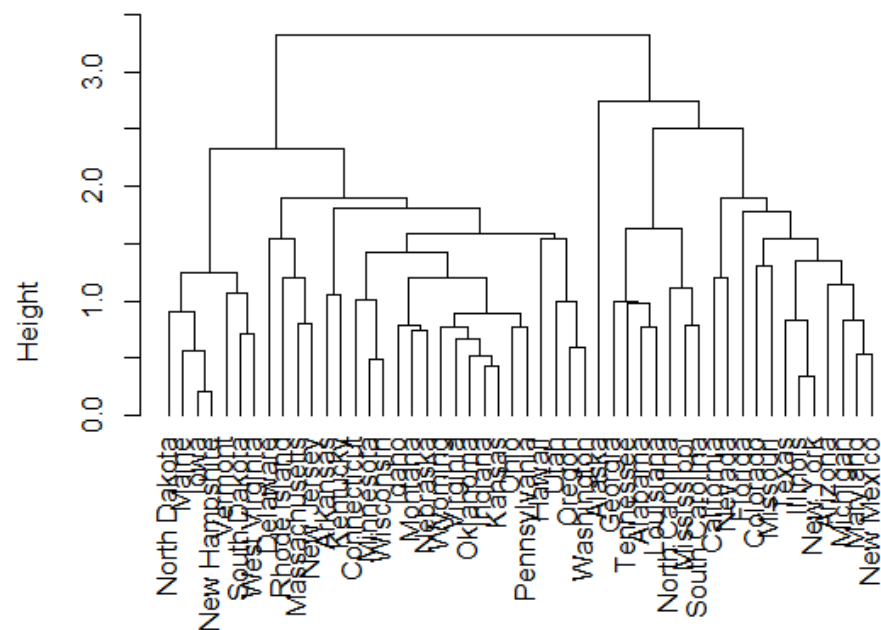


Complete Linkage Method



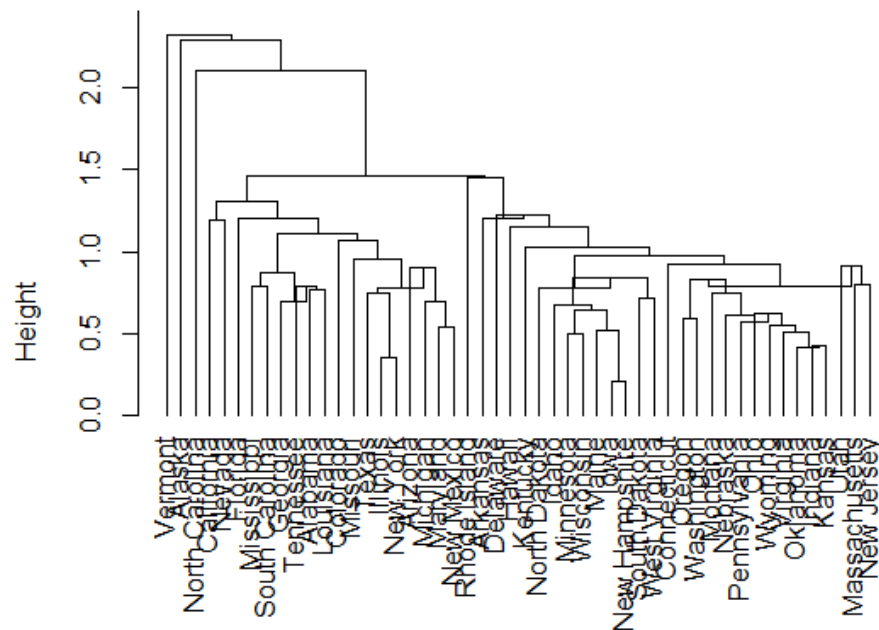
2. 위계적 군집분석 방법

Average Linkage Method



USArrestsDist
hclust(*, "average")

Centroid Linkage Method



USArrestsDist
hclust(*, "centroid")



2. 위계적 군집분석 방법

```
*****
* Among all indices:
* 8 proposed 2 as the best number of clusters
* 2 proposed 3 as the best number of clusters
* 3 proposed 4 as the best number of clusters
* 6 proposed 5 as the best number of clusters
* 1 proposed 8 as the best number of clusters
* 1 proposed 14 as the best number of clusters
* 2 proposed 15 as the best number of clusters

***** Conclusion *****

* According to the majority rule, the best number of clusters is 2

*****
```



2. 위계적 군집분석 방법

```
# 최종결과 획득(군집 지정)
cl.Num <- 2 # number of clusters
hc1.result = cutree(hc1, k=cl.num)
table(hc1.result)

hc2.result = cutree(hc2, k=cl.num)
table(hc2.result)

hc3.result = cutree(hc3, k=cl.num)
table(hc3.result)

hc4.result = cutree(hc4, k=cl.num)
table(hc4.result)
```

```
> table(hc3.result)
hc3.result
 1  2
20 30
```

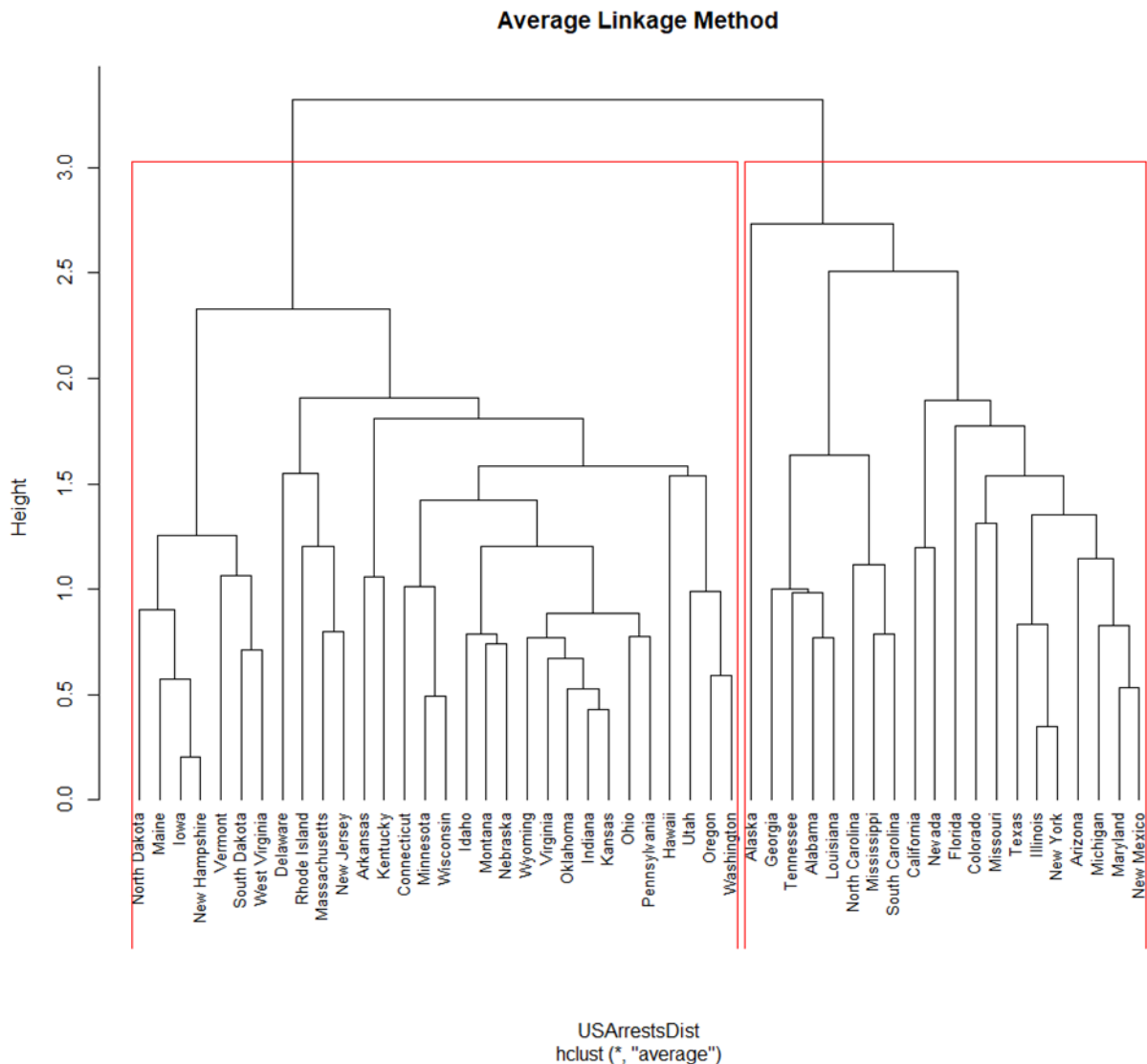
2. 위계적 군집분석 방법

```
#군집 결과 시각화
plot(hc3, hang=-1, cex=.8, main="Average Linkage Method")
rect.hclust(hc3, k=c1.num)

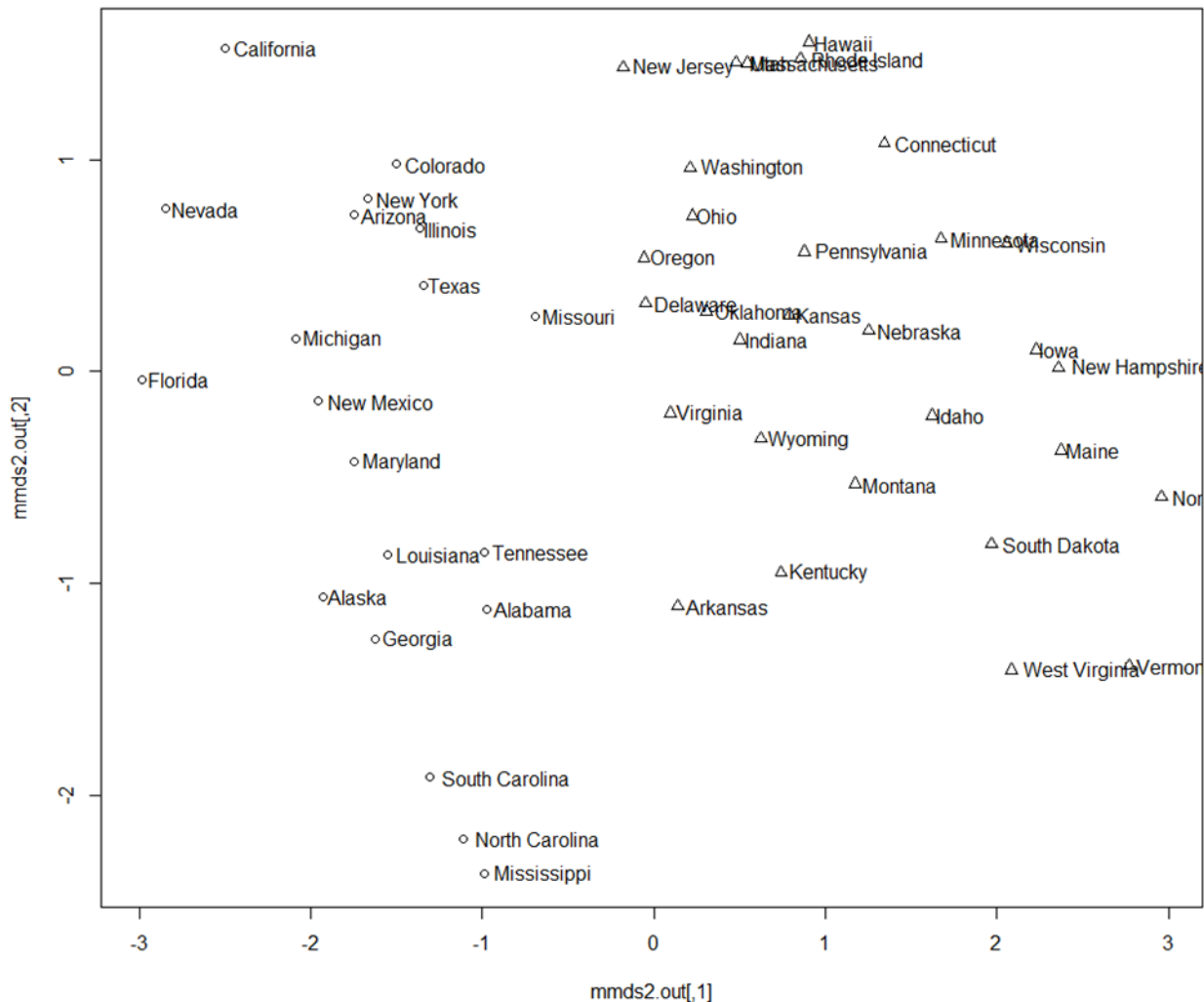
#군집 결과 시각화2 (using MDS)
mmds2.out <- cmdscale(USArrestsDist)
mmds2.out
hc3.result
plot(mmds2.out, pch = hc3.result)
text(mmds2.out, adj = -0.1, rownames(mmds2.out))
plot(mmds2.out, type = "n")
text(mmds2.out, rownames(mmds2.out), col = hc3.result)

#군집의 특징 확인
USArrests[order(USArrests$Assault, decreasing = T),]
USArrests[order(USArrests$UrbanPop, decreasing = T),]
aggregate(USArrests, by=list(cluster=hc3.result), mean)
aggregate(USArrests.scaled, by=list(cluster=hc3.result), mean)
```

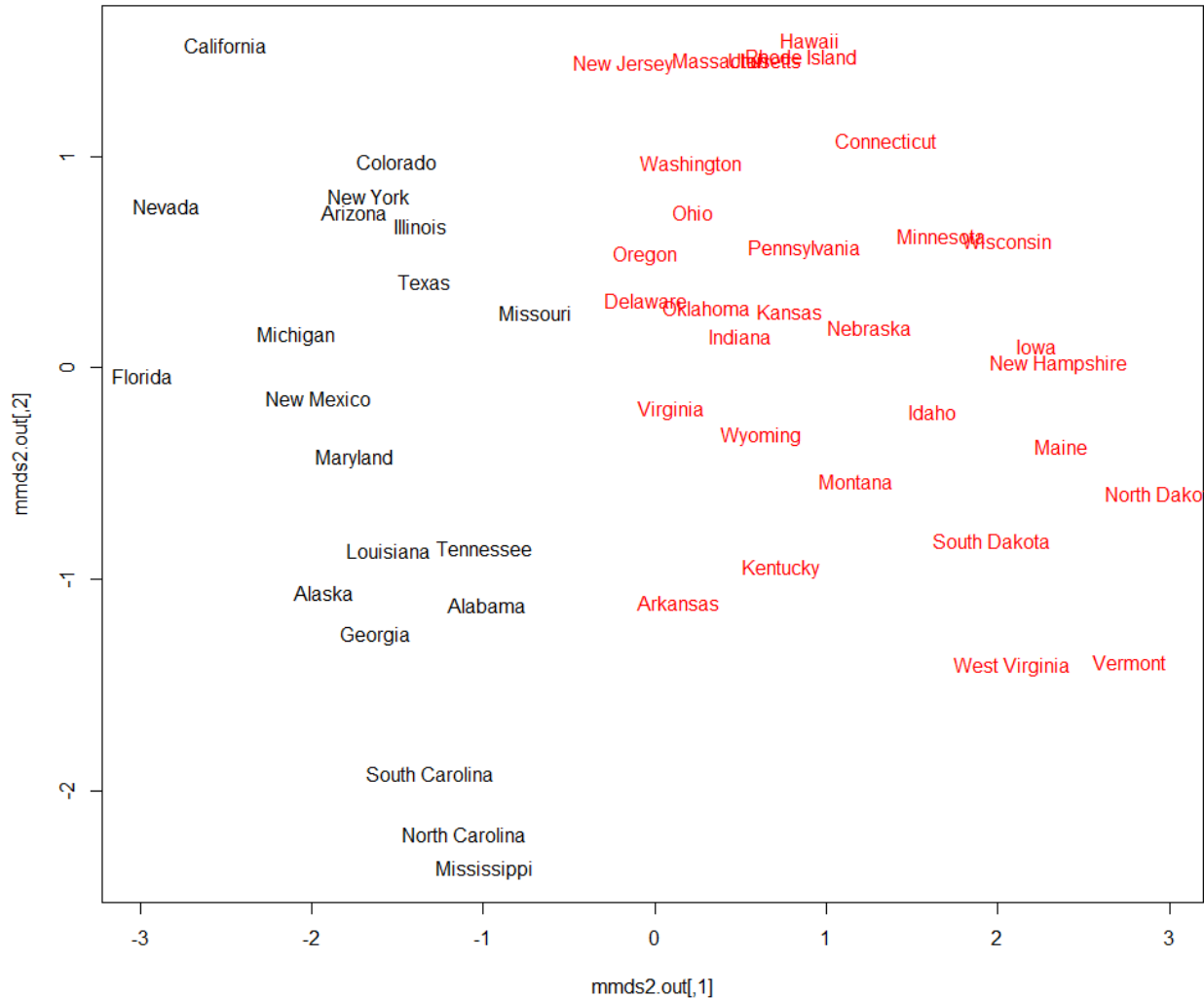
2. 위계적 군집분석 방법



2. 위계적 군집분석 방법



2. 위계적 군집분석 방법



2. 위계적 군집분석 방법

```
> USArrests[order(USArrests$Assault, decreasing = T),]
      Murder Assault UrbanPop Rape
North Carolina  13.0    337      45  16.1
Florida         15.4    335      80  31.9
Maryland        11.3    300      67  27.8
.....
Vermont          2.2     48      32  11.2
Hawaii           5.3     46      83  20.2
North Dakota     0.8     45      44   7.3
> USArrests[order(USArrests$UrbanPop, decreasing = T),]
      Murder Assault UrbanPop Rape
California     9.0    276      91  40.6
New Jersey     7.4    159      89  18.8
Rhode Island    3.4    174      87   8.3
.....
North Dakota    0.8     45      44   7.3
West Virginia   5.7     81      39   9.3
Vermont         2.2     48      32  11.2
> aggregate(USArrests, by=list(cluster=hc3.result), mean)
  cluster Murder Assault UrbanPop Rape
1        1 12.165 255.2500 68.40000 29.16500
2        2  4.870 114.4333 63.63333 15.94333
> aggregate(USArrests.scaled, by=list(cluster=hc3.result), mean)
  cluster Murder Assault UrbanPop Rape
1        1  1.004934  1.0138274  0.1975853  0.8469650
2        2 -0.669956 -0.6758849 -0.1317235 -0.5646433
```



네→너무 건조하다
 넝→너무 가볍다
 념→마찬가지
 넷→군대같음
 넥→스
 넵→적당히 절도있고 깔끔하고
 건조하지도 않고 가볍지도 않고 뭔가
 확실한 인상을 주는 것 같아서 씬 근데
 계속 넵만 하면 성의없어보여서 가끔
 "네"나 "네!", "네~!"도 섞음

2017년 09월 20일 · 5:36 오후

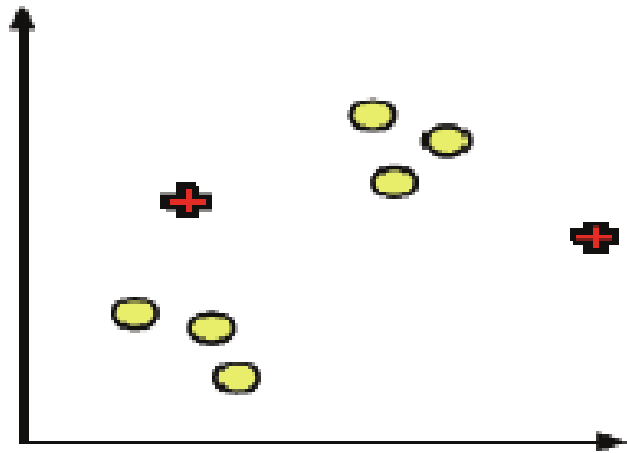
33.3K 리트윗 7,489 마음에 들어요

3. 분할 군집분석 방법

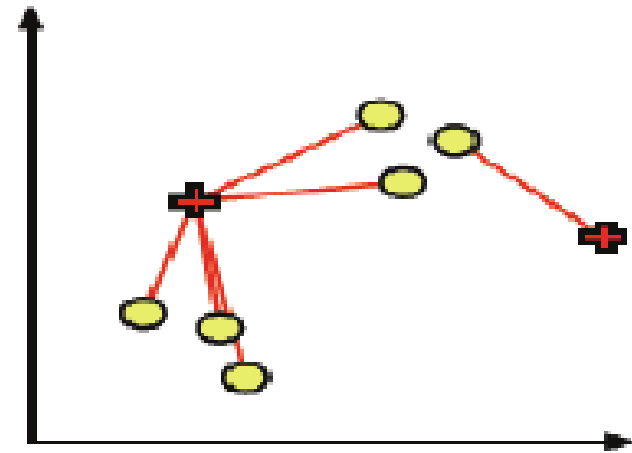
- 분할 군집분석 절차

1. 초기 Seed 역할을 하는 개체 K개를 선택하여 군집의 초기 값으로 활용
 - ✓ 초기 seed는 랜덤하게 선택하거나, 가장 멀리 떨어진 것들로 선택하거나, 최고 거리 이상 갖는 것 등으로 설정 가능
2. 나머지 개체들은 군집의 초기값과 거리를 계산하여 가까운 초기값의 군집에 배치
3. 새로운 군집이 형성되면 초기 seed는 군집 중심(Centroid)으로 대체 (K-means/K-median/K-medoid)
 - ✓ Medoid: most centrally located **point** in the cluster
4. 모든 개체에 대하여 소속 군집 중심과의 거리와 다른 군집 중심과의 거리를 비교하여 소속 군집 중심보다 더 가까운 군집 중심이 있을 경우 군집을 재배치
5. 재배치 후 군집 중심을 재계산 (R에서 default값은 10회)

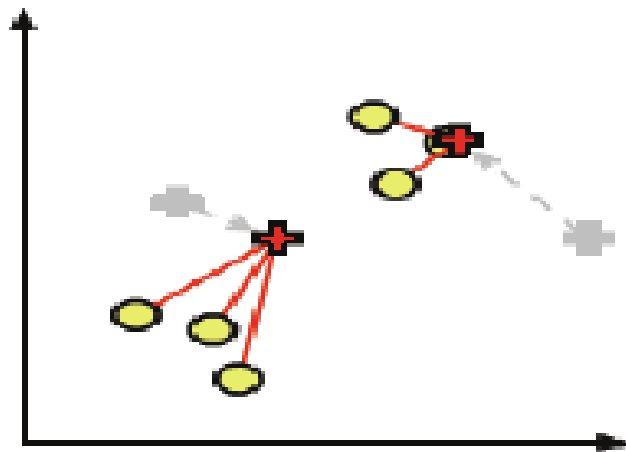
*Example case of K-means clustering



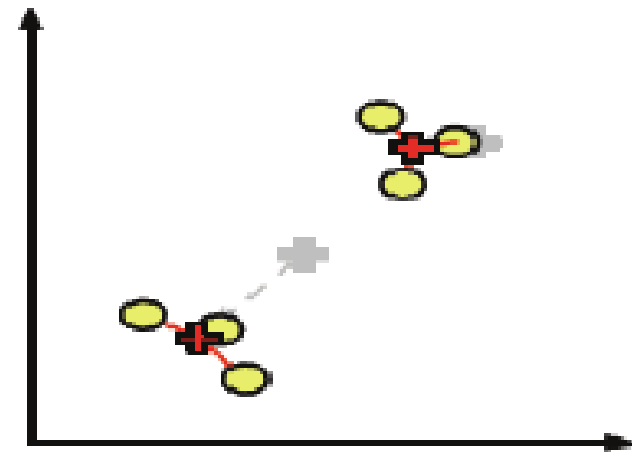
(a)



(b)



(c)



(d)

3. 분할 군집분석 방법

- K-means clustering의 장점
 - 계층적 군집분석에 비해 큰 데이터셋에 사용 가능
 - 관측치가 군집에 **영원히** 할당되는 것이 아니라 최종결과를 개선시키는 방향으로 이동
- K-means clustering의 단점
 - 클러스터의 수를 **미리 지정해** 주어야 함
 - 초기값에 따라 Local optimum에 머물 우려
 - 이상치에 심하게 영향 받음 → 이상치 제거하거나 K-median, K-medoid 사용할 필요
 - Spherical 형태가 아닌 클러스터를 찾는 데는 부적합: 알고리즘이 **중심점으로부터 Spherical 형태로 군집을 구성**

3. 분할 군집분석 방법

1. 탐색적으로 군집의 수(k) 결정

```
#표준화 된 데이터 준비
USArrests.scaled <- scale(USArrests)
USArrestsDist <- dist(USArrests.scaled, method = "euclidean")
#MDS 실시
mmds2.out <- cmdscale(USArrestsDist)
mmds2.out
#그림으로 확인
plot(mmds.out)
plot(mmds.out, type = "n")
text(mmds.out, rownames(mmds.out))
#각종 기준치 활용
library(NbClust)
devAskNewPage(ask = T) #The user is prompted before starting a
new page
nc <- NbClust(USArrests.scaled, distance = "euclidean", min.nc = 2,
max.nc = 15, method = "kmeans")
devAskNewPage(ask = F)
par(mfrow = c(1,1)) # 쪼개진 plot 화면 하나로 복원
table(nc$Best.nc[1,])
```

3. 분할 군집분석 방법

1. 클러스터링 실행 및 결과 확인

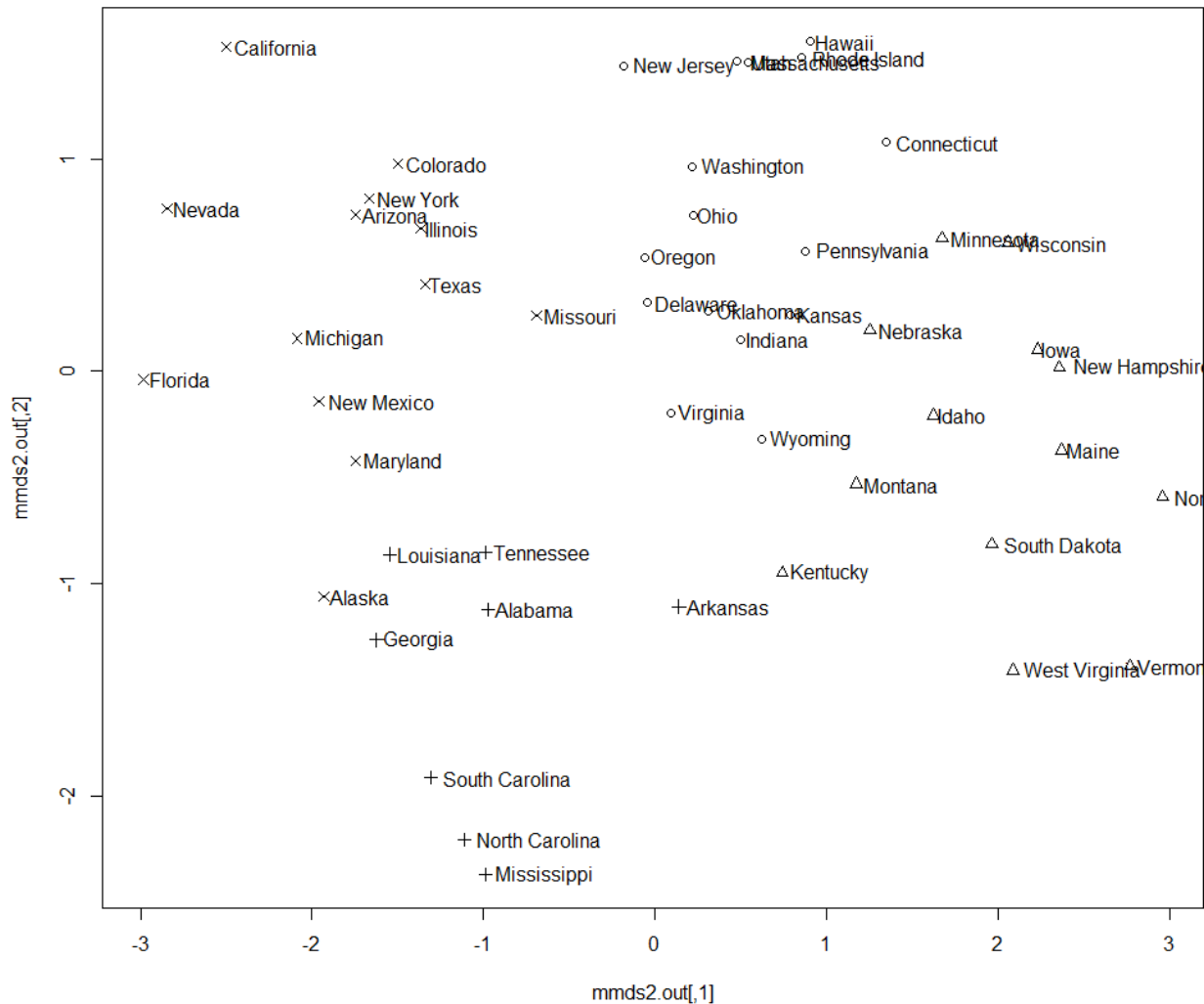
```
#K-means clustering
cl.num <- 4 # number of clusters
pc1 <- kmeans(USArrests.scaled, centers = cl.num, algorithm =
"Hartigan-Wong") #algorithm name can be abbreviated
#algorithm = c("Hartigan-Wong", "Lloyd", "Forgy", "MacQueen")
pc1
table(pc1$cluster)
#Hierarchical Clustering과 Partitional Clustering 비교
cbind(HierClust = hc3.result, PartClust = pc1$cluster)
table(hc3.result, pc1$cluster)

#군집 시각화(using MDS)
#MDS with USArrestsDist
mmds2.out <- cmdscale(USArrestsDist)
mmds2.out
plot(mmds2.out, pch = pc1$cluster)
text(mmds2.out, adj = -0.1, rownames(mmds2.out))
plot(mmds2.out, type = "n")
text(mmds2.out, rownames(mmds2.out), col = pc1$cluster)

#군집의 특징 확인
aggregate(USArrests, by=list(cluster=pc1$cluster), mean)
aggregate(USArrests.scaled, by=list(cluster=pc1$cluster), mean)
pc1$centers
```

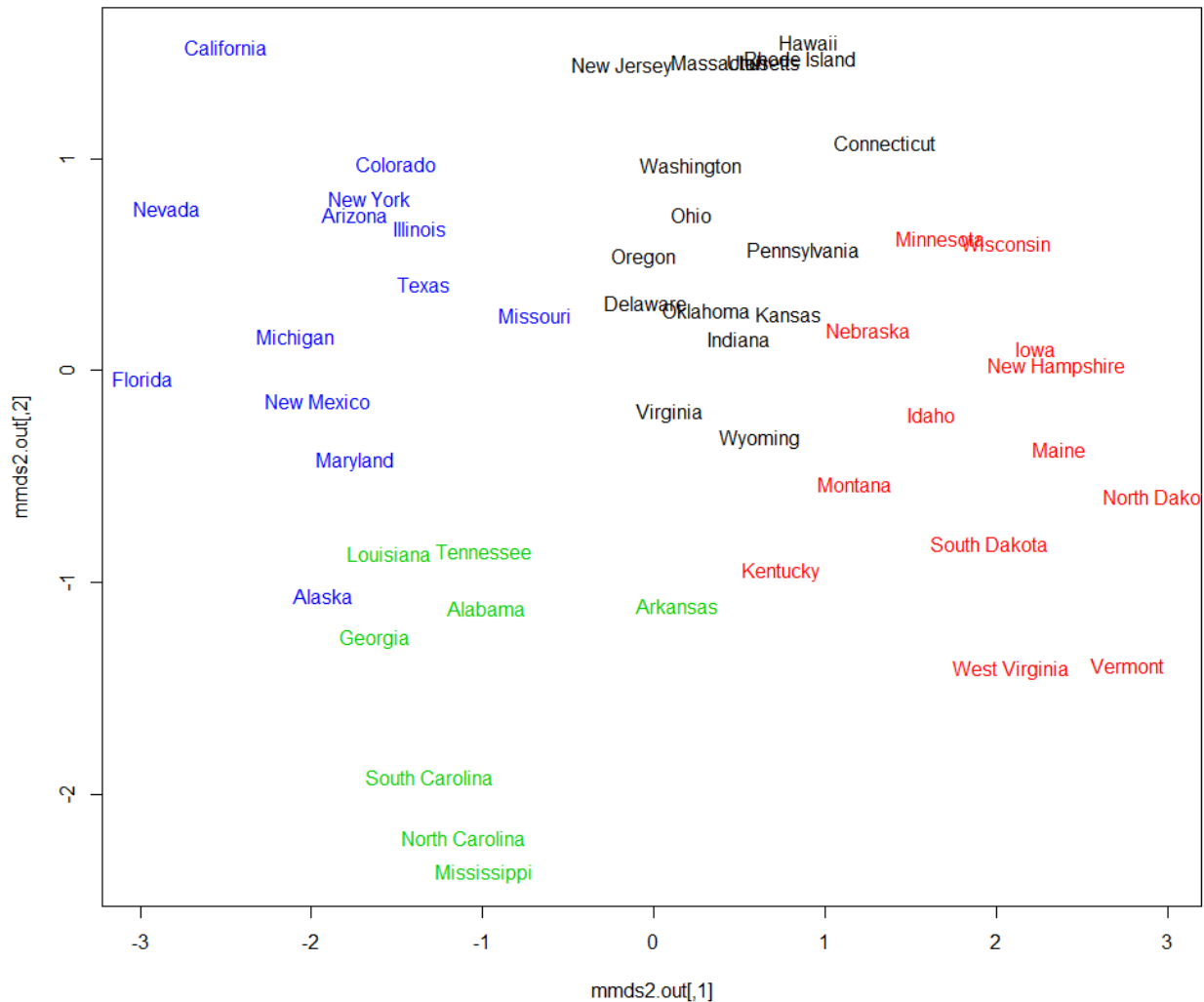



3. 분할 군집분석 방법





3. 분할 군집분석 방법





3. 분할 군집분석 방법

```
> aggregate(USArrests, by=list(cluster=pc1$cluster), mean)
  cluster  Murder  Assault UrbanPop  Rape
1        1  5.65625 138.87500  73.87500 18.78125
2        2  3.60000  78.53846  52.07692 12.17692
3        3 13.93750 243.62500  53.75000 21.41250
4        4 10.81538 257.38462  76.00000 33.19231
> aggregate(USArrests.scaled, by=list(cluster=pc1$cluster), mean)
  cluster  Murder  Assault  UrbanPop  Rape
1        1 -0.4894375 -0.3826001  0.5758298 -0.26165379
2        2 -0.9615407 -1.1066010 -0.9301069 -0.96676331
3        3  1.4118898  0.8743346 -0.8145211  0.01927104
4        4  0.6950701  1.0394414  0.7226370  1.27693964
> pc1$centers
  Murder  Assault  UrbanPop  Rape
1 -0.4894375 -0.3826001  0.5758298 -0.26165379
2 -0.9615407 -1.1066010 -0.9301069 -0.96676331
3  1.4118898  0.8743346 -0.8145211  0.01927104
4  0.6950701  1.0394414  0.7226370  1.27693964
```