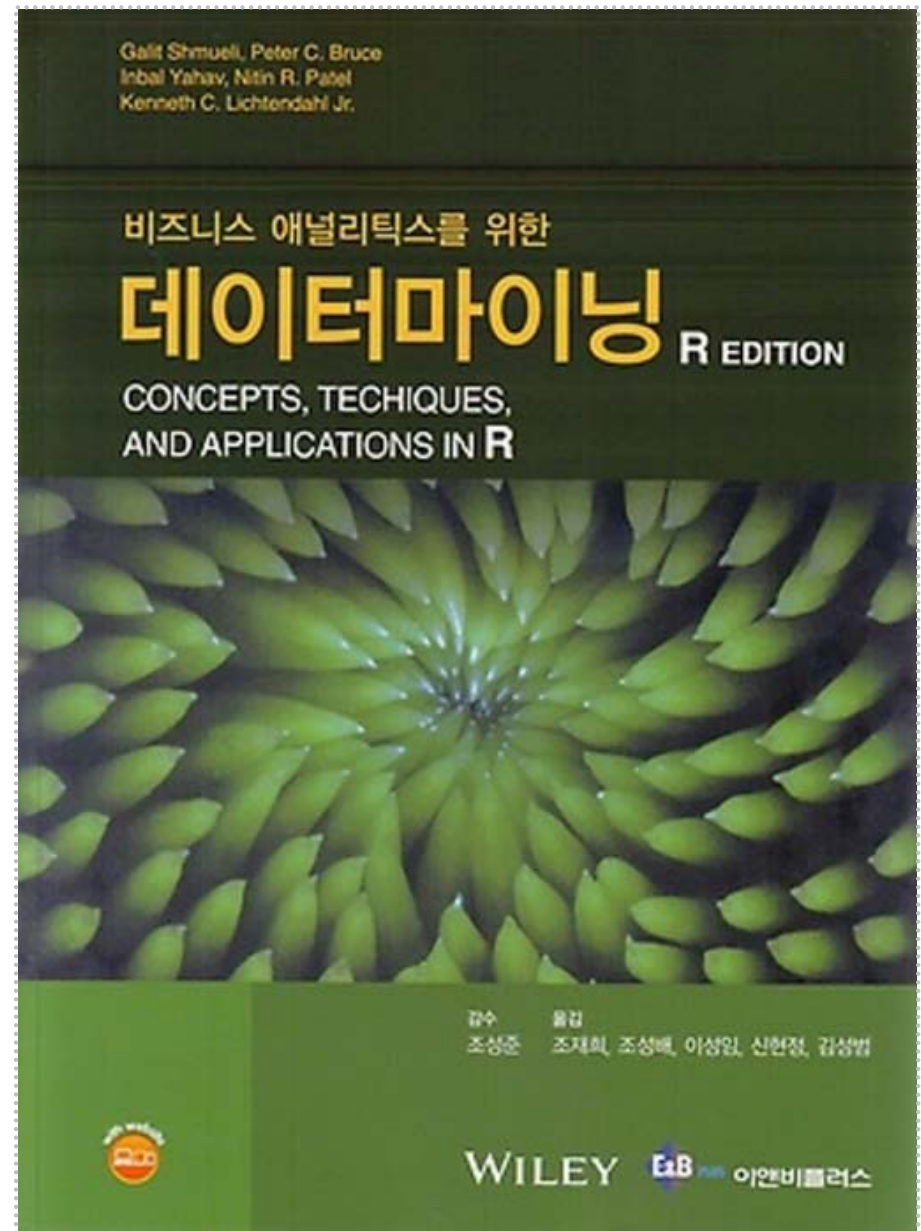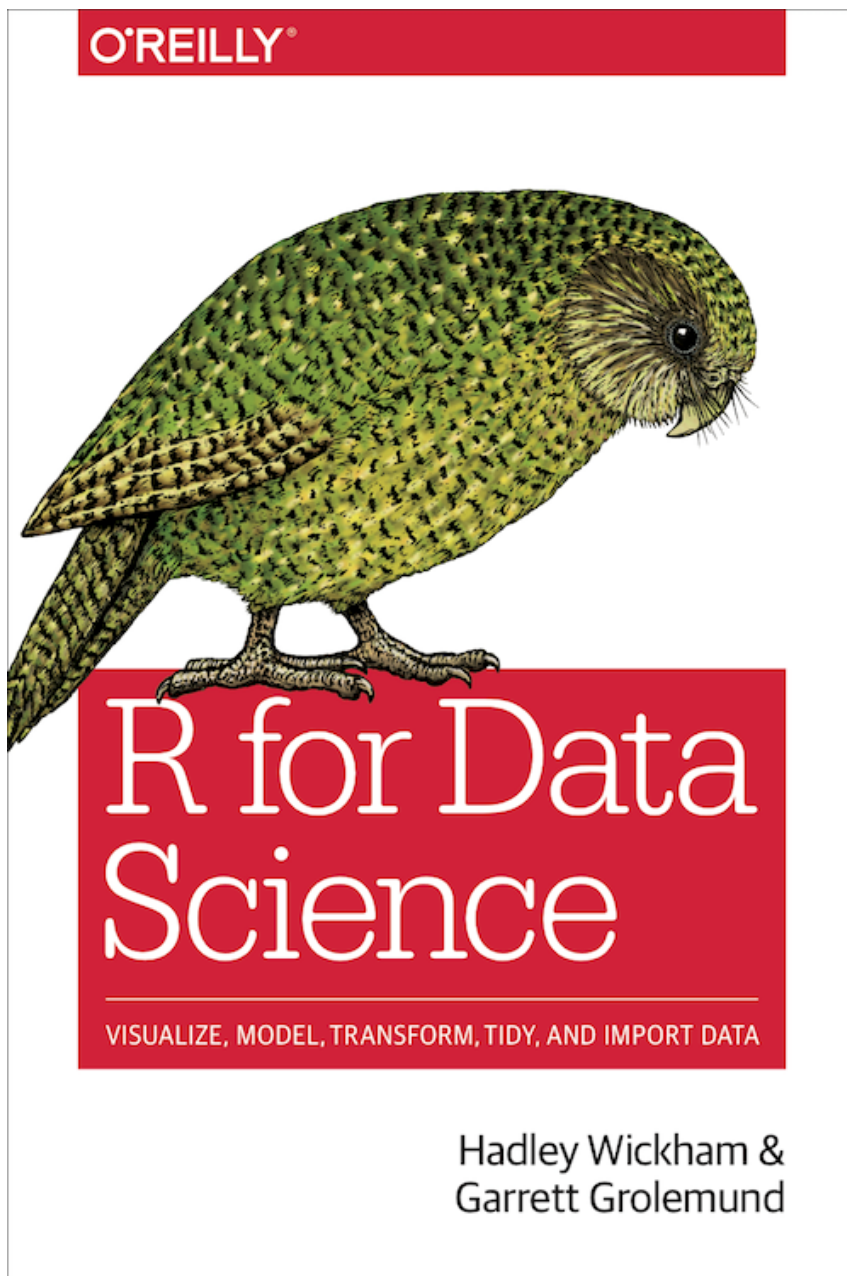# Ch11.군집분석 실습

https://r4ds.had.co.nz/

# 학습목표

- R 프로그래밍 실습

- 실습: K-평균 군집화

# 전력회사 사례

# 데이터형식

- 공공전력회사 사례
  - 22개 전력회사
  - 유사한 전력회사로 군집

| Company | Fixed_charge | RoR | Cost | Load | Δ Demand | Sales | Nuclear | Fuel_Cost |
|---|---|---|---|---|---|---|---|---|
| Arizona | 1.06 | 9.2 | 151 | 54.4 | 1.6 | 9077 | 0 | 0.628 |
| Boston | 0.89 | 10.3 | 202 | 57.9 | 2.2 | 5088 | 25.3 | 1.555 |
| Central | 1.43 | 15.4 | 113 | 53 | 3.4 | 9212 | 0 | 1.058 |
| Commonwealth | 1.02 | 11.2 | 168 | 56 | 0.3 | 6423 | 34.3 | 0.7 |
| Con Ed NY | 1.49 | 8.8 | 192 | 51.2 | 1 | 3300 | 15.6 | 2.044 |
| Florida | 1.32 | 13.5 | 111 | 60 | -2.2 | 11127 | 22.5 | 1.241 |
| Hawaiian | 1.22 | 12.2 | 175 | 67.6 | 2.2 | 7642 | 0 | 1.652 |
| Idaho | 1.1 | 9.2 | 245 | 57 | 3.3 | 13082 | 0 | 0.309 |
| Kentucky | 1.34 | 13 | 168 | 60.4 | 7.2 | 8406 | 0 | 0.862 |
| Madison | 1.12 | 12.4 | 197 | 53 | 2.7 | 6455 | 39.2 | 0.623 |
| Nevada | 0.75 | 7.5 | 173 | 51.5 | 6.5 | 17441 | 0 | 0.768 |
| New England | 1.13 | 10.9 | 178 | 62 | 3.7 | 6154 | 0 | 1.897 |
| Northern | 1.15 | 12.7 | 199 | 53.7 | 6.4 | 7179 | 50.2 | 0.527 |
| Oklahoma | 1.09 | 12 | 96 | 49.8 | 1.4 | 9673 | 0 | 0.588 |
| Pacific | 0.96 | 7.6 | 164 | 62.2 | -0.1 | 6468 | 0.9 | 1.4 |
| Puget | 1.16 | 9.9 | 252 | 56 | 9.2 | 15991 | 0 | 0.62 |
| San Diego | 0.76 | 6.4 | 136 | 61.9 | 9 | 5714 | 8.3 | 1.92 |
| Southern | 1.05 | 12.6 | 150 | 56.7 | 2.7 | 10140 | 0 | 1.108 |
| Texas | 1.16 | 11.7 | 104 | 54 | -2.1 | 13507 | 0 | 0.636 |
| Wisconsin | 1.2 | 11.8 | 148 | 59.9 | 3.5 | 7287 | 41.1 | 0.702 |
| United | 1.04 | 8.6 | 204 | 61 | 3.5 | 6650 | 0 | 2.116 |
| Virginia | 1.07 | 9.3 | 174 | 54.3 | 5.9 | 10093 | 26.6 | 1.306 |

# Package 설치

```
# https://www.tidymodels.org/learn/statistics/k-means/
# 02.K-평균 군집화: Table 15.9


##############################################
# tidyverse: ggplot2, purrr, tibble  3.0.3,          #
#            dplyr, tidyr, stringr, readr, forcats    #
##############################################


# install.packages("tidyverse")
# install.packages("tidymodels")
library(tidyverse)
library(tidymodels)
```

# Package 설치

```
Console   Terminal ×   Jobs ×
C:/Users/leecho/Desktop/R-DM/Ch11.군집분석/
> # install.packages("tidyverse")
> # install.packages("tidymodels")
> library(tidyverse)
-- Attaching packages ------------------------------------------------- tidyverse 1.3.0 --
√ ggplot2 3.3.2     √ purrr   0.3.4
√ tibble  3.0.4     √ dplyr   1.0.2
√ tidyr   1.1.2     √ stringr 1.4.0
√ readr   1.4.0     √ forcats 0.5.0
-- Conflicts --------------------------------------------------------- tidyverse_conflicts() --
x dplyr::filter() masks stats::filter()
x dplyr::lag()    masks stats::lag()
경고메시지(들):
1: 패키지 'tidyverse'는 R 버전 4.0.3에서 작성되었습니다
2: 패키지 'ggplot2'는 R 버전 4.0.3에서 작성되었습니다
3: 패키지 'tibble'는 R 버전 4.0.3에서 작성되었습니다
4: 패키지 'tidyr'는 R 버전 4.0.3에서 작성되었습니다
5: 패키지 'readr'는 R 버전 4.0.3에서 작성되었습니다
6: 패키지 'purrr'는 R 버전 4.0.3에서 작성되었습니다
7: 패키지 'dplyr'는 R 버전 4.0.3에서 작성되었습니다
8: 패키지 'stringr'는 R 버전 4.0.3에서 작성되었습니다
9: 패키지 'forcats'는 R 버전 4.0.3에서 작성되었습니다
>

Files   Plots   Packages   Help   Viewer
```

# Package 설치

# 01.데이터 불러오기

```r
# 01.데이터 불러오기
book_tb <- read_csv('CharlesBookClub.csv',
              col_names = TRUE,
              locale=locale('ko', encoding='euc-kr'),
              na=".") %>% # csv 데이터 읽어오기
  mutate_if(is.character, as.factor)

str(book_tb)
head(book_tb)
```

# 01.데이터 불러오기

```
Console    Terminal ×    Jobs ×
C:/Users/leecho/Desktop/R-DM/Ch11.군집분석/
> utilities_tb <- read_csv('Utilities.csv',
+                          col_names = TRUE,
+                          locale=locale('ko', encoding='euc-kr'),
+                          na=".") %>% # csv 데이터 읽어오기
+   mutate_if(is.character, as.factor)

-- Column specification ------------------------------------------------------
cols(
  Company = col_character(),
  Fixed_charge = col_double(),
  RoR = col_double(),
  Cost = col_double(),
  Load_factor = col_double(),
  Demand_growth = col_double(),
  Sales = col_double(),
  Nuclear = col_double(),
  Fuel_Cost = col_double()
)

> str(utilities_tb)
tibble [22 x 9] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
 $ Company      : Factor w/ 22 levels "Arizona","Boston",..: 1 2 3 4 13 5 6 7 8 9 ...
 $ Fixed_charge : num [1:22] 1.06 0.89 1.43 1.02 1.49 1.32 1.22 1.1 1.34 1.12 ...
 $ RoR          : num [1:22] 9.2 10.3 15.4 11.2 8.8 13.5 12.2 9.2 13 12.4 ...
 $ Cost         : num [1:22] 151 202 113 168 192 111 175 245 168 197 ...
 $ Load_factor  : num [1:22] 54.4 57.9 53 56 51.2 60 67.6 57 60.4 53 ...
 $ Demand_growth: num [1:22] 1.6 2.2 3.4 0.3 1 -2.2 2.2 3.3 7.2 2.7 ...
 $ Sales        : num [1:22] 9077 5088 9212 6423 3300 ...
 $ Nuclear      : num [1:22] 0 25.3 0 34.3 15.6 22.5 0 0 0 39.2 ...
 $ Fuel_Cost    : num [1:22] 0.628 1.555 1.058 0.7 2.044 ...
 - attr(*, "spec")=
  .. cols(
  ..    Company = col_character(),
  ..    Fixed_charge = col_double(),
  ..    RoR = col_double(),
  ..    Cost = col_double(),
```
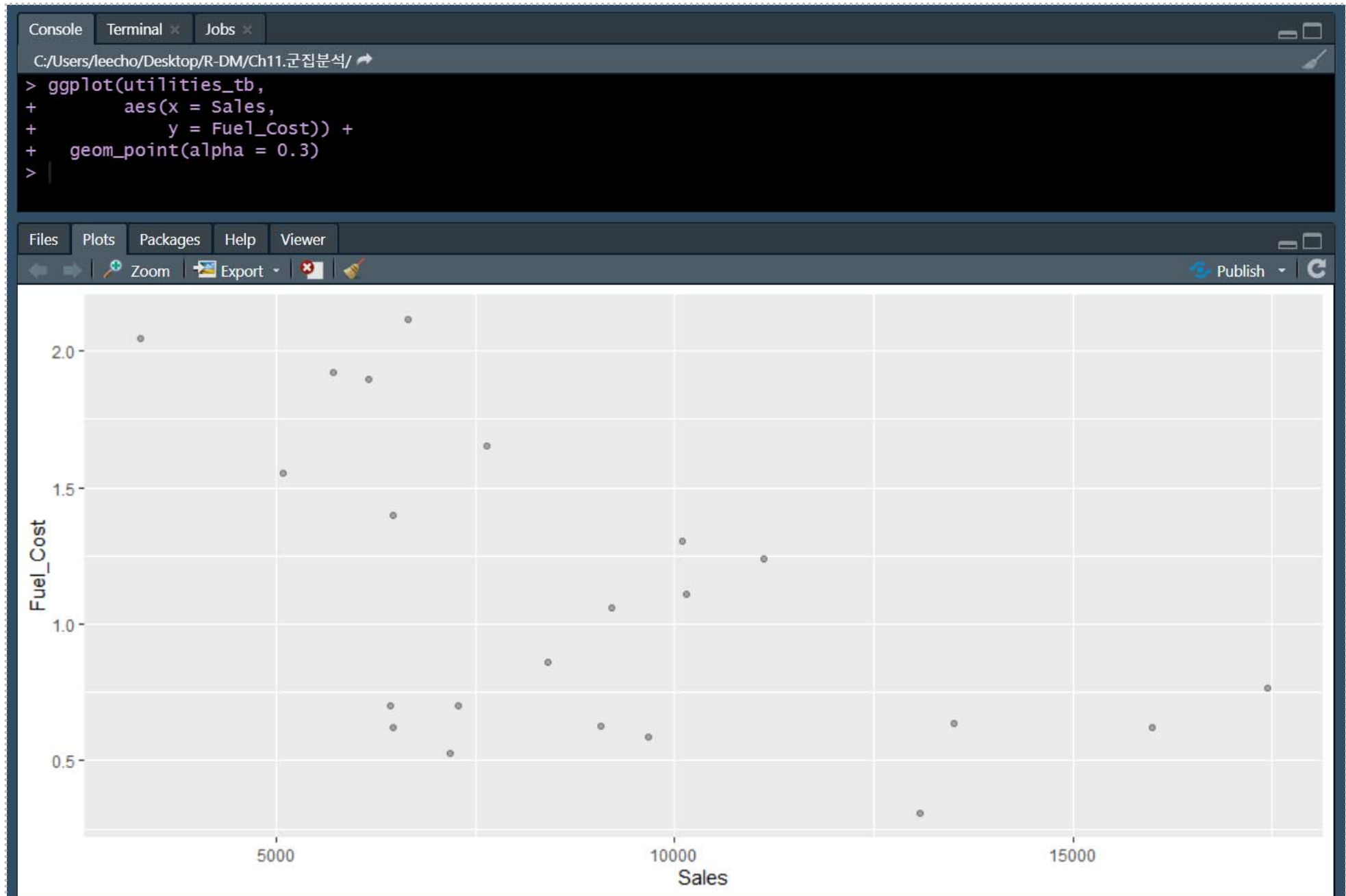
Files    Plots    Packages    Help    Viewer

# 01.데이터 불러오기

```
> str(utilities_tb)
tibble [22 x 9] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
 $ Company      : Factor w/ 22 levels "Arizona","Boston",..: 1 2 3 4 13 5 6 7 8 9 ...
 $ Fixed_charge : num [1:22] 1.06 0.89 1.43 1.02 1.49 1.32 1.22 1.1 1.34 1.12 ...
 $ RoR          : num [1:22] 9.2 10.3 15.4 11.2 8.8 13.5 12.2 9.2 13 12.4 ...
 $ Cost         : num [1:22] 151 202 113 168 192 111 175 245 168 197 ...
 $ Load_factor  : num [1:22] 54.4 57.9 53 56 51.2 60 67.6 57 60.4 53 ...
 $ Demand_growth: num [1:22] 1.6 2.2 3.4 0.3 1 -2.2 2.2 3.3 7.2 2.7 ...
 $ Sales        : num [1:22] 9077 5088 9212 6423 3300 ...
 $ Nuclear      : num [1:22] 0 25.3 0 34.3 15.6 22.5 0 0 0 39.2 ...
 $ Fuel_Cost    : num [1:22] 0.628 1.555 1.058 0.7 2.044 ...
 - attr(*, "spec")=
  .. cols(
  ..    Company = col_character(),
  ..    Fixed_charge = col_double(),
  ..    RoR = col_double(),
  ..    Cost = col_double(),
  ..    Load_factor = col_double(),
  ..    Demand_growth = col_double(),
  ..    Sales = col_double(),
  ..    Nuclear = col_double(),
  ..    Fuel_Cost = col_double()
  .. )
> head(utilities_tb)
# A tibble: 6 x 9
  Company      Fixed_charge  RoR  Cost Load_factor Demand_growth Sales Nuclear Fuel_Cost
  <fct>               <dbl> <dbl> <dbl>      <dbl>         <dbl> <dbl>   <dbl>     <dbl>
1 Arizona              1.06   9.2   151       54.4           1.6  9077     0       0.628
2 Boston               0.89  10.3   202       57.9           2.2  5088    25.3     1.56
3 Central              1.43  15.4   113       53             3.4  9212     0       1.06
4 Commonwealth         1.02  11.2   168       56             0.3  6423    34.3     0.7
5 NY                   1.49   8.8   192       51.2           1    3300    15.6     2.04
6 Florida              1.32  13.5   111       60            -2.2 11127    22.5     1.24
>
```

# 01.데이터 불러오기

```
# 데이터 분포 확인

ggplot(utilities_tb,
      aes(x = Sales,
          y = Fuel_Cost)) +
  geom_point(alpha = 0.3)
```

# 02.데이터 정규화

```
# 02.데이터 정규화
# 데이터 정규화: mutate_if, 수치형 변수만 정규화
# 회사이름을 row 이름으로 변경

utilities_tb <-
  utilities_tb %>%
  mutate_if(is.numeric, funs(scale(.))) %>%
  column_to_rownames(var = "Company")

utilities_tb
```

```
Console  Terminal ×  Jobs ×
C:/Users/leecho/Desktop/R-DM/Ch11.군집분석/

  # Simple named list:
  list(mean = mean, median = median)

  # Auto named with `tibble::lst()`:
  tibble::lst(mean, median)

  # Using lambdas
  list(~ mean(., trim = .2), ~ median(., na.rm = TRUE))
This warning is displayed once every 8 hours.
Call `lifecycle::last_warnings()` to see where this warning was generated.
> utilities_tb
              Fixed_charge          RoR         Cost Load_factor Demand_growth        Sales     Nuclear    Fuel_Cost
Arizona         -0.29315791  -0.68463896  -0.417122002  -0.57771516    -0.52622751   0.04590290  -0.7146294  -0.85367545
Boston          -1.21451134  -0.19445367   0.821002037   0.20683629    -0.33381191  -1.07776413   0.7920476   0.81329670
Central          1.71214073   2.07822360  -1.339645796  -0.89153574     0.05101929   0.08393124  -0.7146294  -0.08043055
Commonwealth    -0.50994695   0.20660702  -0.004413989  -0.21906307    -0.94312798  -0.70170610   1.3280197  -0.72420189
NY               2.03732429  -0.86288816   0.578232617  -1.29501935    -0.71864311  -1.58142837   0.2143888   1.69263800
Florida          1.11597086   1.23153991  -1.388199680   0.67756716    -1.74485965   0.62337028   0.6253007   0.24864810
Hawaiian         0.57399826   0.65223002   0.165524604   2.38116460    -0.33381191  -0.35832428  -0.7146294   0.98772637
Idaho           -0.07636887  -0.68463896   1.864910540   0.00509449     0.01895002   1.17407698  -0.7146294  -1.42731528
Kentucky         1.22436538   1.00872841  -0.004413989   0.76723019     1.26965142  -0.14311204  -0.7146294  -0.43288637
Madison          0.03202565   0.74135462   0.699617327  -0.89153574    -0.17346558  -0.69269198   1.6198267  -0.86266667
Nevada          -1.97327298  -1.44219805   0.116970720  -1.22777208     1.04516655   2.40196983  -0.7146294  -0.60192130
New England      0.08622291   0.07292013   0.238355430   1.12588228     0.14722709  -0.77748109  -0.7146294   1.42829614
Northern         0.19461744   0.87504152   0.748171211  -0.73462545     1.01309729  -0.48874740   2.2749037  -1.03529809
Oklahoma        -0.13056613   0.56310542  -1.752353809  -1.60883993    -0.59036605   0.21379097  -0.7146294  -0.92560521
Pacific         -0.83513051  -1.39763576  -0.101521757   1.17071379    -1.07140505  -0.68902999  -0.6610322   0.53456889
Puget            0.24881470  -0.37270287   2.034849134  -0.21906307     1.91103676   1.99351729  -0.7146294  -0.86806140
San Diego       -1.91907572  -1.93238335  -0.781276132   1.10346652     1.84689822  -0.90142531  -0.2203441   1.46965575
Southern        -0.34735517   0.83047922  -0.441398944  -0.06215278    -0.17346558   0.34534086  -0.7146294   0.00948165
Texas            0.24881470   0.42941852  -1.558138274  -0.66737818    -1.71279038   1.29379583  -0.7146294  -0.83928950
Wisconsin        0.46560374   0.47398082  -0.489952828   0.65515141     0.08308855  -0.45832473   1.7329764  -0.72060540
United          -0.40155243  -0.95201276   0.869555920   0.90172472     0.08308855  -0.63776215  -0.7146294   1.82211157
Virginia        -0.23896065  -0.64007666   0.141247662  -0.60013092     0.85275095   0.33210137   0.8694658   0.36553395
>
```

Files  Plots  Packages  Help  Viewer

# 03.최적 군집수 찾기

```
# 03.최적 군집수 찾기
# 최적 군집수를 찾는 엘보우(Elbow) 챠트


# 군집 9개

kclusts <-
  tibble(k = 1:9) %>%
  mutate(
    kclust = map(k, ~kmeans(utilities_tb, .x)),
    tidied = map(kclust, tidy),
    glanced = map(kclust, glance),
    augmented = map(kclust, augment, utilities_tb)
  )


kclusts
```

# 03.최적 군집수 찾기

```
Console   Terminal ×   Jobs ×
C:/Users/leecho/Desktop/R-DM/Ch11.군집분석/ 
> kclusts <-
+   tibble(k = 1:9) %>%
+   mutate(
+     kclust = map(k, ~kmeans(utilities_tb, .x)),
+     tidied = map(kclust, tidy),
+     glanced = map(kclust, glance),
+     augmented = map(kclust, augment, utilities_tb)
+   )
> kclusts
# A tibble: 9 x 5
      k kclust    tidied            glanced            augmented
  <int> <list>    <list>            <list>             <list>
1     1 <kmeans> <tibble [1 x 11]> <tibble [1 x 4]> <tibble [22 x 10]>
2     2 <kmeans> <tibble [2 x 11]> <tibble [1 x 4]> <tibble [22 x 10]>
3     3 <kmeans> <tibble [3 x 11]> <tibble [1 x 4]> <tibble [22 x 10]>
4     4 <kmeans> <tibble [4 x 11]> <tibble [1 x 4]> <tibble [22 x 10]>
5     5 <kmeans> <tibble [5 x 11]> <tibble [1 x 4]> <tibble [22 x 10]>
6     6 <kmeans> <tibble [6 x 11]> <tibble [1 x 4]> <tibble [22 x 10]>
7     7 <kmeans> <tibble [7 x 11]> <tibble [1 x 4]> <tibble [22 x 10]>
8     8 <kmeans> <tibble [8 x 11]> <tibble [1 x 4]> <tibble [22 x 10]>
9     9 <kmeans> <tibble [9 x 11]> <tibble [1 x 4]> <tibble [22 x 10]>
> 

Files   Plots   Packages   Help   Viewer
```

```
clusters <-
  kclusts %>%
  unnest(cols = c(tidied))

assignments <-
  kclusts %>%
  unnest(cols = c(augmented))

clusterings <-
  kclusts %>%
  unnest(cols = c(glanced))
```
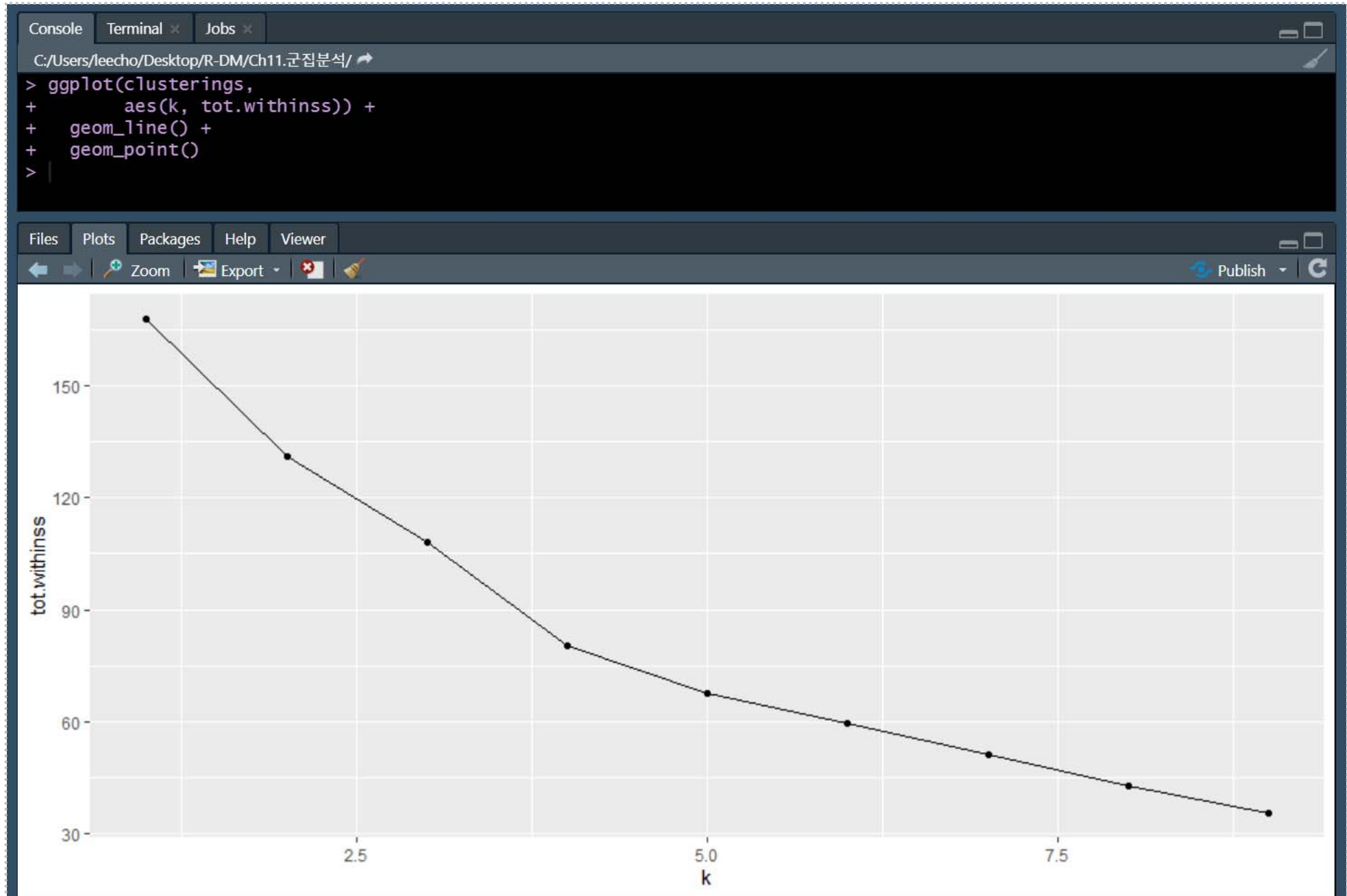
```
# 엘보우(Elbow) 챠트

ggplot(clusterings,
      aes(k, tot.withinss)) +
  geom_line() +
  geom_point()

# 군집별 그래프

ggplot(assignments,
      aes(x = Sales,
        y = Fuel_Cost)) +
  geom_point(aes(color = .cluster),
        alpha = 0.8) +
  facet_wrap(~ k) +
  geom_point(data = clusters,
        size = 5,
        shape = "x")
```

```
#  best model 구축

set.seed(123)

kclust_best <-
  kmeans(utilities_tb,
        centers = 3)


# 군집분석 결과 확인

tidy(kclust_best)
```

```
Console   Terminal ×   Jobs ×
C:/Users/leecho/Desktop/R-DM/Ch11.군집분석/
> set.seed(123)
> kclust_best <-
+   kmeans(utilities_tb,
+         centers = 3)
> tidy(kclust_best)
# A tibble: 3 x 11
  Fixed_charge     RoR    Cost Load_factor Demand_growth   Sales Nuclear Fuel_Cost  size withinss cluster
         <dbl>   <dbl>   <dbl>       <dbl>         <dbl>   <dbl>   <dbl>     <dbl> <int>    <dbl> <fct>
1      -0.239  -0.659   0.256       0.799       -0.0544  -0.860  -0.288      1.25     7     34.2 1
2       0.520   1.03   -1.30       -0.510       -0.834    0.512  -0.447     -0.317     5     15.2 2
3      -0.0926 -0.0519  0.469      -0.304        0.455    0.346   0.425     -0.716    10     57.5 3
> |
```
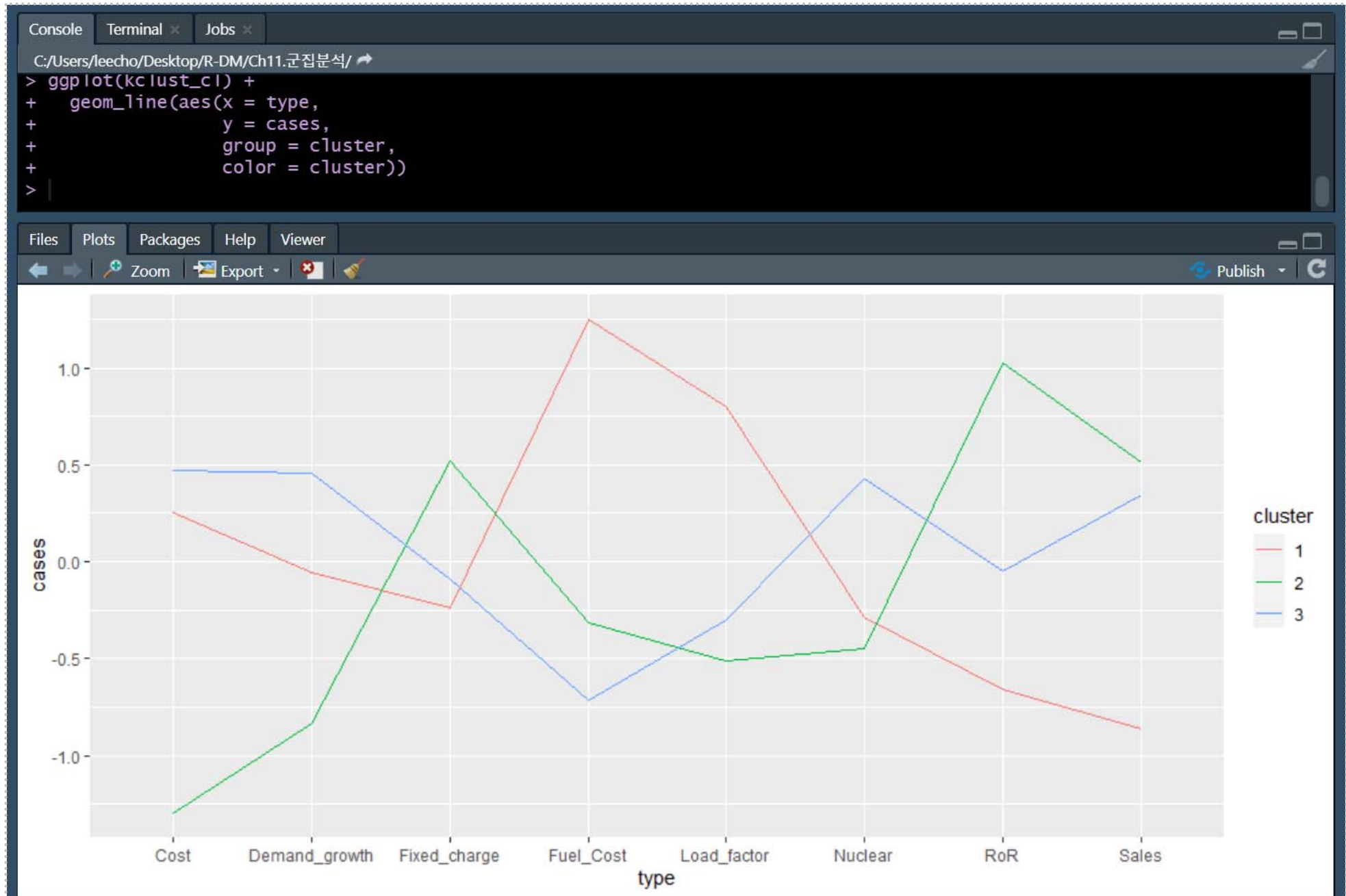
Files   Plots   Packages   Help   Viewer

```
kclust_cl <-
  tidy(kclust_best) %>%
  select(-c(size, withinss))

kclust_cl <-
  kclust_cl %>%
  pivot_longer(c("Fixed_charge", #c("1999, 2000")에러남
            "Cost",
            "RoR",
            "Load_factor",
            "Demand_growth",
            "Sales",
            "Nuclear",
            "Fuel_Cost"),
          names_to = "type",
          values_to = "cases")

ggplot(kclust_cl) +
  geom_line(aes(x = type,
```

28

# 참고자료

# 동영상 및 참고교재

- International
  - Tidymodels, https://www.tidymodels.org/
  - R for Data Science, https://r4ds.had.co.nz/
  - Statistical Inference via Data Science, https://www.tidymodels.org/books/moderndive/
  - Data Mining for Business Analytics_ Concepts, Techniques, and Applications in R, Shmueli et al., WILEY
- Domestic
  - K-MOOC
    - 경영데이터마이닝, 김종우, 한양대학교
    - 빅데이터의 세계, 원리와 응용, 신경식, 이화여자대학교
    - 빅데이터와 텍스트마이닝, 이신행, 세종대학교
  - 데이터마이닝 R edition, 조재희 외, 이앤비플러스
  - 김종우, 김선태, 경영을 위한 데이터마이닝, 한경사, 2009.
  - 경영정보시스템, 김우주 외, 시그마프레스