

Lecture 2

*dplyr & ggplot2 & rmd
(Carseat Sales)*

러닝스푼즈

2018년 4월

프로젝트 개요 (카시트 판매량)

1. 수집 - *from package*
2. 전처리 - *dplyr*
3. 분석, 시각화 - *ggplot2*
4. 문서화 - *rmd*
5. 공유

Module 0

패키지 & 파일 & 작업 경로

패키지 (package, library)

- R을 처음에 설치하면 깔려있는 패키지는 *base*라고 함
- *Base*외의 패키지를 필요에 따라서 *install*하여 사용함.

```
install.packages("package_name")
```

- 소스코드에서 사용할 것을 선언해 주어야 함.

```
library(package_name) .
```

파일 (files)

- 소스 파일
 - 명령어를 기록해둔 파일
 - 한번에 실행이 가능한 단위로 구성
 - `.r` : R 소스 파일
 - `.rmd` : R Markdown 소스파일
- 데이터 파일
 - 데이터가 기록된 파일 – 일반적으로 서식이 없이 문자 자체로 되어 있음.
 - `.csv` : 컴마로 구분된 데이터 파일 (Comma Separated Values)
 - `.txt` : 일반적인 데이터 파일 (Tab, `\n` 등의 기호로 구분되어 있음)
 - `.xls,.xlsx` :
 - 서식 등이 포함되어 있는 경우에는 [다른 이름으로 저장 – `csv`로 저장]이 바람직
 - 바로 불러오기도 가능은 함
- Rdata 파일
 - `.Rdata` : R 작업시에 메모리의 상태를 저장할 수 있어서 작업을 이어서 하는데 유용
 - `.Rda` : 변수 1개를 저장하는 파일

데이터 파일 불러오기/저장하기

	A	B	C
1	이름		
2		나이	
3			—

```
# csv
dataset <- 
  read.csv("filename.csv", header = TRUE, stringsAsFactors = FALSE)
  # if the first line is "header" (default)
dataset <-
  read.csv("filename.csv", header = FALSE, stringsAsFactors = FALSE)
  # if no "header" in the data file
write.csv(dataset, "filename.csv")
```

```
# txt
# read.csv 대신에 read.table
# write.csv 대신에 write.table

dataset <- read.table(나머지 문법은 read.csv 와 똑같고 sep 옵션 추가)
# ex) sep = ","           "comma"
#      sep = " "          "space"
#      sep = "-"          "hyphen"
```

데이터 파일 불러오기/저장하기

```
# xls, xlsx
install.packages("readxl") # 최초 사용시에만 설치 필요
library(readxl)
dataset <- read_excel("filename.xlsx")
# if first Line is "header"
dataset <- read_excel("filename.xlsx", col_names=FALSE)
# if no "header" in the data file
```

```
# rda - 데이터셋 한 개 저장
load("filename.rda")
save(dataset, "filename.rda")

# rdata - 현재 메모리에 있는 모든 변수 저장
load("2018-04-17.rdata")
save(dataset, "2018-04-17.rdata")
```

경로 (Directory)

- 새파일로 R파일을 만들었을 경우에는 시스템 작업 디렉토리
- 존재하는 R파일을 열었을 경우에는 R 파일이 있는 디렉토리가 현재 작업 디렉토리

Working Directory

`getwd()`

Find the current working directory (where inputs are found and outputs are sent).

`setwd('C://file/path')`

Change the current working directory.

Use projects in RStudio to set the working directory to the folder you are working in.

- 작업 디렉토리
- 현재 작업 디렉토리를 조회
(파일 저장시 여기에 저장됨)
- 작업 디렉토리를 설정
- Rstudio에서는 프로젝트 단위로의 설정도 제공 (진행하는 작업이 많아지면 유용함)

Module 1

`dplyr`

전처리 (preprocessing)
(modern and nice way to treat a dataset)

Hello Carseat?

```
install.packages("ISLR")  
library(ISLR)  
class(Carseats)  
  
head(Carseats)      # 처음 6 개 관찰값  
.tail(Carseats, 5) # 처음 6 개 관찰값 마지막 5개  
  
View(Carseats)     # "Viewer"에 보여줌  
dim(Carseats)      # 차원 (dimension)  
str(Carseats)       # 구조 (structure)  
  
summary(Carseats)  # 기초 통계량 (descriptive stat)
```



James, G., Witten, D., Hastie, T., and Tibshirani, R. (2013) *An Introduction to Statistical Learning with applications in R*,
<www.StatLearning.com>, Springer-Verlag, New York

Hello Carseat?

```
library(ISLR)
class(Carseats)

## [1] "data.frame"

head(Carseats)      # 처음 6 개 관찰값

##   Sales CompPrice Income Advertising Population Price ShelveLoc Age
## 1  9.50        138     73          11       276    120      Bad  42
## 2 11.22        111     48          16       260     83      Good  65
## 3 10.06        113     35          10       269     80  Medium  59
## 4  7.40        117    100           4       466     97  Medium  55
## 5  4.15        141     64           3       340    128      Bad  38
## 6 10.81        124    113          13       501     72      Bad  78
##   Education Urban US
## 1            17 Yes Yes
## 2            10 Yes Yes
## 3            12 Yes Yes
## 4            14 Yes Yes
## 5            13 Yes  No
## 6            16  No Yes
```

Hello Carseat?

```
tail(Carseats, 5) # 처음 6 개 관찰값
```

```
##      Sales CompPrice Income Advertising Population Price ShelveLoc Age
## 396 12.57        138     108          17       203    128      Good   33 +
## 397  6.14        139      23           3       37    120    Medium   55 +
## 398  7.41        162      26          12       368    159    Medium   40 +
## 399  5.94        100      79           7       284     95      Bad    50 +
## 400  9.71        134     37           0       27    120      Good   49 +  
##          Education Urban US
## 396            14 Yes Yes +
## 397            11 No  Yes +
## 398            18 Yes Yes +
## 399            12 Yes Yes +
## 400            16 Yes Yes +
```

```
View(Carseats) # "Viewer"에 보여줌
```

```
dim(Carseats) # 차원 (dimension)
```

```
## [1] 400 11 +
```

Hello Carseat?

```
str(Carseats) # 구조 (structure)+  
  
## 'data.frame': 400 obs. of 11 variables:  
## $ Sales : num 9.5 11.22 10.06 7.4 4.15 ...+  
## $ CompPrice : num 138 111 113 117 141 124 115 136 132 132 ....+  
## $ Income : num 73 48 35 100 64 113 105 81 110 113 ...+  
## $ Advertising: num 11 16 10 4 3 13 0 15 0 0 ...+  
## $ Population : num 276 260 269 466 340 501 45 425 108 131 ...+  
## $ Price : num 120 83 80 97 128 72 108 120 124 124 ...+  
## $ ShelveLoc : Factor w/ 3 levels "Bad", "Good", "Medium": 1 2 3 3 1 1 3 2  
3 3 ...+  
## $ Age : num 42 65 59 55 38 78 71 67 76 76 ...+  
## $ Education : num 17 10 12 14 13 16 15 10 10 17 ...+  
## $ Urban : Factor w/ 2 levels "No", "Yes": 2 2 2 2 2 1 2 2 1 1 ...+  
## $ US : Factor w/ 2 levels "No", "Yes": 2 2 2 2 1 2 1 2 1 2 ...+
```

boolen 은 뭘까요.

Hello Carseat?

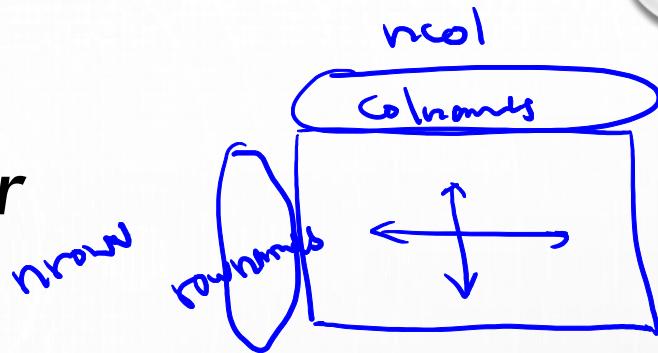
`summary(Carseats) # 기초 통계량 (descriptive stat)`

```
##      Sales          CompPrice        Income       Advertising +  
## Min. : 0.000    Min.   : 77    Min.   : 21.00    Min.   : 0.000 +  
## 1st Qu.: 5.390   1st Qu.:115   1st Qu.: 42.75    1st Qu.: 0.000 +  
## Median : 7.490   Median  :125   Median  : 69.00    Median  : 5.000 +  
## Mean   : 7.496   Mean    :125   Mean    : 68.66    Mean   : 6.635 +  
## 3rd Qu.: 9.320   3rd Qu.:135   3rd Qu.: 91.00    3rd Qu.:12.000 +  
## Max.  :16.270   Max.   :175   Max.   :120.00    Max.  :29.000 +  
##      Population        Price       ShelveLoc        Age      +  
## Min.   : 10.0    Min.   : 24.0    Bad    : 96    Min.   :25.00 +  
## 1st Qu.:139.0   1st Qu.:100.0   Good   : 85    1st Qu.:39.75 +  
## Median :272.0   Median  :117.0   Medium :219   Median :54.50 +  
## Mean   :264.8   Mean    :115.8           Mean   :53.32 +  
## 3rd Qu.:398.5   3rd Qu.:131.0           3rd Qu.:66.00 +  
## Max.  :509.0   Max.   :191.0           Max.  :80.00 +  
##      Education        Urban        US      +  
## Min.   :10.0    No   :118    No  :142    +  
## 1st Qu.:12.0   Yes  :282    Yes :258    +  
## Median :14.0           +  
## Mean   :13.9           +  
## 3rd Qu.:16.0           +  
## Max.  :18.0+  
##
```

dplyr

- 빠르고 직관적인 데이터를 다루는 패키지
 - 가장 빠른 언어인 C를 기반으로 만들어서 빠름
 - 데이터 처리에 있어서 가장 직관적인 SQL (*Standard Query Language*)과 유사하게 만들어져서 직관적임!
 - 직관적이어서 코드 가독성도 높음
 - 그러나 *base* 명령어도 어느정도 같이 알아 두면 타인의 소스코드 참조와 파이썬, SQL등 다른 언어를 배울 때 도움이 됨
 - *Play with “tidy” data*
- 제작자
 - Hadley Wickham, Ph.D., Head Scientist, Rstudio
 - 통계학 박사 후 R에서 다수의 사용하기 좋은 패키지 개발
 - youtube.com에 keynote등 좋은 동영상 많아요
 - 누나도 통계학 전공 교수

dplyr



- Tidy data?

dplyr functions work with pipes and expect **tidy data**. In tidy data:

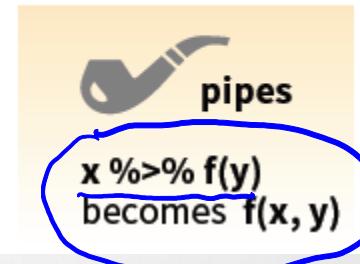


&



Each **variable** is in
its own **column**

Each **observation**, or
case, is in its own **row**



- 각 variable은 column에 대응
- 각 observations은 row에 대응
- $f(x, y)$ 는 $x \%>\% f(y)$ 로 적을 수 있어서 가독성이 좋음
(x 를 f함수와 y 로 처리한다고 읽을 수 있음)
- $f(x)$ 는 $x \%>\% f()$

$f(x, y)$
 $x \%>\% f(y)$

콘솔에 띄워 놓고 보면서!

```
head(Carseats)
```

```
##   Sales CompPrice Income Advertising Population Price ShelveLoc Age+
## 1  9.50      138     73          11       276    120      Bad    42+
## 2 11.22      111     48          16       260     83      Good    65+
## 3 10.06      113     35          10       269     80  Medium    59+
## 4  7.40      117    100           4       466     97  Medium    55+
## 5  4.15      141     64           3       340    128      Bad    38+
## 6 10.81      124    113          13       501     72      Bad    78+
##   Education Urban US+
## 1            17 Yes Yes+
## 2            10 Yes Yes+
## 3            12 Yes Yes+
## 4            14 Yes Yes+
## 5            13 Yes  No+
## 6            16  No Yes+
```

Basic Treatments

```
# install.packages("dplyr")  
library(dplyr)  
  
# rename (이름 바꾸기) ↓  
Carseats <- rename(Carseats, Sales = Revenue) + modern  
names(Carseats)[names(Carseats) == "Sales"] <- "Revenue" + classic  
↓  
  
# filter (관찰값 추출, Row 추출) ↓  
temp <- filter(Carseats, Income > 100) ↓  
temp <- Carseats %>% filter(Income > 100) ↓  
temp <- Carseats[Carseats$Income > 100,] ↓  
↓  
temp <- filter(Carseats, Age >= 30 & Age < 40) ↓  
temp <- Carseats %>% filter(Age >= 30 & Age < 40) ↓  
temp <- Carseats[((Carseats$Age >= 30) & (Carseats$Age < 40)),] ↓  
  
# select (변수 추출, Column 선택) ↓  
temp <- select(Carseats, Income, Population) ↓  
temp <- Carseats %>% select(Income, Population) ↓  
temp <- Carseats[, c("Income", "Population")] ↓
```

And : or

```
# arrange (정렬) +
Carseats <- arrange(Carseats, Price) +
Carseats <- Carseats %>% arrange(Price) +
Carseats <- Carseats[order(Carseats$Price), ] +
+
Carseats <- arrange(Carseats, desc(Price)) +
Carseats <- Carseats %>% arrange(desc(Price)) +
Carseats <- Carseats[order(Carseats$Price, decreasing = TRUE), ] +
```

```
# mutate (새로운 변수) ↓
Carseats <- mutate(Carseats, +
                     AdvPerCapita = Advertising/Population, +
                     RevPerCapita = Revenue/Population) +
Carseats <- Carseats %>% +
  mutate(AdvPerCapita = Advertising/Population, +
        RevPerCapita = Revenue/Population) +
Carseats$AdvPerCapita <- CarSeat$Advertising/Carseats$Population +
Carseats$RevPerCapita <- CarSeat$Revenue/Carseats$Population +
↓
Carseats <- mutate(Carseats, +
                     AgeClass = ifelse(Age>=60, "Silver", "non-Silver")) +
Carseats <- Carseats %>% +
  mutate(AgeClass = ifelse(Age>=60, "Silver", "non-Silver")) +
Carseats$AgeClass <- ↓
  ifelse(Carseats$Age >= 60, "Silver", "non-Silver") +
```

Successive Treatments

- 해당 Carseat 회사의 주요 구매층은 아이를 낳아서 키울 나이인 30대 연령층입니다.
- 소득이 높으면서 도시의 평균 연령이 30대인 도시에 충분한 광고비를 지출하고 있는지 검토해보고 싶습니다.
- *We all can think in what we can speak of!*
- *Think and speak in logical, sequential, (and more often, reverse-sequential), and codable ways!*
- *Learning how to program makes your brain more logical.*

Successive Treatments

```
# successive treatments↓
focusCity <- Carseats %>% ↓
  filter(Income > 100) %>%↓
  filter(Age >= 30 & Age < 40) %>%↓
  mutate(AdvPerCapita = Advertising/Population) %>%↓
  select(Revenue, Income, Age, Population, Education, AdvPerCapita) %>%↓
  arrange(Revenue) ↓
print(focusCity) ↓

##      Revenue Income Age Population Education AdvPerCapita ↓
## 1      5.04    114  34        298       16  0.00000000 ↓
## 2      5.32    116  39        170       16  0.00000000 ↓
## 3      6.80    117  38        337       10  0.01483680 ↓
## 4      7.49    119  35        178       13  0.03370787 ↓
## 5      7.67    117  36        400       10  0.02000000 ↓
## 6      8.55    111  36        480       16  0.04791667 ↓
## 7      8.97    107  33        144       13  0.00000000 ↓
## 8      9.03    102  35        123       16  0.10569106 ↓
## 9      9.39    118  32        445       15  0.03146067 ↓
## 10     9.58    104  37        353       17  0.06515581 ↓
## 11    10.36    105  34        428       12  0.04205607 ↓
## 12    10.59    120  30        262       10  0.05725191 ↓
## 13    12.57    108  33        203       14  0.08374384 ↓
```

Grouping & Summarizing Treatment

- 도시의 평균 나이가 20대, 30대, 40대 이상인 경우에 Revenue에 차이가 있을까요?

```
Carseats %>%  
  mutate(AgeClass =  
         ifelse(Age < 30, "Twenties",  
                ifelse(Age < 40, "Thirties", "FourtyAbove")))) %>%  
  group_by(AgeClass) %>%  
  summarise(avgRevenue = mean(Revenue))  
  
## # A tibble: 3 x 2  
##   AgeClass     avgRevenue  
##   <chr>          <dbl>  
## 1 FourtyAbove    7.30  
## 2 Thirties       8.26  
## 3 Twenties       7.76
```

Discussion

- 평균 나이에 따른 평균 Revenue의 차이는 무엇을 말해주나요?
- 평균 나이가 포함하지 못하고 있는 정보는 어떤 것이 있나요?
- 어떤 데이터가 있으면 더 좋을까요?
- 그 데이터가 있으면 어떻게 하실 건가요?

30세 이상 비중

Discussion

modern	classic
dplyr	base
Everyday Language	Classic Programming language
사용성	타 언어와 범용적
SQL	Pandas package in Python

- SQL
 - 대용량 데이터 베이스와 통신하는 언어
 - R에서도 sqldf라는 패키지를 사용해서 SQL명령어로 R에서 작업할 수 있음
 - R을 할 줄 아는 사람이 SQL을 배우는데 걸리는 시간 < 1 day
 - (얇은 책 하나만 사서 보시면 됩니다.)
 - 데이터를 보는 눈을 키워 줍니다.
- !데이터를 전처리하는 작업의 소요시간은 전체 프로젝트에서 80% 이상입니다.

Data Transformation with dplyr :: CHEAT SHEET



dplyr functions work with pipes and expect **tidy data**. In tidy data:



Each variable is in its own column



Each observation, or case, is in its own row



`x %>% f(y)` becomes `f(x, y)`

Summarise Cases

These apply **summary functions** to columns to create a new table of summary statistics. Summary functions take vectors as input and return one value (see back).

summary function

- `summarise(.data, ...)`
Compute table of summaries.
`summarise(mtcars, avg = mean(mpg))`
- `count(x, ..., wt = NULL, sort = FALSE)`
Count number of rows in each group defined by the variables in ... Also `tally()`.
`count(iris, Species)`

VARIATIONS

- `summarise_all()` - Apply funs to every column.
- `summarise_at()` - Apply funs to specific columns.
- `summarise_if()` - Apply funs to all cols of one type.

Group Cases

Use `group_by()` to create a "grouped" copy of a table. dplyr functions will manipulate each "group" separately and then combine the results.

- → `mtcars %>% group_by(cyl) %>% summarise(avg = mean(mpg))`

`group_by(.data, ..., add = FALSE)`
Returns copy of table grouped by ...
`g_iris <- group_by(iris, Species)`

`ungroup(x, ...)`
Returns ungrouped copy of table.
`ungroup(g_iris)`

Manipulate Cases

EXTRACT CASES

Row functions return a subset of rows as a new table.

- `filter(.data, ...)` Extract rows that meet logical criteria.
`filter(iris, Sepal.Length > 7)`
- `distinct(.data, ..., .keep_all = FALSE)` Remove rows with duplicate values.
`distinct(iris, Species)`
- `sample_frac(tbl, size = 1, replace = FALSE, weight = NULL, .env = parent.frame())` Randomly select fraction of rows.
`sample_frac(iris, 0.5, replace = TRUE)`
- `sample_n(tbl, size, replace = FALSE, weight = NULL, .env = parent.frame())` Randomly select size rows.
`sample_n(iris, 10, replace = TRUE)`
- `slice(.data, ...)` Select rows by position.
`slice(iris, 10:15)`
- `top_n(x, n, wt)` Select and order top n entries (by group if grouped data).
`top_n(iris, 5, Sepal.Width)`

Logical and boolean operators to use with filter()

<	<=	is.na()	%in%		xor()
>	>=	is.na()	!	&	

See `?base::logic` and `?Comparison` for help.

ARRANGE CASES

- `arrange(.data, ...)` Order rows by values of a column or columns (low to high), use with `desc()` to order from high to low.
`arrange(mtcars, mpg)`
`arrange(mtcars, desc(mpg))`

ADD CASES

- `add_row(.data, ..., .before = NULL, .after = NULL)` Add one or more rows to a table.
`add_row(faithful, eruptions = 1, waiting = 1)`

Manipulate Variables

EXTRACT VARIABLES

Column functions return a set of columns as a new vector or table.

- `pull(.data, var = -1)` Extract column values as a vector. Choose by name or index.
`pull(iris, Sepal.Length)`
- `select(.data, ...)` Extract columns as a table. Also `select_if()`.
`select(iris, Sepal.Length, Species)`

Use these helpers with `select()`, e.g. `select(iris, starts_with("Sepal"))`

<code>contains(match)</code>	<code>num_range(prefix, range)</code>	:	e.g. <code>mpg:cyl</code>
<code>ends_with(match)</code>	<code>one_of(...)</code>	-	e.g. <code>-Species</code>
<code>matches(match)</code>	<code>starts_with(match)</code>		

MAKE NEW VARIABLES

These apply **vectorized functions** to columns. Vectorized funs take vectors as input and return vectors of the same length as output (see back).

vectorized function

- `mutate(.data, ...)` Compute new column(s).
`mutate(mtcars, gpm = 1/mpg)`
- `transmute(.data, ...)` Compute new column(s), drop others.
`transmute(mtcars, gpm = 1/mpg)`
- `mutate_all(.tbl, .funs, ...)` Apply funs to every column. Use with `funs()`. Also `mutate_if()`.
`mutate_all(faithful, funs(log(.), log2(.)))`
`mutate_if(iris, is.numeric, funs(log(.)))`
- `mutate_at(.tbl, .cols, .funs, ...)` Apply funs to specific columns. Use with `funs()`, `vars()` and the helper functions for `select()`.
`mutate_at(iris, vars(-Species), funs(log(.)))`
- `add_column(.data, ..., .before = NULL, .after = NULL)` Add new column(s). Also `add_count()`, `add_tally()`.
`add_column(mtcars, new = 1:32)`
- `rename(.data, ...)` Rename columns.
`rename(iris, Length = Sepal.Length)`



Vector Functions

TO USE WITH MUTATE()

mutate() and **transmute()** apply vectorized functions to columns to create new columns. Vectorized functions take vectors as input and return vectors of the same length as output.

vectorized function

OFFSETS

dplyr::lag() - Offset elements by 1
dplyr::lead() - Offset elements by -1

CUMULATIVE AGGREGATES

dplyr::cumall() - Cumulative all()
dplyr::cumany() - Cumulative any()
dplyr::cummax() - Cumulative max()
dplyr::cummean() - Cumulative mean()
 cummin() - Cumulative min()
 cumprod() - Cumulative prod()
 cumsum() - Cumulative sum()

RANKINGS

dplyr::cume_dist() - Proportion of all values <=
dplyr::dense_rank() - rank with ties = min, no gaps
dplyr::min_rank() - rank with ties = min
dplyr::ntile() - bins into n bins
dplyr::percent_rank() - min_rank scaled to [0,1]
dplyr::row_number() - rank with ties = "first"

MATH

+, -, *, /, ^, %/%, %% - arithmetic ops
log(), log2(), log10() - logs
<, <=, >, >=, !=, == - logical comparisons
dplyr::between() - x >= left & x <= right
dplyr::near() - safe == for floating point numbers

MISC

dplyr::case_when() - multi-case if_else()
dplyr::coalesce() - first non-NA values by element across a set of vectors
dplyr::if_else() - element-wise if() + else()
dplyr::na_if() - replace specific values with NA
 pmax() - element-wise max()
 pmin() - element-wise min()
dplyr::recode() - Vectorized switch()
dplyr::recode_factor() - Vectorized switch() for factors

Summary Functions

TO USE WITH SUMMARISE()

summarise() applies summary functions to columns to create a new table. Summary functions take vectors as input and return single values as output.

summary function

COUNTS

dplyr::n() - number of values/rows
dplyr::n_distinct() - # of uniques
sum(is.na()) - # of non-NA's

LOCATION

mean() - mean, also mean(!is.na())
median() - median

LOGICALS

mean() - Proportion of TRUE's
sum() - # of TRUE's

POSITION/ORDER

dplyr::first() - first value
dplyr::last() - last value
dplyr::nth() - value in nth location of vector

RANK

quantile() - nth quantile
min() - minimum value
max() - maximum value

SPREAD

IQR() - Inter-Quartile Range
mad() - median absolute deviation
sd() - standard deviation
var() - variance

Row Names

Tidy data does not use rownames, which store a variable outside of the columns. To work with the rownames, first move them into a column.

rownames_to_column()
Move row names into col.
a <- rownames_to_column(iris, var = "C")

column_to_rownames()
Move col in row names.
column_to_rownames(a, var = "C")

Also **has_rownames()**, **remove_rownames()**

Combine Tables

COMBINE VARIABLES

X Y
+ =
Z

Use **bind_cols()** to paste tables beside each other as they are.

bind_cols(...) Returns tables placed side by side as a single table.
BE SURE THAT ROWS ALIGN.

Use a "**Mutating Join**" to join one table to columns from another, matching values with the rows that they correspond to. Each join retains a different combination of values from the tables.

left_join(x, y, by = NULL,
copy = FALSE, suffix = c("x", "y"), ...)
Join matching values from y to x.

right_join(x, y, by = NULL, copy =
FALSE, suffix = c("x", "y"), ...)
Join matching values from x to y.

inner_join(x, y, by = NULL, copy =
FALSE, suffix = c("x", "y"), ...)
Join data. Retain only rows with matches.

full_join(x, y, by = NULL,
copy = FALSE, suffix = c("x", "y"), ...)
Join data. Retain all values, all rows.

Use **by = c("col1", "col2")** to specify the column(s) to match on.
left_join(x, y, by = "A")

Use a named vector, by = c("col1" = "col2"), to match on columns with different names in each data set.
left_join(x, y, by = c("C" = "D"))

Use suffix to specify suffix to give to duplicate column names.
left_join(x, y, by = c("C" = "D"), suffix = c("1", "2"))

COMBINE CASES

X Y
+ =
Z

Use **bind_rows()** to paste tables below each other as they are.

bind_rows(..., .id = NULL)
Returns tables one on top of the other as a single table. Set .id to a column name to add a column of the original table names (as pictured)

intersect(x, y, ...)
Rows that appear in both x and y.

setdiff(x, y, ...)
Rows that appear in x but not y.

union(x, y, ...)
Rows that appear in x or y.
(Duplicates removed). union_all()
retains duplicates.

Use **setequal()** to test whether two data sets contain the exact same rows (in any order).

EXTRACT ROWS

X Y
+ =
Z

Use a "**Filtering Join**" to filter one table against the rows of another.

semi_join(x, y, by = NULL, ...)
Return rows of x that have a match in y.
USEFUL TO SEE WHAT WILL BE JOINED.

anti_join(x, y, by = NULL, ...)
Return rows of x that do not have a match in y. USEFUL TO SEE WHAT WILL NOT BE JOINED.



Module 2

ggplot2

(Grammar of Graphics)

“The simple graph has brought more information to the data analyst's mind than any other device.”

(간단한 그래프는 다른 어떤 장치보다 데이터 분석가의 마음에 더 많은 정보를 제공합니다.)

– John Tukey

"Visualization is defined as the interpretation of visualized information and the formation of a mental model of the information (Tan and Steinbach, 2006). Visual representation has gained and is also continuously gaining its popularity more and more due to its summarizing power on the huge-sized data with complicated and elusive attributes. Considering the nature of the studies that utilizes both spatial and temporal data, the effective visual presentation alone is well worth it. Effective visualization not only helps better understanding of data itself and the results of the analysis, but it can also potentially provide intuitive monitoring tools in actual operations."

(시각화는 시각화 된 정보의 해석과 정보의 정신적 모델의 형성으로 정의됩니다. 시각적 표현은 복잡하고 애매한 속성이 있는 거대한 크기의 데이터에 대한 요약 기능으로 인해 점점 더 인기를 얻어왔으며 지속적으로 더 많은 인기를 얻고 있습니다. 공간 데이터와 시간 데이터를 모두 사용하는 연구의 본질을 고려할 때 효과적인 시각적 표현만으로 충분한 가치가 있습니다. 효과적인 시각화는 데이터 자체와 분석 결과를 더 잘 이해할 수 있을 뿐 아니라 잠재적으로 실제 작업에서 직관적인 모니터링 도구를 제공 할 수 있습니다.)

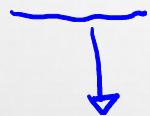
mpg

- Ggplot2 패키지에 내장된 데이터셋
- 1999년 부터 2008년 까지의 38개 차종 연비 데이터

변수 이름	변수 설명
manufacturer	
model	차종 Model name
displ	엔진 크기 engine displacement, in litres 2000 cc
year	연식 year of manufacture
cyl	기통 number of cylinders
trans	트랜스미션 type of transmission
drv	전륜/후륜/사륜 f = front-wheel drive, r = rear wheel drive, 4 = 4wd
cty	도심마일리지 city miles per gallon
hwy	고속도로마일리지 highway miles per gallon
fl	연료 종류 fuel type
class	(세단, SUV...) "type" of car

```
library(ggplot2) +  
? mpg +  
## starting httpd help server ... done +  
  
mpg +  
## # A tibble: 234 x 11  
##   manufacturer model      displ  year   cyl trans  drv   cty   hwy fl  
##   <chr>        <chr>     <dbl> <int> <int> <chr>  <chr> <int> <int> <int> <chr>  
## 1 audi         a4       1.80  1999     4 auto(l~ f      18    29 p  
## 2 audi         a4       1.80  1999     4 manual~ f     21    29 p  
## 3 audi         a4       2.00  2008     4 manual~ f     20    31 p  
## 4 audi         a4       2.00  2008     4 auto(a~ f      21    30 p  
## 5 audi         a4       2.80  1999     6 auto(l~ f      16    26 p  
## 6 audi         a4       2.80  1999     6 manual~ f     18    26 p  
## 7 audi         a4       3.10  2008     6 auto(a~ f      18    27 p  
## 8 audi         a4 quat~  1.80  1999     4 manual~ 4     18    26 p  
## 9 audi         a4 quat~  1.80  1999     4 auto(l~ 4      16    25 p  
## 10 audi        a4 quat~  2.00  2008     4 manual~ 4     20    28 p  
## # ... with 224 more rows, and 1 more variable: class <chr> +  
  
class(mpg) +  
## [1] "tbl_df"     "tbl"        "data.frame" +  
  
# mpg has attributes for multiple classes! +  
# tbl is from package 'tibble' and can be handy +
```

엔진이 크면 연비가 안 좋을까요?



dsppl

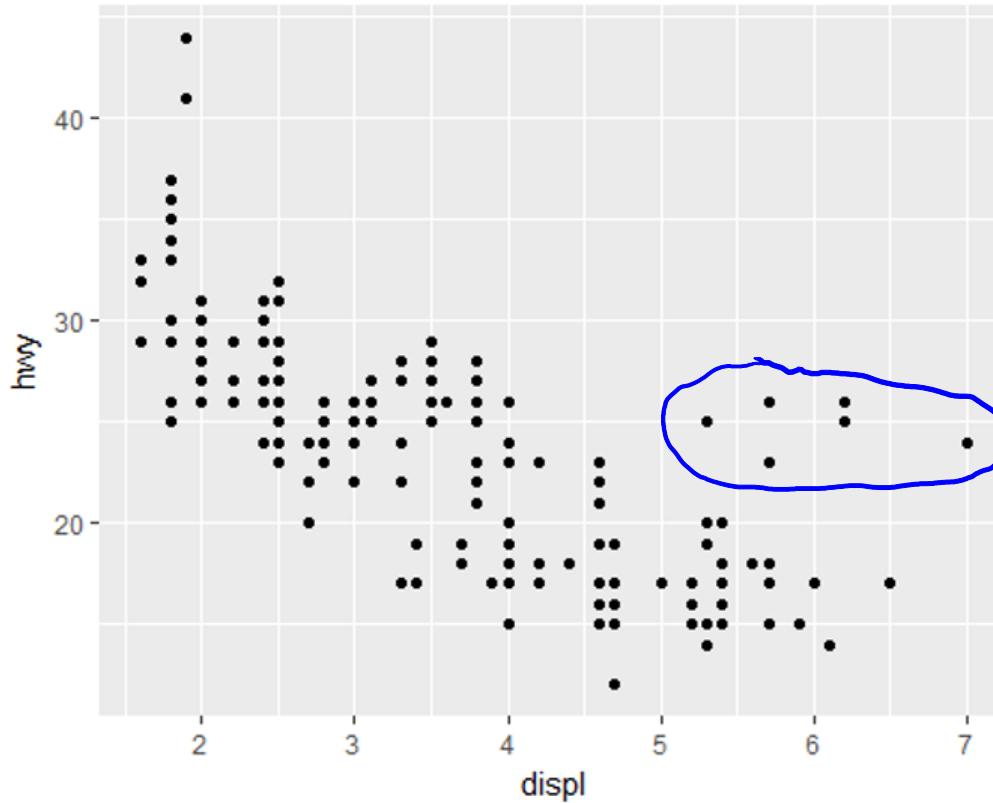


hwy

- mpg 의 dsppl 이 hwy 의 36제
- dsppl → hwy
- X → Y

engine size +
ggplot(data=mpg) +
geom_point(mapping = aes(x = displ, y = hwy)) +

%)>%>

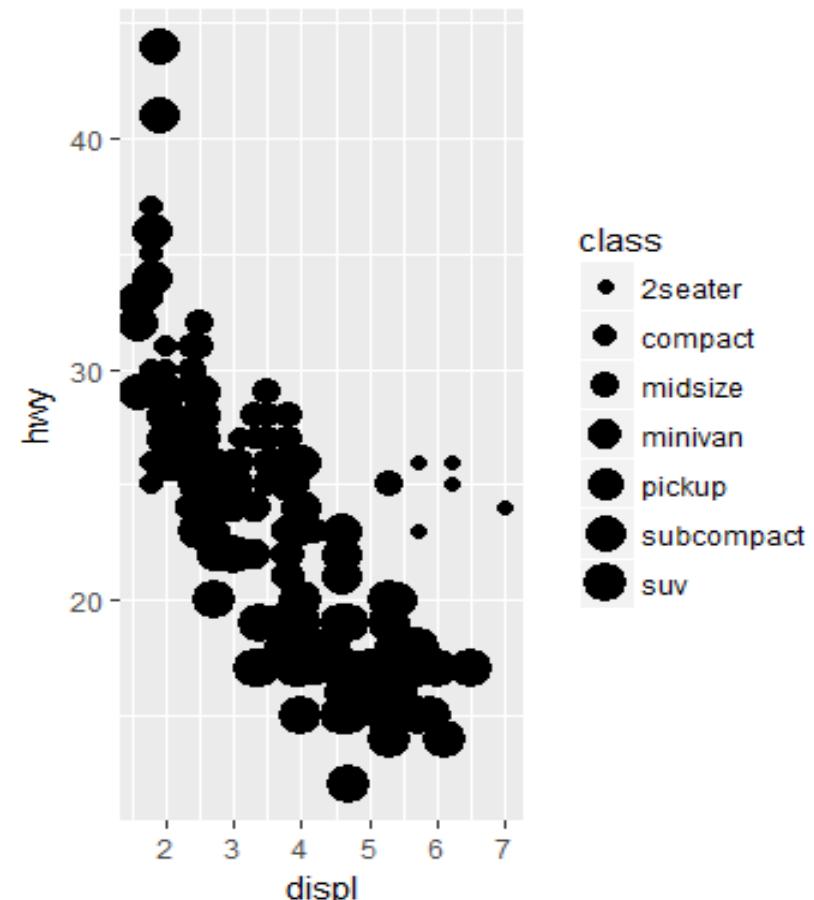
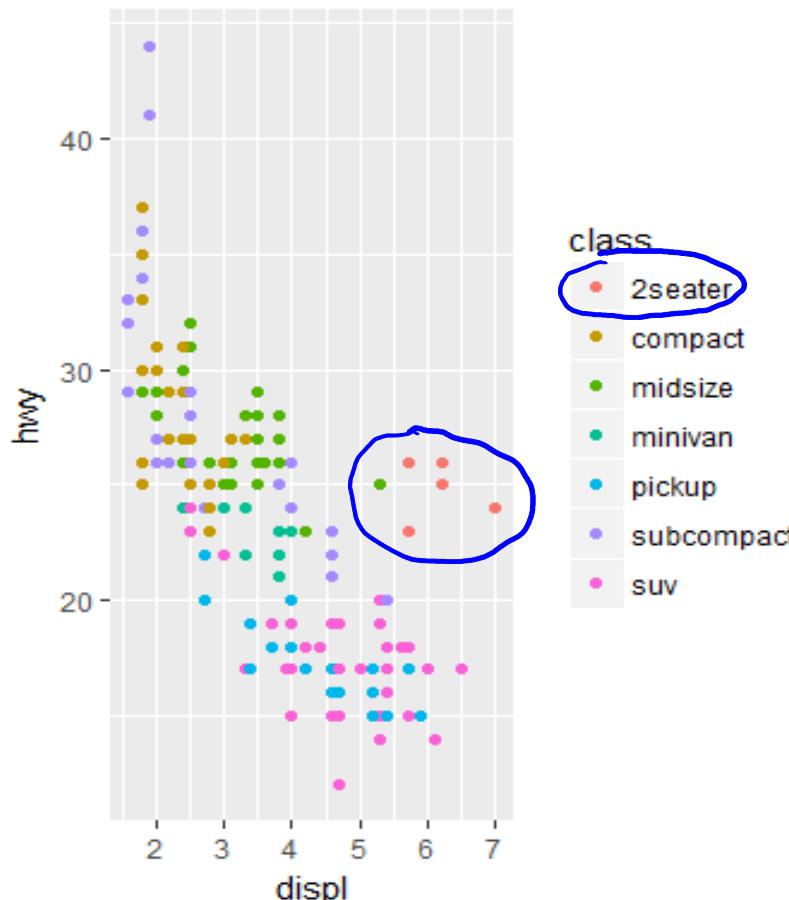


[L2.M2.Fig 1]

template +
ggplot(data = <DATA>) +
<GEOM_FUNCTION>(mapping = aes(< MAPPING >)) +

```
library(gridExtra)  
a <- ggplot(data = mpg) +  
  geom_point(mapping = aes(x = displ, y = hwy, color = class))  
b <- ggplot(data = mpg) +  
  geom_point(mapping = aes(x = displ, y = hwy, size = class))  
grid.arrange(a, b, nrow=1, ncol=2)
```

Warning: Using size for a discrete variable is not advised.



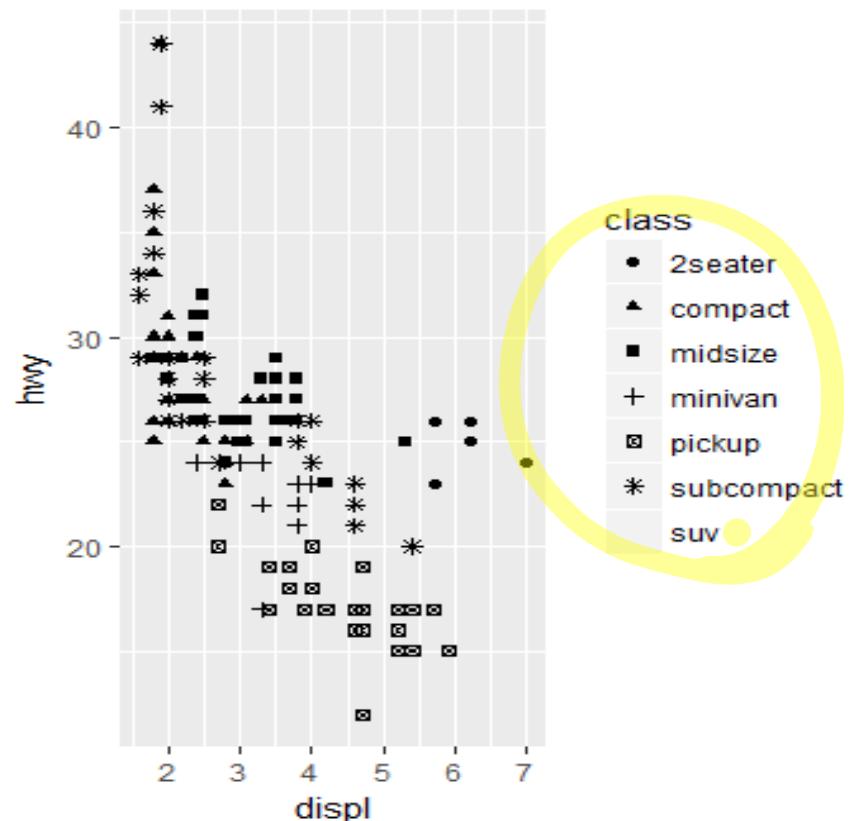
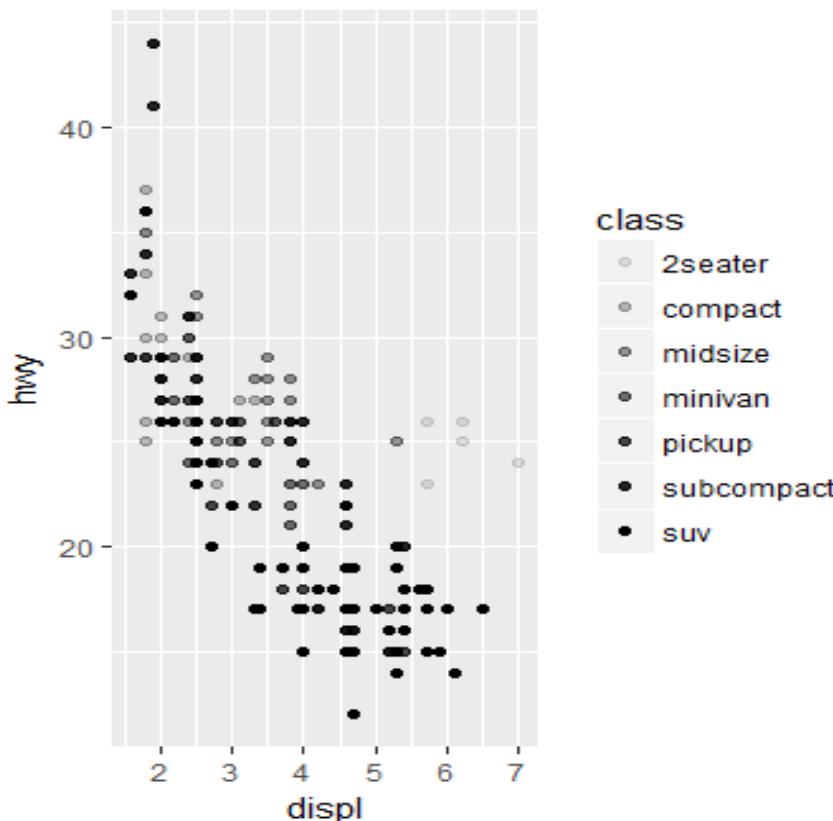
```

a <- ggplot(data = mpg) +
  geom_point(mapping = aes(x = displ, y = hwy, alpha = class))|
b <- ggplot(data = mpg) +
  geom_point(mapping = aes(x = displ, y = hwy, shape = class))|
grid.arrange(a, b, nrow=1, ncol=2)

## Warning: The shape palette can deal with a maximum of 6 discrete values +
## because more than 6 becomes difficult to discriminate; you have 7. +
## Consider specifying shapes manually if you must have them. +

## Warning: Removed 62 rows containing missing values (geom_point). +

```



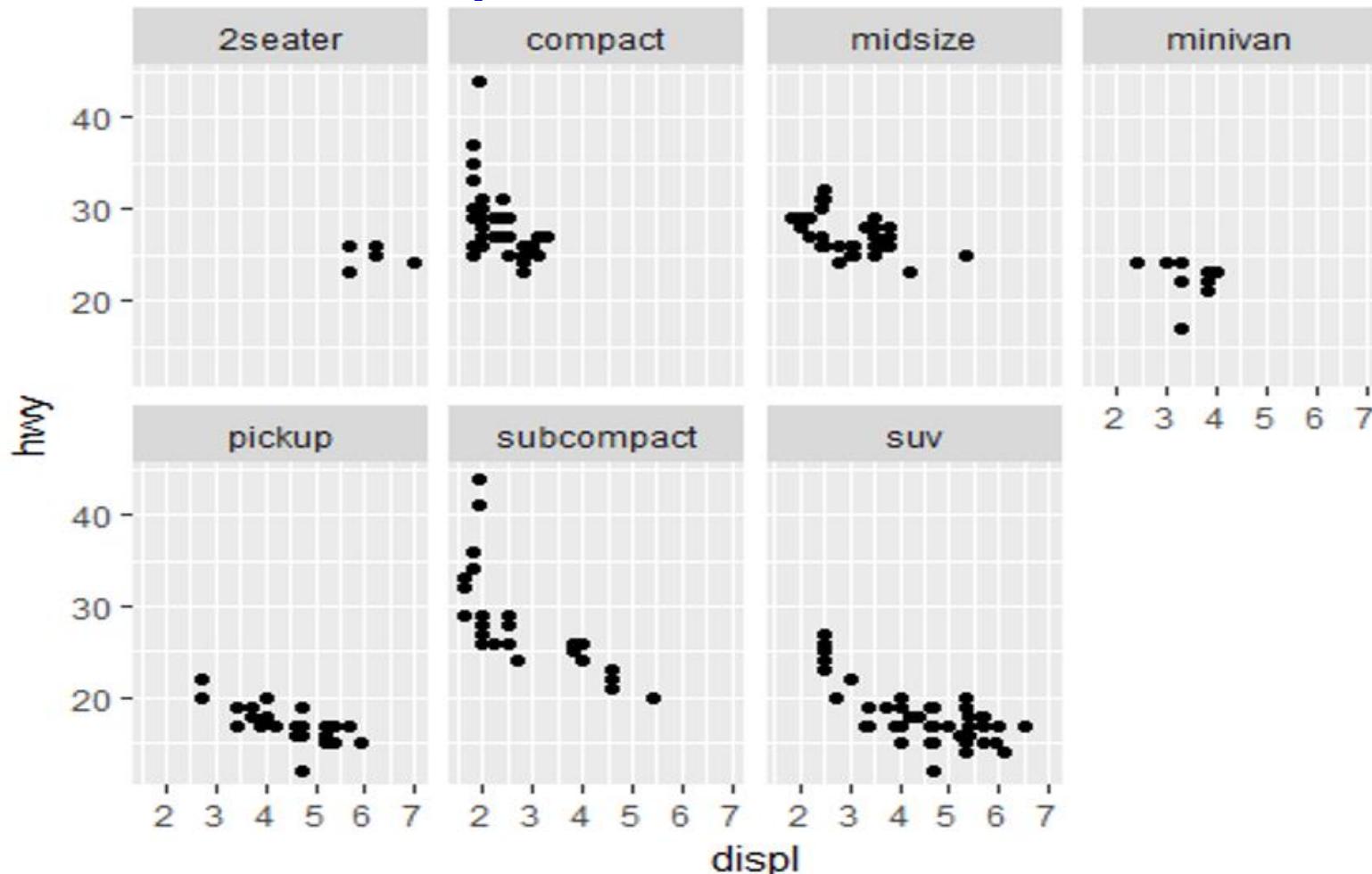
Aesthetics (aes) – size, alpha, color, shape

class

numeric (quantitative) ($\text{Age}, \text{Height}$)	size alpha	continuous or discrete	grayscale 33-75%
factor (qualitative) (Gender) (categorical) (group)	color (≤10) shape (≤5)	discrete	grayscale 75-100%

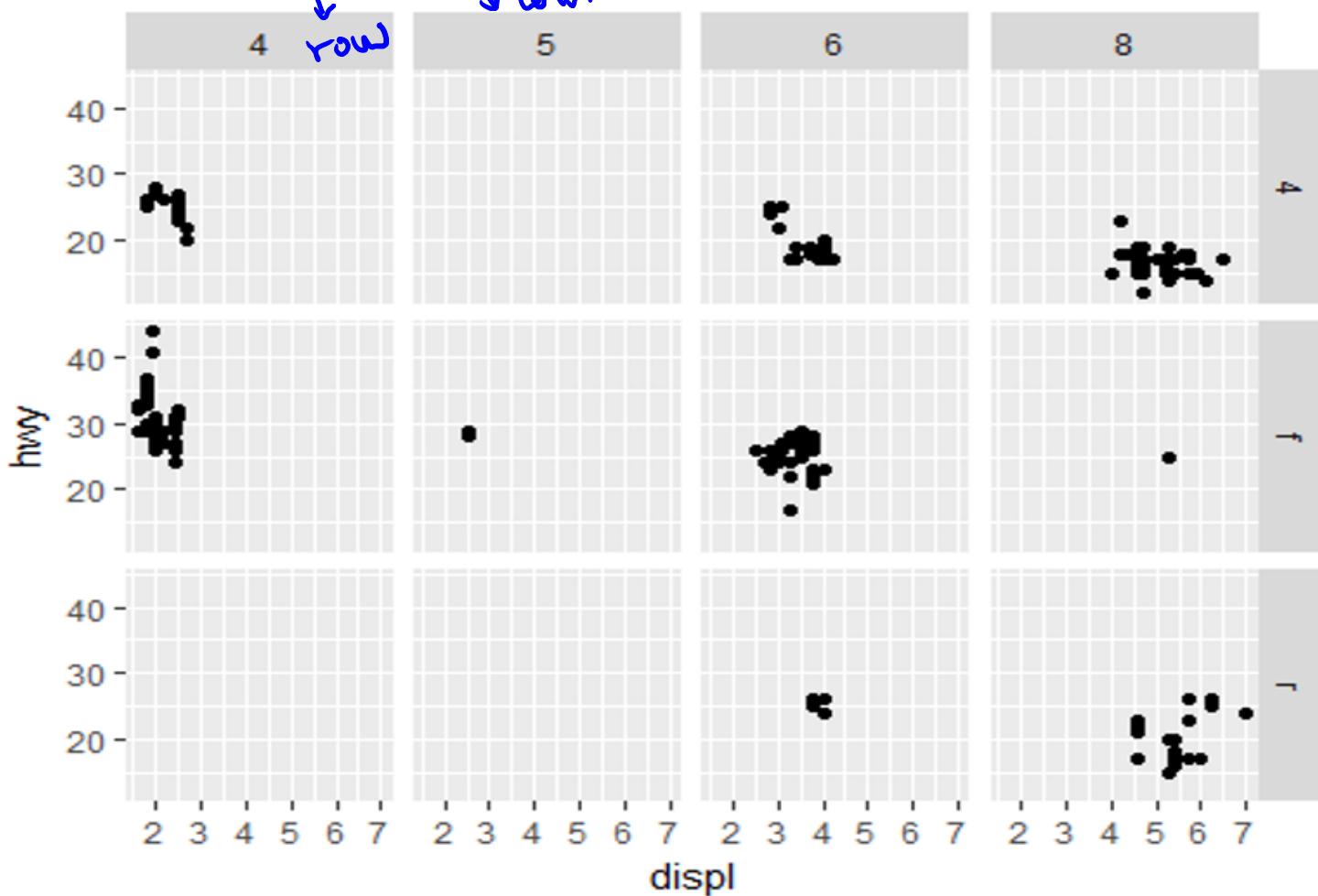
FACETS (1 var)

```
ggplot(data = mpg) +  
  geom_point(mapping = aes(x = displ, y = hwy)) +  
  facet_wrap(~ class, nrow = 2)  
    ↪ column
```



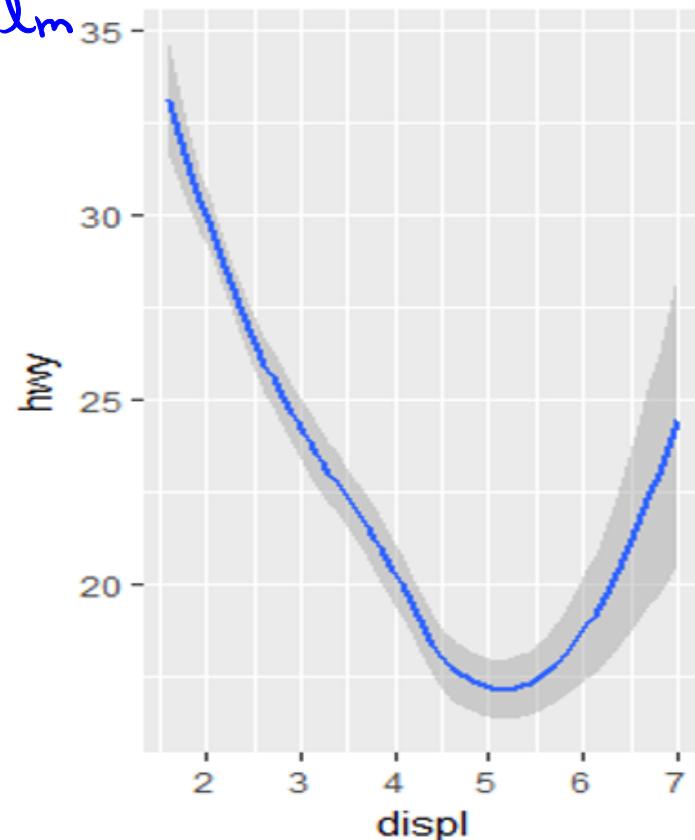
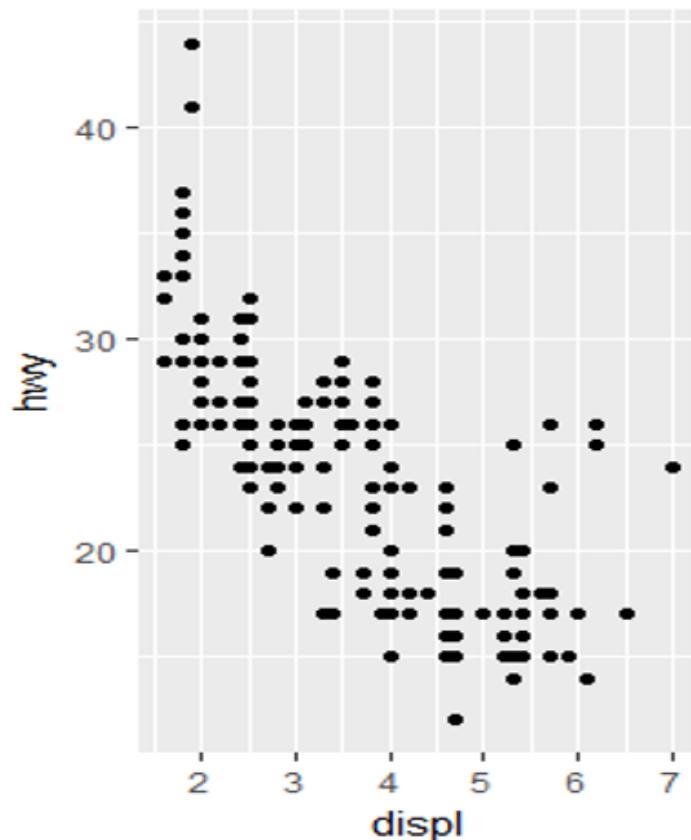
FACETS (2 vars)

```
ggplot(data = mpg) +  
  geom_point(mapping = aes(x = displ, y = hwy)) +  
  facet_grid(drv ~ cyl)
```



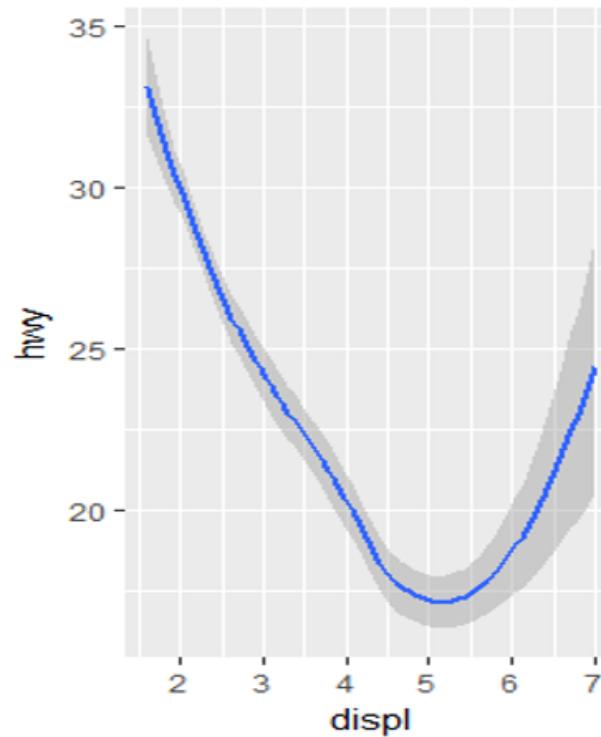
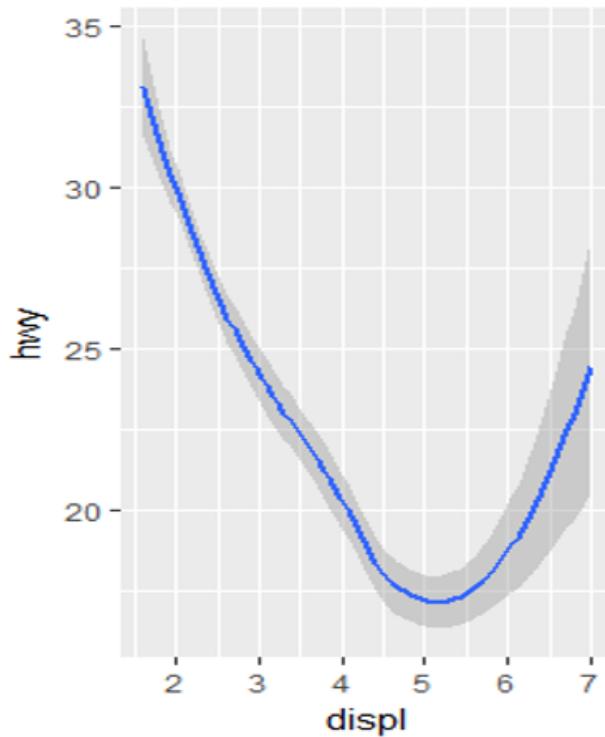
POINT VS SMOOTH

```
# point vs smooth↓  
a <- ggplot(data = mpg) +  
  geom_point(mapping = aes(x = displ, y = hwy))↓  
b <- ggplot(data = mpg) +  
  geom_smooth(mapping = aes(x = displ, y = hwy))↓  
grid.arrange(a, b, nrow=1, ncol=2)↓  
## `geom_smooth()` using method = 'loess'  
          ^ loess  
          ^ lm
```



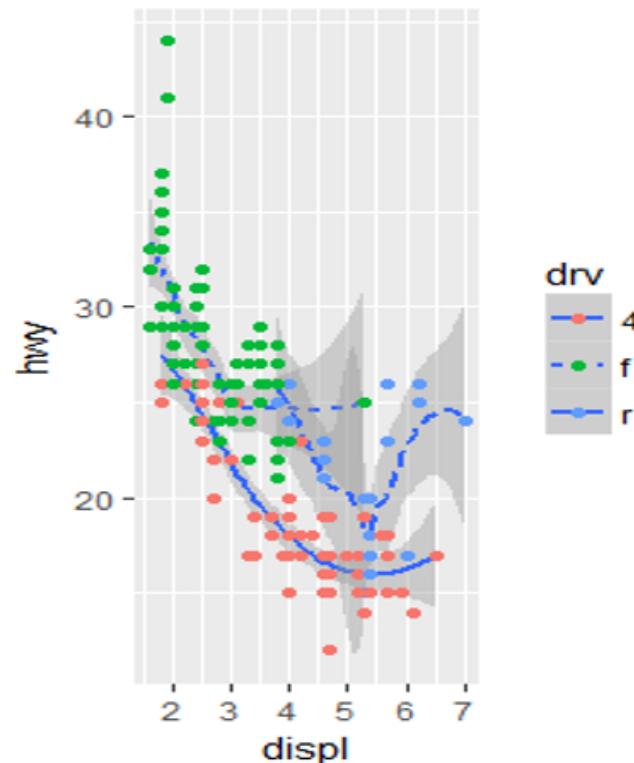
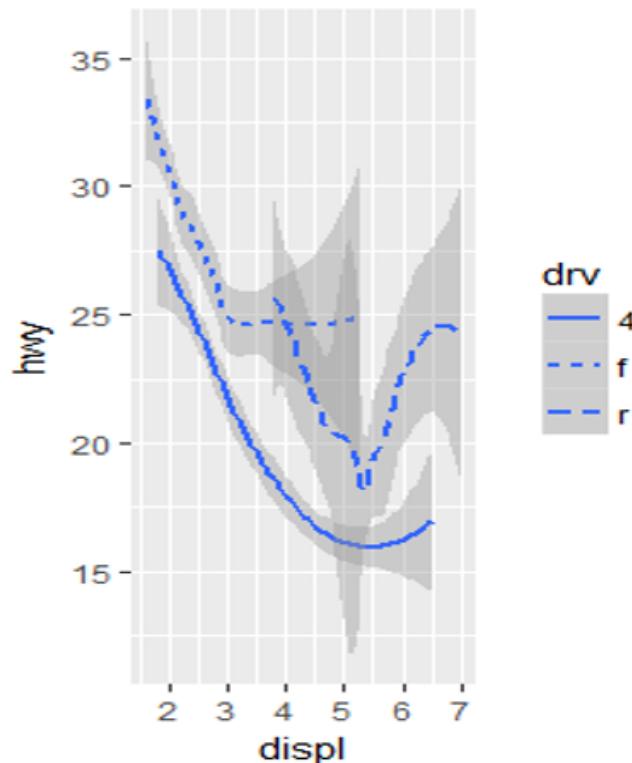
AES MAPPING의 중첩

```
# aes extends or overwrites ↓  
a <- ggplot(data = mpg) +  
  geom_smooth(mapping = aes(x = displ, y = hwy)) ✓  
b <- ggplot(data = mpg, mapping = aes(x = displ, y = hwy)) +  
  geom_smooth()  
grid.arrange(a, b, nrow=1, ncol=2)  
  
## `geom_smooth()` using method = 'loess' ↓  
## `geom_smooth()` using method = 'loess' ↓
```



FINAL

```
# point + smooth ↓  
a <- ggplot(data = mpg, mapping = aes(x = displ, y = hwy)) +  
  geom_smooth(aes(linetype = drv)) +  
b <- ggplot(data = mpg, mapping = aes(x = displ, y = hwy)) +  
  geom_smooth(aes(linetype = drv)) +  
  geom_point(aes(color = drv)) # Line type ignored ↓  
grid.arrange(a, b, nrow=1, ncol=2) +  
## `geom_smooth()` using method = 'loess' ↓  
## `geom_smooth()` using method = 'loess' ↓
```



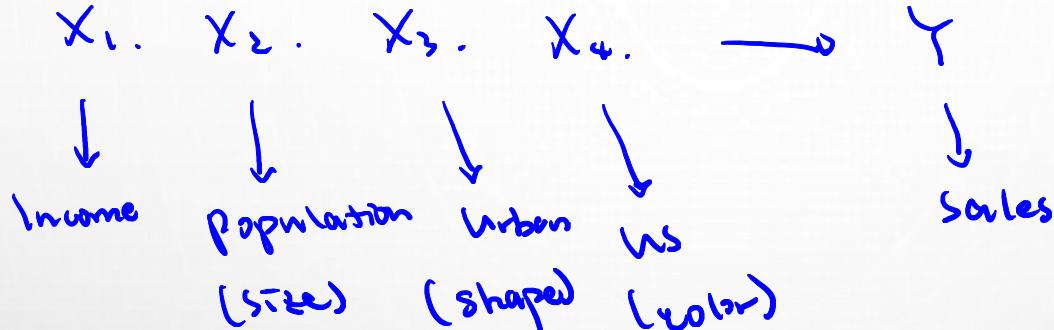
SUMMARY

- *Data, Aesthetic, Geometric Object*를 명시하여 그림을 *rendering*.

- ```
template↓
ggplot(data = <DATA>) +↓
<GEOM_FUNCTION>(mapping = aes(<MAPPING>)) +↓
```

- *grid.arrange* 함수를 이용하여 *grid* 형태로 *positioning* 시킬 수 있음
- *Size > Alpha > Color >= Shape* 의 특성을 이해
- 그림의 목적을 정하고
  - *X (explanatory variable – 설명 변수)*의 개수, 각 변수의 중요도와 순서
  - *Y (dependent variable)*를 명시하여 *ggplot* 객체를 *rendering*

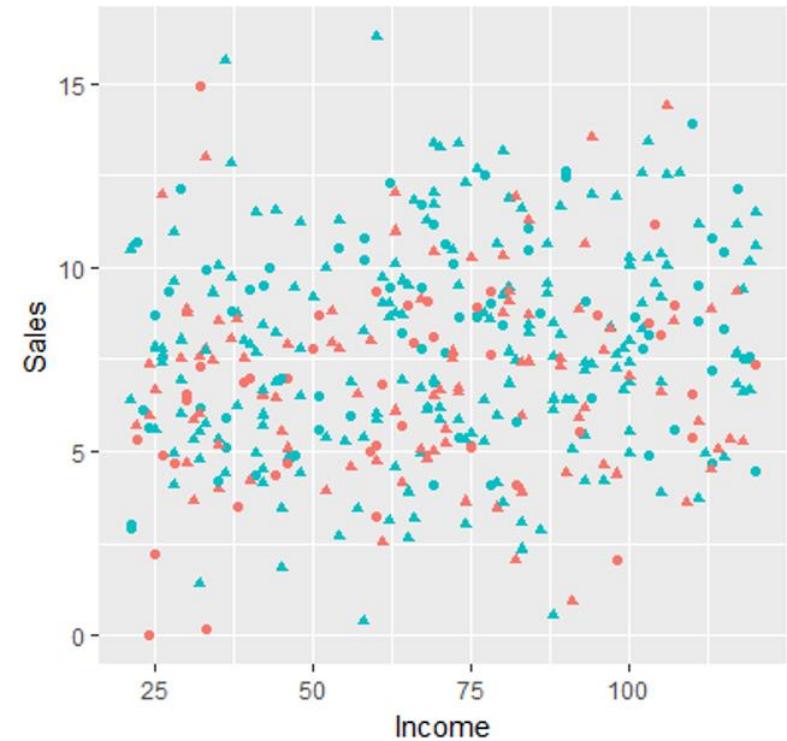
# back to Carseat



```
library(ggplot2)
library(ISLR)
str(Carseats)

'data.frame': 400 obs. of 11 variables:
$ Sales : num 9.5 11.22 10.06 7.4 4.15 ...
$ CompPrice : num 138 111 113 117 141 124 115 136 132 132 ...
$ Income : num 73 48 35 100 64 113 105 81 110 113 ...
$ Advertising: num 11 16 10 4 3 13 0 15 0 0 ...
$ Population : num 276 260 269 466 340 501 45 425 108 131 ...
$ Price : num 120 83 80 97 128 72 108 120 124 124 ...
$ ShelveLoc : Factor w/ 3 levels "Bad","Good","Medium": 1 2 3 3 1 1 3 2
$ Age : num 42 65 59 55 38 78 71 67 76 76 ...
$ Education : num 17 10 12 14 13 16 15 10 10 17 ...
$ Urban : Factor w/ 2 levels "No","Yes": 2 2 2 2 2 1 2 2 1 1 ...
$ US : Factor w/ 2 levels "No","Yes": 2 2 2 2 1 2 1 2 1 2 ...
```

```
a <- ggplot(data = Carseats, aes(x = Income, y = Sales)) +
 geom_point(aes(shape = Urban, color = US)) +
 print(a)
```

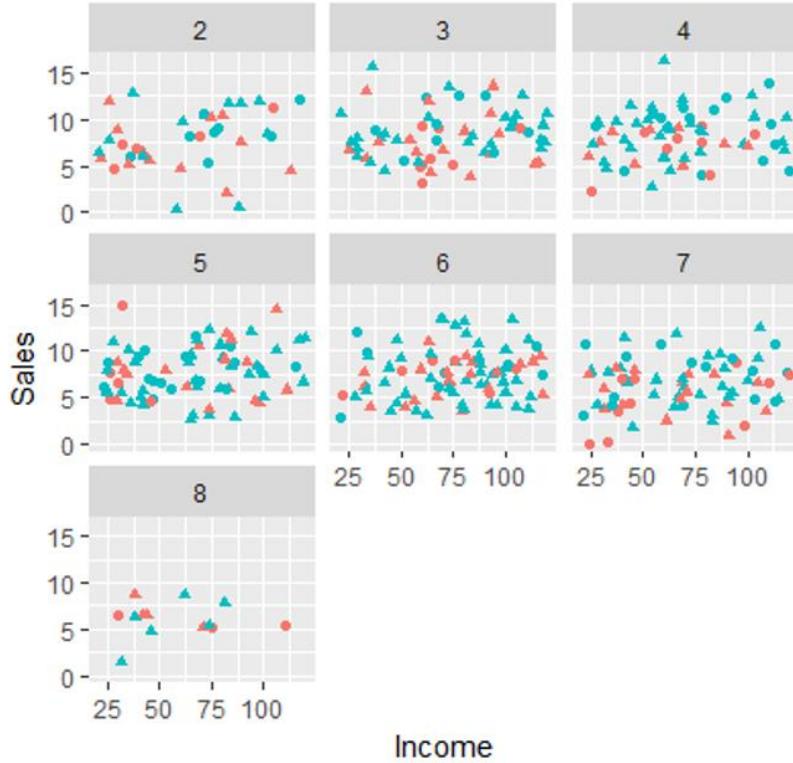


size = Population

```
a <- a + facet_wrap(~ floor(Age/10)) +
print(a)
```

US  
• No  
• Yes

Urban  
• No  
▲ Yes



US  
• No  
• Yes

Urban  
• No  
▲ Yes

```
doFacetWrap <- TRUE +
a <- ggplot(data = Carseats, aes(x = Income, y = Sales)) +
 geom_point(aes(shape = Urban, color = US)) +
if (doFacetWrap) {
 a <- a + facet_wrap(~ floor(Age/10)) +
}
print(a)
```

# preview

Rich Country = Live Longer??

Final Project: [ggplot](#) + [API](#) + [rmarkdown](#) + [shiny](#) + [flexdashboard](#)

[/](#) Source Code

OECD

OECD

non-OECD

Continents

Asia

Europe

Africa

North America

South America

Oceania

Range of Country GDP (USD)



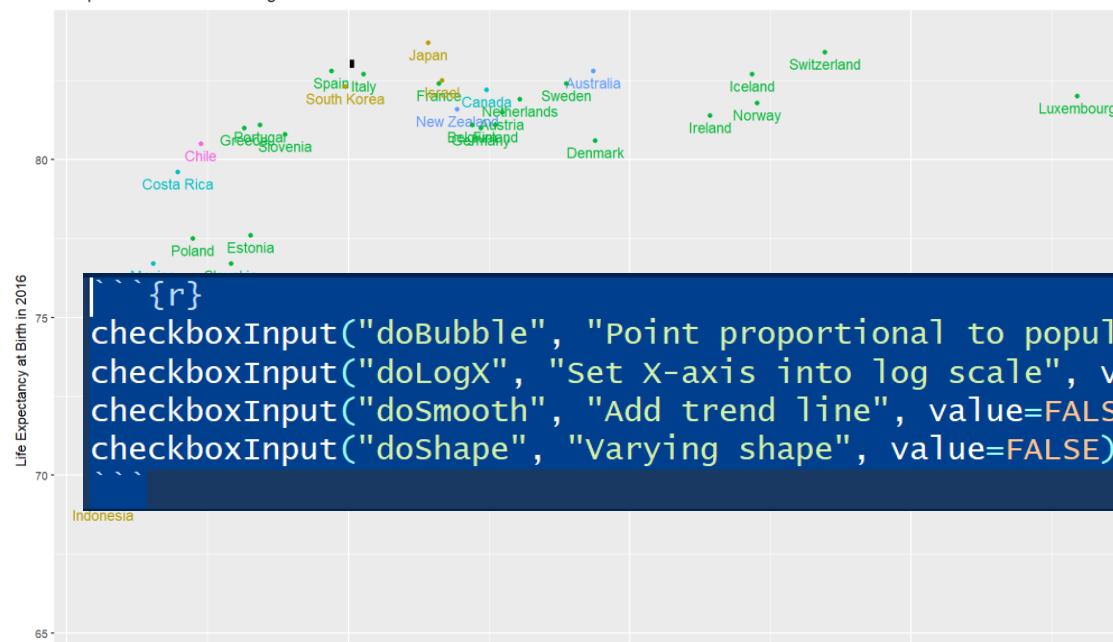
Point proportional to population

Set X-axis into log scale

Add trend line

Varying shape

2D space cannot limit our imagination!



```
| myGgplot <- ggplot(lifeCountry,
| aes(x=GDP_per_Capita, y=Life.Expectancy, color=Continent)) +
| geom_text(aes(label=Country), size=4, vjust=1.5) +
| geom_point() +
| xlab("GDP per Capita in 2017 (in US Dollars)") +
| ylab("Life Expectancy at Birth in 2016") +
| ggtitle("2D space cannot limit our imagination!")
| if (input$doBubble) { myGgplot <- myGgplot + geom_point(aes(size=Population)) }
| if (input$doLogX) { myGgplot <- myGgplot + scale_x_log10() + xlab("GDP per Capita (log-scale)")}
| if (input$doSmooth) { myGgplot <- myGgplot + geom_smooth(method='lm', se=TRUE, size=1) }
| if (input$doShape) { myGgplot <- myGgplot + geom_point(aes(shape=Continent))}
```





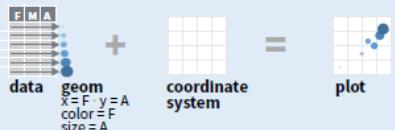
# Data Visualization with ggplot2 :: CHEAT SHEET

## Basics

ggplot2 is based on the **grammar of graphics**, the idea that you can build every graph from the same components: a **data set**, a **coordinate system**, and **geoms**—visual marks that represent data points.



To display values, map variables in the data to visual properties of the geom (**aesthetics**) like **size**, **color**, and **x** and **y** locations.



Complete the template below to build a graph.

```
ggplot (data = <DATA>) +
 <GEOM_FUNCTION> (mapping = aes(<MAPPINGS>),
 stat = <STAT>, position = <POSITION>) +
 <COORDINATE_FUNCTION> +
 <FACET_FUNCTION> +
 <SCALE_FUNCTION> +
 <THEME_FUNCTION>
```

required

Not required, sensible defaults supplied

ggplot(data = mpg, aes(x = cyl, y = hwy)) Begins a plot that you finish by adding layers to. Add one geom function per layer.

**aesthetic mappings**   **data**   **geom**

qplot(x = cyl, y = hwy, data = mpg, geom = "point") Creates a complete plot with given data, geom, and mappings. Supplies many useful defaults.

last\_plot() Returns the last plot

ggsave("plot.png", width = 5, height = 5) Saves last plot as 5' x 5' file named "plot.png" in working directory. Matches file type to file extension.



## Geoms

Use a geom function to represent data points, use the geom's aesthetic properties to represent variables. Each function returns a layer.

### GRAPHICAL PRIMITIVES

```
a <- ggplot(economics, aes(date, unemploy))
b <- ggplot(seals, aes(x = long, y = lat))

a + geom_blank()
 (Useful for expanding limits)

b + geom_curve(aes(yend = lat + 1,
 xend = long + 1, curvature = z)) -> x, yend, y, yend,
 alpha, angle, color, curvature, linetype, size

a + geom_path(lineend = "butt", linejoin = "round",
 linemetre = 1)
 x, y, alpha, color, group, linetype, size

a + geom_polygon(aes(group = group))
 x, y, alpha, color, fill, group, linetype, size

b + geom_rect(aes(xmin = long, ymin = lat, xmax =
 long + 1, ymax = lat + 1)) -> xmax, xmin, ymax,
 ymin, alpha, color, fill, linetype, size

a + geom_ribbon(aes(ymin = unemploy - 900,
 ymax = unemploy + 900)) -> x, ymax, ymin,
 alpha, color, fill, group, linetype, size
```

### LINE SEGMENTS

common aesthetics: x, y, alpha, color, linetype, size

```
b + geom_abline(aes(intercept = 0, slope = 1))
b + geom_hline(aes(yintercept = lat))
b + geom_vline(aes(xintercept = long))

b + geom_segment(aes(yend = lat + 1, xend = long + 1))
b + geom_spoke(aes(angle = 1:155, radius = 1))
```

### ONE VARIABLE continuous

```
c <- ggplot(mpg, aes(hwy)); c2 <- ggplot(mpg)

c + geom_area(stat = "bin")
 x, y, alpha, color, fill, linetype, size

c + geom_density(kernel = "gaussian")
 x, y, alpha, color, fill, group, linetype, size, weight

c + geom_dotplot()
 x, y, alpha, color, fill

c + geom_freqpoly()
 x, y, alpha, color, group, linetype, size

c + geom_histogram(binwidth = 5)
 x, y, alpha, color, fill, linetype, size, weight

c2 + geom_qq(aes(sample = hwy))
 x, y, alpha, color, fill, linetype, size, weight
```

### discrete

```
d <- ggplot(mpg, aes(f1))

d + geom_bar()
 x, alpha, color, fill, linetype, size, weight
```

### TWO VARIABLES

#### continuous x , continuous y

```
e + geom_label(aes(label = cty), nudge_x = 1,
 nudge_y = 1, check_overlap = TRUE) x, y, label,
 alpha, angle, color, family, fontface, hjust,
 lineheight, size, vjust

e + geom_l jitter(height = 2, width = 2)
 x, y, alpha, color, fill, shape, size

e + geom_point(), x, y, alpha, color, fill, shape,
 size, stroke

e + geom_quantile(), x, y, alpha, color, group,
 linetype, size, weight

e + geom_rug(sides = "bl") x, y, alpha, color,
 linetype, size

e + geom_smooth(method = lm) x, y, alpha,
 color, fill, group, linetype, size, weight

e + geom_text(aes(label = cty), nudge_x = 1,
 nudge_y = 1, check_overlap = TRUE) x, y, label,
 alpha, angle, color, family, fontface, hjust,
 lineheight, size, vjust
```

#### discrete x , continuous y

```
f <- ggplot(mpg, aes(class, hwy))

f + geom_col()
 x, y, alpha, color, fill, group, linetype, size

f + geom_boxplot()
 x, y, lower, middle, upper,
 ymax, ymin, alpha, color, fill, group, linetype,
 shape, size, weight

f + geom_dotplot(binaxis = "y", stackdir =
 "center") x, y, alpha, color, fill, group

f + geom_violin(scale = "area") x, y, alpha, color,
 fill, group, linetype, size, weight
```

#### discrete x , discrete y

```
g <- ggplot(diamonds, aes(cut, color))

g + geom_count()
 x, y, alpha, color, fill, shape,
 size, stroke
```

### THREE VARIABLES

```
seals$z <- with(seals, sqrt(delta_long^2 + delta_lat^2)) l <- ggplot(seals, aes(long, lat))
l + geom_contour(aes(z = z))
x, y, z, alpha, colour, group, linetype,
size, weight

l + geom_raster(aes(fill = z), hjust = 0.5, vjust = 0.5,
 interpolate = FALSE) x, y, alpha, fill

l + geom_tile(aes(fill = z)) x, y, alpha, color, fill,
 linetype, size, width
```

### continuous bivariate distribution

```
h <- ggplot(diamonds, aes(carat, price))

h + geom_bln2d(binwidth = c(0.25, 500))
 x, y, alpha, color, fill, linetype, size, weight

h + geom_density2d()
 x, y, alpha, colour, group, linetype, size

h + geom_hex()
 x, y, alpha, colour, fill, size
```

### continuous function

```
i <- ggplot(economics, aes(date, unemploy))

i + geom_area()
 x, y, alpha, color, fill, linetype, size

i + geom_line()
 x, y, alpha, color, group, linetype, size

i + geom_step(direction = "hv")
 x, y, alpha, color, group, linetype, size
```

### visualizing error

```
df <- data.frame(grp = c("A", "B"), fit = 4:5, se = 1:2)
j <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))

j + geom_crossbar(fatten = 2)
 x, y, ymax, ymin, alpha, color, fill, group, linetype,
 size

j + geom_errorbar()
 x, y, max, min, alpha, color, group, linetype, size
 width (also
 geom_errorbarh())

j + geom_linerange()
 x, ymin, ymax, alpha, color, group, linetype, size

j + geom_pointrange()
 x, y, min, max, alpha, color, fill, group, linetype,
 size, weight
```

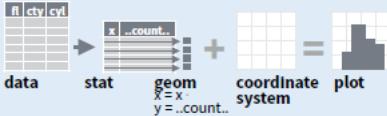
### maps

```
data <- data.frame(murder = USArrests$Murder,
 state = tolower(rownames(USArrests)))
map <- map_data("state")
k <- ggplot(data, aes(fill = murder))

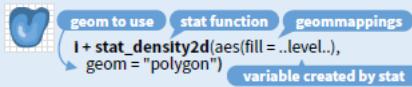
k + geom_map(aes(map_id = state), map = map)
 + expand_limits(x = map$long, y = map$lat),
 map_id, alpha, color, fill, linetype, size
```

## Stats An alternative way to build a layer

A stat builds new variables to plot (e.g., count, prop).



Visualize a stat by changing the default stat of a geom function, `geom_bar(stat="count")` or by using a stat function, `stat_count(geom="bar")`, which calls a default geom to make a layer (equivalent to a geom function). Use `..name..` syntax to map stat variables to aesthetics.



```

c + stat_bin(binwidth = 1, origin = 10)
x, y | ..count., ..ncount., ..density., ..ndensity..
c + stat_count(width = 1) x, y, | ..count., ..prop..
c + stat_density(adjust = 1, kernel = "gaussian")
x, y | ..count., ..density., ..scaled..
e + stat_bin_2d(bins = 30, drop = T)
x, y, fill | ..count., ..density..
e + stat_bin_hex(bins = 30) x, y, fill | ..count., ..density..
e + stat_density_2d(contour = TRUE, n = 100)
x, y, color, size | ..level..
e + stat_ellipse(level = 0.95, segments = 51, type = "t")

```

```

l + stat_contour(aes(z = z)) x, y, z, order | ..level..
l + stat_summary_hex(aes(z = z), bins = 30, fun = max)
x, y, z, fill | ..value..
l + stat_summary_2d(aes(z = z), bins = 30, fun = mean)
x, y, z, fill | ..value..
f + stat_boxplot(coef = 1.5) x, y | ..lower.,.
..middle.,..upper.,..width.,..ymin.,..ymax..
f + stat_ydensity(kernel = "gaussian", scale = "area") x, y |
..density.,..scaled.,..count.,..n.,..violinwidth.,..width..

```

```

e + stat_ecdf(n = 40) x, y | ..x.,..y..
e + stat_quantile(quartiles = c(0.1, 0.9), formula = y ~ log(x), method = "rq") x, y | ..quantile..
e + stat_smooth(method = "lm", formula = y ~ x, se = T, level = 0.95) x, y | ..se.,..x.,..y.,..ymin.,..ymax..

```

```

ggplot() + stat_function(aes(x = -3:3), n = 99, fun = dnorm, args = list(sd = 0.5)) x | ..x.,..y..
e + stat_identity(na.rm = TRUE)
ggplot() + stat_qq(aes(sample = 1:100), dist = qt, dparam = list(df = 5)) sample, x, y | ..sample.,..theoretical..
e + stat_sum() x, y, size | ..n.,..prop..
e + stat_summary(fun.data = "mean_cl_boot")
h + stat_summary_bin(fun.y = "mean", geom = "bar")
e + stat_unique()

```

## Scales

Scales map data values to the visual values of an aesthetic. To change a mapping, add a new scale.



### GENERAL PURPOSE SCALES

Use with most aesthetics

```

scale_*_continuous() - map cont' values to visual ones
scale_*_discrete() - map discrete values to visual ones
scale_*_identity() - use data values as visual ones
scale_*_manual(values = c()) - map discrete values to manually chosen visual ones
scale_*_date(date_labels = "%m/%d", date_breaks = "2 weeks") - treat data values as dates.
scale_*_datetime() - treat data x values as date times. Use same arguments as scale_x_date(). See ?strptime for label formats.

```

### X & Y LOCATION SCALES

Use with x or y aesthetics (x shown here)

```

scale_x_log10() - Plot x on log10 scale
scale_x_reverse() - Reverse direction of x axis
scale_x_sqrt() - Plot x on square root scale

```

### COLOR AND FILL SCALES (DISCRETE)

```

n <- d + geom_bar(aes(fill = fl))
n + scale_fill_brewer(palette = "Blues")
For palette choices: RColorBrewer::display.brewer.all()
n + scale_fill_grey(start = 0.2, end = 0.8, na.value = "red")

```

### COLOR AND FILL SCALES (CONTINUOUS)

```

o <- c + geom_dotplot(aes(fill = ..x..))
o + scale_fill_distiller(palette = "Blues")
o + scale_fill_gradient(low = "red", high = "yellow")
o + scale_fill_gradient2(low = "red", high = "blue",
mid = "white", midpoint = 25)
o + scale_fill_gradientn(colours = topo.colors(6))
Also: rainbow(), heat.colors(), terrain.colors(),
cm.colors(), RColorBrewer::brewer.pal()

```

### SHAPE AND SIZE SCALES

```

p <- e + geom_point(aes(shape = fl, size = cyl))
p + scale_shape() + scale_size()
p + scale_shape_manual(values = c(3:7))
0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25
p + scale_radius(range = c(1,6))
p + scale_size_area(max_size = 6)

```

## Coordinate Systems

`r <- d + geom_bar()`

```

r + coord_cartesian(xlim = c(0, 5))
xlm, ylim
The default cartesian coordinate system
r + coord_fixed(ratio = 1/2)
ratio, xlim, ylim
Cartésian coordinates with fixed aspect ratio between x and y units
r + coord_flip()
xlim, ylim
Flipped Cartesian coordinates
r + coord_polar(theta = "x", direction = 1)
theta, start, direction
Polar coordinates
r + coord_trans(xtrans = "sqrt")
xtrans, ytrans, xlim, ylim
Transformed cartesian coordinates. Set xtrans and ytrans to the name of a window function.

```

`r + coord_quickmap()`

```

r + coord_map(projection = "ortho"
orientation = c(41, -74, 0)) projection, orientation,
xlim, ylim
Map projections from the mapproj package
(mercator (default), aezqualaea, lagrange, etc.)

```

## Position Adjustments

Position adjustments determine how to arrange geoms that would otherwise occupy the same space.

```

s <- ggplot(mpg, aes(fl, fill = drv))
s + geom_bar(position = "dodge")
Arrange elements side by side
s + geom_bar(position = "fill")
Stack elements on top of one another, normalize height
e + geom_point(position = "jitter")
Add random noise to X and Y position of each element to avoid overplotting
e + geom_label(position = "nudge")
Nudge labels away from points
s + geom_bar(position = "stack")
Stack elements on top of one another

```

Each position adjustment can be recast as a function with manual width and height arguments

`s + geom_bar(position = position_dodge(width = 1))`



## Faceting

Facets divide a plot into subplots based on the values of one or more discrete variables.

`t <- ggplot(mpg, aes(cty, hwy)) + geom_point()`

```

t + facet_grid(~ fl)
facet into columns based on fl
t + facet_grid(year ~)
facet into rows based on year
t + facet_grid(year ~ fl)
facet into both rows and columns
t + facet_wrap(~ fl)
wrap facets into a rectangular layout

```

Set scales to let axis limits vary across facets

```

t + facet_grid(drv ~ fl, scales = "free")
x and y axis limits adjust to individual facets
"free_x" - x axis limits adjust
"free_y" - y axis limits adjust

```

Set labeller to adjust facet labels

|                                                                |            |            |            |            |            |
|----------------------------------------------------------------|------------|------------|------------|------------|------------|
| t + facet_grid(. ~ fl, labeller = label_both)                  | fl:c       | fl:d       | fl:e       | fl:p       | fl:r       |
| t + facet_grid(drv ~ , labeller = label_bquote(alpha ^ .(fl))) | $\alpha^c$ | $\alpha^d$ | $\alpha^e$ | $\alpha^p$ | $\alpha^r$ |
| t + facet_grid(. ~ fl, labeller = label_parsed)                | c          | d          | e          | p          | r          |

## Labels

`t + labs(x = "New x axis label", y = "New y axis label", title = "Add a title above the plot", subtitle = "Add a subtitle below title", caption = "Add a caption below plot", <AES> = "New <AES> legend title")`

Use scale functions to update legend labels

`geom to place` manual values for geom's aesthetics

## Legends

`n + theme(legend.position = "bottom")`  
Place legend at "bottom", "top", "left", or "right"

`n + guides(fill = "none")`  
Set legend type for each aesthetic: colorbar, legend, or none (no legend)

`n + scale_fill_discrete(name = "Title", labels = c("A", "B", "C", "D", "E"))`  
Set legend title and labels with a scale function.

## Zooming

Without clipping (preferred)  
`t + coord_cartesian(xlim = c(0, 100), ylim = c(10, 20))`

With clipping (removes unseen data points)  
`t + xlim(0, 100) + ylim(10, 20)`

`t + scale_x_continuous(limits = c(0, 100)) + scale_y_continuous(limits = c(0, 100))`



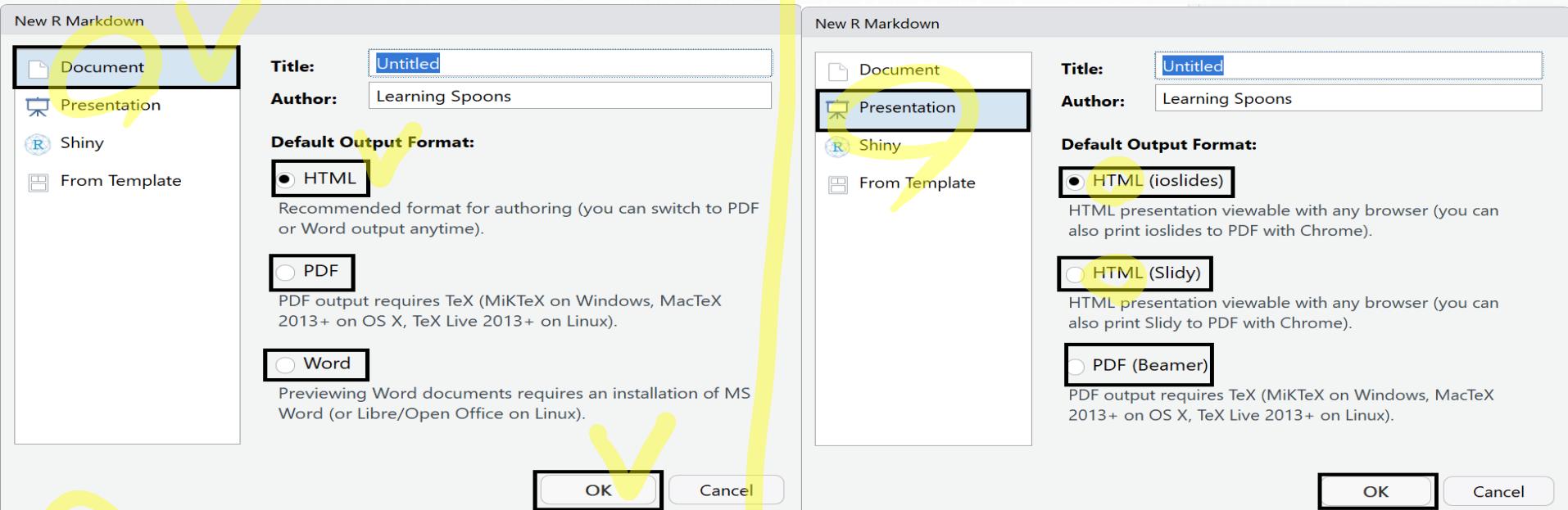


# Module 3

*rmarkdown*

(New Gen. of Computer Programming)  
(Literate Programming)

# 파일 → 새파일 → R Markdown



- html
  - Interactive Feature 가능
- PDF
  - Professional 문서
  - Texlive가 설치되어 있어야 함.
  - 한글을 위해서 별도의 template 사용
- Word
  - Office 없이도 문서 만들 수 있음
- ioslide, slidy
  - Interactive Feature 가능
- PDF (beamer)
  - Scientific Presentation
  - Texlive가 설치되어 있어야 함.
  - 한글을 위해서 별도의 template 사용

# html

R Studio

File Edit Code View Plots Session Build Debug Profile Tools Help

Untitled1 x a.Rmd x

Knit

Untitled

Learning Spoons  
2018-04-17

R Markdown

For more details on using R Markdown see <http://rmarkdown.rstudio.com>. When you click the **Knit** button a document will be generated.

summary(cars)

|            | speed | dist           |
|------------|-------|----------------|
| ## Min.    | 4.0   | 2.00           |
| ## 1st Qu. | 12.0  | 1st Qu.: 26.00 |
| ## Median  | 15.0  | Median : 36.00 |
| ## Mean    | 15.4  | Mean : 42.98   |
| ## 3rd Qu. | 19.0  | 3rd Qu.: 56.00 |
| ## Max.    | 25.0  | Max. : 120.00  |

인블럭 R코드는 `r ncol(cars)` 와 같이 사용.

## Including Plots

You can also embed plots, for example:

```{r pressure, echo=FALSE}  
plot(pressure)

pdf

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Untitled1 x a.Rmd x rmd_pdfRmd x

Knit Insert Run

```
1 ---  
2 title: "rmd 한글 template"  
3 author: "Learning Spoons"  
4 date: "`r Sys.Date()`"  
5 output:  
6   pdf_document:  
7     latex_engine: xelatex  
8     keep_tex: true  
9     # pandoc_args: [  
10       # "-v", "classoption=twocolumn"  
11     # ]  
12     smaller: true  
13     mainfont: NanumGothic  
14     classoption: a4paper  
15 ---  
16 ``{r setup, include=FALSE}  
17 knitr::opts_chunk$set(echo = TRUE)  
18 ``  
19 ## R Markdown 한글  
20  
21 ## R Markdown 한글  
22  
23 Code의 9,10,11 번째 라인의 pound sign(#)을 제거하고 위의 줄과 indent를  
24 맞추면 2 컬럼의 문서가 render 됩니다. <http://rmarkdown.rstudio.com>.  
25 when you click the **Knit** button a document will be generated.  
26 ``{r cars}  
27 summary(cars)  
28 ``  
29 ## Including Plots  
30 ``{r pressure, echo=FALSE, fig.height = 2}  
11:7 rmd한글 template
```

Console Terminal R Markdown

Adobe Acrobat Reader DC

파일 편집 보기(V) 창(W) 도움말(H)

도구 rmd_pdf (1단).pdf x

로그인

1 / 1 50%

rmd 한글 template
Learning Spoons
2018-04-17

R Markdown 한글

Code의 9,10,11 번째 라인의 pound sign(#)을 제거하고 위의 줄과 indent를 맞추면 2 컬럼의 문서가 render 됩니다. <http://rmarkdown.rstudio.com>. When you click the Knit button a document will be generated.

summary(cars)

| | speed | dist |
|------------|-------|----------------|
| ## Min. | 4.0 | 2.00 |
| ## 1st Qu. | 12.0 | 1st Qu.: 26.00 |
| ## Median | 15.0 | Median : 36.00 |
| ## Mean | 15.4 | Mean : 42.98 |
| ## 3rd Qu. | 19.0 | 3rd Qu.: 56.00 |
| ## Max. | 25.0 | Max. :120.00 |

Including Plots

rmd_pdf.pdf

페이지: 1 / 1

rmd 한글 template
Learning Spoons
2018-04-17

R Markdown 한글

Code의 9,10,11 번째 라인의 pound sign(#)을 제거하고 위의 줄과 indent를 맞추면 2 컬럼의 문서가 render 됩니다. <http://rmarkdown.rstudio.com>. When you click the Knit button a document will be generated.

summary(cars)

| | speed | dist |
|------------|-------|----------------|
| ## Min. | 4.0 | 2.00 |
| ## 1st Qu. | 12.0 | 1st Qu.: 26.00 |
| ## Median | 15.0 | Median : 36.00 |
| ## Mean | 15.4 | Mean : 42.98 |
| ## 3rd Qu. | 19.0 | 3rd Qu.: 56.00 |
| ## Max. | 25.0 | Max. :120.00 |

docx

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

a.Rmd x

1 ---
2 title: "Untitled"
3 author: "Learning spoons"
4 date: "2018년 4월 17일"
5 output: word_document

7
8 ````{r setup, include=FALSE}
9 knitr::opts_chunk\$set(echo = TRUE)
```  
11  
12 ## R Markdown  
13  
14 This is an R Markdown document. Markdown is a simple formatting syntax  
for authoring HTML, PDF, and MS Word documents. For more details on  
using R Markdown see <<http://rmarkdown.rstudio.com>>.  
15  
16 When you click the **Knit** button a document will be generated that  
includes both content as well as the output of any embedded R code  
chunks within the document. You can embed an R code chunk like this:  
17  
18 ````{r cars}  
19 summary(cars)  
20```  
21  
22 ## Including Plots  
23  
24 You can also embed plots, for example:  
25  
26 ````{r pressure, echo=FALSE}  
27 plot(pressure)  
28```  
29  
2:1 # Untitled

R Markdown

Console Terminal R Markdown

자동 저장 (켜기) 파일 흠 삽입 그리기 디자인 레이아웃 참조 편지 검토 보기 개발 도구 Acrobat 입력하세요

로그인 저장됨 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44

## Untitled

Learning Spoons

2018년 4월 17일

### R Markdown

This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://rmarkdown.rstudio.com>.

When you click the Knit button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this:

```
summary(cars)
```

##	speed	dist
##	Min. : 4.0	Min. : 2.00
##	1st Qu.:12.0	1st Qu.: 26.00
##	Median :15.0	Median : 36.00
##	Mean :15.4	Mean : 42.98
##	3rd Qu.:19.0	3rd Qu.: 56.00
##	Max. :25.0	Max. :120.00

1/2 페이지 151개 단어 영어(미국)

# ioslides

```

```

```
title: "Untitled"
author: "Learning Spoons"
date: "2018년 4월 17일"
output: ioslides_presentation
```

```
--
```

```
```{r setup, include=FALSE}
knitr::opts_chunk$set(echo = FALSE)
```
```

## ## R Markdown

This is an R Markdown presentation.  
<http://rmarkdown.rstudio.com>.

when you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document.

## ## slide with Bullets

- Bullet 1
- Bullet 2
- Bullet 3

## ## slide with R output

```
```{r cars, echo = TRUE}
summary(cars)
```
```

## ## slide with Plot

### Slide with Bullets

- Bullet 1
- Bullet 2
- Bullet 3

### R Markdown

This is an R Markdown presentation.  
<http://rmarkdown.rstudio.com>.

When you click the Knit button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document.

### Untitled

Learning Spoons  
2018년 4월 17일

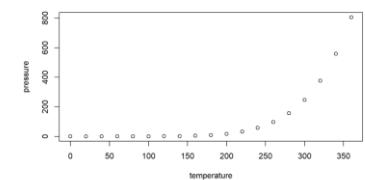


### Slide with R Output

```
summary(cars)
```

```
speed dist
Min. :4.0 Min. : 2.00
1st Qu.:12.0 1st Qu.:26.00
Median :15.0 Median :36.00
Mean :15.4 Mean :42.98
3rd Qu.:19.0 3rd Qu.:56.00
Max. :25.0 Max. :120.00
```

### Slide with Plot



# slidy

```

```

```
title: "untitled"
author: "Learning Spoons"
date: "`r sys.Date()`"
output:
 slidy_presentation: default
 ioslides_presentation: default

```

```
```{r setup, include=FALSE}
knitr::opts_chunk$set(echo = FALSE)
```

R Markdown

This is an R Markdown presentation.
<http://rmarkdown.rstudio.com>.

When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document.

slide with Bullets

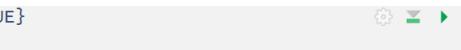
- Bullet 1
- Bullet 2
- Bullet 3

slide with R output

```
```{r cars, echo = TRUE}
summary(cars)
```
```

Slide with Bullets

- Bullet 1
- Bullet 2
- Bullet 3



Untitled

Learning Spoons

2018-04-17

R Markdown

This is an R Markdown presentation. <http://rmarkdown.rstudio.com>.

When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document.

Contents

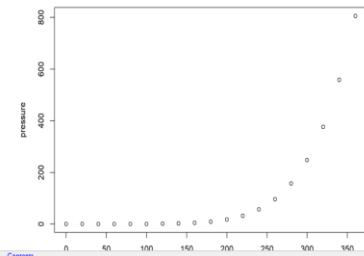
slide 2/5

Slide with R Output

```
summary(cars)
```

```
##   speed      dist
## Min. :4.0  Min. :  8.00
## 1st Qu.:12.0 1st Qu.: 24.00
## Median:15.0 Median : 36.00
## Mean  :15.4 Mean  : 42.98
## 3rd Qu.:19.0 3rd Qu.: 56.00
## Max. :25.0  Max. :120.00
```

Slide with Plot



Contents

slide 3/5

beamer (pdf)

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

+ a.Rmd x a.Rmd x rmd_beamer.Rmd x

1 ---
2 title: "한글로 제목"
3 author: "러닝스푼즈"
4 output:
5 beamer_presentation:
6 latex_engine: xelatex
7 # keep_tex: true
8 # template: myTemplate.tex
9 includes:
10 in_header: myRmdBeamerStyle/latex-topmatter.tex
11 classoption: t
12 mainfont: NanumGothic
13 ---
14
15 ```{r setup, include=FALSE}
16 knitr::opts_chunk\$set(echo = FALSE)
17 ```
18
19 ## R Markdown
20
21 알 마크다운
22
23 ## slide with Bullets
24
25 - Bullet 1
26 - Bullet 2
27 - Bullet 3
28
29 ## slide with R output
30
31 ```{r cars, echo = TRUE}
32 summary(cars)

4:8 # 한글로 제목

Console Terminal R Markdown

페이지: 1 (1 / 4)

beamer

파일 흡 공유 보기

한글로 제목 러닝스푼즈

클립보드 구성 새로 만들기 열기 선택

이름 myRmdBeamerStyle rmd_beamer rmd_beamer

Slide with Bullets

- Bullet 1
- Bullet 2
- Bullet 3

Slide with R Output

```
summary(cars)
##      speed      dist
##  Min.   :4.0   Min.   : 2.00
##  1st Qu.:12.0  1st Qu.: 26.00
##  Median :15.0  Median : 36.00
##  Mean   :15.4  Mean   : 42.98
##  3rd Qu.:19.0  3rd Qu.: 86.00
##  Max.   :25.0  Max.   :120.00
```

3개 항목

Example – 1주차 숙제!

The image shows a dual-screen setup. On the left screen, the RStudio interface is visible, displaying an R Markdown file named 'Review-Week1.Rmd'. The code includes YAML front matter and several R code chunks. One chunk defines variable 'a' as 'Hello', and another defines 'a' and 'b' as 'Hello' and 'world' respectively, then pastes them together. The right screen shows a Beamer presentation slide titled 'Review - Week1' with sections for 'Module 1 - Hello World' and four problems labeled 1 through 4.

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

a.Rmd x a.Rmd x rmd_beamer.Rmd x Review-Week1.Rmd*

1 ---
2 title: "Review - week1"
3 author: "learningSpoonsR"
4 date: "2018년 4월 15일"
5 output: word_document
6 ---
7
8 ```{r setup, include=FALSE}
9 knitr::opts_chunk\$set(echo = TRUE)
10 knitr::opts_chunk\$set(results = 'hide')
11
12
13 ## Module 1 - Hello world
14
15 `r problem <- 1`
16
17 **Problem** `r problem`. `r problem <- problem + 1`
18 ```{r}
19 a <- "Hello"
20 a
21
22 ans:
23
24
25 **Problem** `r problem`. `r problem <- problem + 1`
26 ```{r}
27 a <- "Hello"
28 b <- "world"
29 paste(a,b)
30
31 ans:
32

18:7 Chunk 2 R Markdown

자동 저장 (●) Review... - 저장됨 로그인 파일 험 삽입 그리기 디자인 레이아 참조 편지 검토 보기 개발 도 Acrobat 입력하세요

1 | 4 | 2 | 1 | 2 | 4 | 6 | 8 | 10 | 12 | 14 | 16 | 18 | 20 | 22 | 24 | 26 | 28 | 30 | 32 | 34 | 36 | 38 | 40 | 42 | 44

Review - Week1

learningSpoonsR
2018년 4월 15일

Module 1 - Hello World

Problem 1.
a <- "Hello"
a

ans;

Problem 2.
a <- "Hello"
b <- "World"
paste(a,b)

ans;

Problem 3.
paste0(a,b)

ans;

Problem 4.
paste(a,b,sep="-")

ans;

Example – 1주차 숙제!

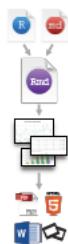
```
**Problem** `r problem`.  
`r problem <- problem + 1`  
```{r}  
grep("ui", font)
```  
```{r, echo=FALSE}  
ans0 <- "ANS: The function grep somehow returns 0 or 1, which is
equivalent to TRUE or FALSE!"
ans0
```  
ans:
```

Problem 3.

```
grep("ui", font)  
## [1] 1  
## [1] "ANS: The function grep somehow returns 0 or 1, which is equivalent to  
TRUE or FALSE!"  
ans:
```

R Markdown :: CHEAT SHEET

What is R Markdown?

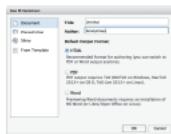


.Rmd files - An R Markdown (.Rmd) file is a record of your research. It contains the code that a scientist needs to reproduce your work along with the narration that a reader needs to understand your work.

Reproducible Research - At the click of a button, or the type of a command, you can rerun the code in an R Markdown file to reproduce your work and export the results as a finished report.

Dynamic Documents - You can choose to export the finished report in a variety of formats, including html, pdf, MS Word, or RTF documents; html or pdf based slides, Notebooks, and more.

Workflow



① Open a new .Rmd file at File ▶ New File ▶ R Markdown. Use the wizard that opens to pre-populate the file with a template

② Write document by editing template

③ Knit document to create report; use knit button or render() to knit

④ Preview Output in IDE window

⑤ Publish (optional) to web server

⑥ Examine build log in R Markdown console

⑦ Use output file that is saved along side .Rmd

Embed code with knitr syntax

INLINE CODE

Insert with `r<code>`. Results appear as text without code.

Built with `r getRversion()` ➔ Built with 3.2.3

IMPORTANT CHUNK OPTIONS

cache - cache results for future knits (default = FALSE)

cache.path - directory to save cached results in (default = "cache/")

child - file(s) to knit and then include (default = NULL)

collapse - collapse all output into single block (default = FALSE)

comment - prefix for each line of results (default = '#')

```

1 -+
2 title: "R Markdown"
3 author: "RStudio"
4 output: 2
5 html_document:
6 toc: TRUE
7 -
8
9 ``{r setup, include=FALSE}
10 knitr::opts_chunk$set(echo = TRUE)
11 -
12
13 ## R Markdown
14
15 This is an R Markdown document.
16 Markdown is a simple formatting
17 syntax for authoring HTML, PDF,
18 and MS Word documents.
19
20 ``{r cars}
21 summary(cars)
22
23
24 For more details on using R Markdown
25 see http://rmarkdown.rstudio.com.
  
```

This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents.

| summary(cars) |
|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| ## speed dist
Min. : 4.0 Min. : 2.00
1st Qu.:12.0 1st Qu.: 26.00
Median :15.0 Median : 36.00
Mean :15.4 Mean : 42.98
3rd Qu.:19.0 3rd Qu.: 56.00
Max. :25.0 Max. :120.0 |

For more details on using R Markdown see <http://rmarkdown.rstudio.com>.

render

Use rmarkdown::render() to render/knit at cmd line. Important args:

input - file to render

output_options - List of render options (as in YAML)

output_file

output_dir

params - list of params to use

envir - environment to evaluate code chunks in

encoding - of input file

CODE CHUNKS

One or more lines surrounded with ``{r} and ```. Place chunk options within curly braces, after r. Insert with

``{r echo=TRUE}

getRversion()

```

getRversion()
#> [1] '3.2.3'
  
```

GLOBAL OPTIONS

Set with knitr::opts_chunk\$set(), e.g.

``{r include=FALSE}

knitr::opts_chunk\$set(echo = TRUE)

````

**message** - display code messages in document (default = TRUE)

**results** (default = 'markup')

'asis' - passthrough results

'hide' - do not display results

'hold' - put all results below all code

**tidy** - tidy code for display (default = FALSE)

**warning** - display code warnings in document (default = TRUE)

## Interactive Documents

Turn your report into an interactive Shiny document in 4 steps

1. Add runtime: shiny to the YAML header.
2. Call Shiny input functions to embed input objects.
3. Call Shiny render functions to embed reactive output.
4. Render with rmarkdown::run or click Run Document in RStudio IDE

--  
output: html\_document  
runtime: shiny  
--

``{r, echo = FALSE}  
numericInput("n",  
"How many cars?", 5)

renderTable({  
 head(cars, input\$n)  
})  
,,



How many cars?

5

speed dist

1 4.0 2.00

2 4.0 10.00

3 7.00 4.00

4 7.00 22.00

5 8.00 16.00

Embed a complete app into your document with shiny::shinyAppDir()

NOTE: Your report will be rendered as a Shiny app, which means you must choose an html output format, like html\_document, and serve it with an active R Session.





# Pandoc's Markdown

Write with syntax on the left to create effect on right (after render)

Plain text  
End a line with two spaces to start a new paragraph.  
"Italics" and "bold"  
verbatim code  
sub/superscript<sup>2</sup>  
~~strikethrough~~  
escaped: `^` \\\  
endash: --, emdash: ---  
equation: \$A = \pi r^2\\$  
equation block:

$\$SE = mc^2\$$

> block quote

# Header1 [#anchor]

## Header 2 [#css\_id]

### Header 3 [.css\_class]

#### Header 4

##### Header 5

##### Header 6

<!--Text comment-->

\textbf{Text ignored in HTML}

<em>HTML ignored in pdfs</em>

<http://www.rstudio.com>

[link](www.rstudio.com)

Jump to [Header 1](#anchor)

image:

Plain text  
End a line with two spaces to start a new paragraph.  
"Italics" and "bold"  
verbatim code  
sub/superscript<sup>2</sup>  
~~strikethrough~~  
escaped: `^` \\\  
endash: --, emdash: ---  
equation: \$A = \pi r^2\\$  
equation block:

$E = mc^2$

block quote

**Header1**

**Header 2**

**Header 3**

**Header 4**

**Header 5**

**Header 6**

HTML (ignored in pdfs)

http://www.rstudio.com

link

Jump to Header 1

image:



Caption

- \* unordered list
  - + sub-item 1
    - o sub-item 1
    - o sub-item 2
    - sub-sub-item 1
  - item 2
    - Continued (indent 4 spaces)

- 1. ordered list
  1. item 1
  2. item 2
    - i. sub-item 1
      - A. sub-sub-item 1

(@) A list whose numbering continues after

continues after

2. an interruption

Term 1

Definition 1

Right	Left	Default	Center
12	12	12	12
123	123	123	123
1	1	1	1

- slide bullet 1

- slide bullet 2

(> to have bullets appear on click)

horizontal rule/slide break:

...

A footnote [\*1]

[\*1]: Here is the footnote.

- > to have bullets appear on click
- horizontal rule/slide break:
- ...

A footnote [\*1]

1. Here is the footnote.<sup>2</sup>

# Set render options with YAML

When you render, R Markdown

1. runs the R code, embeds results and text into .md file with knitr
2. then converts the .md file into the finished format with pandoc



Set a document's default output format in the YAML header:

```
-- output: html_document
Body
```

**output value**

**creates**

html_document	html
pdf_document	pdf (requires Tex)
word_document	Microsoft Word (.docx)
odt_document	OpenDocument Text
rtf_document	Rich Text Format
md_document	Markdown
github_document	Github compatible markdown
ioslides_presentation	ioslides HTML slides
slidy_presentation	slidy HTML slides
beamer_presentation	Beamer pdf slides (requires Tex)

Customize output with sub-options (listed to the right):

```
-- Indent 2 spaces
-- output: html_document: code_folding: hide
-- toc_float: TRUE
Body
```

**html tabsets**

Use tablet css class to place sub-headers into tabs

```
Tabset {.tabset.tabset-fade.tabset-pills}
```



Tabset

Tab 1 Tab 2

text 1

End tabset

## Create a Reusable Template

1. Create a new package with a `inst/rmarkdown/templates` directory

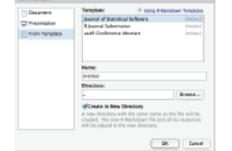
2. In the directory, Place a folder that contains: `template.yaml` (see below) `skeleton.Rmd` (contents of the template)

any supporting files

3. Install the package

4. Access template in wizard at File ▶ New File ▶ R Markdown template.yaml

-- name: My Template



**sub-option**

**description**

`citation_package` The LaTeX package to process citations, natbib, biblatex or none

`code_folding` Let readers to toggle the display of R code, "none", "hide", or "show"

`colortheme` Beamer color theme to use

`css` CSS file to use to style document

`dev` Graphics device to use for figure output (e.g. "png")

`duration` Add a countdown timer (in minutes) to footer of slides

`fig_caption` Should figures be rendered with captions?

`fig_height, fig_width` Default figure height and width (in inches) for document

`highlight` Syntax highlighting: "tango", "pygments", "kate", "zenburn", "textmate"

`includes` File of content to place in document (in\_header, before\_body, after\_body)

`incremental` Should bullets appear one at a time (on presenter mouse clicks)?

`keep_md` Save a copy of .md file that contains knitr output

`keep_tex` Save a copy of .tex file that contains knitr output

`latex_engine` Engine to render latex, "pdflatex", "xelatex", or "lualatex"

`lib_dir` Directory of dependency files to use (Bootstrap, MathJax, etc.)

`mathjax` Set to local or a URL to use a local/URL version of MathJax to render equations

`md_extensions` Markdown extensions to add to default definition or R Markdown

`number_sections` Add section numbering to headers

`pandoc_args` Additional arguments to pass to Pandoc

`preserve_yaml` Preserve YAML front matter in final document?

`reference_docx` docx file whose styles should be copied when producing docx output

`selfContained` Embed dependencies into the doc

`slide_level` The lowest heading level that defines individual slides

`smaller` Use the smaller font size in the presentation?

`smart` Convert straight quotes to curly, dashes to em-dashes, ... to ellipses, etc.

`template` Pandoc template to use when rendering file quarterly\_report.html

`theme` Bootswatch or Beamer theme to use for page

`toc` Add a table of contents at start of document

`toc_depth` The lowest level of headings to add to table of contents

`toc_float` Float the table of contents to the left of the main content

html	X
pdf	X
word	
odt	
rtf	X
ind	X
gbub	
isides	
slidy	X
beamer	

html	
pdf	X
word	X
odt	X
rft	
ind	
gbub	
isides	
slidy	X
beamer	

html	X
pdf	X
word	X
odt	X
rft	X
ind	X
gbub	
isides	
slidy	X
beamer	

html	X
pdf	X
word	X
odt	X
rft	X
ind	X
gbub	
isides	
slidy	X
beamer	

html	X
pdf	X
word	X
odt	X
rft	X
ind	X
gbub	
isides	
slidy	X
beamer	

html	X
pdf	X
word	X
odt	X
rft	X
ind	X
gbub	
isides	
slidy	X
beamer	

html	X
pdf	X
word	X
odt	X
rft	X
ind	X
gbub	
isides	
slidy	X
beamer	

html	X
pdf	X
word	X
odt	X
rft	X
ind	X
gbub	
isides	
slidy	X
beamer	

html	X
pdf	X
word	X
odt	X
rft	X
ind	X
gbub	
isides	
slidy	X
beamer	

html	X
pdf	X
word	X
odt	X
rft	X
ind	X
gbub	
isides	
slidy	X
beamer	

html	X
pdf	X
word	X
odt	X
rft	X
ind	X
gbub	
isides	
slidy	X
beamer	

html	X
pdf	X
word	X
odt	X
rft	X
ind	X
gbub	
isides	
slidy	X
beamer	

html	X
pdf	X
word	X
odt	X
rft	X
ind	X
gbub	
isides	
slidy	X
beamer	

html	X
pdf	X
word	X
odt	X
rft	X
ind	X
gbub	
isides	
slidy	X
beamer	

html	X
pdf	X
word	X
odt	X
rft	X
ind	X
gbub	
isides	
slidy	X
beamer	

html	X
pdf	X
word	X
odt	X
rft	X
ind	X
gbub	
isides	
slidy	X
beamer	

html	X
pdf	X
word	X
odt	X
rft	X
ind	X
gbub	
isides	
slidy	X
beamer	

html	X
pdf	X
word	X
odt	X
rft	X
ind	X
gbub	
isides	
slidy	X
beamer	

html	X
pdf	X
word	X
odt	X
rft	X
ind	X
gbub	
isides	
slidy	X
beamer	

html	X
pdf	X
word	X
odt	X
rft	X
ind	X
gbub	
isides	
slidy	X
beamer	

html	X
pdf	X
word	X
odt	X
rft	X
ind	X
gbub	
isides	
slidy	X
beamer	

html	X
pdf	X
word	X
odt	X
rft	X
ind	X
gbub	
isides	
slidy	X
beamer	

html	X
pdf	X
word	X
odt	X
rft	X
ind	X
gbub	
isides	
slidy	X
beamer	

html	X
pdf	X
word	X
odt	X
rft	X
ind	X
gbub	
isides	
slidy	X
beamer	

html	X
pdf	X
word	X
odt	X
rft	X
ind	X
gbub	
isides	
slidy	X
beamer	

html	X
pdf	X
word	X
odt	X
rft	X
ind	X
gbub	
isides	
slidy	X
beamer	

html	X
pdf	X
word	X
odt	X
rft	X
ind	X
gbub	
isides	
slidy	X
beamer	

html	X
pdf	X
word	X
odt	X
rft	X
ind	X
gbub	
isides	
slidy	X
beamer	

html	X
pdf	X
word	X
odt	X
rft	X
ind	X
gbub	
isides	
slidy	X
beamer	

html	X
pdf	X
word	X
odt	X
rft	X
ind	X
gbub	
isides	
slidy	X
beamer	

html	X
pdf	X
word	X
odt	X
rft</	

# R 마크다운

컨닝쪽지

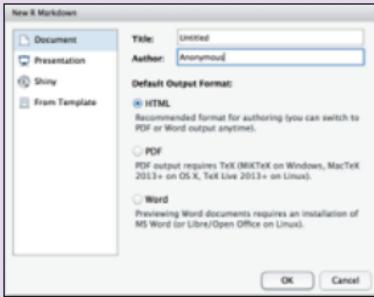
추가 학습 정보 [rmarkdown.rstudio.com](http://rmarkdown.rstudio.com)

rmarkdown 0.2.50 최종갱신일: 8/14



2. 파일 열기 .Rmd 확장자를 갖는 텍스트 파일로 저장해서 시작하거나, Studio Rmd 템플릿을 열어 시작한다.

- 메뉴막대에서, 다음순으로 클릭한다.  
**File ▶ New File ▶ R Markdown...**
- 윈도우가 열리면, .Rmd 파일로 작성하려는 출력유형을 선택한다.
- 라디오 버튼으로 출력형식을 선택한다 (나중에 출력형식은 변경할 수 있다)
- OK 버튼을 클릭한다.



4. 출력형식 설정 R 마크다운 파일에서 생성할 문서 유형을 기술하는 YAML 헤더정보를 작성한다.

## YAML

YAML 헤더는 키(key) 집합: 파일 시작지점에 나오는 키-값 쌍. 헤더 시작과 끝은 3개 대쉬를 갖는 라인 (- - -)

```

```

title: "xwMOOC 보고서"  
author: "무명씨"  
output: html\_document  
---

보고서의 시작지점으로, YAML 헤더 정보를 작성해 놓는다.

RStudio 템플릿것이  
기준으로 YAML 헤더  
정보를 작성해 놓는다.

출력값이 .Rmd 파일에서 어떤 형식 파일을 생성할 것인지 결정한다(단계 6)

**output: html\_document** ..... html 파일 (웹페이지)

**output: pdf\_document** ..... pdf 문서

**output: word\_document** ..... MS 워드문서 .docx

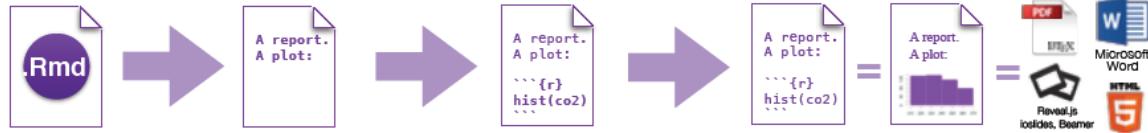
**output: beamer\_presentation** ..... Beamer 발표자료 (pdf)

**output: ioslides\_presentation** ..... 맥 발표자료 (html)



1. 작업흐름 R 마크다운은 R로 재현가능하고, 동적인 보고서를 작성하는 서식이다. R 마크다운을 사용해서 R 코드와 실행결과를 발표자료, pdf, html, 워드 문서 등에 삽입할 수 있다. 보고서를 작성하려면:

- 파일열기 - .Rmd 확장자를 갖는 파일을 연다.
- 작성하기 - 본문을 작성하기 쉬운 R 마크다운 구문을 사용해서 작성한다.
- 내장하기 - 리포트에 포함될 출력 결과를 생성하는 R 코드를 내장한다.
- 렌더링(Render) - R 코드를 출력형식으로 치환하고 보고서를 발표자료, pdf, html, MS 워드 파일 형식으로 변환한다.



3. 마크다운 다음으로, 일반 텍스트로 보고서를 작성한다. 마크다운 구문을 사용해서 최종 보고서에 적용할 텍스트 서식을 기술한다.

## 입력 구문

### 일반 텍스트

새로운 단락을 시작하려면 줄 마지막을 공백 2개로 끝낸다.

\*기술인 글씨\* and \_기술인 글씨\_

\*\*굵은 글씨\*\* and \_\_굵은 글씨\_\_

원첨자`2`

~~취소선~~

[링크](www.rstudio.com)

# 제목 1

## 제목 2

### 제목 3

#### 제목 4

##### 제목 5

###### 제목 6

N자 크기 대시 부호: --

m자 크기 대시 부호: ---

생략: ...

즉시 처리하는 수식: \$A = \pi r^2\$

이미지: 

수평선(혹은 슬라이드 멤버):

---

### > 인용 블록

\* 순서없는 목록

\* 항목 2

+ 하위 항목 1  
+ 하위 항목 2

1. 순서있는 목록

2. 항목 2

+ 하위 항목 1  
+ 하위 항목 2

표 제목 | 두번째 제목

표간 | 칸 2

칸 3 | 칸 4

## 출력 결과

### 일반 텍스트

새로운 단락을 시작하려면 줄 마지막을 공백 2개로 끝낸다.

\*기술인 글씨\* and \_기술인 글씨\_

굵은 글씨\*\* and \_\_굵은 글씨\_\_

윗첨자`2`

취소선~~

[링크]

## 제목 1

## 제목 2

## 제목 3

제목 4

제목 5

제목 6

N자 크기 대시 부호: --

m자 크기 대시 부호: ---

생략: ...

즉시 처리하는 수식:  $A = \pi r^2$

이미지: 

수평선(혹은 슬라이드 멤버):

## 인용 블록

\* 순서없는 목록

\* 항목 2

+ 하위 항목 1  
+ 하위 항목 2

1. 순서있는 목록

2. 항목 2

+ 하위 항목 1  
+ 하위 항목 2

표 제목

두번째 제목

표간

칸 2

칸 3

칸 4

## 5. 코드내장하기

knitr 구문을 사용해서 R 코드를 보고서에 내장한다.  
R이 코드를 실행하고, 보고서를 렌더링할 때 결과를 포함시킨다.

### 인라인 코드

r 코드를 백틱(`)으로 감싼다.  
R이 인라인 코드를 실행된 결과로 대체한다.

2 더하기 2는 `r 2`  
2와 같다.

Two plus two  
equals 4.

실행결과는 다음과 같다  
```{r}  
dim(iris)
```

Here's some code  
dim(iris)  
## [1] 150 5

### 화면 출력 선택옵션

knitr 선택옵션을 사용해서 코드 덩어리 출력 스타일을 적용한다.  
코드 상단 괄호 내부에 선택옵션을 지정한다.

Here's some code  
```{r eval=FALSE}  
dim(iris)

Here's some code
dim(iris)

Here's some code
```{r echo=FALSE}  
dim(iris)

Here's some code  
## [1] 150 5

### 선택옵션 기본설정 효과

| 선택옵션       | 기본설정     | 효과                               |
|------------|----------|----------------------------------|
| eval       | TRUE     | 코드를 평가하고 실행결과를 포함한다.             |
| echo       | TRUE     | 실행결과와 함께 코드를 출력한다.               |
| warning    | TRUE     | 경고메시지를 출력한다.                     |
| error      | FALSE    | 오류메시지를 출력한다.                     |
| message    | TRUE     | 메시지를 출력한다.                       |
| tidy       | FALSE    | 깔끔한 방식으로 코드 형태를 변형한다.            |
| results    | "markup" | "markup", "asis", "hold", "hide" |
| cache      | FALSE    | 결과값을 캐싱해서 향후 실행시 건너뛰게 설정한다.      |
| comment    | "##"     | 주석문자로 출력결과에 서두를 붙인다.             |
| fig.width  | 7        | 덩어리로 생성되는 그래프에 대한 폭을 인치로 지정한다.   |
| fig.height | 7        | 덩어리로 생성되는 그래프에 대한 높이를 인치로 지정한다.  |

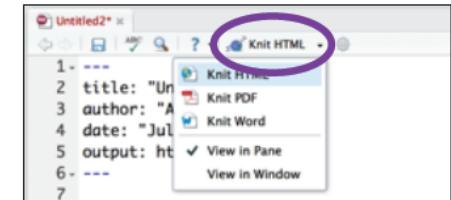
보다 자세한 사항은 웹사이트를 참조: [yihui.name/knitr/](http://yihui.name/knitr/)

## 6. 렌더링

최종보고서를 생성하는데 .Rmd 파일을 사용하여 청사진을 제작한다.

두 가지 방식으로 보고서를 렌더링한다.

1. `rmarkdown::render("〈파일 경로〉")` 명령어를 실행한다.
2. RStudio 스크립트 작성창 상단에 **knit HTML** 버튼을 클릭한다.



- レン더링 명령을 실행시키면, R은 다음을 수행한다
- 내장된 코드 덩어리를 각각 실행시키고, 실행결과를 보고서에 삽입한다.
  - 출력 파일형식에 맞춰 신규 보고서를 생성한다.
  - 미리보기로 뷰어창에 출력파일을 연다.
  - 작업디렉토리에 출력파일을 저장한다.

## 7. 인터랙티브 문서

작성한 보고서를 3단계를 거쳐 인터랙티브 Shiny 문서변환.

- 1 YAML 헤더에 `runtime: shiny` 을 추가한다.

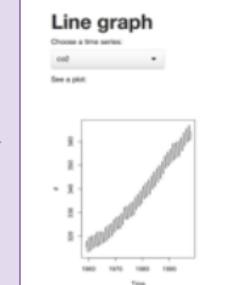
```
title: "Line graph"
output: html_document
runtime: shiny
```

- 2 코드 덩어리에, 위젯을 내장하는 Shiny `input` 함수를 추가한다. Shiny `render` 함수를 추가해서 반응형 출력결과를 내장한다.

```
title: "Line graph"
output: html_document
runtime: shiny

Choose a time series:
```{r echo = FALSE}
selectInput("data", "", 
  c("co2", "lh"))
```
See a plot:
```{r echo = FALSE}
renderPlot({
  d <- get(input$data)
  plot(d)
})
```

- 3 `rmarkdown::run` 명령어로 렌더링하거나 RStudio Run Document 버튼을 클릭한다.



- * 주목: 보고서는 Shiny 앱이 된다. 따라서, (인터랙티브 보고서를 위해) `html_document` 혹은 (인터랙티브 발표자료) `ioslides_presentation` 출력형식을 선택한다.

8. Publish

온라인으로 접속하는 사용자와 보고서를 공유한다.

Rpubs.com

RStudio 무료 R 마크다운 게시 사이트를 통해 정적 문서를 공유한다.

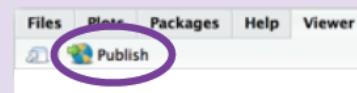
www.rpubs.com

ShinyApps.io

Studio 서버에 인터랙티브 문서를 올려 호스팅한다. 무료와 유료 선택옵션이 있다.

www.shinyapps.io

RStudio 미리보기 창에 "Publish" 버튼을 클릭하여 rpubs.com 사이트에 버튼 한번 클릭으로 바로 올린다.



9. 추가 학습

문서와 예제 - rmarkdown.rstudio.com

추가 기사 - shiny.rstudio.com/articles

- blog.rstudio.com

- [@rstudio](https://twitter.com/rstudio)



RStudio® and Shiny™ are trademarks of RStudio, Inc.
[CC BY RStudio info@rstudio.com](http://info@rstudio.com)
844-448-1212 rstudio.com

Going back to Carseats!

(시연)

Carseat_pmd.Rmd x ABC Knit

```
1 ---  
2 title: 'carseat 판매량'  
3 author: "Learning Spoons"  
4 date: `r Sys.Date()`  
5 output:  
6   pdf_document:  
7     latex_engine: xelatex  
8     keep_tex: true  
9     # pandoc_args: [  
10       # "-V", "classoption=twocolumn"  
11     # ]  
12     smaller: true  
13     mainfont: NanumGothic  
14     classoption: a4paper|  
15 ---  
16  
17 ```{r setup, include=FALSE}  
18 knitr::opts_chunk$set(echo = TRUE)  
19 ````  
20  
21 ## carseat 소개  
22  
23 ```{r}  
24 library(ISLR)  
25 library(dplyr)  
26 library(ggplot2)  
27 str(Carseats)  
28  
29  
30 ## Focus City  
31  
32 ```{r}  
33 focusCity <- Carseats %>%  
34   filter(Income > 100) %>%  
35   filter(Age >= 30 & Age < 40) %>%  
36   mutate(AdvPerCapita = Advertising/Population) %>%  
37   select(Sales, Income, Age, Population, Education, AdvPerCapita) %>%  
38   arrange(Sales)  
39   print(focusCity)  
40  
41
```

```
41   ## Income vs Sales  
42  
43   ```{r}  
44   doFacetWrap <- FALSE  
45   a <- ggplot(data = Carseats, aes(x = Income, y = sales)) +  
46     geom_point(aes(shape = Urban, color = US))  
47   if (doFacetWrap) {  
48     a <- a + facet_wrap(~ floor(Age/10))  
49   }  
50   print(a)  
51   ```  
52  
53   Your comment!  
54  
55   ## Income vs Sales  
56  
57   ```{r}  
58   doFacetWrap <- TRUE  
59   a <- ggplot(data = Carseats, aes(x = Income, y = sales)) +  
60     geom_point(aes(shape = Urban, color = US))  
61   if (doFacetWrap) {  
62     a <- a + facet_wrap(~ floor(Age/10))  
63   }  
64   print(a)  
65   ```  
66
```

Carseat 판매량

Learning Spoons
2018-04-21

Carseat 소개

```
library(ISLR)
library(dplyr)

##
## Attaching package: 'dplyr'
## The following objects are masked from 'package:stats':
##
##     filter, lag
## The following objects are masked from 'package:base':
##
##     intersect, setdiff, setequal, union
library(ggplot2)
str(Carsesats)

## 'data.frame':  400 obs. of  11 variables:
##   $ Sales    : num  9.5 11.22 10.06 7.4 4.15 ...
##   $ CompPrice: num  138 111 113 117 141 124 115 136 132 132 ...
##   $ Income   : num  73 48 35 100 64 113 105 81 110 113 ...
##   $ Advertising: num  11 16 10 4 3 13 0 15 0 0 ...
##   $ Population: num  276 260 269 466 340 501 45 425 108 131 ...
##   $ Price    : num  120 83 80 97 128 72 108 120 124 124 ...
##   $ ShelveLoc: Factor w/ 3 levels "Bad","Good","Medium": 1 2 3 3 1 1 3 2 3
##   $ Age      : num  42 65 59 55 38 78 71 67 76 76 ...
##   $ Education: num  17 10 12 14 13 16 15 10 17 ...
##   $ Urban    : Factor w/ 2 levels "No","Yes": 2 2 2 2 2 1 2 2 1 1 ...
##   $ US       : Factor w/ 2 levels "No","Yes": 2 2 2 2 1 2 1 2 1 2 ...

Focus City
```

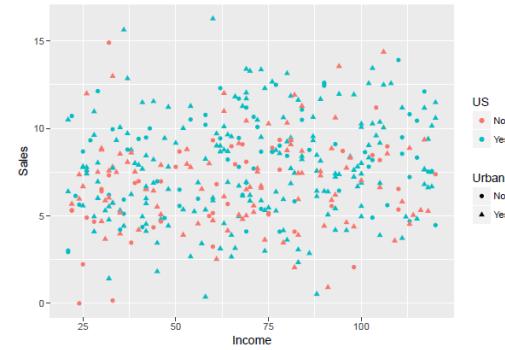
```
focusCity <- Carsesats %>%
  filter(Income > 100) %>%
  filter(Age >= 30 & Age < 40) %>%
  mutate(AdvPerCapita = Advertising/Population) %>%
  select(Sales, Income, Age, Population, Education, AdvPerCapita) %>%
  arrange(Sales)
print(focusCity)

##   Sales Income Age Population Education AdvPerCapita
## 1  5.04    114  34      298        16  0.00000000
## 2  5.32    116  39      170        16  0.00000000
## 3  6.80    117  38      337        10  0.01483680
## 4  7.49    119  35      178        13  0.03370787
## 5  7.67    117  36      400        10  0.02000000
## 6  8.55    111  36      480        16  0.04791667
## 7  8.97    107  33      144        13  0.00000000
```

```
##  9  9.39    118  32      445        15  0.03146067
## 10 9.58    104  37      353        17  0.06515581
## 11 10.36    105  34      428        12  0.04205607
## 12 10.59    120  30      262        10  0.05725191
## 13 12.57    108  33      203        14  0.08374384
```

Income vs Sales

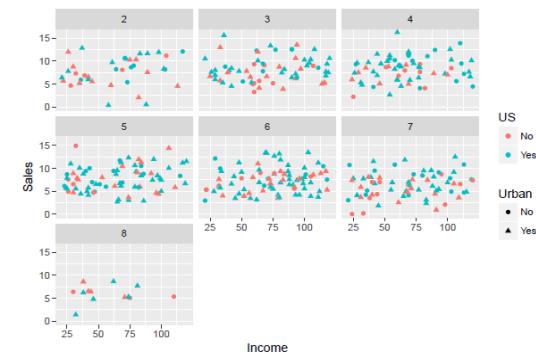
```
doFacetWrap <- FALSE
a <- ggplot(data = Carsesats, aes(x = Income, y = Sales)) +
  geom_point(aes(shape = Urban, color = US))
if (doFacetWrap) {
  a <- a + facet_wrap(~ floor(Age/10))
}
print(a)
```



Your comment!

Income vs Sales

```
doFacetWrap <- TRUE
a <- ggplot(data = Carsesats, aes(x = Income, y = Sales)) +
  geom_point(aes(shape = Urban, color = US))
if (doFacetWrap) {
  a <- a + facet_wrap(~ floor(Age/10))
}
print(a)
```



Discussion - rmarkdown

- Markdown이라는 Markup Language로 변환하여 사용성 높음
- 변환 과정
 - Rmd → .md → .html, .docx
 - Rmd → .md → .tex → .pdf *texlive*
- 약간의 수정으로 여러가지 형식의 문서화가 가능함 ✓
- 자연스럽게 문법에 맞는 Color Coding
- R output에 가장 비슷한 quality의 output을 얻을 수 있음

Discussion – Literature Programming

- *Points*

- 프로그래밍이 아닌 글쓰기가 초점이 되는 차세대 트렌드
- 데이터 분석 업무에 연관성이 높음!
- 글을 쓴다는 것과 프로그래밍을 하는 것은 자기 자신을 표현하는 가장 높은 수준의 지적인 활동인데 이것을 한꺼번에 할 수 있게 해주는 도구 (#Express Yourself)

- *Advantages*

- 코드와 문서가 하나의 파일이라서 관리가 편함
- 코드에서 주석을 조금만 달아도 됨
- 코드만 있는 것에 비해서 의사소통이 용이함
- 데이터나 분석의 결과가 달라지는 것이 문서에 즉각적으로 반영됨
- 반복적으로 작성하는 문서 작업과 엑셀 작업을 안해도 됨
- Colin Powell: “You don’t know what you will can get away with until you try”

Appendix – pdf 파일을 위한 준비

- *Texlive 2017 설치 (수백메가)*
 - <<http://www.ktug.org/xe/index.php?mid=install>>
- *한글 폰트 설치*
 - google "nanumgothic download"
 - google "nanummyeongjo download"
 - google "nanumgothiccoding download"
- *Reference: google "latex beamer에서 한글 쓰기!"*

Appendix – 다른 Programming Language

```
# C, C++ +
for(int i=1; i <=10, i=i+1) {
    printf("value of i: %d\n", i);
}

#
# MATLAB +
for i=1:10 {
    disp("value of i: %d", i)
}

#
# Python +
for i in range(1, 11):
    print "value of i: %d" % (i)

##
## R +
for (i in 1:10) {
    print(paste("value of i:", i))
}
```