

Orientation

2025-09-02

오늘 목표

- 빅데이터분석 강의 설명
- 프로그램 설치
- 주요 툴 설명, 실습

강의 설명

강의 목표

- (의료) 데이터를 효과적으로 처리하고 분석하기 위해 **R 등의 도구**를 잘 활용할 수 있도록 한다.
 - 데이터 분석의 과정에는 Communication이 포함
- 데이터 기반 인사이트를 도출하는 **실무 기술 능력**을 배양한다.

강의 방향

데이터 분석의 흐름



- Wrangle = Preprocess
- 데이터 수집: 데이터 엔지니어 / 실무진
- 데이터 모델링: 도메인 경험이 필요. (기술적 분석, 통계에 갇힘)

Why R?

- 대용량 데이터 >> Excel
- 분석 결과의 재현성: `reproducible`
- 통계 분석의 효율성: `effective`
- Academic / Open Source
- !확장성
- !성능

Good to know

- Markdown
- SQL
- SAS
- Excel
- Python
- Web (HTML / CSS / JS)
- Statistics ([Biostatistics](#))
- Computer science
([Bioinformatics](#))

다루지 않는 것

- 고급 모델링 / 통계 이론
- 빅데이터 엔지니어링: Hadoop, Spark
- 다른 프로그래밍 언어: Python, SQL, SAS, Matlab, Julia
- 의료 데이터
- **기초 통계에 대한 이론적 설명**



강의 진행 방법

- 오프라인, 실습 + 이론 (+ 트렌드)
- $(45 + 15) * 3$
- 절대 평가
- 결석 4회 이상 => F

강의 계획 1 (기초)

- OT & 툴 사용
- R 기초
- 데이터 로드
- 데이터 전처리
- 데이터 시각화
- 트렌드(세미나)

강의 계획 2 (응용)

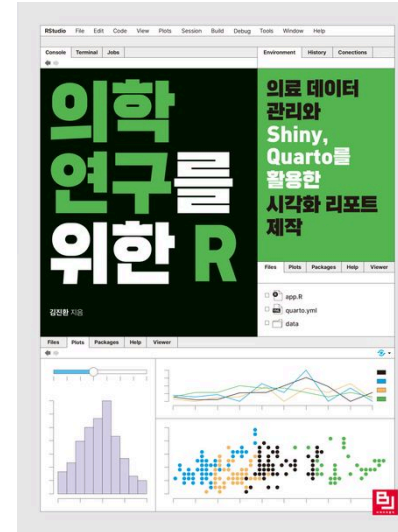
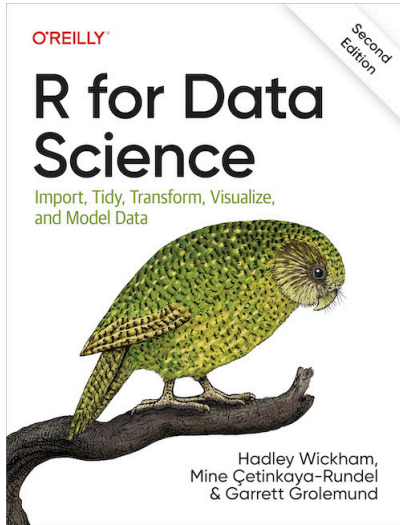
- Quarto
- Shiny
- R Package
- 프로젝트
- 트렌드(세미나)

성적 평가 방법

- 출석 10%
- 중간 / 기말고사 각 35%
- (프로그래밍) 과제, 프로젝트 20%
- 오픈북 시험

교재 (필수 아님)

- R for Data Science (2e) [링크](#)
- 의학 연구를 위한 R [링크](#)



✨ 생성형 AI

- 한국은행 리포트
- ChatGPT (Edu)
- GitHub Copilot (Edu)
- **사용 기록을 답변에 포함 (프롬프트)**

✨ 생성형 AI

- 중간고사 / 기말고사 **택일**, 선착순
- Cheatsheet

Cheatsheet

Data visualization with ggplot2 : : CHEATSHEET

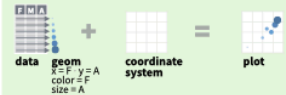


Basics

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>
```

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

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.

Aes

Common aesthetic values.

color and fill - string ("red", "#RRGGBB")

linetype - integer or string (0 = "blank", 1 = "solid", 2 = "dashed", 3 = "dotted", 4 = "dotdash", 5 = "longdash", 6 = "twodash")

size - integer (in mm for size of points and text)

linewidth - integer (in mm for widths of lines)

shape - integer/shape name or a single character ("a")



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, unemployment))  
b <- ggplot(seals, aes(x = long, y = lat))
```

a + geom_blank() and **a + expand_limits()**
Ensure limits include values across all plots.

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

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

a + geom_polygon() (aes(alpha = 50)) - x, y, alpha, color, fill, group, subgroup, 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 = unemployment - 900, ymax = unemployment + 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 = 1:1155, 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(f))
```

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

TWO VARIABLES both continuous

```
e <- ggplot(mpg, aes(cty, hwy))
```

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

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, group, 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) - x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust

one discrete, one continuous

```
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

both discrete

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

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

g + geom_jitter() (height = 2, width = 2)
x, y, alpha, color, fill, shape, size

THREE VARIABLES

```
sealsSz <- 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, color, group, linetype, size, weight

l + geom_contour_filled() (aes(fill = z))
x, y, alpha, color, fill, group, linetype, size, subgroup

continuous bivariate distribution

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

h + geom_bin2d() (binwidth = c(0.25, 500))
x, y, alpha, color, fill, linetype, size, weight

h + geom_density_2d()
x, y, alpha, color, group, linetype, size

h + geom_hex()
x, y, alpha, color, fill, size

continuous function

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

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, ymax, ymin, 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, ymin, ymax, alpha, color, fill, group, linetype, shape, size

maps

Draw the appropriate geometric object depending on the simple features present in the data. aes() arguments: map_id, alpha, color, fill, linetype, linewidth.

```
nc <- sf::read(system.file("shape/nc.shp", package = "sf"))
```

```
ggplot(nc) +  
  geom_sf(aes(fill = AREA))
```