

~\OneDrive\3학년 (2025)\빅 데이터\study\_0528.R

```

1  setwd("C:/r_workdate")
2  Sys.setlocale("LC_ALL", "Korean")
3
4  library(lubridate)
5  library(dplyr)
6  library(ggplot2)
7  library(googleVis)
8
9  library(plyr)
10 # =====
11 # ggplot2 패키지 : 다양한 형태의 그래프를 쉽게 표현
12 # =====
13 # =====
14 # ■ ggplot(): 그래프 틀
15 # =====
16 # - plot() 함수의 확장 버전
17 # - 다양한 레이어와 옵션을 조합하여 고급 그래프 구현 가능
18 # - 형식: ggplot(데이터프레임, aes(x = x축 데이터, y = y축 데이터)) + geom_*() 함수
19
20 # =====
21 # ■ geom 함수 (geom_*) : 실제 그래프 형태
22 # =====
23 # - 그래프의 종류와 모양을 지정하는 함수
24 # - 항상 ggplot() 뒤에 + 기호로 연결
25
26 # =====
27 # 1. geom 함수의 stat 옵션
28 # =====
29 # - stat: 주어진 데이터에서 geom에 필요한 데이터를 생성한다.
30 # - stat_bin: 아래와 같은 데이터를 갖는 dataframe을 출력
31 #   1. count : 항목별 빈도수 (기본 막대그래프)
32 #   2. density : 항목별 밀도 (히스토그램 분석용)
33 #   3. ncount : count를 0~1 범위로 정규화 (상대 비교용)
34 #   4. ndensity : density를 0~1 범위로 정규화 (밀도 비교용)
35
36 # =====
37 # 2. geom 함수의 종류
38 # =====
39
40 # -----
41 # 2-1. 산점도 (geom_point)
42 # -----
43 kor = read.table("학생별국어성적_new.txt", header=TRUE, sep=',')
44 ggplot(kor, aes(x=이름, y=점수)) + geom_point()
45
46 # -----
47 # 2-2. 막대그래프 (geom_bar): barplot와 비슷
48 # -----
49 ggplot(kor, aes(x=이름, y=점수)) + geom_bar(stat='identity')
50 # stat = "identity": y 값을 그대로 막대 높이로 사용
51

```

```

52 # ● 막대 테두리/채우기 색상 지정
53 gg1 = ggplot(kor, aes(x=이름, y=점수)) +
54     geom_bar(stat='identity', color='red', fill='green')
55
56 # ● x축 이름 색상, 크기 변경
57 gg1 + theme(axis.text.x = element_text(angle=45, hjust=1, vjust=1, colour='blue', size=8))
58
59
60 # -----
61 # 2-3. 누적 막대그래프 (geom_bar + geom_text)
62 # -----
63 kem = read.csv("학생별과목별성적_국영수_new.csv")
64 skem = arrange(kem, 이름, 과목) #오름차순 정렬
65
66
67 # 하나의 막대그래프에 국/영/수 성적을 표현
68 skem2 = ddply(skem, '이름', transform, 누적합계 = cumsum(점수)) # 누적합계 열 생성: 사람 기
69     준 점수 누적합
70
71 # 각 점수 영역 중간에 점수를 표현
72 skem3 = ddply(skem2, '이름', transform,
73     누적합계 = cumsum(점수),
74     label = cumsum(점수) - 0.5 * 점수) # 점수 영역의 50% 위치 지정
75
76 # 누적 막대그래프 + 텍스트 + 범례 반전
77 gg2 = ggplot(skem3, aes(x=이름, y=점수, fill=과목)) +
78     geom_bar(stat='identity', position=position_stack(reverse=T)) + # 쌓는 순서 역전
79     geom_text(aes(y=label, label=paste(점수, '점')), color='black', size=4) + # 텍스트 표
80     시
81     guides(fill=guide_legend(reverse=T)) # 범례 순서도 역전
82     # position=position_stack(reverse=T)가 없으면 데이터가 반대로 나온다
83
84 # x축 이름 회전 및 정렬
85 gg2 + theme(axis.text.x = element_text(
86     angle = 45,          # 글자를 45도 기울임
87     hjust = 1,          # 수평 정렬 (horizontal justification)
88     vjust = 1,          # 수직 정렬 (vertical justification)
89     color = 'black',    # 글자 색상
90     size = 8            # 글자 크기
91 ))
92
93 # -----
94 # 2-4. 롤리팝 그래프 (geom_segment + geom_point)
95 # -----
96 install.packages("gridExtra")
97 library(gridExtra)
98 mt = mtcars
99
100 # (1) x축 기준 롤리팝
101 # x축에서 수직선으로 연결된 그래프: geom_segment()의 aes(xend=변수x, yend=0) 설정
102 ggplot(mt, aes(x = mpg, y = disp)) +
103     geom_segment(aes(xend = mpg, yend = 0, color = gear), size = 1.3) + # y = 0부터 y =
104     disp 까지 수직선(위)
105     geom_point(aes(color = gear), size = 6) + # 점 찍기

```

```

103     scale_color_continuous(type = "gradient") +           # 색상
104     theme_minimal()                                       # 배경(깔끔)
105
106 # (2) y축 기준 롤리팝
107 # y축에서 수직선으로 연결된 그래프: geom_segment()의 aes(xend=0, yend=변수y) 설정
108 ggplot(mt, aes(x = mpg, y = disp)) +
109     geom_segment(aes(xend = 0, yend = disp, color = gear), size = 1.3) + # x = 0부터 x = mpg
    까지 수평선 (오른쪽으로)
110     geom_point(aes(color = gear), size = 6) +             # 점 찍기
111     scale_color_continuous(type = "viridis") +           # 색상
112     theme_minimal()                                       # 배경(깔끔)
113
114 # -----
115 # 2-5. geom_point()의 다양한 옵션
116 # -----
117 g1=ggplot(mt, aes(x=hp, y=mpg))
118
119 # ● 단순 점
120 g1+geom_point()
121
122 # ● 색상 지정
123 g2= g1+geom_point(color = 'blue')
124 g2
125
126 # ● 그룹별 색상 (am: 0 = 자동, 1 = 수동 (두 가지 색))
127 g3= g1+geom_point(aes(color =factor(am)))
128 g3
129
130 # ● 사이즈 변경
131 g4 = g1+geom_point(size=7)
132 g4
133
134 # ● 각각 사이즈 조절
135 g5=g1+geom_point(aes(size=wt))
136 g5
137
138 # ● 모양 변경
139 g6=g1+geom_point(aes(size=wt, shape=factor(am)))
140 g6
141
142 # ● 색상을 다르게
143 g7 = g1+geom_point(aes(size=wt, shape=factor(am), color=factor(am)))
144 g7
145
146 # ● 원하는 색으로
147 g8 = g1+geom_point(aes(size=wt, shape=factor(am), color=factor(am)))+
148     scale_color_manual(values=c('red', 'green'))
149 g8
150
151 # ● 라인 추가(점 + 선)
152 g9 = g1+geom_point(aes(size=wt, shape=factor(am), color=factor(am)))+
153     scale_color_manual(values=c('red', 'green'))+
154     geom_line()
155 g9

```

```

156
157 # ● 축 이름 변경
158 g10 = g1 + geom_point(aes(size=wt, shape=factor(am), color=factor(am))) +
159     scale_color_manual(values=c('red', 'green')) +
160     geom_line() +
161     labs(x='마력', y='연비')
162 g10
163
164
165 # -----
166 # 2-6. 선 그래프 (geom_line)
167 # -----
168 th = read.csv("학생별과목별성적_3기_3명.csv")
169 ss = arrange(th, 이름, 과목)
170
171 # 학생별 과목 점수 이름 기준으로 선 연결 + 점 표시
172 # - 한 학생이 여러 과목을 가지고 있기 때문에(6개) 표현도 그만큼 해줘야함 6개를 이름으로 묶어서
  표현
173 ggplot(ss, aes(x=과목, y=점수, group=이름, color=이름)) +
174     geom_line() +
175     geom_point(size=6, shape=22) #shape : 0~25
176
177 # =====
178 # R 데이터분석을 위한 패키지 설치
179 # =====
180 # 1. 다국어 처리를 위한 멀티링구얼 패키지
181 install.packages("multilinguer")
182 library(multilinguer)
183
184 # 2. KoNLP 설치 전 필수 패키지 (문자열 처리, 사전, 데이터 처리용)
185 install.packages(c('string', 'hash', 'tau', 'Sejong', 'RSQLite', 'devtools'), type='binary')
186
187 # 3. GitHub에서 KoNLP 설치 (공식 지원 종료 → 수동 설치)
188 install.packages("remotes")
189 remotes::install_github("haven-jeon/KoNLP", upgrade='never', INSTALL_opts=c("--no-
multiarch"), force=TRUE)
190
191 # 4. KoNLP 기본 설정
192 library(KoNLP)
193 useSejongDic() # 세종 사전 사용
194
195 # =====
196 # 워드클라우드 : 텍스트 데이터를 시각적으로 표현
197 # =====
198 # - 텍스트 데이터에서 단어를 추출하고, 단어의 빈도수에 따라 글자 크기를 다르게 표시
199 # - 핵심 패키지: KoNLP (명사 추출), wordcloud (시각화), RColorBrewer (색상 팔레트)
200
201 # -----
202 # 0. 시각화에 필요한 필수 패키지 설치 및 로딩
203 # -----
204 install.packages("stringr")
205 install.packages("wordcloud")
206 install.packages("RColorBrewer")
207

```

```

208 library(stringr)                # 문자열 처리 (단어 전처리)
209 library(wordcloud)              # 워드클라우드 생성
210 library(RColorBrewer)           # 색상 팔레트 제공
211
212 # -----
213 # 1. 데이터에서 단어만 추출
214 # -----
215 d1 = readLines("BTS유엔연설_국문.txt")      # 텍스트 파일 한 줄씩 읽기
216 d2 = sapply(d1, extractNoun, USE.NAMES = FALSE) # extractNoun로 명사를 추출 (벡터형 반환)
217
218 # -----
219 # 2. 단어 집합 생성 (unlist 이용)
220 # -----
221 d3 = unlist(d2)                          # 리스트 → 벡터 변환
222
223 # -----
224 # 3. 단어 필터링 (2글자 이상 단어만 추출)
225 # -----
226 d3=Filter(function(x){
227             nchar(x) >=2
228             }, d3)
229
230 # -----
231 # 4. 단어 핸들링 (불필요한 단어 제거 - 생략 가능)
232 # -----
233 # 불필요한 단어를 일일이 삭제하기 힘들어 미리 목록작성 후 반복문으로 제거
234 # 남산 이라는 단어가 있으면 남/산으로 각각 쪼개지기 때문에 쪼개지지 않게 설정
235
236 # -----
237 # 5. 텍스트 파일로 저장 후 table로 다시 불러오며 공백제거
238 # -----
239 write(unlist(d3), "BTS_kor.txt")          # 단어 벡터를 파일로 저장
240 d4 = read.table("BTS_kor.txt")           # 다시 불러옴→ 데이터프레임 형태(공백
제거)
241
242 # -----
243 # 6. 단어 빈도수 저장 (table() 함수)
244 # -----
245 wc = table(d4)
246
247 # -----
248 # 7. 워드클라우드 출력
249 # -----
250 pal = brewer.pal(9, "Set3")             # 색상 팔레트 정의
251
252 wordcloud(words = names(wc),             # 단어 목록 (table의 이름)
253           freq = wc,                    # 빈도수
254           scale = c(5, 1),             # 단어 크기 (가장 큰 단어: 5, 작은 단어: 1)
255           rot.per = 0.25,              # 회전될 단어 비율 (25%)
256           min.freq = 2,                # 최소 빈도수: 2번 이상 등장만 표시
257           random.color = TRUE,         # 색상 랜덤 적용
258           random.order = FALSE,       # 빈도 높은 단어 중심으로
259           colors = pal)               # 색상 팔레트 적용

```