

eda

eda를 위한 trainset 데이터 셋 정리

```
setwd("C:/Users/ChangYong/Desktop/나노디그리/1. 정규강의 학습자료/1차 프로젝트/소상공인/데이터")
load("dataset.rda")
summary(smallbz_total)
smallbz_total_dummy <- smallbz_total %>% filter(점포수 !=0)
train_dummy <- trainset %>% filter(점포수 != 0)
test_dummy <- testset %>% filter(점포수 !=0)

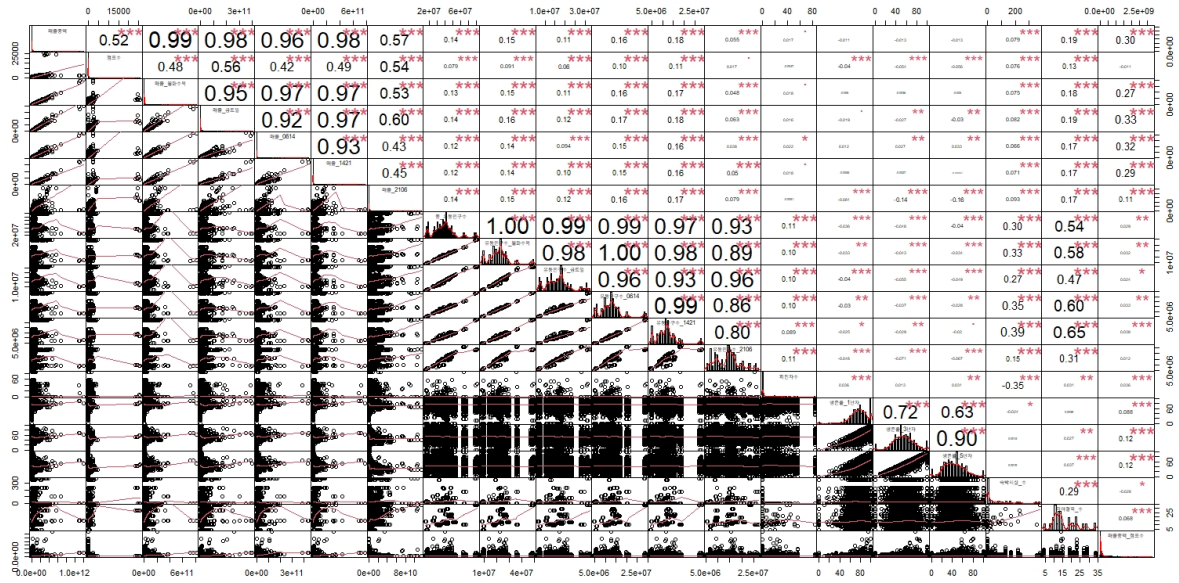
#거리&상권별 값이 다른 컬럼은 sum / 거리&상권에 동일한 값을 적용한 컬럼은 mean
train_dummy <- train_dummy %>%
  group_by(년도, 분기, 행정구역, 대분류, 소분류, 년분기) %>%
  summarise(매출총액 = sum(매출총액),
            점포수 = sum(점포수),
            매출_월화수목 = sum(매출_월화수목),
            매출_금토일 = sum(매출_금토일),
            매출_0614 = sum(매출_0614),
            매출_1421 = sum(매출_1421),
            매출_2106 = sum(매출_2106),
            총_유동인구수 = mean(총_유동인구수),
            유동인구수_월화수목 = mean(유동인구수_월화수목),
            유동인구수_금토일 = mean(유동인구수_금토일),
            유동인구수_0614 = mean(유동인구수_0614),
            유동인구수_1421 = mean(유동인구수_1421),
            유동인구수_2106 = mean(유동인구수_2106),
            확진자수 = mean(확진자수),
            생존률_1년차 = mean(생존률_1년차),
            생존률_3년차 = mean(생존률_3년차),
            생존률_5년차 = mean(생존률_5년차),
            숙박시설_수 = mean(숙박시설_수),
            지하철역_수 = mean(지하철역_수)) %>%
  mutate(매출총액_점포수 = 매출총액/점포수) %>%
  as.data.frame()

#매출 및 유동인구수 scale 조정 및 데이터 정규분포화를 위한 자연로그 적용
vars <- c(7,9:19,26)
train_dummy_log <- train_dummy
train_dummy_log[,vars] <- log(train_dummy_log[,vars])
for(i in vars){
  train_dummy_log[,i] <- ifelse(is.infinite(train_dummy_log[,i])==T,0,train_dummy_log[,i])
}
```

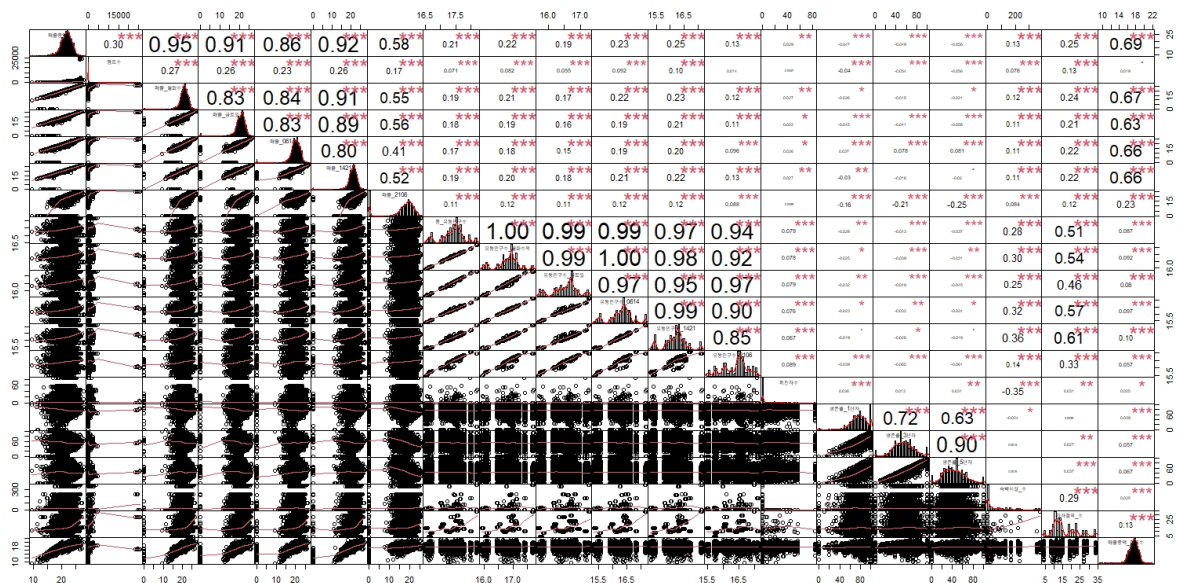
log 적용 전 값이 0인 경우 자연로그를 취하면 무한대가 되는데 이걸 0으로 처리해도 무방한지..

train_dummy 자연로그 적용&미적용 데이터 상관관계 분석

자연로그 미적용



*자연로그 적용



매출 데이터의 경우 분기 및 업태별 매출액이 다르게 표기되어 있으나, 유동인구 & 코로나 확진자수 & 생존율 & 숙박시설 & 지하철 개수 데이터는 분기별로만 데이터가 있어서 상대적으로 데이터가 부족한데, 상관관계가 작아도 사용해도 괜찮은지..

행정구역&년도&분기&업종별 업태 점포 평균 매출액 top_3 누적횟수

This bar chart displays the frequency of the top three search results across various industries and categories. The y-axis represents the count of top-3 hits, ranging from 0 to 150. The x-axis lists 36 different categories. A legend on the right identifies each category by color.

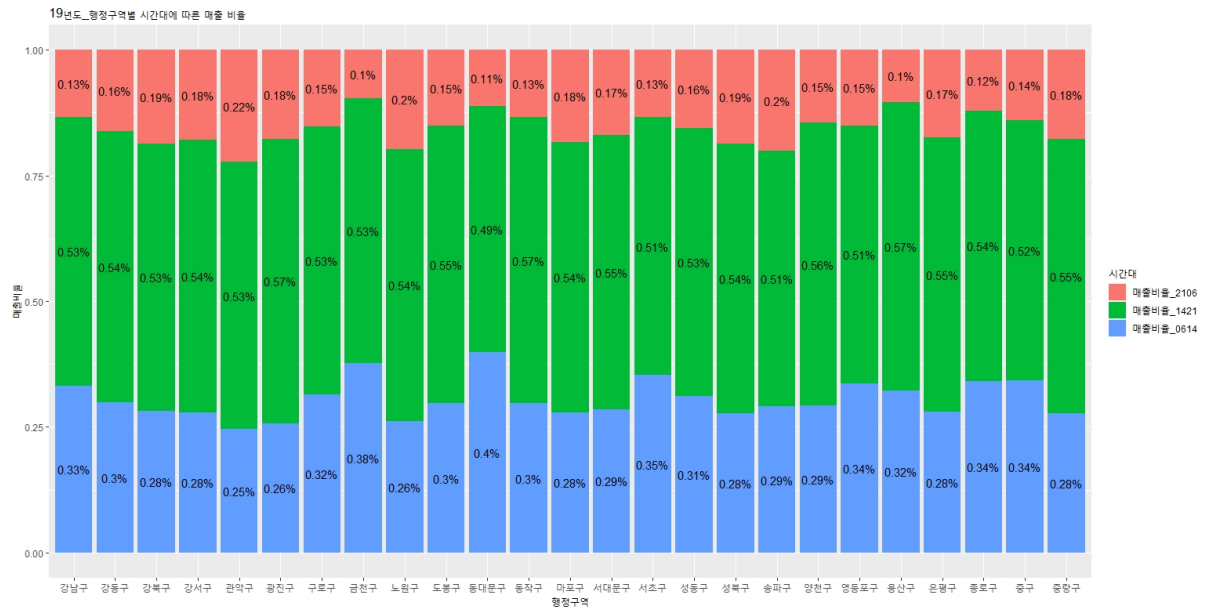
업태	부록 top_3 횟수
pc방	48
가수	17
가방	2
가전제품	62
가전제품수리	2
고시원	22
골프연습장	12
문구	9
미용관매	20
반찬가게	11
저택	11
숙박공간매	13
스노우볼판매	25
양식음식점	48
여관	4
완구	18
외국어배움	10
운동경기용품	7
육류판매	7
의복기	32
의료기관	48
월반의원	145
왕치마	52
자전거 및 기타운송장비	2
자전거	34
조명장치	120
조명장동	5
종이책	49
종교상	21
차량리워	135
커피숍로	2
컴퓨터및전자기기매	12
패션잡화	52
편의점	48
한식음식점	80
한의원	2
호프간이주점	3

행정구역별 시간대에 따른 매출 비율

eda

```
gather(시간대, 매출비율, 매출비율_0614:매출비율_2106) %>%
  arrange(행정구역, 시간대, 매출비율) %>%
  group_by(행정구역) %>%
  mutate(시간대 = factor(x = 시간대, levels = c("매출비율_2106", "매출비율_1421", "매출비율_0614")),
         매출비율_cumsum = cumsum(매출비율) - (0.5*매출비율)) %>%
  ggplot(aes(x = 행정구역, y = 매출비율, fill = 시간대)) +
  geom_bar(stat = "identity") +
  geom_text(aes(label = paste0(round(매출비율, digits=2L), "%"), y=매출비율_cumsum))
```

2019년



2020년

