

## R 데이터(텍스트) 마이닝 및 시각화 분석

# Outline

분산분석(ANOVA, analysis of variance)

ggplot2를 이용한 시각화 분석

## 분산분석(ANOVA, analysis of variance)

▶ n개의 집단을 비교하는 통계적 분포 (단,  $n > 2$ )



## 네카라쿠배당토



# 데이터 생성

```
> #네카라쿠배당토
> coupang <- sample(c(6000:9000),100,replace = TRUE)
> kakao <- sample(c(4000:9000),100,replace = TRUE)
> naver <- sample(c(4000:8000),100,replace = TRUE)
> companyName <- c("coupang","kakao","naver")
> companies <- NULL
> income <- NULL
> employee <- NULL
> for(i in 1:3){
+   for (j in 1:100){
+     companies <- c(companies, i)
+     employee <- c(employee,j)
+     if (i==1){
+       income <- c(income, coupang[j])
+     } else if (i==2){
+       income <- c(income, kakao[j])
+     } else {
+       income <- c(income, naver[j])
+     }
+   }
+ }
```

## R 데이터(텍스트) 마이닝 및 시각화 분석

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ▶ ↺ 🔍 ↻





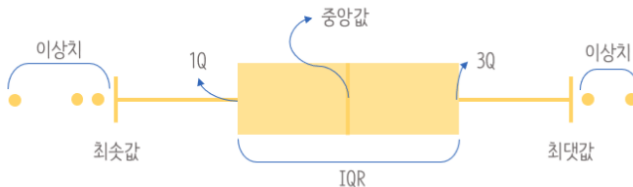




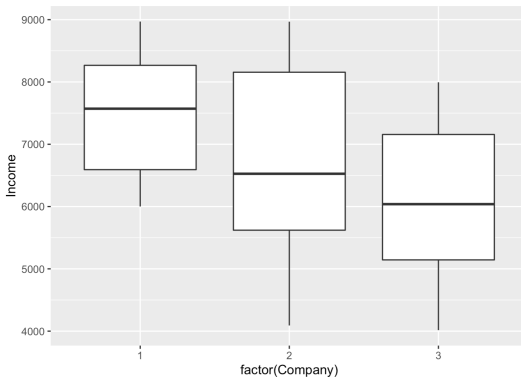


## ggplot2를 이용한 시각화 분석

# 상자그림



# 상자 그림 그리기



```
> library(ggplot2)
> ggplot(overall, aes(x=factor(Company),y=Income)) + geom_boxplot()
```

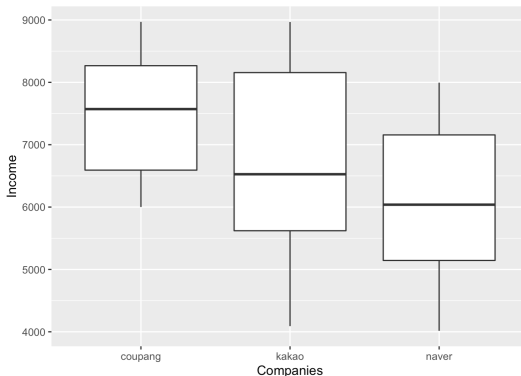
# 문자형 데이터 추가하기

```
> library(dplyr)
> overall_2 <- overall %>% mutate(Companies = ifelse(Company==1, "coupang", ifelse(Company==2, "kaka
ao", "naver")))
> summary(overall_2)
```

Company		Employee		Income		Companies
Min.	:1	Min.	: 1.00	Min.	:4016	Length:300
1st Qu.:	:1	1st Qu.:	25.75	1st Qu.:	5774	Class :character
Median	:2	Median	: 50.50	Median	:6777	Mode :character
Mean	:2	Mean	: 50.50	Mean	:6738	
3rd Qu.:	:3	3rd Qu.:	75.25	3rd Qu.:	7811	
Max.	:3	Max.	:100.00	Max.	:8970	

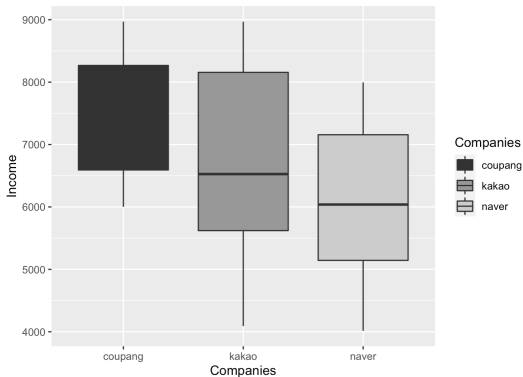


# 문자 데이터를 사용한 상자 그림



```
> #문자 데이터를 사용한 상자 그림; 회사 이름을 사용해서 다시 상자그림 그리기  
> ggplot(overall_2, aes(x=Companies,y=Income)) + geom_boxplot()
```

# 색상 지정하기



```
> # 상자 그림 회색 계열로 색상지정 (논문)
> ggplot(overall_2, aes(x=Companies,y=Income, fill=Companies)) + geom_boxplot() + scale_fill_grey()
```