

## < 데이터 2 단계 정리 >

0 데이터파일 읽기 ← read\_csv("file.csv") : file를 데이터파일 이름으로 저장하고 불러오기

☆ 우선 library(readr) 설치하기

0 str(데이터 아 데이터 \$변수) : 데이터 또는 데이터의 변수에 대해 다양한 정보 (class, 개수, etc...)

0 as\_tibble(데이터) : 데이터의 다양한 정보를 보며 as\_tibble과 조금 같음

☆ 우선 library(dplyr) 설치하기

0 summary(데이터) : 기본통계량을 보여줌, 즉 문맥형식도 제아하고 계량형식도 (int, num...)의 수서적인 것

결측치의 개수 확인 (NA)

0 freq(데이터 \$ 변수) : 해당변수의 빈도와 막대 그래프

☆ 우선 library(desc) 설치하기

0 데이터 %>% group\_by(변수) %>% summarise(변수명 = mean or var or sd... (변수))

: 변수별로 그룹을 나고 기본통계량을 조건에 맞게 제시

☆ 우선 library(dplyr) 설치하기

0 mean/var/sd(데이터 \$ 변수, na.rm=T) : 데이터 변수에서 결측치 (NA)를 제외하고 계산

0 boxplot(데이터 \$ 변수) : 해당변수의 상자 그래프, 단 변수 계량형식도 이상치 check 가능

0 데이터 \$ 변수 ← as.factor(데이터 \$ 변수) : 해당변수 class를 범주형으로 변경

0 데이터 \$ 변수 ← as.POSIXct(데이터 \$ 변수) : 해당변수 class를 시간형으로 변경

0 table(데이터 \$ 변수) : 해당변수의 빈도

0 데이터 \$ 변수 ← factor(데이터 \$ 변수, levels=c("a", "b", "c"...))

: a, b, c 순으로 변수를 배열하고 class를 범주형으로 바뀐 해당데이터의 변수에 지정

0 데이터 \$ 변수 ← ifelse(데이터 \$ 변수 == "a", "b", 데이터 \$ 변수) :

만약 데이터의 변수값이 a면 b로 바꾸고 아니면 데이터의 원래 변수값 그대로 두어라

0 is.na(데이터 \$ 변수) : 데이터의 변수가 결측치 (NA)이다.

0 데이터 ← 데이터 %>% drop\_na() : 결측치 (NA)가 있는 행들을 다 제거

☆ 우선 library(tidy) 설치하기

0 descriptive ← describe(데이터) : 해당데이터들을 계량적으로 계산해 descriptive에 해당 (mean, sd, min, max...)

☆ 우선 library(psych) 설치하기

0 descriptive ← descriptive %>% mutate(LL = mean - 3 \* sd, UL = mean + 3 \* sd)

0 데이터 ← 데이터 %>% filter(변수 < (K) - UL)

: 위에서 생성한 UL 값보다 작은 값들만 filter 이용해서 추출 (UL보다 큰 값은 이상치)

0 데이터 ← 데이터 %>% filter(변수 > (K) + UL)



## 1. 다중회귀분석

→ 완벽하게 파악하지 (결측치 제거)

0 `corr.test`(데이터[,k:n], method="pearson", alpha=0.05, use="pairwise.complete.obs")  
(k가 0이면 데이터의 전체 행을 의미함) `data`에서 `C(7:14)` → 7월~14월까지의 온도 변화

※ `library(psych)` 설치하기

※ 상관계수를 알아보는 데 데이터에서 k행 n행끼리, 방법은 피어슨 상관행

→ '조각'을 나눠줌

0 데이터의 잔차 = `lm(DV ~ IV, data=데이터)` : 다중회귀수립하기

0 `plot()` : 5종한 4개의 그래프로 산형성, 정규성, 등분산성 확인

0 `ks.test`(데이터\$변수, pnorm, mean=mean(데이터\$DV), sd=sd(데이터\$DV)) :

데이터의 종속변수에서 정규분포를 2가지 검정, 평균과 표준편차 이용

0 `shapiro.test`(데이터\$변수) : N이 작을 때의 데이터의 종속변수의 정규성 검토

0 `hist`(데이터\$변수, breaks=seq(M1, M2, k)) : 데이터의 분포를 범위를 M1 ~ M2, 간격을 k

안정하는 히스토그램 만들기 → 실제 10을 취하면 상대적으로 정규분포 형태를 띈다.

0 `durbinWatsonTest`(데이터) : 오차의 자기상관 확인

#### 다중회귀분석 실습####

### 1) 상관분석 실행###

`corr.test(bicycle[,c(7:14)], method="pearson",  
alpha=0.05, use="pairwise.complete.obs")`

#bicycle 데이터프레임에서 모든 행에서 7월~14월까지

#----> 결과: 2p-value가 알파인 0.05와 비교했을 때  
대부분의 값이 0이므로 모든 변수들 간에 상관관계와  
상관계수 추정치가 통계적으로 유의하다. (2p-value < 0.05)  
#가장 강한 양의 상관관계: temp & atemp(0.99)

### 2) 연구가설 수립###

# IV: temp, atemp, humidity, windspeed, difference & DV  
: total

#H1: temp -> total(+) #temp는 total에 양의 인과관계를  
미친다/두 변수간에는 양의 인과관계를 갖고 있다./temp가  
높아질 수록 total도 높아진다.

#H2: atemp -> total(+)

#H3: humidity -> total(-)

#H4: windspeed -> total(+)

#H5: difference -> total(+)

### 3) 다중회귀식 수립###

`lm1 <- lm(total ~ temp + atemp + humidity + windspeed`

`+ difference, data=bicycle)` `lm1 <- lm(종속변수 ~ 독립변수)`

# $\hat{Y}_i = a + b_1X_{1i} + b_2X_{2i} + \dots + b_5X_{5i}$

# $X_{1i}$ =temp..... $X_{5i}$ =difference

### 4) 다중회귀분석 전제조건 확인: plot(m1)

#residuals vs fitted: 빨간색 선이 하얀색 점선을 따라서  
일직선이여야 선형성 조건 만족(그렇지 않음) # 정규성 조건

#normal Q-Q: 점들이 대각선 위에 존재하면 정규성 조건  
만족(그렇지 않음)

#scale-location: 빨간색 선이 일직선이여야 등분산성과  
선형성 조건 만족(그렇지 않음)

#residuals vs leverage: leverage는 사례가 다른  
사례로부터 떨어진 정도로 0에 가까울 수록 바람직함. 일부  
사례가 떨어져 존재하여 이상치로 볼 수도 있음. 빨간선은  
수평의 일직선이 바람직함(등분산성). cook's distance는 0  
.5나 혹은 1을 넘으면 해당 사례가 회귀계수 추정치에  
지나치게 많은 영향을 미침(그런 사례가 없음) ~~있으면 0.5, 1이 넘으면~~  
#6720, 6721, 6722, 8915, 8917, 8918을 제거하면 선형성  
/정규성/등분산성이 개선됨(각 case 번호는 6728, 6729..)

```
bicycle <- bicycle %>% filter(case != 6728, case != 6729,
case != 6730, case != 9004, case != 9006, case != 9007)
lm1 <- lm(total~temp+atemp+humidity+windspeed
+difference, data=bicycle) #6개 사례 제외한 bicycle로
업데이트
```

### 5) 정규성 조건 확인 ### ~~종속변수를 기준으로 check~~

```
ks.test(bicycle$total, pnorm,
mean=mean(bicycle$total), sd=sd(bicycle$total))
```

#Kolmogorov-Smirnov test: n이 클 때 정규성 검토. 2p  
-value=0 < 0.05 즉 유의하기에 정규성 조건 만족 못함.

0.2(P-value) > 0.05(alpha) ~~정규성 만족~~   
shapiro.test(bicycle\$total)   
L> 2표준, 통계적으로 정상이

#shapiro.test: n이 작을 때 정규성 검토(현재는 n이  
10723이기에 실행하면 오류가 뜸)

```
hist(bicycle$total, breaks = seq(0,1000,10))
```

```
hist(log(bicycle$total), breaks = seq(0,10,0.1))
```

#log: 밑이 e(무리수)인 자연로그

### 6) 독립성(오차의 자기상관) 검토###

```
library(car)
```

```
durbinwatsonTest(lm1)
```

# 오차는 서로 양의(0.9153073) 자기상관 존재(p-value=0  
혹은 rho !=0, 즉 상관관계가 있고 독립성이 아니다.)

#DW가 0.1693747로 2보다 많이 멀리 떨어져 있기에 오차의  
자기 상관이 존재한다(2에 가까우면 존재) 0에 가까우면 음의 자기 상관이  
존재한다. 4에 가까우면 음의 자기 상관 존재

~~오차의 자기상관이 존재하지 않아야 독립성조건에 만족함~~  
(즉 DW값이 2에 가까워야 함)



## 다중회귀분석 Real 실험

### ① 다중회귀 분석 결과를 통한 가설 검증

summary(lm1) : summary 함수에 의해 만든 회귀식을 통해서 회귀계수 추정치(Estimate)와 유의성(P-value) 검토하기 (변수명, P-value vs alpha 형태, R에서 나온 P-value 값은 0.05보다 작으면 0.1과 비교)

### ② 수정된 다중회귀식 구하기

$\hat{Y}_i(\text{Chat}) = 44.630 (\text{Intercept의 estimate}) + 1.125 \times 1T + 4.129 \times 2T$

(0에서 나온 P-value 값과  $\alpha=0.1$ 을 비교해서 유의한 변수들의 estimate를 이용해서 수정 다중회귀식 작성)

(Multiple R: 0.9612, Adjusted R: 0.9761) 큰 수이기 때문에 적합도, 설명력이 높은 회귀식이다

(P-value < 2.2e-16으로 0에 가깝고 유의 즉, 모든 회귀계수가 0이 아닌 귀무가설 기각하고 적어도 하나의 회귀계수는 0이 아닌 대립가설 채택)

### ③ 모형 적합도 제고를 위한 다중회귀 분석 방법

lm1 <- stepAIC(lm1, direction="forward") : lm1 데이터로 'forward' 다중회귀 분석  
 lm1 <- stepAIC(lm1, direction="backward") : lm1 데이터로 'backward' 다중회귀 분석  
 lm1 <- stepAIC(lm1, direction="both") : lm1 데이터로 'forward', 'backward' 혼합  
 summary(lm1 <- stepAIC(lm1, direction="both")) : summary 함수를 이용해서 위에 데이터 형상할 때 나온 star와 AIC=k  
 의미를 각각 비교해서 작문값이 나온 방법을 택함 (결과가 3가지 방식이야! 그중 유의성 높은 변수를 제거하는 게 다함)

### ④ 다중공선성 확인

VIF(lm1) : VIF 함수를 이용해서 다중공선성을 확인

주요 변수 library(CAR) 설치하기

아무 VIF의 값이 5.3을 넘고 1을 제일 값이 큰 변수를 제외하고 앞선 과정 반복하기 (Atemp 변수 제외)

lm2 <- lm(total ~ temp + humidity + windspeed + difference, data=bicycle)

VIF(lm2) : Atemp 변수 제외하고 만든 회귀식의 다중공선성 확인

summary(lm2) : 새로운 회귀식의 요약치 검토(P-value, estimate, ...)

동일한 방법으로 temp 변수도 제외한 데이터 lm3을 만들고 VIF, summary로 검토

### ⑤ IV의 중요도 (표준화 회귀계수 추정치의 절대값 크기)

lm.beta(lm1) : 표준화 회귀계수 추정치들을 절대값 기준으로 크기 비교해서 중요도 확인

### ⑥ 새로운 변수(IV) 추가 타당성 검토

lm4 <- lm(total ~ humidity + windspeed + difference + working, data=bicycle) : working 추가

summary(lm4) : R-squared 와 adjusted R-squared의 값을 확인해서 추가 타당성 검토

anova(lm3, lm4) : R2의 변화량인 0.045가 통계적으로 유의한지 확인

이후 summary를 이용해서 lm4에 대한 추정결과 자세히 확인

lm5 <- lm(total ~ humidity + windspeed + difference + working + season, data=bicycle)

두 번째 추가 IV로 season 추가하기



0 summary(lm5)  
 0 anova(lm4, lm5):  $R^2$ 의 변화량인 0.0348이 통계적으로 유의한지 확인

⊕ 이후 summary를 이용하여 lm5에 대한 추정결과 자세히 확인

0 vif(lm5): 다중공선성 확인

\* 우선 (tbl\_dty(car) 인차해)

G VIF  $\wedge (1/(2 \times DF))$ 의 값이 2보다 작으면 다중공선성은 문제없음

(7) 새로운 독립변수 atemp 추가 문제 확인

0 lm6 <- lm(total ~ humidity + windspeed + difference + working + season + atemp, data =

! atemp 변수 추가

0 summary(lm6)

0 anova(lm5, lm6):  $R^2$ 의 변화량인 0.032가 통계적으로 유의한지 확인

0 vif(lm6): 다중공선성 확인

0 lm.beta(lm6): 각 독립변수의 상대적 중요도를 정량화 기준으로 확인

⊕ 가장 설명력이 높은 lm5에 기반하여 IV속성값이 존재할 때 종속변수(매출액) 예측

: lm5 회귀계수가 부여한 속성값 안에서 매출액 구함

→ 새로운 독립변수를 추가했을 때 모형설명력과 적합도가 좋아졌다고 무조건 이 변수를

추가하는 것이 애매, 상대적으로 중요도가 떨어지는 변수들에 의해 왜곡이 발생한 것 (season full year

⑨ 조절효과 확인하기

0 bicycle\$working <- as.numeric(bicycle\$working): working 변수를 수치형으로 변경

0 bicycle <- bicycle %>% mutate(inter = humidity \* working): 조절효과에 따른 새로운

0 lm7 <- lm(total ~ humidity + windspeed + difference + working + inter, data = bicycle)

: inter를 추가한 회귀식 실행 (lm4에서 inter 추가)

0 summary(lm7)

0 anova(lm4, lm7):  $R^2$ 의 변화량인 0.0122가 통계적으로 유의한지 확인

\* 이후 조절효과를 상세히 확인



## <정태분 2 기말 정리> <3장 과제 완성>

데이터 파일 읽기 ← `read_csv("파일.csv")` : 데이터 파일 불러오기

↳ `library(readr)` 설치하기

• `summary(데이터)` : 값의 NA 존재 확인하기

• `str(데이터)` : 해당 데이터의 구조 살펴보기 (ex) num, chr, ...)

• `desc <- describe(데이터[A:B])` : 데이터에서 A와 B 사이의 변수의 기술 통계량 확인

• `desc <- desc %>% mutate(U = mean + 3 * sd, L = mean - 3 * sd)` : 이상치 기준 추가

↳ `library(dplyr)` / `library(psych)` 설치하기

• `logit1 <- glm(종속변수 ~ 독립변수, data=데이터이름, family=binomial())` : 로짓 회귀 모델

• `outlierTest(logit1)` : 이상치 확인하기 (위에 로짓 회귀식을 바탕으로) → `library(car)` 설치하기

• `티어택 <- 데이터 %>% filter(변수 > 2인것)` : 변수가 2인 것만 선택해서 데이터에 할당

↳ `library(dplyr)` 설치하기

• `summary(logit1)` : 로짓 회귀에서 중요한 결과들을 설명

• `hoslem.test(logit1$y, logit1$fitted.values)` : 로짓 회귀의 모형 적합도 확인

↳ `library(resourceSelection)` 설치하기

• `pseudor2(logit1, which=c("CoxSnell", "Nagelkerke"))` : 로짓 회귀의 모형 적합도 확인

↳ `library(descTools)` 설치하기

→ 이후 로짓 회귀에 `intercept` 추가한 로짓 회귀 2인 것 → 이상치 확인과 제1회 로짓 회귀에 대안됨

→ `summary`로 각각의 2개를 비교 해서 로짓 회귀 2가 더 영향을 줬는지 확인 후 증명

• `difference <- logit1$deviance - logit2$deviance`

• `dof <- logit1$df.residual - logit2$df.residual`

• `1 - pchisq(difference, dof)`

↳ 나오는 값이 0.3 정도 확인

→ 이후 로짓 회귀를 대상으로 `hoslem.test` 함수와 `pseudor2` 함수를 이용해서 모형 적합도 확인

• `prediction <- predict(logit2, newdata=employ)` : 로짓 value 도출

• `prediction <- ifelse(prediction < 0, 0, 1)` : 위에 값을 0을 기준으로 0, 1로 변경

• `employ$result <- as.factor(employ$result)` : result의 척도를 범주형으로 변경

(prediction도 동일하게 변경)

• `ConfusionMatrix(employ$result, prediction)` : 결과를 통해 hit ratio 구하기 (= accuracy)

• `cl001 <- data.frame(id=1001, GPA=3.58, ...)` : 문제에 맞게 만들기

• `predict(logit2, cl001)` : 로짓 회귀를 이용해서 합격 여부 예측

↳ ~~생성된 모델에~~ 예측을 원하는 변수들



## <군집분류> 1- 계층군집분류 데이터 이름

```

distance-hr <- dist(x(HRA-hr, method="power"))
# 사례간의 거리 측정하여 객체 생성
# 원의 library(cluster) 불러오기
HRA-CA-hr <- hclust(distance-hr, method="ward.D2")
# 객체 hclust와 함께
plot(HRA-CA-hr, col="steelblue", main="HRA") # 덴드로그램 만들기
HRA-hr <- as.numeric(HRA-hr$merge) # 계량형 척도로 변경
set.seed(k) # 실행할 때마다 결과가 달라지는 것을 막기 위해 설정 (k는 아무 숫자 아)
HRA-hr-NC1 <- Nbclust(HRA-hr, distance="euclidean", min.nc=
, max.nc=b, method="averager") # 최소 a 최대 b를 설정하고 최적의 군집 개수
HRA-hr-HCA <- cutree(HRA-CA-hr, q) # 덴드로그램에서 q개만큼 잘라내기
# HRA-hr-NC1에서 나온 숫자 불러
(군집별 번호)
result-hr <- aggregate(HRA-hr, by=list(cluster=HRA-hr-HCA), mean)
# 군집별 차이와 특징 도출
library(Nbclust) 불러오기

```

## 2- 비계층군집분류

```

HRA-hr-k <- scale(HRA-hr-k) # 변수의 표준화한 데이터 만들기
HRA-hr-NC2 <- Nbclust(HRA-hr-k, distance="euclidean", min.nc=
max.nc=b, method="kmeans") # 최적의 군집 개수 찾기
HRA-hr-kCA <- kmeans(HRA-hr-k, centers=k, nstart=25)
# 위에서 구한 최적의 군집 개수 k를 적용하여 비계층적 군집분석 실행
result-hr2 <- aggregate(HRA-hr, by=list(cluster=HRA-hr-kCA$
), mean) # 군집별 차이와 특징 도출
library(CMclust) 불러오기

```

## <KNN>

```

cancer-z <- as.data.frame(scale(cancer[2:31])) # [-1] : 1행을 제외
# 제외하고 나머지 IV변들을 표준화해서 데이터 생성
ind <- sample(2, nrow(cancer-z), replace=T, prob=c(0.7, 0.3))
# test와 train을 구분하는 ind 객체 생성
cancer-train <- cancer-z[ind=1,]
cancer-test <- cancer-z[ind=2,]
# train과 test 내에서 생성

```



0 gird 1 ← expand.grid(k=3:10) : k범위를 3~10에서 설정

↳ library(caret) 먼저하기

0 Control ← trainControl(method="repeated cv", number=10,

repeats=5) : 학습 방법 선택

0 knn.train ← train(diagnost\$N, data=cancer.train, method="knn",

fitControl=Control, tuneGrid=gird 1) : 최적의 K값 선정

↳ library(caret)

0 ValIMP(knn.train, scale=F) : 모델 검증 확인

↳ 앞서 scale 적용했기에

0 pred.test ← predict(knn.train, newdata=cancer.test)

생성한 knn.train을 test에 데이터에 적용해서 성능 판단

0 ConfusionMatrix(pred.test, cancer.test\$diagnost) : 실제값 vs 예측값

↳ library(caret)

↳ 예측값

↳ 실제값

0 kkn.train ← train(kknn(diagnost\$N, data=cancer.train,

kmax=25, (distance=2) kernel=c("다양한 10가지 방법 작성"))

: 최적의 K값과 방법 찾기

↳ 유클리드 거리

0 kkn.pred.test ← predict(kkn.train, newdata=cancer.test)

: 개선된 모델과 실제값 적용해서 비교

0 ConfusionMatrix(kkn.pred.test, cancer.test\$diagnost)

: 개선된 모델을 실제값 vs 예측값

## <SVM실습>

0 ltnear.svm ← tune.svm(diagnost\$N, data=cancer.svm.train,

kernel="ltnear" cost=c(0.1, 0.25, 0.5, 0.75, 1, 2, 3, 4, 5, 7, 10))

: 선형(ltnear) 방법으로 train 훈련

↳ 옵션 (파라미터) 값 지정

↳ library(e1071) 먼저하기

0 ltnear.test ← predict(ltnear.svm\$best.model, newdata=cancer.svm

test) : 생성한 ltnear의 최적의 모델을 test에 적용해서 성능 판단

0 ConfusionMatrix(ltnear.test, cancer.svm.test\$diagnost)

: Accuracy, kappa 값은 맞는지 확인 → 예측값 vs 실제값

↳ library(caret) 먼저하기

=> 나머지 3가지 방법 증명하기 진행