

# 분석프로그래밍 I.12

## 고급 통계 분석

국민대학교 경영학부 빅데이터경영통계학 전공

2018. 6. 04(월)

### 1 일반 선형 모형 General Linear Model

- 분산분석ANOVA: ANalysis Of VAriance와 다중회귀Multiple Regression을 모두 포함하는 개념

```
1 data(tips, package='reshape2')
2 ## ANOVA
3 tipAnova <- aov(tip ~ day - 1, tips)
4 tipIntercept <- aov(tip ~ day, tips)
5
6 # 계수 비교
7 coef(tipAnova)
8 coef(tipIntercept)
9
10 # 계수 검정
11 library(lmtest) # coeftest
12 coeftest(tipAnova)
13 coeftest(tipIntercept)
14
15 ## 다중회귀
16 tipLm <- lm(tip ~ day - 1, tips)
17 tipLmIntercept <- lm(tip ~ day, tips)
18
19 ## 코딩 방법 바꾸기: 코딩 방법에 따라 계수의 의미가 달라진다.
20 tips2 <- tips
```

```

21 levels(tips2$day)
22 tips2$day <- relevel(tips2$day, ref="Sun"); # new reference,
    Sunday
23 coef(lm(tip ~ day, tips2))
24 coef(lm(tip ~ day, tips2, contrasts = list(day = contr.sum)))
25 coef(lm(tip ~ day, tips2, contrasts = list(day = MASS::contr.sdif
    )))

```

Listing 1: aov와 lm의 분석 결과 비교

## 2 다중 선형 회귀 Multiple Linear Regression

- 여러 개의 예측 변수predictors로 하나의 결과 변수outcome variable을 예측하고자 한다.
- 결과 변수  $y$ 의 조건부 평균은 계수  $\beta_0, \beta_1, \dots, \beta_k$ 의 선형 결합linear combination으로 나타낼 수 있다.

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik} + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2)$$

$$E[y_i | x_{i1}, x_{i2}, \dots, x_{ik}] = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik}$$

이를 행렬로 나타내면  $Y = X\beta + \epsilon$ 이 되고,  $\hat{\beta} = (X^T X)^{-1} X^T Y$ 를 구할 수 있다.

- 기본적인 가정 : LINE(Linear, Independent, Normal, Equal variance)
- lm 함수를 사용한다. `lm(formula, data)`에서 *formula*를 통해 모형을 정한다. 이때 *formula*에 쓰이는  $+$ ,  $-$ ,  $*$ ,  $/$ 는 산술적인 가감승제가 아니라 Wilkinson-Rogers 표기법을 따른다.

```

1 dat = data.frame(y= , x1= , x2= , x3= )
2 lm(y ~ x1 + x2, data=dat)
3 lm(y ~ x1 + x2 + 1, data=dat)
4 lm(y ~ x1 + x2 - 1, data=dat)
5 lm(y ~ I(x1 + x2), data=dat)
6 lm(exp(y) ~ x1 + x2, data=dat)
7 lm(y ~ exp(x1) + x2, data=dat)
8 lm(y ~ x1 * x2, data=dat) # x1 + x2 + x1:x2

```

```

9  lm(y ~ (x1 + x2 + x3)^2 - x1:x2, data=dat) # x1 + x2 + x3 + x1:
    x3 + x1:x2
10  lm(y ~ x1 / x2, data=dat) # x1 + x2:x2, x1 + x2 %in% x1
11  fit <- lm(y ~ ., data=dat)
12  update(fit, formula = . ~ . - x1)

```

Listing 2: 함수 `lm`의 *formula* 예

- 교재의 housing 자료에 대한 분석의 예

```

1  # 교재 p.356
2  housing <- read.table('https://www.jaredlander.com/data/housing.
    csv',
3
4                      sep = ",",
5                      header = TRUE,
6                      stringsAsFactors = FALSE)
7  head(housing)
8  names(housing) <-c("Neighborhood", "Class", "Units", "YearBuilt",
9                    "SqFt", "Income", "IncomePerSqFt", "Expense",
10                   "ExpensePerSqFt", "NetIncome", "value",
11                   "ValuePerSqFt", "Boro")
12 #p364
13 house1 <- lm(ValuePerSqFt ~ Units + SqFt + Boro, data=housing)
14 fHouse1 <- formula(ValuePerSqFt ~ Units + SqFt + Boro)
15 summary(house1)
16 house1$coefficients
17 coef(house1)
18 coefficients(house1)
19
20 library(coefplot) #install.packages('coefplot')
21 coefplot(house1)
22
23 # coding
24 housing2 <- housing # reference: Bronx
25 levels(housing2$Boro)

```

```
26 housing2$Boro <- relevel(housing2$Boro, ref="Manhattan"); # new
    reference, Manhattan
27 coef(lm(fHouse1, housing2))
28 coef(lm(fHouse1, housing2, contrasts = list(Boro = contr.sum)))
29 coef(lm(fHouse1, housing2, contrasts = list(Boro = MASS::contr.
    sdif)))
30
31 # coefficients
32 coef(house1)
33 confint(house1) # confidence intervals of the parameter estimates
34 lmtest::coefci(house1)
35 # compare the speed of confint and coefci for glm.
36 lmtest::coeftest(house1)
37 vcov(house1) # estimated variance-covariance matrix for parameter
    estimates
38
39 # fitted values, new predictions
40 fitted(house1) # same as predict(house1, housing) or predict(
    house1, house1$data)
41 # fitted values
42 predict(house1, data.frame(Units=45, SqFt = 554000, Boro = "
    Manhattan")) # prediction for new data
43
44 # residuals
45 residuals(house1)
46 resid(house1)
47 car::residualPlots(house1) # car(Companion to Applied Regression)
48
49 # fit
50 deviance(house1) # residual sum of squares(RSS) for linear models
51 logLik(house1)
52 AIC(house1)
53 BIC(house1)
54
55 # variance inflation factors
56 car::vif(house1)
```

```
57
58 # diagnostics
59 plot(house1, which=1:6)
60 car::qqPlot(house1)
61 car::inverseResponsePlot(house1)
62
63 # outliers
64 car::outlierTest(house1)
65 car::influenceIndexPlot(house1)
66 car::influenceIndexPlot(house1, id = list(n=5)) # 5 most
    influential obs.
67
68 # New models
69 house2 <- lm(ValuePerSqFt ~ Units * SqFt + Boro, data=housing) #
    Nesting Model
70 house3 <- lm(ValuePerSqFt ~ Units : SqFt + Boro, data=housing) #
    Non-Nesting Model
71
72 # model comparison
73 anova(house1, house2)
74 lmtest::waldtest(house1, house2)
75 lmtest::encomptest(house1, house3)
76
77 multiplot(house1, house2)
78 car::compareCoefs(house1, house2, se=F)
79 car::compareCoefs(house1, house2)
80
81 AIC(house1, house2, house3)
82 BIC(house1, house2, house3)
83
84 # Stepwise variable selection p.415
85 nullModel <- lm(ValuePerSqFt ~ 1, data = housing)
86 fullModel <- lm(ValuePerSqFt ~ Units + SqFt * Boro + Boro * Class
    , data =housing)
87
88 houseStep <- step(nullModel, scope = list(lower = nullModel,
```

```
upper = fullModel), direction = 'both')
```

Listing 3: 교재의 housing 자료에 대한 분석 예

### 3 일반화 선형 모형 Generalized Linear Model

- 결과 변수가  $-\infty \sim \infty$ 의 연속값을 갖지 않는 경우(예. 개체수, 사망/생존 등)에 사용할 수 있다.
- 결과 변수의 기댓값은 계수의 선형 결합에 선형/비선형 함수를 적용한 결과로 결정된다.<sup>1</sup>

$$\mathbb{E}[Y|X] = g^{-1}(X\beta)$$

로 표현할 수도 있다.

- 교재의 예

```
1 # 20. Generalized linear model. p376
2 acs <- read.table('http://jaredlander.com/data/acs_ny.csv', sep='
  ', header=T, stringsAsFactors=F)
3 head(acs)
4 acs$Income <- with(acs, FamilyIncome >= 150000)
5 library(ggplot2)
6 library(useful)
7 ggplot(acs, aes(x=FamilyIncome)) +
8   geom_density(fill = 'grey', color='grey') +
9   geom_vline(xintercept=150000) +
10  scale_x_continuous(labels=multiple.dollar, limits=c(0,1000000))
11
12 income1 <- glm(Income ~ HouseCosts + NumWorkers + OwnRent +
  NumBedrooms + FamilyType, data = acs, family=binomial(link='
  logit'))
13 summary(income1)
14
15 invlogit <- function(x) {
16   1/(1+exp(-x))
```

<sup>1</sup> $g[\mathbb{E}[Y|X]] = X\beta$

```

17 }
18
19 invlogit(income1$coefficients)
20
21 # p380
22 ggplot(acs, aes(x=NumChildren)) + geom_histogram(binwidth=1)
23 children1 <- glm(NumChildren ~ FamilyIncome + FamilyType +
24   OwnRent, data = acs, family = poisson(link="log"))
25 summary(children1)
26 coefplot(children1)
27
28 z <- (acs$NumChildren - children1$fitted.values) / sqrt(children1
29   $fitted.values)
30 sum(z^2)/children1$df.residual
31 pchisq(sum(z^2), children1$df.residual)
32
33 children2 <- glm(NumChildren ~ FamilyIncome + FamilyType +
34   OwnRent, data = acs, family = quasipoisson(link='log'))
35 summary(children2)
36 coefplot::multiplot(children1, children2)
37
38 deviance(children1) #
39 logLik(children1)

```

Listing 4: 교재의 housing 자료에 대한 분석 예

## 4 조절 효과 Moderation Effect

- 조절 효과의 예: 남녀에 따른 키와 체중의 관계

```

1 gender <- rep(c("M","F"), each=100)
2 height <- c(rnorm(100, 175, 13), rnorm(100, 160, 10))
3 weight <- rep(NA, 200)
4 weight[1:100] <- 0.4*height[1:100] + 8 + rnorm(100)*10
5 weight[101:200] <- 0.25*height[101:200] + 18 + rnorm(100)*10
6
7 plot(weight ~ height, col=factor(gender))

```

```

8
9 fit1 <- lm(weight ~ height)
10 fit2 <- lm(weight ~ gender)
11 fit3 <- lm(weight ~ height : gender)
12 fit4 <- lm(weight ~ height * gender)
13
14 car::compareCoefs(fit1, fit2, fit3, fit4)
15 # fit1
16 plot(weight ~ height, col=factor(gender))
17 abline(fit1, col="black")
18 # fit2
19 plot(weight ~ height, col=factor(gender))
20 abline(h = coef(fit2)["(Intercept)"])
21 abline(h = sum(coef(fit2)), col="red")
22 # fit3
23 plot(weight ~ height, col=factor(gender))
24 abline(a = coef(fit3)["(Intercept)"], b= coef(fit3)["height:
    genderF"])
25 abline(a = coef(fit3)["(Intercept)"], b= coef(fit3)["height:
    genderM"], col='red')
26 # fit4
27 plot(weight ~ height, col=factor(gender))
28 abline(a = coef(fit4)["(Intercept)"], b= coef(fit4)["height"])
29 abline(a = sum(coef(fit4)[c("(Intercept)", "genderM")]),
30       b = sum(coef(fit4)[c("height", "height:genderM")]), col='
    red')
31
32 coplot(weight ~ height | gender)

```

Listing 5: 조절 효과의 예

## 5 선형 혼합 효과 모형 Linear Mixed Effect Model

- 집단별로 다른 선형 회귀선을 고려하여 분석하는 방법이다.
- 집단별로 다른 선형 계수가 특정한 분포(예. 정규분포)를 따른다고 가정하고 분석한다.



## 6 요인 분석Factor Analysis과 문항 반응 이론Item Response Theory

- 여러 변수에 나타나는 상관관계를 설명할 수 있는 잠재 변수를 가정한다.
- 변수가 연속이면 요인 분석, 변수가 이항 또는 다항이라면 문항 반응 이론으로 분석할 수 있다.

```

1 library(mirt)
2 library(psych)
3 data(Science)
4
5 # EFA
6 fit0 <- fa(Science)
7 print(fit0)
8 summary(fit0)
9
10 # IRT(Graded Response Model)
11 fit1 <- mirt(Science, 1)
12 summary(fit1)
13 fscores(fit1)
14 personfit(fit1)
15 itemfit(fit1)
16 itemplot(fit1, 4)
17 testinfo(fit1, Theta=0)
18
19 ##
20 library(lavaan)
21 data(HolzingerSwineford1939)
22 dat <- HolzingerSwineford1939[,paste0("x",1:9)]
23
24 # determining the number of factors
25 fa.parallel(dat, fm='ml')
26 fit0 <- fa(dat, nfactors=3)
27 summary(fit0)
28 diagram(fit0)

```

Listing 6: 요인 분석과 문항 반응 분석의 예

## 7 경로 분석 Path Analysis

- 선형 모형에서 계수는 모형 (설명 변수의 종류)에 따라 달라진다.
- 예측 모형은 인과 관계를 나타내지 않는다.
- 경로 분석은 여러 변수 사이의 인과 관계를 분석하기 위해 사용된다.
- 예. Simpson's paradox, Berkinson's paradox

## 8 구조 방정식 모형 Structural Equation Model

- 요인 분석 모형과 경로 분석 모형을 합쳐 놓은 것이다.
- 여러 변수에 나타나는 상관관계를 설명할 수 있는 잠재 변수를 가정한다.
- 변수가 연속이면 요인 분석, 변수가 이항 또는 다항이라면 문항 반응 이론으로 분석할 수 있다.

```

1 library(lavaan)
2 model <- '
3 # measurement model
4   ind60 =~ x1 + x2 + x3
5   dem60 =~ y1 + y2 + y3 + y4
6   dem65 =~ y5 + y6 + y7 + y8
7 # regression model
8   dem60 ~ ind60
9   dem65 ~ ind60 + dem60
10 # residual correlations
11   y1 ~~ y5
12   y2 ~~ y4 + y6
13   y3 ~~ y7
14   y4 ~~ y8
15   y6 ~~ y8'
16 fit <- sem(model, data=PoliticalDemocracy)
17 summary(fit, standardized=T)
18 semPlot::semPaths(fit)

```

Listing 7: SEM 분석의 예