# 회귀분석



# 1. 단순선형회귀 모형 적합

• 반응변수 Y와 설명변수 X 사이의 선형관계 가정

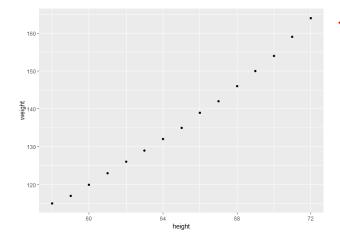
$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i, \qquad i = 1, ..., n$$

- 일차적인 관심: 회귀계수  $\beta_0$ 와  $\beta_1$ 의 추정
- 오차항  $\varepsilon_i$ 에 대한 가정: 서로 독립, 동일 분포 N 0, $\sigma^2$



### 예제: 데이터 프레임 women

- 변수 height와 weight의 관계 탐색
- 첫 번째 작업: 산점도 작성
  - > library(ggplot2)
- > ggplot(women, aes(x=height, y=weight)) +
   geom\_point()



#### 선형관계가 있는 것으로 보임

# R

#### 선형회귀모형 적합: 함수 1m()

```
> fit <- lm(weight ~ height, women)
> fit

Call:
lm(formula = weight ~ height, data = women)
```

```
> names(fit)
 [1] "coefficients" "residuals" "effects" "rank"
 [5] "fitted.values" "assign" "qr" "df.residual"
 [9] "xlevels" "call" "terms" "model"
```

- 사용자마다 필요한 정보가 서로 다를 수 있음
- 필요한 정보를 각자 선택해서 추출
- 모든 결과를 한번에 출력하는 SAS, SPSS와는 다른 접근 방식

# R

#### 선형회귀모형 적합: 함수 1m()

```
> fit <- lm(weight ~ height, women)
> fit

Call:
lm(formula = weight ~ height, data = women)
```

```
> names(fit)
 [1] "coefficients" "residuals" "effects" "rank"
 [5] "fitted.values" "assign" "qr" "df.residual"
 [9] "xlevels" "call" "terms" "model"
```

- 사용자마다 필요한 정보가 서로 다를 수 있음
- 필요한 정보를 각자 선택해서 추출
- 모든 결과를 한번에 출력하는 SAS, SPSS와는 다른 접근 방식



# 2. 다중선형회귀모형 적합

• 반응변수 Y와 설명변수  $X_1, X_2, ..., X_k$ 사이에 선형 관계 가정

$$Y_i = \beta_0 + \beta_1 X_{1i} + \dots + \beta_k X_{ki} + \varepsilon_i, \qquad i = 1, \dots, n$$

• 오차항  $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$  가정:  $N 0, \sigma^2,$  서로 독립

- 선형 및 오차항 가정
  - 회귀모형의 추정 및 추론의 정당성 보장
  - 가정 위반 시 추론 결과에 대한 신뢰성 저하



#### 함수 1m()의 기본적인 사용법

lm(formula, data, subset, ... )

- formula: 회귀모형 설정을 위한 R 공식
- data: 데이터 프레임
- subset: 데이터의 일부분만을 이용하는 경우 처음 100개 자료만 이용: lm(y ~ x, subset=1:100) 변수 z가 0 이상인 자료만 이용: lm(y ~ x, subet= z>=0)



#### 함수 1m()의 기본적인 사용법

lm(formula, data, subset, ... )

- formula: 회귀모형 설정을 위한 R 공식
- data: 데이터 프레임
- subset: 데이터의 일부분만을 이용하는 경우 처음 100개 자료만 이용: lm(y ~ x, subset=1:100) 변수 z가 0 이상인 자료만 이용: lm(y ~ x, subet= z>=0)



### R 공식에 사용되는 기호

- 1) 물결표(~): 반응변수 ~ 설명변수
- 2) 플러스(+): 설명변수 구분. y ~ x1 + x2 + x3
- 3) 콜론(:): 설명변수 사이의 상호작용. y ~ x1 + x2 + x1:x2
- 4) 점(.): 반응변수를 제외한 데이터 프레임에 있는 모든 변수. 데이터 프레임에 y, x1, x2, x3가 있다면 y ~ . → y ~ x1 + x2 + x3
- 5) 마이너스(-): 모형에서 제외되는 변수
- 6) 1 또는 + 0: 절편제거
- 7) I(): 괄호 안의 연산자를 수학 연산자로 인식.  $y \sim I(x1+x2) \rightarrow Y = \beta_0 + \beta_1 X_1 + X_2$

- R
- 예제 1: 행렬 state.x77
  - 미국 50개 주와 관련된 8개 변수로 구성된 행렬
  - 반응변수: Murder
- 행렬을 데이터 프레임으로 전환



#### 모형에 포함될 변수들의 관계 탐색

- 상관계수
- 산점도 행렬
- 상관계수 계산: 함수 cor()

cor(x, y=NULL, use="everything")

- x, y: 벡터, 행렬, 데이터 프레임 x만 있는 경우: x에 있는 모든 변수들 사이의 상관계수 계산 x와 y가 있는 경우: x에 있는 변수와 y에 있는 변수를 하나씩 짝을 지어 상관계수 계산
- use: 결측값 처리 방식.
  - "everything": 결측값이 있으면 NA
  - "pairwise": 상관계수가 계산되는 변수만을 대상으로 NA가 있는 케이스 제거





```
> cor(states)
            Population
                                    Illiteracy
                                                  Life_Exp
                           Income
Population
            1.00000000
                         0.2082276
                                    0.10762237
                                                -0.06805195
            0.20822756
                         1.0000000 -0.43707519
                                                0.34025534
Income
Illiteracy 0.10762237
                        -0.4370752
                                    1.00000000
                                                -0.58847793
           -0.06805195
                         0.3402553 - 0.58847793
Life Exp
                                                 1.00000000
                        -0.2300776
Murder
           0.34364275
                                    0.70297520
                                                -0.78084575
           -0.09848975
                         0.6199323 -0.65718861
                                                0.58221620
HS Grad
Frost
           -0.33215245
                         0.2262822 -0.67194697
                                                0.26206801
            0.02254384
                         0.3633154
                                    0.07726113
                                                -0.10733194
Area
               Murder
                           HS_Grad
                                        Frost
                                                     Area
Population |
            0.3436428 -0.09848975 -0.3321525
                                               0.02254384
Income
           -0.2300776
                        0.61993232
                                    0.2262822
                                               0.36331544
Illiteracy 0.7029752 -0.65718861 -0.6719470
                                               0.07726113
                                    0.2620680 - 0.10733194
Life Exp
           -0.7808458
                        0.58221620
Murder
            1.0000000 -0.48797102 -0.5388834
                                               0.22839021
           -0.4879710
                        1.00000000 0.3667797
                                               0.33354187
HS Grad
Frost
           -0.5388834
                        0.36677970
                                    1.0000000
                                               0.05922910
            0.2283902
                        0.33354187
                                    0.0592291
                                               1.00000000
Area
```

- 상관계수 행렬: 변수의 개수가 많아지면 변수 사이 관계 파악이 어려움
- 상관계수 행렬을 그래프로 표현: 패키지 GGally의 함수 ggcorr()

# R

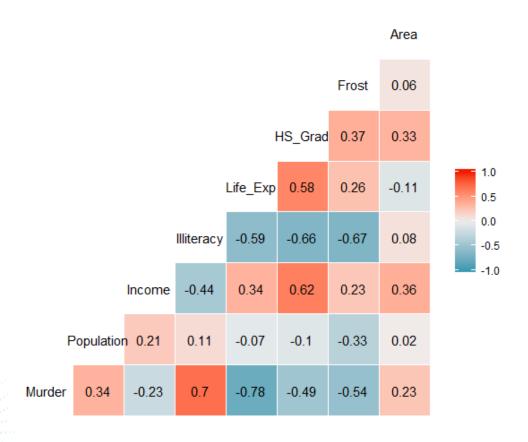
### 패키지 GGally의 함수 ggcorr()

- label: 그래프에 상관계수 표시 여부
- label\_round: 상관계수 반올림 자릿수



#### states 변수들의 상관계수 그래프

- > library(GGally)
- > states <- select(states, Murder, everything())</pre>
- > ggcorr(states, label=TRUE, label\_round=2)





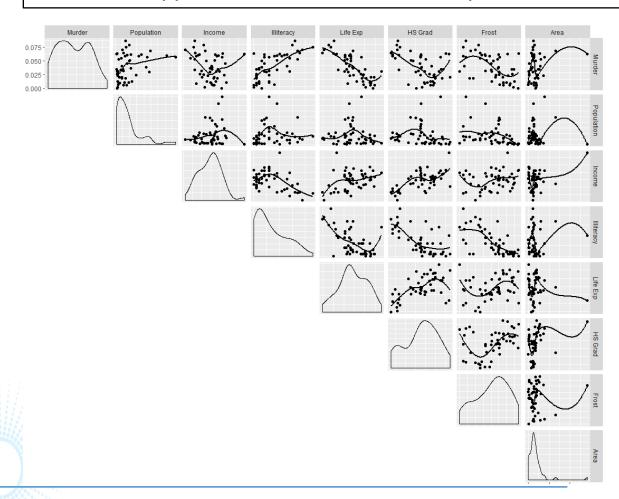
#### 산점도 행렬

- 여러 변수로 이루어진 자료에서 두 변수끼리 짝을 지어 작성된 산점도를 행렬 형태로 배열
- 회귀분석에서 필수적인 그래프



### 패키지 GGally의 함수 ggpairs()

- > library(GGally)



R

#### 예제 1 계속: states에 대한 회귀모형 적합

```
> fit <- lm(Murder ~ ., states)</pre>
> fit
call:
lm(formula = Murder ~ ., data = states)
Coefficients: (Intercep
t)
     Population
                                   Illiteracy
                           Income
 1.222e+02 1.880e-04 -1.592e-04
                                    1.373e+00
  Life_Exp Hs_Grad
                                        Area
                            Frost
 -1.655e+00 3.234e-02 -1.288e-02
                                    5.967e-06
```

- 함수 1m()으로 생성된 객체(회귀분석 결과)의 내용 확인을 위한 함수
  - anova(): 분산분석표
  - coefficients(): 추정된 회귀계수, coef()도 가능
  - confint(): 회귀계수 신뢰구간.
  - fitted(): 반응변수 적합값
  - residuals(): 잔차. resid()도 가능
  - summary(): 중요한 적합 결과 요약



#### 예제 2: women

- 데이터 프레임 women의 변수 weight와 height의 관계
- 선형보다는 2차가 더 적합한 것으로 보임
- 다항회귀모형

$$Y_i = \beta_0 + \beta_1 X_i + \beta_2 X^2 + \dots + \beta_p X^p + \varepsilon_i$$

- 차수 p를 너무 높이면 다중공선성의 문제가 발생할 수 있음
- 3차를 넘지 않는 것이 일반적

R

#### 반응변수 weight에 대한 height의 2차 다항회귀모형 적합

모형식:  $\Psi = 261.87 - 7.34X_i + 0.083X_i^2$ 



#### 예제 3: 질적 변수를 설명변수로 사용

- 회귀모형에서 사용되는 변수 형태

반응변수: 연속형(정규분포 가정 필요)

설명변수: 연속형(정규분포 가정은 필요 없으나, 가능한 좌우대칭)

범주형(가변수 필요)

- 가변수 회귀모형: 2개 범주(yes, no) → 1개 가변수

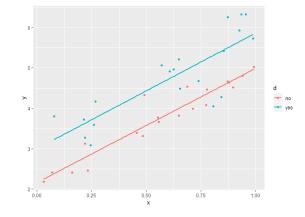
$$Y = \beta_0 + \beta_1 X + \beta_2 D + \varepsilon$$
,  $D = \varphi_1^0$  no yes

- D=0인 범주: 기준 범주
- 회귀계수  $\beta_2$ : yes 범주와 기준 범주의 차이





- → 절편 제거하면 추정 가능
- → 회귀계수의 해석이 달라짐(해당 범주의 효과)
- → 두 개 이상의 범주형 변수가 포함되는 경우에는 적용이 어려움





## 3. 회귀모형의 추론

- 회귀모형:  $Y = \beta_0 + \beta_1 X_1 + \cdots + \beta_k X_k + \varepsilon$
- 회귀계수에 대한 가설
  - 1)  $H_0: \beta_1 = \dots = \beta_k = 0$
  - 2)  $H_0: \beta_q = \beta_{q+1} = \dots = \beta_r = 0, \quad q < r \le k$
  - 3)  $H_0: \beta_i = 0, H_1: \beta_i \neq 0$
- 회귀계수의 신뢰구간
- 회귀모형 적합 정도에 대한 통계량
  - 결정계수, 수정된 결정계수
  - AIC, BIC



#### 적합된 회귀모형 추론을 위한 함수

함수 summary()

```
> fit <- lm(Murder ~ ., states)</pre>
```

> summary(fit)

```
Residuals:
```

```
Min 1Q Median 3Q Max -3.4452 -1.1016 -0.0598 1.1758 3.2355
```

#### Coefficients:

```
Estimate Std. Error t value Pr(>|t|)
                     1.789e+01
                                  6.831 2.54e-08 ***
(Intercept)
            1.222e+02
Population
                      6.474e-05 2.905 0.00584 **
           1.880e-04
Income
           -1.592e-04
                      5.725e-04 -0.278 0.78232
Illiteracv
                     8.322e-01 1.650 0.10641
           1.373e+00
Life_Exp
           -1.655e+00
                      2.562e-01
                                 -6.459 8.68e-08 ***
           3.234e-02
                      5.725e-02
                                 0.565 0.57519
Hs_Grad
           -1.288e-02 7.392e-03
                                 -1.743 0.08867 .
Frost
            5.967e-06
                      3.801e-06
                                 1.570
                                         0.12391
Area
```

Residual standard error: 1.746 on 42 degrees of freedom Multiple R-squared: 0.8083, Adjusted R-squared: 0.7763

F-statistic: 25.29 on 7 and 42 DF, p-value: 3.872e-13

- 개별 회귀계수 추 정 및 검정
- $\sqrt{MSE}$
- 결정계수 및 수정 된 결정계수
- 모든 회귀계수의 유의성 검정

R

### 두 회귀모형의 비교

- 1) 확장모형( $\Omega$ ):  $Y = \beta_0 + \beta_1 X_1 + \cdots + \beta_k X_k + \varepsilon$
- 2) 축소모형( $\omega$ ): 다음의 귀무가설이 사실인 모형

$$H_0$$
:  $\beta_q = \beta_{q+1} = \cdots = \beta_r = 0$ ,  $q < r \le k$ 

 $RSS_{\Omega}$ : 확장모형의 잔차제곱합  $RSS_{\omega}$ : 축소모형의 잔차제곱합

- 만일  $RSS_{\omega} RSS_{\Omega}$ 가 적다면, 축소모형이 확장모형만큼 좋다는 의미
- 모수절약의 원칙에 따라 축소모형 선택 가능
- 검정통계량

$$F = rac{(RSS_{\omega} - RSS_{\Omega})/두$$
 모형의 모수 차이  $RSS_{\Omega}/n - k - 1$ 

```
R
```

#### 함수 anova()

- 두 회귀모형의 비교

anova(축소모형, 확장모형)

귀무가설의 기각이 어려움



### 회귀모형 적합 정도에 대한 통계량

- 결정계수 $(R^2)$ : 반응변수의 변량 중 회귀모형으로 설명되는 변량의 비율
  - 모형에 포함된 설명변수의 개수가 증가하면 증가하는 특성이 있음
  - 설명변수의 개수가 같은 모형 비교에는 의미가 있는 통계량
- 수정 결정계수(adj.  $R^2$ ): 추가된 설명변수가 모형 적합도에 도움이 되는 경우에만 증가
- AIC & BIC: 설명변수의 개수가 p인 모형

$$- AIC = n\log\left(\frac{SSE}{n}\right) + 2p$$

- 
$$BIC = n\log\left(\frac{SSE}{n}\right) + p\log(n)$$

- AIC, BIC가 작은 모형이 더 적합도가 높은 모형



### 4. 변수 선택

- 반응변수의 변동을 설명할 수 있는 많은 설명변수 중 '최적'의 변수를 선택하여 모형에 포함시키는 절차
- 검정에 의한 방법
  - 변수의 유의성 검정을 이용하여 단계적으로 모형 선택
  - 후진소거법, 전진선택법, 단계별 선택법
- 모형선택 기준에 의한 방법
  - 모형의 적합도 등을 측정하는 통계량을 기반으로 모형 선택
  - 결정계수, 수정결정계수, 잔차제곱합,  $C_p$  통계량, AIC, BIC 등등
- 어떤 모형이 '최적' 모형인가?



### 모형선택 기준에 의한 방법

- 모형 수립 목적을 고려한 변수 선택 방법
- 모형의 적합도 등을 나타내는 통계량을 선택 기준으로 사용
- 사용되는 통계량
  - 수정 결정계수(adj. R<sup>2</sup>)
  - AIC, BIC
- 선택 방법
  - 모든 가능한 회귀(All possible regression)
  - 단계별 선택법



### 모형선택 기준에 의한 방법

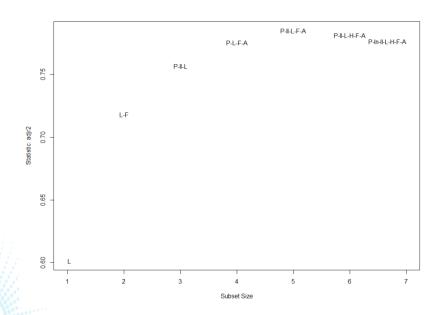
- 모형 수립 목적을 고려한 변수 선택 방법
- 모형의 적합도 등을 나타내는 통계량을 선택 기준으로 사용
- 사용되는 통계량
  - 수정 결정계수(adj. R<sup>2</sup>)
  - AIC, BIC
- 선택 방법
  - 모든 가능한 회귀(All possible regression)
  - 단계별 선택법



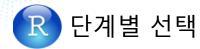
### 함수 car::subsets()에 의한 확인

- 옵션 statistics: 디폴트 bic
- 옵션 legend: 범례의 위치
   legend="interactive": 디폴트, 마우스로 위치 지정 가능
   legend=FALSE: Console 창에 범례 출력





#### - P-I1-L-F-A 모형 선택



- AIC 혹은 BIC에 의한 단계별 선택
- 변수의 개수가 많은 경우에 적용 가능
- MASS::stepAIC()로 실시

stepAIC( object, scope, k=2)

- object: 함수 glm()으로 생성된 객체
- scope: 모든 설명변수가 포함된 full 모형의 formula. 생략되면 object에 설정된 모형이 full 모형.
- k: 탐색에 사용되는 IC. k=2는 AIC, k=log(n)은 BIC.

### 전진선택법에 의한 단계별 선택

- > library(MASS)
  > fit\_full <- lm(Murder ~ ., states)
  > fit <- lm(Murder ~ 1, states)
  > stepAIC(fit, scope=formula(fit\_full))

Start: AIC=131.59 Step: AIC=65.15 Murder ~ 1 Murder ~ Life\_Exp + Frost + Population Df Sum of Sq RSS AIC Df Sum of Sq RSS AIC + Life\_Exp 1 407.14 260.61 86.550 1 19.040 137.75 60.672 + Area + Illiteracy 1 329.98 337.76 99.516 + Illiteracy 1 11.826 144.97 63.225 1 + Frost 193.91 473.84 116.442 156.79 65.146 <none> + Hs\_Grad 159.00 508.75 119.996 1.821 154.97 66.561 + Hs\_Grad + Population 1 78.85 588.89 127.311 + Income 1 0.739 156.06 66,909 1 35.35 632.40 130.875 + Income - Population 1 23.710 180.50 70.187 1 34.83 632.91 130.916 + Area - Frost 1 47.198 203.99 76.303 667.75 131.594 <none> Life\_Exp 1 296,694 453,49 116,247 Step: AIC=86.55 Step: AIC=60.67 Murder ~ Life\_Exp Murder ~ Life\_Exp + Frost + Population + Area Df Sum of Sq RSS AIC Df Sum of Sq RSS AIC 1 80.10 180.50 70.187 + Frost + Illiteracy 1 8.723 129.03 59,402 + Illiteracy 1 60.55 200.06 75.329 <none> 137.75 60.672 + Population 1 76.303 56.62 203.99 1.241 136.51 + Income 62.220 1 14.12 246.49 85.764 + Area 1 0.771 136.98 62.392 + Hs Grad 86.550 260.61 <none> 19.040 156.79 65.146 - Area + Hs\_Grad 1.12 259.48 88.334 - Population 1 21.666 159.42 65.976 1 0.96 259.65 88.366 + Income - Frost 52.970 190.72 74.940 Life\_Exp 407.14 667.75 131.594 Life\_Exp 1 272,927 410,68 113,290 Step: AIC=70.19 Step: AIC=59.4 Murder ~ Life\_Exp + Frost Murder ~ Life\_Exp + Frost + Population + Area + Illiteracy Df Sum of Sq RSS AIC Df Sum of Sq RSS AIC + Population 1 23.710 156.79 65.146 <none> 129.03 59.402 1 65.976 + Area 21.084 159.42 - Illiteracy 1 8.723 137.75 60.672 70.187 180.50 <none> + Hs\_Grad 0.763 128.27 61.105 6.066 174.44 70.477 + Illiteracy 1 + Income 0.026 129.01 61.392 1 5.560 174.94 70.622 + Income Frost 11.030 140.06 61.503 + Hs\_Grad 1 2.068 178.44 71.610 1 - Area 15.937 144.97 63.225 Frost 80.104 260.61 86, 550 - Population 1 26.415 155.45 66.714 Life\_Exp 293.331 473.84 116.442 Life\_Exp 140.391 269.42 94.213

#### - 후진소거법에 의한 단계별 선택

6.804e-06

```
> stepAIC(fit_full, trace=FALSE)

Call:
lm(formula = Murder ~ Population + Illiteracy + Life_Exp + Frost +
    Area, data = states)

Coefficients:
(Intercept) Population Illiteracy Life_Exp Frost
    1.202e+02    1.780e-04    1.173e+00    -1.608e+00    -1.373e-02
    Area
```

#### - BIC에 의한 단계별 선택

AIC에 의한 단계별 선택과는 다른 결과





#### 주요 이력

現) ㈜비즈스프링 웹로그분석 및 DP사업 진행 中

煎) 하이호금속 회계팀

煎) ㈜벽산 회계팀

煎) K문고 CRM VIP 군집전략 CRM프로젝트 보조연구원

煎) L백화점 CRM Alert 전략 CRM프로젝트 보조연구원

BSL(스위스 로잔 비즈니스 스쿨) MBA 진행 中 ASSIST 빅데이터 MBA 진행 中

국가공인 ADSP(빅데이터 준전문가) 현 코리아IT아카데미 빅데이터 R 강사 현 코리아IT아카데미 빅데이터 기초 파이썬 강사 현 코리아IT아카데미 빅데이터 기초통계 전담강사

[박영식] 완성에 이르기까지