

# 第六次作业

1907402030 熊雄

## 1 QUESTION.

假设样本  $(x_1, \dots, x_p) \sim N(0, I_p)$ , 建立回归模型  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \varepsilon$ , 其中回归系数  $\beta_0, \beta_1, \beta_2, \beta_3 \sim Unif(1, 2)$  为随机数, 误差项  $\varepsilon \sim N(0, 1)$ , 并且  $p \in \{10, 20\}$

请重复下述过程  $n$  次, 记录正确(即选择  $x_1, x_2, x_3$ ), 多选, 少选(即不全是  $x_1, x_2, x_3$ ), 错选的次数(即无  $x_1, x_2, x_3$ ), 其中  $n \in \{200, 500, 1000\}$ .

1. 生成  $\beta_0, \beta_1, \beta_2, \beta_3$ ;
2. 生成  $\varepsilon$  和  $X$  数据;
3. 生成  $Y$  数据;
4. 利用  $AIC$  准则,  $BIC$  准则,  $R_{adj}$ ,  $C_p$  统计量寻找最优子集.

## 2 ANSWER.

### 2.1 CODE.

输入以下代码:

```

1  n <- 200
2  p <- 10
3  correct = 0 # 记录正确的个数
4  add = 0 # 记录多选的个数
5  less = 0 # 少选的个数
6  error = 0 # 记录错误的个数
7  for (i in 1:200) # 与 n = 200 一起进行动态调整
8  {
9    xmean <- matrix(c(runif(p, 0, 0)), ncol = 1)
10   xsigma <- diag(p)
11   library(MASS)
12   x <- mvrnorm(n, xmean, xsigma)
13   beta <- matrix(runif(3, 1, 2))
14   beta0 <- matrix(rep(runif(1, 1, 2), n))

```

```

15 e <- matrix(rnorm(n, 0, 1))
16 y <- x[,c(1:3)]%*%beta + beta0 + e
17 dx <- data.frame(x)
18 library(leaps)
19 sub.fit <- regsubsets(y~.,data = dx)
20 best.summary <- summary(sub.fit)
21 k <- which.max(best.summary$adjr2) # 调整R方
22 k <- which.min(best.summary$bic) # BIC
23 k <- which.min(best.summary$cp) # c_p
24 best.summary$aic = best.summary$bic - p * log(n)+2 * p # AIC
25 k <- which.min(best.summary$aic)
26 jg <- coef(sub.fit,k)
27 xz <- matrix(rownames(data.frame(jg)),ncol = 1)
28 if('x1'%in%xz & 'x2'%in%xz & 'x3'%in%xz )
29 {
30     if(length(xz) == 4)
31         correct = correct + 1
32     else
33         add = add + 1
34 }
35 else if('x4'%in%xz | 'x5'%in%xz | 'x6'%in%xz | 'x7'%in%xz |
'x8'%in%xz | 'x9'%in%xz | 'x10'%in%xz)
36     error = error + 1
37 else
38     less = less + 1
39 }

```

## 2.2 RESULT and ANALYSIS.

由上述代码知,  $n$  为200时的运行结果为:

```

1 add = 24 # 记录多选个数为70
2 correct = 176 # 记录正确的个数为430
3 less = 0 # 少选的个数为0
4 error = 0 # 记录错误的个数为0

```

将  $n$  调整为500, 运行结果为:

```
1 add = 34 # 记录多选个数为34
2 correct = 466 # 记录正确的个数为466
3 less = 0 # 少选个数为0
4 error = 0 # 记录错误的个数为0
```

将  $n$  调整为1000, 运行结果为:

```
1 add = 34 # 记录多选个数为42
2 correct = 466 # 记录正确的个数为958
3 less = 0 # 少选个数为0
4 error = 0 # 记录错误的个数为0
```

故大致可以看出, 当  $n$  进行改变时, 利用各个准则求得的最优子集记录正确的个数比记录多选的个数多得多, 并且记录少选或错选的情况几乎不可能发生.