

biostatistical methods homework 5

```
library(knitr)
library(tidyverse)
library(faraway)
library(broom)
library(leaps)
library(boot)
library(modelr)
library(caret)
```

R dataset 'state.x77' from library(faraway) contains information on 50 states from 1970s collected by US Census Bureau. The goal is to predict 'life expectancy' using a combination of remaining variables.

```
life_data = as.data.frame(state.x77) %>%
  janitor::clean_names()
```

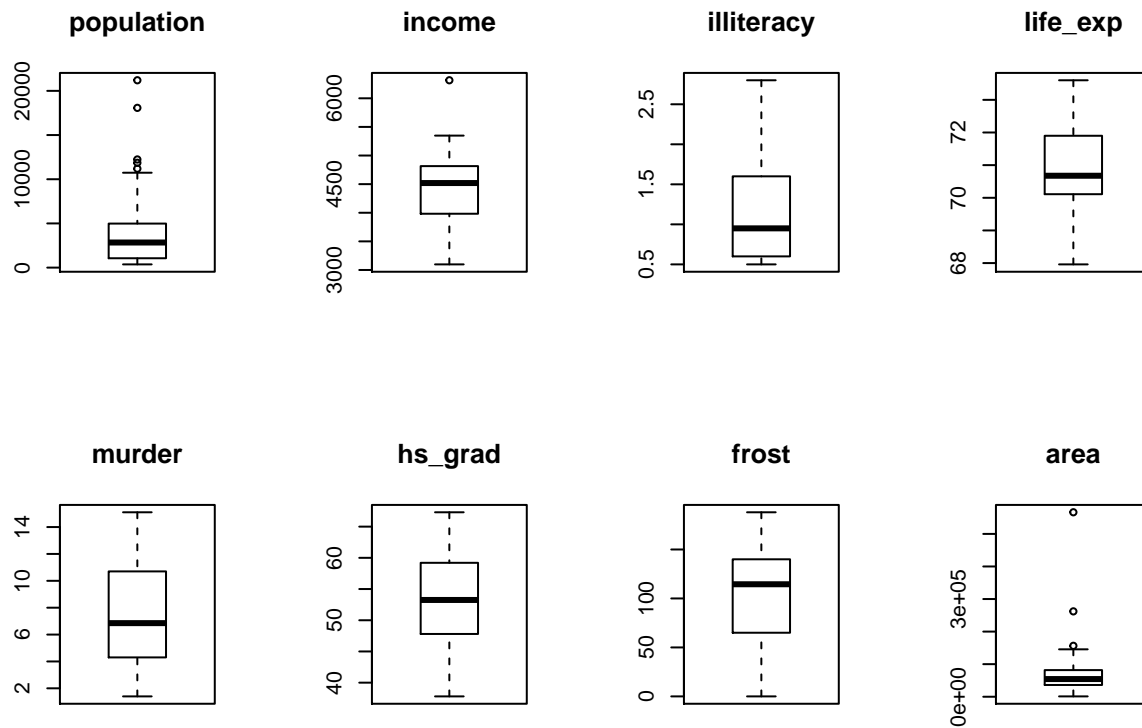
1. Explore the dataset and generate appropriate descriptive statistics and relevant graphs

```
mean_and_sd = function(x) {
  if (!is.numeric(x)) {
    stop("Argument x should be numeric")
  } else if (length(x) == 1) {
    stop("Cannot be computed for length 1 vectors")
  }

  mean_x = mean(x)
  sd_x = sd(x)
  tibble(
    mean = mean_x,
    sd = sd_x
  )
}
```

```
attach(life_data)
```

```
par(mfrow = c(2, 4))
boxplot(population, main = 'population')
boxplot(income, main = 'income')
boxplot(illiteracy, main = 'illiteracy')
boxplot(life_exp, main = 'life_exp')
boxplot(murder, main = 'murder')
boxplot(hs_grad, main = 'hs_grad')
boxplot(frost, main = 'frost')
boxplot(area, main = 'area')
```



Population

```
summary(population)
```

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	365	1080	2838	4246	4968	21198

Income

```
summary(income)
```

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	3098	3993	4519	4436	4814	6315

Illiteracy

```
summary(illiteracy)
```

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	0.500	0.625	0.950	1.170	1.575	2.800

Life Exp

```
summary(life_exp)
```

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	67.96	70.12	70.67	70.88	71.89	73.60

Murder

```
summary(murder)
```

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	1.400	4.350	6.850	7.378	10.675	15.100

HS Grad

```
summary(hs_grad)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##    37.80   48.05   53.25   53.11   59.15   67.30
```

Frost

```
summary(frost)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      0.00   66.25  114.50  104.46  139.75  188.00
```

Area

```
summary(area)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##     1049   36985   54277   70736   81162  566432
```

2. Use automatic procedures to find a ‘best subset’ of the full model. Present the results and comment on the following:

Backward elimination

```
backward_fit <- lm(life_exp ~ ., data=life_data)
step(backward_fit, direction='backward') %>%
  summary()
```

```
## Start:  AIC=-22.18
## life_exp ~ population + income + illiteracy + murder + hs_grad +
##      frost + area
##
##              Df Sum of Sq  RSS    AIC
## - area         1     0.0011 23.298 -24.182
## - income        1     0.0044 23.302 -24.175
## - illiteracy    1     0.0047 23.302 -24.174
## <none>                          23.297 -22.185
## - population   1     1.7472 25.044 -20.569
## - frost        1     1.8466 25.144 -20.371
## - hs_grad      1     2.4413 25.738 -19.202
## - murder       1    23.1411 46.438  10.305
##
## Step:  AIC=-24.18
## life_exp ~ population + income + illiteracy + murder + hs_grad +
##      frost
##
##              Df Sum of Sq  RSS    AIC
## - illiteracy    1     0.0038 23.302 -26.174
## - income        1     0.0059 23.304 -26.170
## <none>                          23.298 -24.182
## - population   1     1.7599 25.058 -22.541
## - frost        1     2.0488 25.347 -21.968
## - hs_grad      1     2.9804 26.279 -20.163
## - murder       1    26.2721 49.570  11.569
##
## Step:  AIC=-26.17
```

```
## life_exp ~ population + income + murder + hs_grad + frost
##
##           Df Sum of Sq   RSS   AIC
## - income      1      0.006 23.308 -28.161
## <none>                23.302 -26.174
## - population  1      1.887 25.189 -24.280
## - frost       1      3.037 26.339 -22.048
## - hs_grad     1      3.495 26.797 -21.187
## - murder      1     34.739 58.041  17.456
##
## Step: AIC=-28.16
## life_exp ~ population + murder + hs_grad + frost
##
##           Df Sum of Sq   RSS   AIC
## <none>                23.308 -28.161
## - population  1      2.064 25.372 -25.920
## - frost       1      3.122 26.430 -23.877
## - hs_grad     1      5.112 28.420 -20.246
## - murder      1     34.816 58.124  15.528
##
## Call:
## lm(formula = life_exp ~ population + murder + hs_grad + frost,
##     data = life_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.47095 -0.53464 -0.03701  0.57621  1.50683
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  7.103e+01  9.529e-01  74.542 < 2e-16 ***
## population    5.014e-05  2.512e-05   1.996  0.05201 .
## murder       -3.001e-01  3.661e-02  -8.199 1.77e-10 ***
## hs_grad       4.658e-02  1.483e-02   3.142  0.00297 **
## frost        -5.943e-03  2.421e-03  -2.455  0.01802 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.7197 on 45 degrees of freedom
## Multiple R-squared:  0.736, Adjusted R-squared:  0.7126
## F-statistic: 31.37 on 4 and 45 DF, p-value: 1.696e-12
```

Forward elimination

```
fit1 <- lm(life_exp ~ population, data=life_data)
tidy(fit1)
```

```
## # A tibble: 2 x 5
##   term          estimate std.error statistic  p.value
##   <chr>          <dbl>     <dbl>     <dbl>    <dbl>
## 1 (Intercept)    71.0         0.265      267.    7.90e-78
## 2 population    -0.0000205 0.0000433   -0.473  6.39e- 1
```

```
fit2 <- lm(life_exp ~ income, data=life_data)
tidy(fit2)
```

```
## # A tibble: 2 x 5
##   term          estimate std.error statistic  p.value
##   <chr>          <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept)  67.6        1.33      50.9  1.98e-43
## 2 income       0.000743  0.000297    2.51  1.56e- 2
```

```
fit3 <- lm(life_exp ~ illiteracy, data=life_data)
tidy(fit3)
```

```
## # A tibble: 2 x 5
##   term          estimate std.error statistic  p.value
##   <chr>          <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept)   72.4      0.338    214.   3.47e-73
## 2 illiteracy    -1.30     0.257    -5.04  6.97e- 6
```

```
fit4 <- lm(life_exp ~ murder, data=life_data)
tidy(fit4)
```

```
## # A tibble: 2 x 5
##   term          estimate std.error statistic  p.value
##   <chr>          <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept)   73.0      0.270    270.   4.72e-78
## 2 murder        -0.284    0.0328   -8.66  2.26e-11
```

```
fit5 <- lm(life_exp ~ hs_grad, data=life_data)
tidy(fit5)
```

```
## # A tibble: 2 x 5
##   term          estimate std.error statistic  p.value
##   <chr>          <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept)   65.7      1.05     62.8  9.92e-48
## 2 hs_grad        0.0968    0.0195    4.96  9.20e- 6
```

```
fit6 <- lm(life_exp ~ frost, data=life_data)
tidy(fit6)
```

```
## # A tibble: 2 x 5
##   term          estimate std.error statistic  p.value
##   <chr>          <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept)   70.2      0.419    168.   4.33e-68
## 2 frost          0.00677  0.00360    1.88  6.60e- 2
```

```
fit7 <- lm(life_exp ~ area, data=life_data)
tidy(fit7)
```

```
## # A tibble: 2 x 5
##   term          estimate std.error statistic  p.value
##   <chr>          <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept)   71.0      0.249    285.   3.46e-79
## 2 area          -0.00000169 0.00000226  -0.748 4.58e- 1
```

```
forward1 <- lm(life_exp ~ murder, data = life_data)
tidy(forward1)
```

```
## # A tibble: 2 x 5
##   term          estimate std.error statistic  p.value
##   <chr>          <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept)   73.0      0.270    270.   4.72e-78
```

```
## 2 murder      -0.284    0.0328    -8.66 2.26e-11
```

```
fit1 <- update(forward1, . ~ . +population)
tidy(fit1)
```

```
## # A tibble: 3 x 5
##   term          estimate std.error statistic  p.value
##   <chr>          <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept)  72.9      0.258      282.  1.55e-77
## 2 murder      -0.312     0.0332     -9.42 2.15e-12
## 3 population   0.0000683 0.0000274    2.49 1.64e- 2
```

```
fit2 <- update(forward1, . ~ . +income)
tidy(fit2)
```

```
## # A tibble: 3 x 5
##   term          estimate std.error statistic  p.value
##   <chr>          <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept)  71.2      0.967      73.6  3.32e-50
## 2 murder      -0.270     0.0328     -8.21 1.22e-10
## 3 income        0.000370 0.000197    1.88 6.66e- 2
```

```
fit3 <- update(forward1, . ~ . +illiteracy)
tidy(fit3)
```

```
## # A tibble: 3 x 5
##   term          estimate std.error statistic  p.value
##   <chr>          <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept)  73.0      0.286     256.  1.56e-75
## 2 murder      -0.264     0.0464     -5.69 7.96e- 7
## 3 illiteracy   -0.172     0.281     -0.613 5.43e- 1
```

```
fit4 <- update(forward1, . ~ . +hs_grad)
tidy(fit4)
```

```
## # A tibble: 3 x 5
##   term          estimate std.error statistic  p.value
##   <chr>          <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept)  70.3      1.02      69.2  5.91e-49
## 2 murder      -0.237     0.0353     -6.72 2.18e- 8
## 3 hs_grad       0.0439     0.0161     2.72 9.09e- 3
```

```
fit5 <- update(forward1, . ~ . +frost)
tidy(fit5)
```

```
## # A tibble: 3 x 5
##   term          estimate std.error statistic  p.value
##   <chr>          <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept)  73.9      0.500     148.  2.36e-64
## 2 murder      -0.328     0.0375     -8.74 2.05e-11
## 3 frost        -0.00578 0.00266     -2.17 3.52e- 2
```

```
fit6 <- update(forward1, . ~ . +area)
tidy(fit6)
```

```
## # A tibble: 3 x 5
##   term          estimate std.error statistic  p.value
##   <chr>          <dbl>    <dbl>    <dbl>    <dbl>
```

```
## 1 (Intercept) 72.9      0.275      265.    2.73e-76
## 2 murder      -0.290     0.0338     -8.58  3.47e-11
## 3 area         0.00000118 0.00000146    0.806 4.24e- 1
```

```
forward2 <- update(forward1, . ~ . + hs_grad)
tidy(forward2)
```

```
## # A tibble: 3 x 5
##   term      estimate std.error statistic  p.value
##   <chr>      <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept) 70.3      1.02      69.2  5.91e-49
## 2 murder     -0.237     0.0353    -6.72  2.18e- 8
## 3 hs_grad      0.0439    0.0161     2.72  9.09e- 3
```

```
fit1 <- update(forward2, . ~ . + population)
tidy(fit1)
```

```
## # A tibble: 4 x 5
##   term      estimate std.error statistic  p.value
##   <chr>      <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept) 70.4      0.969     72.7  3.95e-49
## 2 murder     -0.266     0.0357    -7.45  1.91e- 9
## 3 hs_grad      0.0407    0.0154     2.64  1.12e- 2
## 4 population   0.0000625 0.0000259    2.41  1.99e- 2
```

```
fit2 <- update(forward2, . ~ . + income)
tidy(fit2)
```

```
## # A tibble: 4 x 5
##   term      estimate std.error statistic  p.value
##   <chr>      <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept) 70.1      1.10     64.0  1.33e-46
## 2 murder     -0.239     0.0358    -6.66  2.92e- 8
## 3 hs_grad      0.0391    0.0203     1.92  6.05e- 2
## 4 income       0.0000953 0.000239    0.398 6.92e- 1
```

```
fit3 <- update(forward2, . ~ . + illiteracy)
tidy(fit3)
```

```
## # A tibble: 4 x 5
##   term      estimate std.error statistic  p.value
##   <chr>      <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept) 69.7      1.22     57.1  2.41e-44
## 2 murder     -0.258     0.0435    -5.93  3.63e- 7
## 3 hs_grad      0.0518    0.0188     2.76  8.25e- 3
## 4 illiteracy   0.254     0.305     0.833 4.09e- 1
```

```
fit4 <- update(forward2, . ~ . + frost)
tidy(fit4)
```

```
## # A tibble: 4 x 5
##   term      estimate std.error statistic  p.value
##   <chr>      <dbl>    <dbl>    <dbl>    <dbl>
## 1 (Intercept) 71.0      0.983     72.2  5.25e-49
## 2 murder     -0.283     0.0367    -7.71  8.04e-10
## 3 hs_grad      0.0499    0.0152     3.29  1.95e- 3
## 4 frost       -0.00691    0.00245    -2.82  6.99e- 3
```

```
fit5 <- update(forward2, . ~ . +area)
tidy(fit5)
```

```
## # A tibble: 4 x 5
##   term          estimate std.error statistic  p.value
##   <chr>          <dbl>     <dbl>     <dbl>    <dbl>
## 1 (Intercept)  69.9         1.16        60.1  2.30e-45
## 2 murder      -0.224        0.0404       -5.56  1.30e- 6
## 3 hs_grad      0.0504        0.0190        2.65  1.10e- 2
## 4 area        -0.00000106 0.00000162    -0.658 5.14e- 1
```

```
forward3 <- update(forward2, . ~ . + frost)
tidy(forward3)
```

```
## # A tibble: 4 x 5
##   term          estimate std.error statistic  p.value
##   <chr>          <dbl>     <dbl>     <dbl>    <dbl>
## 1 (Intercept)  71.0         0.983       72.2  5.25e-49
## 2 murder      -0.283        0.0367      -7.71  8.04e-10
## 3 hs_grad      0.0499        0.0152        3.29  1.95e- 3
## 4 frost        -0.00691     0.00245     -2.82  6.99e- 3
```

```
fit1 <- update(forward3, . ~ . +population)
tidy(fit1)
```

```
## # A tibble: 5 x 5
##   term          estimate std.error statistic  p.value
##   <chr>          <dbl>     <dbl>     <dbl>    <dbl>
## 1 (Intercept)  71.0         0.953       74.5  8.61e-49
## 2 murder      -0.300        0.0366      -8.20  1.77e-10
## 3 hs_grad      0.0466        0.0148        3.14  2.97e- 3
## 4 frost        -0.00594     0.00242     -2.46  1.80e- 2
## 5 population    0.0000501 0.0000251     2.00  5.20e- 2
```

```
fit2 <- update(forward3, . ~ . +income)
tidy(fit2)
```

```
## # A tibble: 5 x 5
##   term          estimate std.error statistic  p.value
##   <chr>          <dbl>     <dbl>     <dbl>    <dbl>
## 1 (Intercept)  70.8         1.05        67.4  7.53e-47
## 2 murder      -0.286        0.0373      -7.66  1.07e- 9
## 3 hs_grad      0.0436        0.0190        2.30  2.64e- 2
## 4 frost        -0.00698     0.00247     -2.83  6.96e- 3
## 5 income        0.000127 0.000223     0.571 5.71e- 1
```

```
fit3 <- update(forward3, . ~ . +illiteracy)
tidy(fit3)
```

```
## # A tibble: 5 x 5
##   term          estimate std.error statistic  p.value
##   <chr>          <dbl>     <dbl>     <dbl>    <dbl>
## 1 (Intercept)  71.5         1.32        54.2  1.28e-42
## 2 murder      -0.273        0.0411      -6.64  3.50e- 8
## 3 hs_grad      0.0450        0.0178        2.53  1.49e- 2
## 4 frost        -0.00768     0.00283     -2.72  9.36e- 3
## 5 illiteracy   -0.182        0.328       -0.554 5.82e- 1
```



```
fit4 <- update(forward3, . ~ . +area)
tidy(fit4)
```

```
## # A tibble: 5 x 5
##   term          estimate std.error statistic  p.value
##   <chr>          <dbl>     <dbl>     <dbl>   <dbl>
## 1 (Intercept)    70.9         1.15        61.7 3.92e-45
## 2 murder        -0.279        0.0427       -6.52 5.34e- 8
## 3 hs_grad         0.0519        0.0179        2.91 5.66e- 3
## 4 frost         -0.00682       0.00251       -2.71 9.40e- 3
## 5 area          -0.000000329 0.00000154    -0.214 8.32e- 1
```

```
forward_fit <- lm(life_exp ~ murder + hs_grad + frost, data = life_data) %>%
  summary()
```

Stepwise regression

```
stepwise_fit <- lm(life_exp ~ ., data=life_data)
step(stepwise_fit, direction='both') %>%
  summary()
```

```
## Start:  AIC=-22.18
## life_exp ~ population + income + illiteracy + murder + hs_grad +
##   frost + area
##
##           Df Sum of Sq  RSS    AIC
## - area      1    0.0011 23.298 -24.182
## - income     1    0.0044 23.302 -24.175
## - illiteracy 1    0.0047 23.302 -24.174
## <none>                23.297 -22.185
## - population 1    1.7472 25.044 -20.569
## - frost      1    1.8466 25.144 -20.371
## - hs_grad    1    2.4413 25.738 -19.202
## - murder     1   23.1411 46.438  10.305
##
## Step:  AIC=-24.18
## life_exp ~ population + income + illiteracy + murder + hs_grad +
##   frost
##
##           Df Sum of Sq  RSS    AIC
## - illiteracy 1    0.0038 23.302 -26.174
## - income     1    0.0059 23.304 -26.170
## <none>                23.298 -24.182
## - population 1    1.7599 25.058 -22.541
## + area       1    0.0011 23.297 -22.185
## - frost      1    2.0488 25.347 -21.968
## - hs_grad    1    2.9804 26.279 -20.163
## - murder     1   26.2721 49.570  11.569
##
## Step:  AIC=-26.17
## life_exp ~ population + income + murder + hs_grad + frost
##
##           Df Sum of Sq  RSS    AIC
## - income     1    0.006 23.308 -28.161
## <none>                23.302 -26.174
```

```

## - population 1      1.887 25.189 -24.280
## + illiteracy 1      0.004 23.298 -24.182
## + area       1      0.000 23.302 -24.174
## - frost      1      3.037 26.339 -22.048
## - hs_grad    1      3.495 26.797 -21.187
## - murder     1     34.739 58.041  17.456
##
## Step: AIC=-28.16
## life_exp ~ population + murder + hs_grad + frost
##
##           Df Sum of Sq    RSS    AIC
## <none>                23.308 -28.161
## + income      1      0.006 23.302 -26.174
## + illiteracy  1      0.004 23.304 -26.170
## + area        1      0.001 23.307 -26.163
## - population  1      2.064 25.372 -25.920
## - frost       1      3.122 26.430 -23.877
## - hs_grad     1      5.112 28.420 -20.246
## - murder      1     34.816 58.124  15.528
##
## Call:
## lm(formula = life_exp ~ population + murder + hs_grad + frost,
##     data = life_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.47095 -0.53464 -0.03701  0.57621  1.50683
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  7.103e+01  9.529e-01  74.542  < 2e-16 ***
## population    5.014e-05  2.512e-05   1.996  0.05201 .
## murder       -3.001e-01  3.661e-02  -8.199  1.77e-10 ***
## hs_grad       4.658e-02  1.483e-02   3.142  0.00297 **
## frost        -5.943e-03  2.421e-03  -2.455  0.01802 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.7197 on 45 degrees of freedom
## Multiple R-squared:  0.736, Adjusted R-squared:  0.7126
## F-statistic: 31.37 on 4 and 45 DF, p-value: 1.696e-12

```

a) Do the procedures generate the same model?

Using backward elimination, the model we obtained is: $\text{life_exp} \sim \text{population} + \text{murder} + \text{hs_grad} + \text{frost}$. Using forward elimination, the model we obtained is: $\text{life_exp} \sim \text{murder} + \text{hs_grad} + \text{frost}$. Using stepwise regression, the model we obtained is: $\text{life_exp} \sim \text{population} + \text{murder} + \text{hs_grad} + \text{frost}$. However, we generated the model using different function/code, so the inclusion/exclusion criterion might be different. For instance, using `stepwise()` function in backward elimination and stepwise regression, the variable `population` would be automatically included in the model, but the p-value is not that significant. If we discard `population` in backward and stepwise models, or, we include `population` in forward models, then these 3 procedures generate the same model. We will discuss this problem in the next question.

b) Is there any variable a close call? What was your decision: keep or discard? Provide arguments for your choice. (Note: this question might have more or less relevance depending on the 'subset' you choose).

Using `stepwise()` function, as we talked about before, in backward elimination or stepwise regression, population is a close call variable with p-value of 0.05201.

```
bw_s = lm(life_exp ~ murder + hs_grad + frost, data = life_data)
bw_l = lm(life_exp ~ murder + hs_grad + frost + population, data = life_data)
summary(bw_s)
```

```
##
## Call:
## lm(formula = life_exp ~ murder + hs_grad + frost, data = life_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.5015 -0.5391  0.1014  0.5921  1.2268
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  71.036379   0.983262   72.246 < 2e-16 ***
## murder       -0.283065   0.036731  -7.706 8.04e-10 ***
## hs_grad       0.049949   0.015201   3.286 0.00195 **
## frost        -0.006912   0.002447  -2.824 0.00699 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.7427 on 46 degrees of freedom
## Multiple R-squared:  0.7127, Adjusted R-squared:  0.6939
## F-statistic: 38.03 on 3 and 46 DF,  p-value: 1.634e-12
```

```
summary(bw_l)

##
## Call:
## lm(formula = life_exp ~ murder + hs_grad + frost + population,
##     data = life_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.47095 -0.53464 -0.03701  0.57621  1.50683
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  7.103e+01  9.529e-01  74.542 < 2e-16 ***
## murder       -3.001e-01  3.661e-02  -8.199 1.77e-10 ***
## hs_grad       4.658e-02  1.483e-02   3.142 0.00297 **
## frost        -5.943e-03  2.421e-03  -2.455 0.01802 *
## population    5.014e-05  2.512e-05   1.996 0.05201 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.7197 on 45 degrees of freedom
## Multiple R-squared:  0.736, Adjusted R-squared:  0.7126
```

```
## F-statistic: 31.37 on 4 and 45 DF, p-value: 1.696e-12
```

Judging from the Adjusted R-square, the differences between two models are less than 6%. So according to the principle of parsimony, I choose to discard population.

```
backward_fit = lm(life_exp ~ murder + hs_grad + frost, data = life_data)
```

c) Is there any association between 'Illiteracy' and 'HS graduation rate'? Does your 'subset' contain both?

```
cor.test(illiteracy, hs_grad, method="pearson")
```

```
##
## Pearson's product-moment correlation
##
## data: illiteracy and hs_grad
## t = -6.0408, df = 48, p-value = 2.172e-07
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
## -0.7908657 -0.4636561
## sample estimates:
## cor
## -0.6571886
```

Yes, there is association between illiteracy and hs_grad. The subset we got from forward elimination contains both.

3. Use criterion-based procedures studied in class to guide your selection of the 'best subset'. Summarize your results (tabular or graphical).

```
life_data = life_data %>%
  select(life_exp, everything())
```

```
# Printing the 1 best models of each size, using the Cp criterion:
leaps(x = life_data[,2:8], y = life_data[,1], nbest=1, method="Cp")
```

```
## $which
##      1      2      3      4      5      6      7
## 1 FALSE FALSE FALSE TRUE FALSE FALSE FALSE
## 2 FALSE FALSE FALSE TRUE  TRUE FALSE FALSE
## 3 FALSE FALSE FALSE TRUE  TRUE  TRUE FALSE
## 4  TRUE FALSE FALSE TRUE  TRUE  TRUE FALSE
## 5  TRUE  TRUE FALSE TRUE  TRUE  TRUE FALSE
## 6  TRUE  TRUE  TRUE TRUE  TRUE  TRUE FALSE
## 7  TRUE  TRUE  TRUE TRUE  TRUE  TRUE  TRUE
##
## $label
## [1] "(Intercept)" "1"          "2"          "3"          "4"
## [6] "5"              "6"          "7"
##
## $size
## [1] 2 3 4 5 6 7 8
##
```

```
## $Cp
## [1] 16.126760 9.669894 3.739878 2.019659 4.008737 6.001959 8.000000

# Printing the 1 best models of each size, using the adjusted R2 criterion:
leaps(x = life_data[,2:8], y = life_data[,1], nbest=1, method="adjr2")

## $which
##      1      2      3      4      5      6      7
## 1 FALSE FALSE FALSE TRUE FALSE FALSE FALSE
## 2 FALSE FALSE FALSE TRUE  TRUE FALSE FALSE
## 3 FALSE FALSE FALSE TRUE  TRUE  TRUE FALSE
## 4  TRUE FALSE FALSE TRUE  TRUE  TRUE FALSE
## 5  TRUE  TRUE FALSE TRUE  TRUE  TRUE FALSE
## 6  TRUE  TRUE  TRUE TRUE  TRUE  TRUE FALSE
## 7  TRUE  TRUE  TRUE TRUE  TRUE  TRUE  TRUE
##
## $label
## [1] "(Intercept)" "1"          "2"          "3"          "4"
## [6] "5"          "6"          "7"
##
## $size
## [1] 2 3 4 5 6 7 8
##
## $adjr2
## [1] 0.6015893 0.6484991 0.6939230 0.7125690 0.7061129 0.6993268 0.6921823

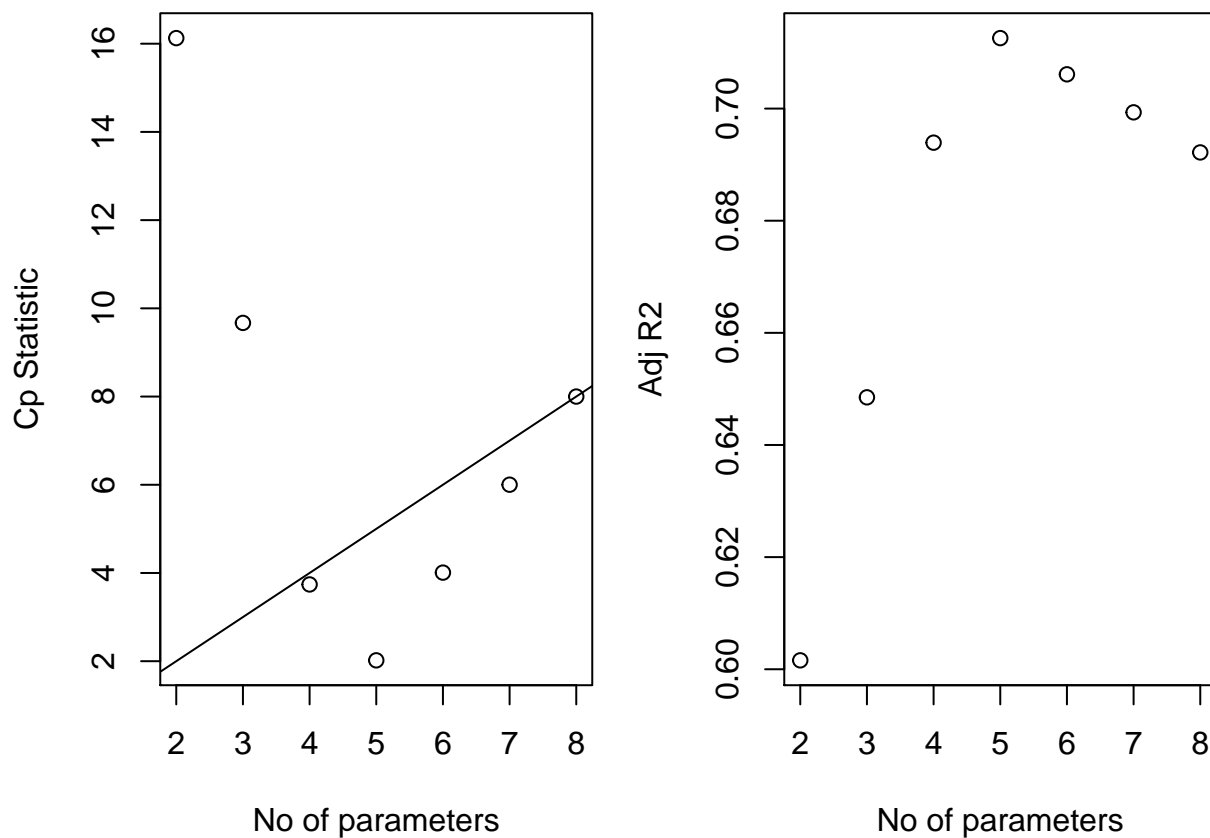
# Summary of models for each size (one model per size)
b<-regsubsets(life_exp ~ ., data=life_data)
(rs<-summary(b))

## Subset selection object
## Call: regsubsets.formula(life_exp ~ ., data = life_data)
## 7 Variables (and intercept)
##      Forced in Forced out
## population      FALSE      FALSE
## income          FALSE      FALSE
## illiteracy      FALSE      FALSE
## murder          FALSE      FALSE
## hs_grad         FALSE      FALSE
## frost          FALSE      FALSE
## area           FALSE      FALSE
## 1 subsets of each size up to 7
## Selection Algorithm: exhaustive
##      population income illiteracy murder hs_grad frost area
## 1 ( 1 ) " "      " "      " "      "*"      " "      " "
## 2 ( 1 ) " "      " "      " "      "*"      "*"      " "
## 3 ( 1 ) " "      " "      " "      "*"      "*"      "*"
## 4 ( 1 ) "*"      " "      " "      "*"      "*"      "*"
## 5 ( 1 ) "*"      "*"      " "      "*"      "*"      "*"
## 6 ( 1 ) "*"      "*"      "*"      "*"      "*"      "*"
## 7 ( 1 ) "*"      "*"      "*"      "*"      "*"      "*"

# Plots of Cp and Adj-R2 as functions of parameters
par(mar=c(4,4,1,1))
par(mfrow=c(1,2))
```

```
plot(2:8, rs$cp, xlab="No of parameters", ylab="Cp Statistic")
abline(0,1)
```

```
plot(2:8, rs$adjr2, xlab="No of parameters", ylab="Adj R2")
```



Judging from the Cp statistics and Adjusted R-square, models with 4-8 parameters are better.

```
# AIC of the 3-predictor model:
```

```
pre_3 <- lm(life_exp ~ murder + hs_grad + frost, data = life_data)
AIC(pre_3)
```

```
## [1] 117.9743
```

```
# BIC
```

```
AIC(pre_3, k = log(length(life_exp)))
```

```
## [1] 127.5344
```

```
# AIC of the 4-predictor model:
```

```
pre_4 <- lm(life_exp ~ murder + hs_grad + frost + population, data = life_data)
AIC(pre_4)
```

```
## [1] 115.7326
```

```
# BIC
```

```
AIC(pre_4, k = log(length(life_exp)))
```

```
## [1] 127.2048
```

```
# AIC of the 5-predictor model:
```

```
pre_5 <- lm(life_exp ~ murder + hs_grad + frost + population + income, data = life_data)
```

```

AIC(pre_5)

## [1] 117.7196
# BIC
AIC(pre_5, k = log(length(life_data$life_exp)))

## [1] 131.1038
# AIC of the 6-predictor model:
pre_6 <- lm(life_exp ~ murder + hs_grad + frost + population + income + illiteracy, data = life_data)
AIC(pre_6)

## [1] 119.7116
# BIC
AIC(pre_6, k = log(length(life_data$life_exp)))

## [1] 135.0077
# AIC of the 7-predictor model:
pre_7 <- lm(life_exp ~ murder + hs_grad + frost + population + income + illiteracy + area, data = life_data)
AIC(pre_7)

## [1] 121.7092
# BIC
AIC(pre_7, k = log(length(life_data$life_exp)))

## [1] 138.9174

```

No of parameter	4	5	6	7	8
Adjusted R-square	0.6939230	0.7125690	0.7061129	0.6993268	0.6921823
Cp	3.7399	2.0197	4.0087	6.0020	8.0000
AIC	117.974	115.733	117.720	119.712	121.709
BIC	127.534	127.205	131.104	135.008	138.917

The model with 5 parameters (4 predictors) has the highest Adjusted R-square and lowest AIC and BIC. So the best model is the one with 5 parameters.

4. Compare the two ‘subsets’ from parts 2 and 3 and recommend a ‘final’ model. Using this ‘final’ model do the following:

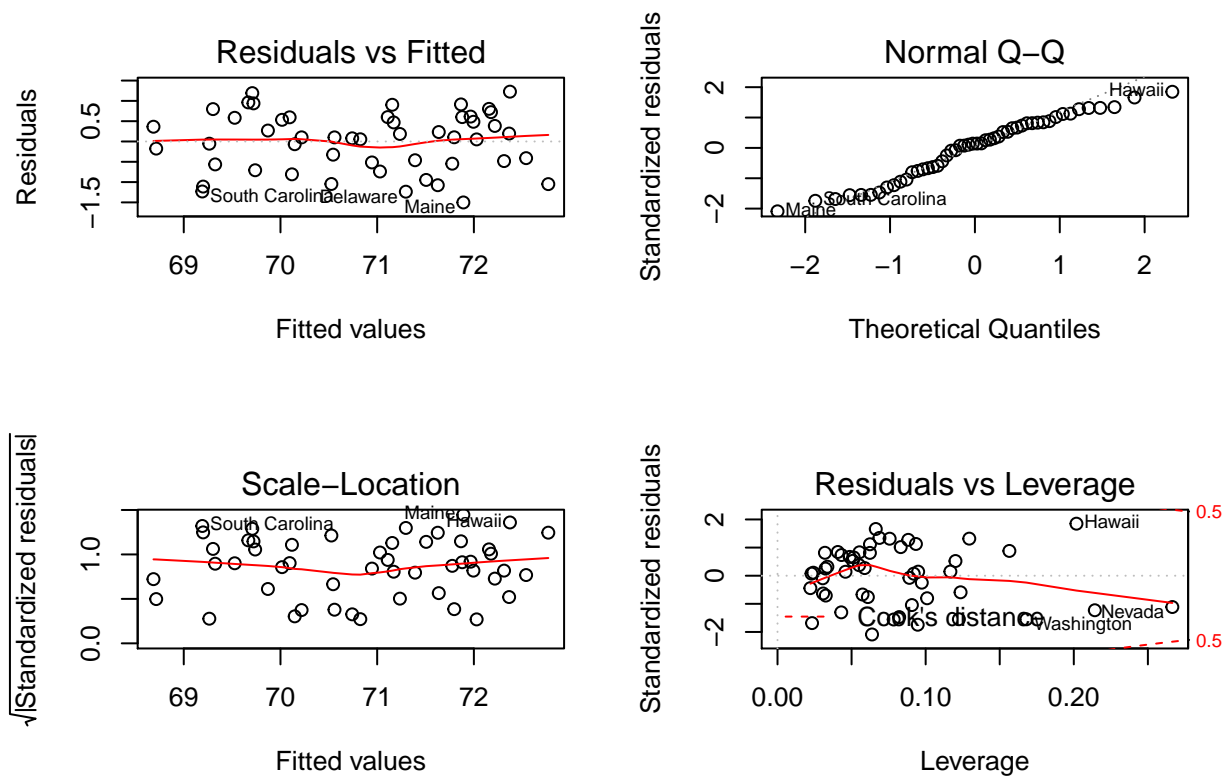
Comparing model with 3 predictors with model with 4 predictors, since the differences between Adjusted R-square, AIC and BIC are pretty small, according to the principle of parsimony, I choose model with 3 predictors, which is `life_exp ~ murder + hs_grad + frost`.

a) Identify any leverage and/or influential points and take appropriate measures.

```

par(mfrow=c(2,2))
plot(pre_3)

```



According to the Residuals vs Leverage plot, there is no leverage or influential points.

b) Check the model assumptions.

Judging from the QQ plot, the residuals are almost normally distributed. Judging from the Residuals vs Fitted values plot and Scale-Location plot, the residuals have constant variance. There is no certain pattern in Residuals vs Fitted values plot, so the residuals are independent.

5. Using the 'final' model chosen in part 4, focus on MSE to test the model predictive ability:

a) Use a 10-fold cross-validation (10 repeats).

```
set.seed(1)
data_train<-trainControl(method="cv", number=10)

model_caret<-train(life_exp ~ murder + hs_grad + frost,
  data=life_data,
  trControl=data_train,
  method='lm',
  na.action=na.pass)

model_caret

## Linear Regression
##
## 50 samples
## 3 predictor
```



```
##
## No pre-processing
## Resampling: Cross-Validated (10 fold)
## Summary of sample sizes: 44, 44, 44, 45, 45, 45, ...
## Resampling results:
##
##      RMSE      Rsquared    MAE
##  0.759794  0.7869101  0.642568
##
## Tuning parameter 'intercept' was held constant at a value of TRUE
```

The RMSE is 0.759794.

b) Experiment a new, but simple bootstrap technique called “residual sampling”.

```
boot_res = lm(life_exp ~ murder + hs_grad + frost, data=life_data)

pred = predict(boot_res)
resid = residuals(boot_res)

res_data = tibble(resid = residuals(boot_res))

boot_sample = function(df) {
  sample_frac(df, replace = TRUE)
}
```

Repeat 10 times

```
set.seed(1)

list = ls()
i = 1

for (i in 1:10){
  res_boot = boot_sample(res_data)
  y_star = res_boot$resid + pred
  life_boot_data = bind_cols(life_data, tibble(y_star))
  boot_res_reg = lm(y_star ~ murder + hs_grad + frost, data=life_boot_data)
  list[i] = rmse(boot_res_reg, life_boot_data)
  i = i + 1
}

repeat_10 = tibble(rMSE = list[1:10])
repeat_10

## # A tibble: 10 x 1
##   rMSE
##   <chr>
## 1 0.704451873543589
## 2 0.581079840881366
## 3 0.638528072155515
## 4 0.637400021055365
## 5 0.728060056647469
## 6 0.71606216374113
```

```
## 7 0.755885751965547
## 8 0.769793897157417
## 9 0.777938189350777
## 10 0.733343726196286
```

```
repeat_10 %>%
  mutate(rMSE = as.numeric(rMSE)) %>%
  summary()
```

```
##      rMSE
## Min.   :0.5811
## 1st Qu.:0.6550
## Median :0.7221
## Mean   :0.7043
## 3rd Qu.:0.7503
## Max.   :0.7779
```

Repeat 1000 times

```
set.seed(1)

list = ls()
i = 1

for (i in 1:1000){
  res_boot = boot_sample(res_data)
  y_star = res_boot$resid + pred
  life_boot_data = bind_cols(life_data, tibble(y_star))
  boot_res_reg = lm(y_star ~ murder + hs_grad + frost, data=life_boot_data)
  list[i] = rmse(boot_res_reg, life_boot_data)
  i = i + 1
}

repeat_1000 = tibble(rMSE = list[1:1000])
repeat_1000
```

```
## # A tibble: 1,000 x 1
##   rMSE
##   <chr>
## 1 0.704451873543589
## 2 0.581079840881366
## 3 0.638528072155515
## 4 0.637400021055365
## 5 0.728060056647469
## 6 0.71606216374113
## 7 0.755885751965547
## 8 0.769793897157417
## 9 0.777938189350777
## 10 0.733343726196286
## # ... with 990 more rows
```

```
repeat_1000 %>%
  mutate(rMSE = as.numeric(rMSE)) %>%
  summary()
```

```
##      rMSE
## Min.   :0.4972
```

```
## 1st Qu.:0.6459
## Median :0.6813
## Mean   :0.6809
## 3rd Qu.:0.7192
## Max.    :0.8220
```

c) In a paragraph, compare the MSE values generated by the two methods a) and b). Briefly comment on the differences and your recommendation for assessing model performance.

Comparing rMSEs generating from different methods, we can see that the rMSEs of CV are higher than the rMSE of residual sampling. The rMSE of 1000 repeat residual sampling is lower than 10 repeat. If the model assumptions are not satisfied, the residual sampling procedure is recommended.