# Linear Regression: A Revisit

Yifei Sun

# Contents

```r
library(tidyverse)
library(summarytools)
library(leaps)
```

# Data

In this example, we assess the association between high density lipoprotein (HDL) cholesterol and body mass index, blood pressure, and other demographic factors (age, gender, race) using the NHANES data (https://wwwn.cdc.gov/nchs/nhanes/ContinuousNhanes/Default.aspx?BeginYear=2001). The data can be downloaded using functions in the package `RNHANES`.

```r
load("L4_data.RData")
```

Summary statistics of the predictors and the response:

```r
st_options(plain.ascii = FALSE,
           style = "rmarkdown",
           dfSummary.silent = TRUE,
           footnote = NA,
           subtitle.emphasis = FALSE)

dfSummary(dat[,-1])
```

**Data Frame Summary**

**dat**
**Dimensions:** 6434 x 6
**Duplicates:** 0

| No | Variable | Stats / Values | Freqs (% of Valid) | Graph | Valid | Missing |
|----|----------|----------------|--------------------|-------|-------|---------|
| 1 | gender [factor] | 1. 1<br>2. 2 | 3108 (48.3%)<br>3326 (51.7%) | IIIIIIIII<br>IIIIIIIIII | 6434<br>(100%) | 0<br>(0%) |
| 2 | race [factor] | 1. 1<br>2. 2<br>3. 3<br>4. 4<br>5. 5 | 1593 (24.8%)<br>262 ( 4.1%)<br>2910 (45.2%)<br>1448 (22.5%)<br>221 ( 3.4%) | IIII<br><br>IIIIIIIII<br>IIII | 6434<br>(100%) | 0<br>(0%) |
| 3 | age [numeric] | Mean (sd) : 35.3 (22.1)<br>min < med < max:<br>5 < 29 < 85<br>IQR (CV) : 36 (0.6) | 79 distinct values | :<br>. :<br>: :<br>: : : : : : . .<br>: : : : : : : : : : | 6434<br>(100%) | 0<br>(0%) |
| 4 | bmi [numeric] | Mean (sd) : 26 (6.5)<br>min < med < max:<br>13.4 < 25.3 < 64.2<br>IQR (CV) : 8.2 (0.2) | 2266 distinct values | . :<br>: :<br>: : :<br>. : : :<br>: : : : : . | 6434<br>(100%) | 0<br>(0%) |

| No | Variable | Stats / Values | Freqs (% of Valid) | Graph | Valid | Missing |
|----|----------|----------------|--------------------|-------|-------|---------|
| 5 | sbp [numeric] | Mean (sd) : 119.5 (20.1) min < med < max: 74 < 116 < 228 IQR (CV) : 22 (0.2) | 73 distinct values | : : : . : : : : : : : . | 6434 (100%) | 0 (0%) |
| 6 | hdl [numeric] | Mean (sd) : 51.6 (14.5) min < med < max: 19 < 49 < 160 IQR (CV) : 17 (0.3) | 102 distinct values | : . : : : : : : : : : : : . | 6434 (100%) | 0 (0%) |

## Multiple linear regression: a small experiment

```
fit1 <- lm(hdl ~ bmi + age + gender + race + sbp,
           data = dat)

fit2 <- lm(bmi ~ age + gender + race + sbp,
           data = dat)

r2 <- fit2$residuals

fit3 <- lm(hdl ~ r2,
           data = dat)

coef(fit1)["bmi"]
```

```
    bmi
```

-0.6649902

```
coef(fit3)["r2"]
```

```
    r2
```

-0.6649902

## Prediction interval vs. confidence interval

```
newdata <- dat[1,]
predict(fit1, newdata, interval = "confidence")
```

```
    fit      lwr      upr
```

1 44.48379 43.83743 45.13016

```
predict(fit1, newdata, interval = "predict") # much wider!
```
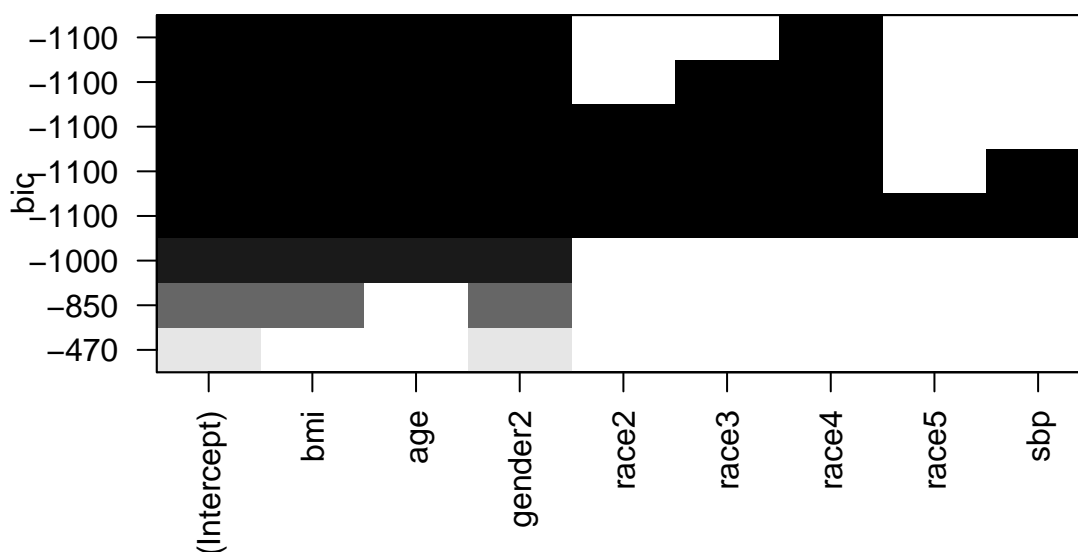
```
     fit      lwr      upr
```

1 44.48379 18.50864 70.45895

# Best subset model selection

```
regsubsetsObj <- regsubsets(hdl ~ bmi + age + gender + race + sbp, data = dat,
                            method = "exhaustive", nbest = 1)

plot(regsubsetsObj, scale = "bic")
```



```
# summary(regsubsetsObj)
```