

Theoretical Solutions

Numerical Computations

$$1) \hat{\beta} = (x^T x)^{-1} (x^T y)$$

$$= \begin{bmatrix} 0.5 & 0 & -0.25 \\ 0 & 0.01 & 0 \\ -0.25 & 0 & 0.25 \end{bmatrix} \begin{bmatrix} 12 \\ -50 \\ 20 \end{bmatrix} = \begin{bmatrix} 6+0-5 \\ 0-0.5+0 \\ -3+0+5 \end{bmatrix} = \begin{bmatrix} 1 \\ -0.5 \\ 2 \end{bmatrix}$$

$$\hat{y} = 1 - 0.5 x_1 + 2 x_2$$

$$x_1 = \begin{bmatrix} 5 \\ 5 \\ -5 \\ -5 \end{bmatrix} \quad x_2 = \begin{bmatrix} 2 \\ 0 \\ 2 \\ 0 \end{bmatrix} \quad \hat{y} = \begin{bmatrix} 1-2.5+4 \\ 1-2.5+0 \\ 1+2.5+4 \\ 1+2.5+0 \end{bmatrix} = \begin{bmatrix} -2.5 \\ -1.5 \\ 7.5 \\ 3.5 \end{bmatrix}$$

Proof for $\hat{\beta}_1$

2) From the normal equations:

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i)}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

We can compute the variance of $\hat{\beta}_1$ as:

$$c_i = \frac{\sum_{i=1}^n (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\text{Var}(\hat{\beta}_1) = \text{Var}\left(\sum_{i=1}^n c_i y_i\right) = \sum_{i=1}^n \text{Var}(c_i y_i)$$

However c_i is non-random, such that

$$\begin{aligned} &= \sum_{i=1}^n c_i^2 \text{Var}(y_i) = \sum_{i=1}^n c_i^2 \sigma^2 = \frac{\sum (x_i - \bar{x})^2}{\sum (x_i - \bar{x})^2} \sigma^2 \\ &= \frac{\sigma^2}{\sum (x_i - \bar{x})^2} \end{aligned}$$

Proof for $\hat{\beta}_0$

$$\begin{aligned} \text{Var}(\hat{\beta}_0) &= \text{Var}(\bar{y} - \hat{\beta}_1 \bar{x}) \\ &= \frac{1}{n^2} [\text{Var}(\sum_i y_i - \hat{\beta}_1 \sum_i x_i)] \\ &= \frac{1}{n^2} [\text{Var}(\sum_i y_i) + \sum_i x_i^2 \text{Var}(\hat{\beta}_1)] \\ &= \frac{1}{n^2} [n \sigma^2 + \sum_i x_i^2 \cdot \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}] \\ &= \sigma^2 \left[\frac{1}{n} + \frac{\sum x_i^2}{n} \cdot \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \right] \\ &= \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right] \end{aligned}$$

Applied Analysis

Part A

The following reduction on the data was performed:

- Samples below 2% and above 40% body fat were excluded. The minimum was obtained through the chart on the ([ACE](#)) website and maximum was set to limit outliers.
- The minimum height accepted was cut to be over 50 inches to remove outliers.
- The maximum weight was also clipped to be below 300 pounds. Although the outlier may have been legitimate, it could have also been an error.

For each individual linear model created between the response of body fat percentage and predictors weight, height and abdomen circumference, the following hypotheses were assumed.

$$\begin{aligned} H_0 : \beta_0 &= \beta_1 = 0 \\ H_1 : \exists \beta_i &\neq 0 \end{aligned}$$

To test these hypotheses, a T-distribution test was performed on each model with the resultant p-values:

- Weight Predictor: $p < 2e-16$
- Height Predictor: $p = 0.805$
- Abdomen Circumference Predictor: $p < 2e-16$

After performing the tests, we cannot reject the null hypothesis for the height as a predictor of body fat percentage. We can however reject the null hypotheses for the weight and abdomen circumference, implying that there is some non-zero slope relationship between the body fat percentage and these predictors.

Based on our regression models, we can obtain the following 99% confidence intervals for β_1 of each model:

- For each pound increase in weight, we are 99% confident that, on average, the body fat percentage will increase between 0.142 and 0.221 percent.
- For each inch increase in height, we are 99% confident that, on average, the body fat percentage will increase between 0.571 and 0.471 percent.
- For each cm increase in abdomen circumference, we are 99% confident that, on average, the body fat percentage will increase between 0.574 and 0.729 percent.

As we can see from the confidence intervals above, not rejecting the null hypothesis for the height as a predictor is further reinforced due to the inclusion of 0 within the interval.

Looking at the coefficients of determination for each model, the values for the weight ($R^2 = 0.3658$) and abdomen circumference ($R^2 = 0.6607$) models are significantly higher than that of the height model ($R^2 = 0.0611$). Again, looking at the confidence intervals, we can see that the abdomen circumference has a lower bound that is much higher (more than twice) than the upper bound for the weight, indicating that abdomen circumference is the strongest predictor out of the three.

Part B

Next, we are interested whether or not there is evidence to indicate that the average increase in body fat percentage for each cm increase of abdomen circumference is actually greater than 0.5%. A t-test was performed on the following hypotheses.

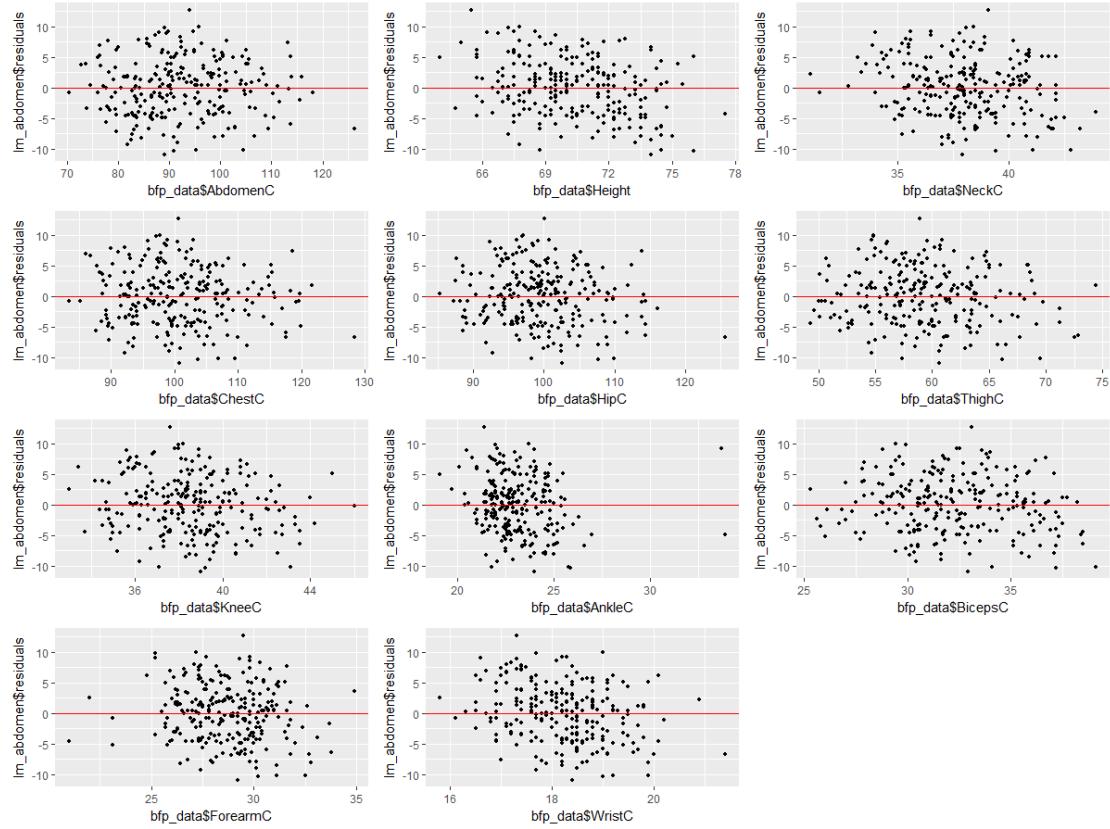
$$H_0 : \beta_1 \leq 0.5$$

$$H_1 : \beta_1 > 0.5$$

The t-value obtained was 5.067856 and an underlying normal distribution was assumed for β_1 . As such, the p-value was calculated using $1 - pnorm(t)$ in R to obtain a p-value of $2.011611 * 10^{-7}$. With the p-value confirming that our confidence interval for the slope is correct, and that the actual value for β_1 is over 0.5.

Part C

The following figure includes the residuals from the abdomen circumference predictor model plotted against weight, height, neck circumference, chest circumference, hip circumference, thigh circumference, knee circumference, ankle circumference, biceps circumference, forearm circumference, and wrist circumference.



From these plots, we can see that there may potentially be a relationship between the body fat percentage and the ankle circumference. It may make sense that wrist circumference could indicate an increase in fat percentage, based on where the wrist is defined to lie. If the measurements were taken further up the forearm where one would wear a watch, as opposed to the actual wrist joint, then the measurement may indicate something about the body fat percentage.

Part D

Next if we fit a linear model (stated below) using all of these variables to predict the body fat percentage, we will be able to compare it to the reduced (abdomen circumference) only model in order to determine whether the other predictors are

actually necessary. The test statistic used for this was an f-distribution (also noted below).

$$BFP = \beta_0 + \beta_1 * AbdomenC + \beta_2 * Weight + \beta_3 * Height + \beta_4 * NeckC + \beta_5 * ChestC + \beta_6 * HipC + \beta_7 * ThighC + \beta_8 * KneeC + \beta_9 * AnkleC + \beta_{10} * BicepsC + \beta_{11} * ForearmC + \beta_{12} * WristC + \epsilon$$

$$\widehat{BFP} = 1.36987591 + 0.98199026 * AbdomenC - 0.04420200 * Weight - 0.26163775 * Height - 0.38320373 * NeckC - 0.10481085 * ChestC - 0.20730987 * HipC + 0.06182966 * ThighC + 0.14670469 * KneeC + 0.12543748 * AnkleC + 0.17961142 * BicepsC + 0.23096461 * ForearmC - 1.40026669 * WristC$$

$$\hat{\sigma}^2 = 17.9776$$

$$R^2 = 0.7304$$

The F-value used in the f-test was computed to be the following value, and using the given parameters:

$$H_0 : \omega : y_i = \beta_0 + \beta_1 * AbdomenC_i$$

$$H_1 : \Omega : \text{The full model stated above where } \exists \beta_i \neq 0, i = [2, 12]$$

$$n = 246, \quad p = 13, \quad q = 2$$

$$F \sim F_{11,233}$$

$$F_{value} = 5.482263$$

$$p-value = 9.011923 * 10^{-8}$$

Based on the results of the F-test and looking at the R-squared values, we cannot reject the null hypothesis that the abdomen circumference model is not strong enough to predict the body fat percentage accurately.

R Output

This section contains the output from the R console for important values and summaries.

```

1 #####
2 ##### Display Results
3 #####
4
5 ## Part A
6 summary(lm_weight)
7 Residuals:
8      Min       1Q    Median       3Q      Max
9 -18.1734   -4.5946    0.0303   4.8440   18.5365
10 Coefficients:
11              Estimate Std. Error t value Pr(>|t|)
12 (Intercept) -13.30646   2.76104 -4.819 2.53e-06 ***
13 Weight        0.18185   0.01533 11.864 < 2e-16 ***
14
15 Signif. codes:  0 *** 0.001 ** 0.01 * 0.05 . 0.1     1
16 Residual standard error: 6.355 on 244 degrees of freedom
17 Multiple R-squared:  0.3658, Adjusted R-squared:  0.3632
18 F-statistic: 140.8 on 1 and 244 DF,  p-value: < 2.2e-16
19
20 summary(lm_height)
21 Residuals:
22      Min       1Q    Median       3Q      Max
23 -16.2242   -6.4851    0.2486   6.0752   20.7394
24 Coefficients:

```

```

25             Estimate Std. Error t value Pr(>|t|)
26 (Intercept) 22.5848     14.1206   1.599    0.111
27 Height      -0.0496     0.2006  -0.247    0.805
28 Residual standard error: 7.979 on 244 degrees of freedom
29 Multiple R-squared:  0.0002504, Adjusted R-squared:  -0.003847
30 F-statistic: 0.06111 on 1 and 244 DF, p-value: 0.805
31
32 summary(lm_abdomen)
33 Residuals:
34     Min      1Q Median      3Q      Max
35 -10.9111 -3.5650  0.1989  3.1421 12.7618
36 Coefficients:
37                 Estimate Std. Error t value Pr(>|t|)
38 (Intercept) -41.03612   2.77484 -14.79 <2e-16 ***
39 AbdomenC     0.65148   0.02989  21.80 <2e-16 ***
40
41 Signif. codes:  0 *** 0.001 ** 0.01 * 0.05 . 0.1
42 Residual standard error: 4.648 on 244 degrees of freedom
43 Multiple R-squared:  0.6607, Adjusted R-squared:  0.6593
44 F-statistic: 475 on 1 and 244 DF, p-value: < 2.2e-16
45
46 lm_weight_conf
47             0.5 %    99.5 %
48 (Intercept) -20.4744889 -6.1384398
49 Weight       0.1420582  0.2216441
50
51 lm_height_conf
52             0.5 %    99.5 %
53 (Intercept) -14.0740731 59.2436692
54 Height       -0.5705262  0.4713198
55
56 lm_abdomen_conf
57             0.5 %    99.5 %
58 (Intercept) -48.2399646 -33.8322748
59 AbdomenC     0.5738825  0.7290846
60
61 # Part B
62 t_abdomen
63 [1] 5.067856
64 p_abdomen
65 [1] 2.011611e-07
66
67
68 # Part D
69 f
70 [1] 5.482263
71 pvalue
72 [1] 9.011923e-08
73
74 # Check results
75 anova(lm_abdomen, lm_full)
76 Analysis of Variance Table
77
78 Model 1: SiriBFperc ~ AbdomenC
79 Model 2: SiriBFperc ~ (Density + AbdomenC + Weight + Height + NeckC +
80                      ChestC + HipC + ThighC + KneeC + AnkleC + BicepsC + ForearmC +
81                      WristC) - Density
82 Res.Df   RSS Df Sum of Sq      F      Pr(>F)
83 1     244  5272.4
84 2     233 4188.4 11      1084 5.4823 9.012e-08 ***
85
86 Signif. codes:  0 *** 0.001 ** 0.01 * 0.05 . 0.1

```