

method	age	weight	height	adiposity	neck	chest	abdomen	hip	thigh	knee	ankle	biceps	forearm	wrist
stepwise (BIC)			•				•							•

stepwise (AIC)	•	•	•	•	•	•	•	•	•	•	•	•	•	•
backward (BIC)	•													•
backward (AIC)	•		•	•	•	•	•	•	•	•	•	•	•	•
forward (BIC)	•	•	•	•	•	•	•	•	•	•	•	•	•	•
forward (AIC)	•	•	•	•	•	•	•	•	•	•	•	•	•	•
best subset		•												•
Lasso	•		•					•						•
Mallow's Cp	•		•				•	•				•		•

Statistical Analysis

A. Best Model Fitting

Here are some models listed to compare R^2 and MSE (mean of least square error)[Table1]. The listed models are either models with variable selection methods or models with small amount of features.

We decide to choose the model with two features: abdomen and wrist. There are mainly three reasons.

1. The results of different variable selection methods above all include abdomen and wrist variables.
2. We find that measurement of Adjusted R^2 and MSE do not improve too much after adding more variables.
3. The model with only two variables has small model complexity and it is easy to interpret each variable.

Our final model is $\text{Bodyfat} = \beta_0 + \beta_1 \cdot \text{abdomen} + \beta_2 \cdot \text{wrist}$.

B. Model Diagnostics

We generally check the model based on three parts: the influential points, the assumption on normality and the assumption on residuals.

1. Checking for Influential Points

	features	Adj R2	MSE
model0_1	abdomen	0.6576	4.3123
model1_0	abdomen+wrist	0.699	4.0433
model1_1	weight+abdomen+wrist	0.715	3.9342
model1_2	age+height+abdomen+wrist	0.7178	3.9145
model1_3	age+height+chest+abdomen+biceps+wrist	0.7206	3.8952
model1_4	weight+height+adiposity+abdomen+biceps+wrist	0.716	3.9273

Table1.1

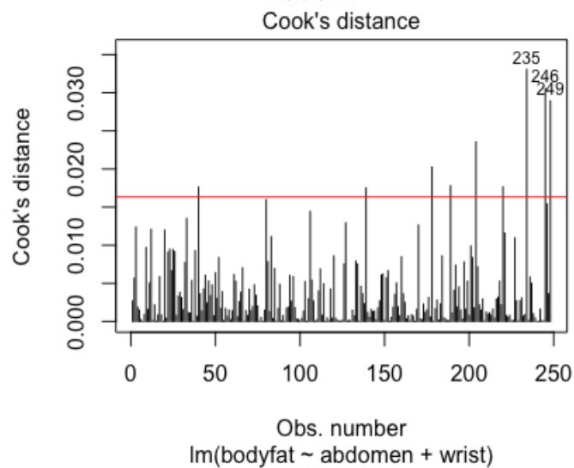


Figure1.2

Residuals vs Fitted

Body Fat vs. 1/Density

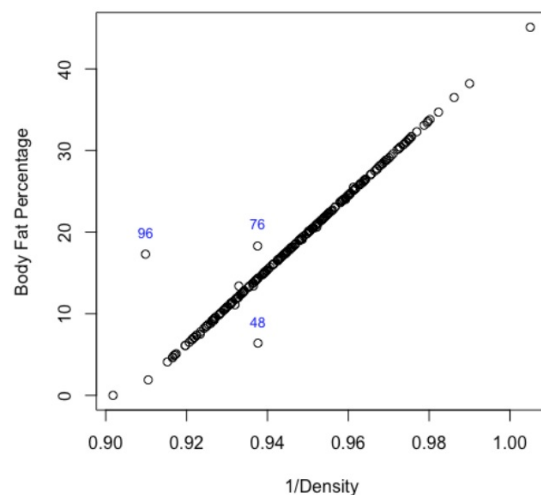
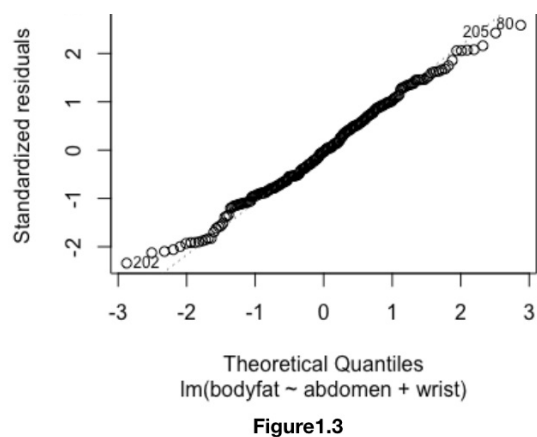
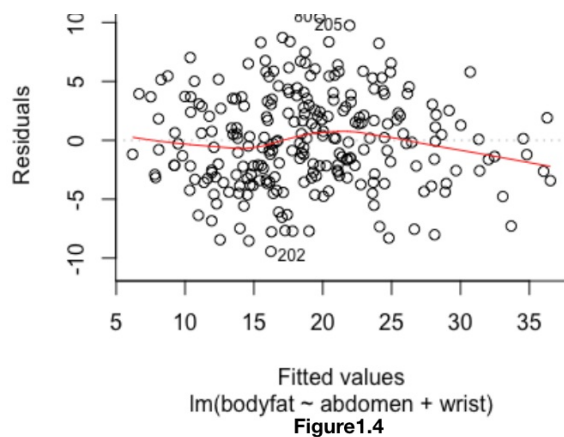


Figure1.1

Normal Q-Q



Based on **Figure1.2**, we find that there are 9 points which Cook's Distance is greater than the threshold value. We check measurements for these observations and think that they all look normal. Therefore we decided to retain all the observations in our model.

2. Checking for Normality

According to the QQ-plot **Figure1.3**, all points are close to the dashed diagonal line, so we believe that the residuals follow the normally distribution.

3. Checking for Residuals

According to **Figure1.4**, the model tends to have a higher variance when fitted values range from 15 to 25. Also, the mean of the residuals become smaller when fitted value is in the range of 30 to 35. So we believe that our model does not satisfy $E(\epsilon) = 0$ and $\text{Var}(\epsilon) = \text{constant}$ perfectly.

C. Hypothesis Testing

We generally do the following 2 hypothesis testing:

1. t-testing on coefficients $H_0: \beta_i = 0$ v.s. $H_1: \beta_i \neq 0$
2. F-testing on the full model we choose $H_0: \beta_1 = \beta_2 = 0$ v.s. $H_1: H_0$ not true

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	-9.26943	5.23360	-1.771	0.0778	.
abdomen	0.71469	0.03216	22.225	< 2e-16	***
wrist	-2.07472	0.35157	-5.901	1.19e-08	***

F-statistic: 287.8 on 2 and 245 DF, p-value: < 2.2e-16

The result is shown above. Based on F-testing, our model has a significant effect for all variables.

Based on t-testing, all of our variables are significant according to p-values and standard errors for each variable is small, with a significant level less than 0.05. The abdomen variable has a positive effect on bodyfat while the wrist variable has a negative effect on bodyfat.

Contributions

1. Bi Qing Teng: Analysing raw data, data cleaning and variable selection
2. Xiaoxiang Hua: Variable selection, Model Diagnostics, Model interpretation and Model strength and weaknesses
3. Yijie Liu: Write the Shiny app code, and make the slides, and jupyter notebook summary.