

module2

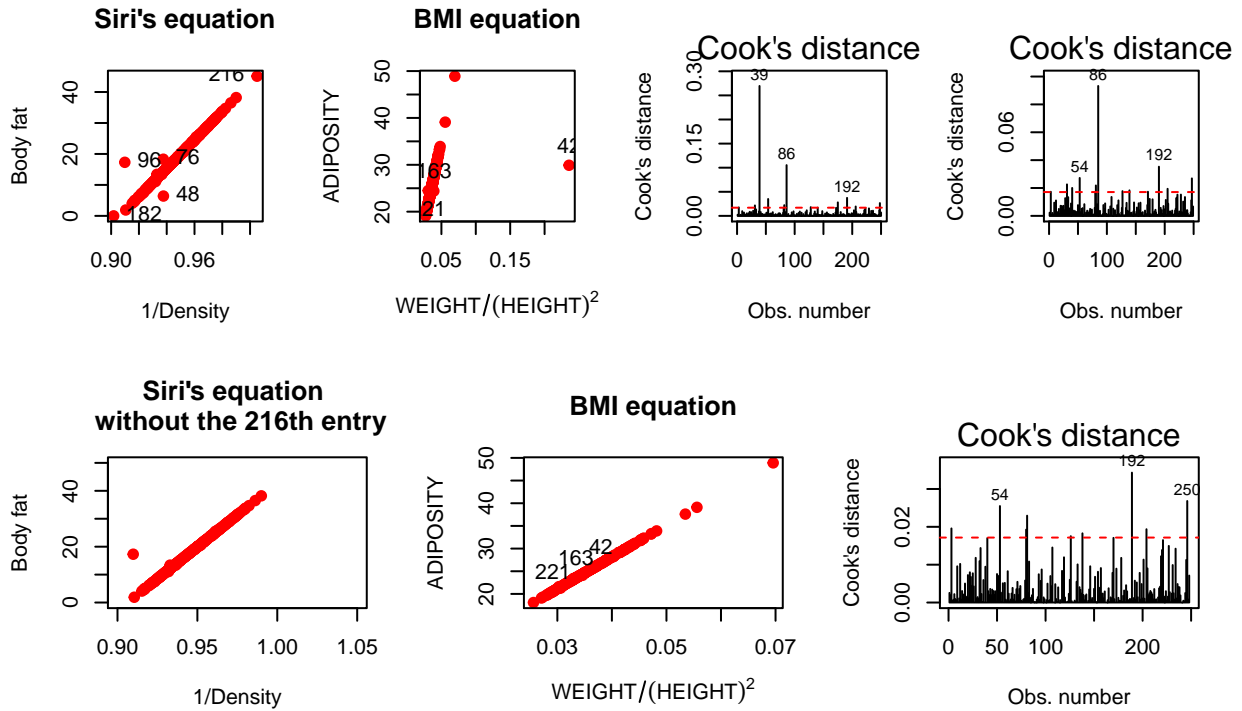
Hao Qin

1.Cleaning data

2.Detection of the outliers and their treatment

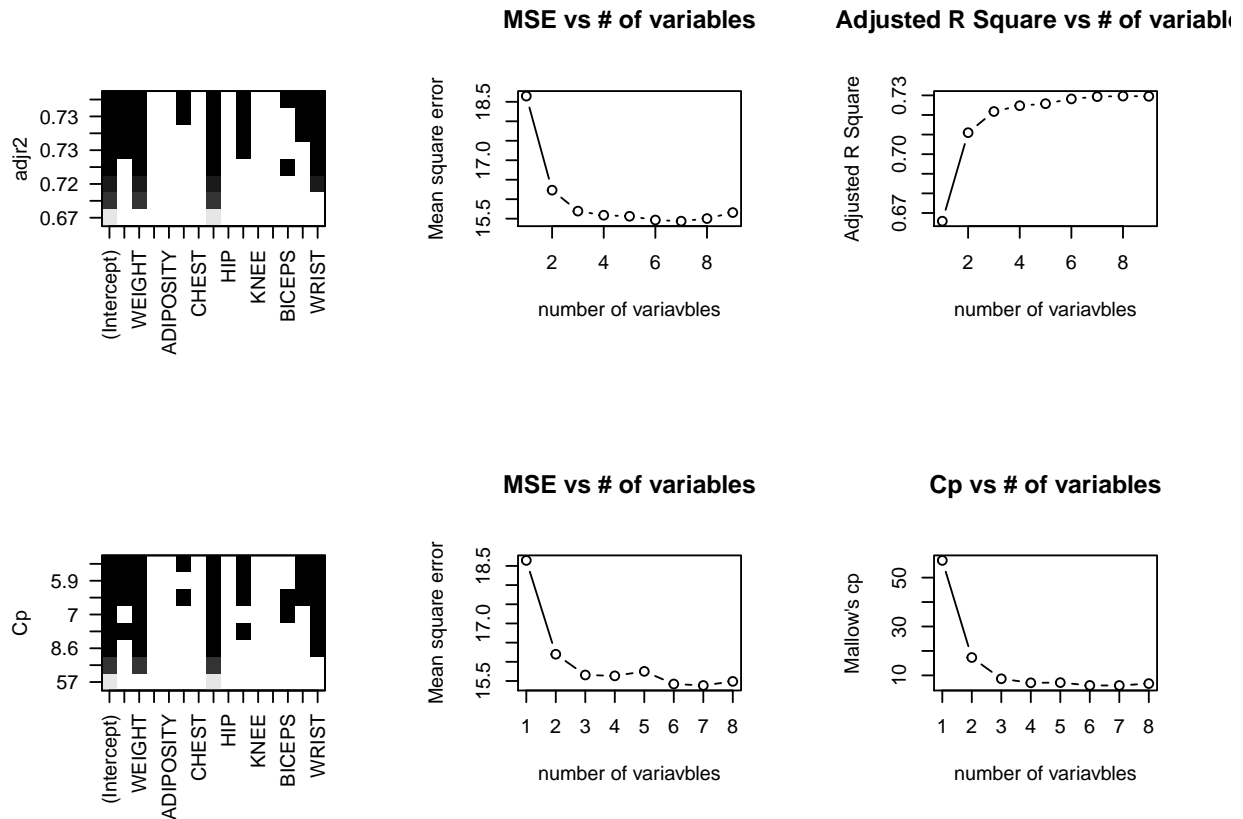
We use three different types of method to detect not only the outliers, but also the influence points in the original dataset, which is Siri's equation, BMI formula and cook's distance. Those points will deviate from the line or the region on which most points are, thus, it is easy for us to detect them through graphics.

Outliers	Detection_method	Abnormal_Component	Treatment
39	Cook's distance	large Cook's distance	deleting
42	BMI formula	too short in height	imputation
48	Siri's equation	abnormal body fat	imputation
76	Siri's equation	abnormal body fat	imputation
86	influence test	deviating from the majority	deleting
96	Siri's equation	abnormal body fat	imputation
163	BMI formula	abnormal BMI index	imputation
182	Siri's equation	zero body fat	deleting
216	Siri's equation	density below than 1	deleting
221	BMI formula	lighter than normal	imputation

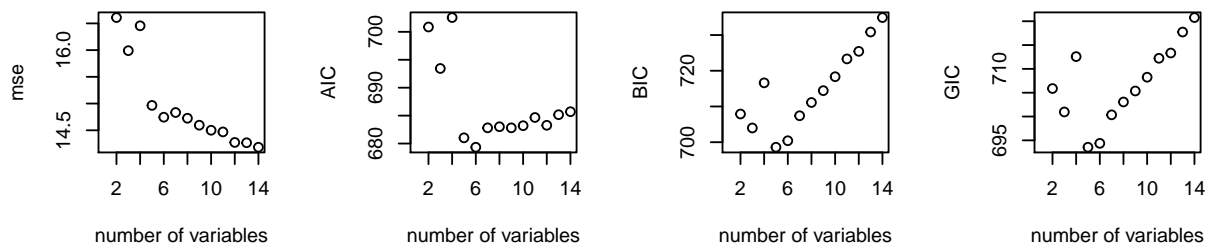


3. Selecting variables

In this part, several methods have been applied, including stepwise aic and bic selection, lasso and group lasso, mallow's cp and Bess, which is a new proposed way to selecting variables. After applying these methods, since the total sample size is not large, we can use cross validation to measure the performance of each method, to decide which component should be treated as the independent variable in the final model. To be specific, the final model should be as a rule of thumb, which requires that the quantity of the independent variables should not exceed four.



#lasso and group lasso



model	Mean_Sq	variables
12 Adjusted R square	15.58826	WEIGHT + ABDOMEN + BICEPS + WRIST

	model	Mean_Sq	variables
8	Mallow's Cp	15.63376	WEIGHT + ABDOMEN + BICEPS + WRIST
7	Mallow's Cp	15.65693	WEIGHT + ABDOMEN + WRIST
1	BIC	15.68525	WEIGHT + ABDOMEN + WRIST
11	Adjusted R square	15.69427	WEIGHT + ABDOMEN + WRIST
15	BeSS	15.96599	WEIGHT + ADIPOSITIVITY + CHEST + ABDOMEN + WRIST
2	BIC	16.19068	WEIGHT + ABDOMEN
6	Mallow's Cp	16.19761	WEIGHT + ABDOMEN
10	Adjusted R square	16.23208	WEIGHT + ABDOMEN
14	BeSS	16.49337	ADIPOSITIVITY + CHEST + ABDOMEN + HIP
13	BeSS	16.74293	HEIGHT + CHEST + ABDOMEN
5	Mallow's Cp	18.64326	ABDOMEN
9	Adjusted R square	18.64655	ABDOMEN
4	BIC	18.65362	ABDOMEN
3	BIC	35.19864	WEIGHT

Model building