# Body Fat Measurement Model
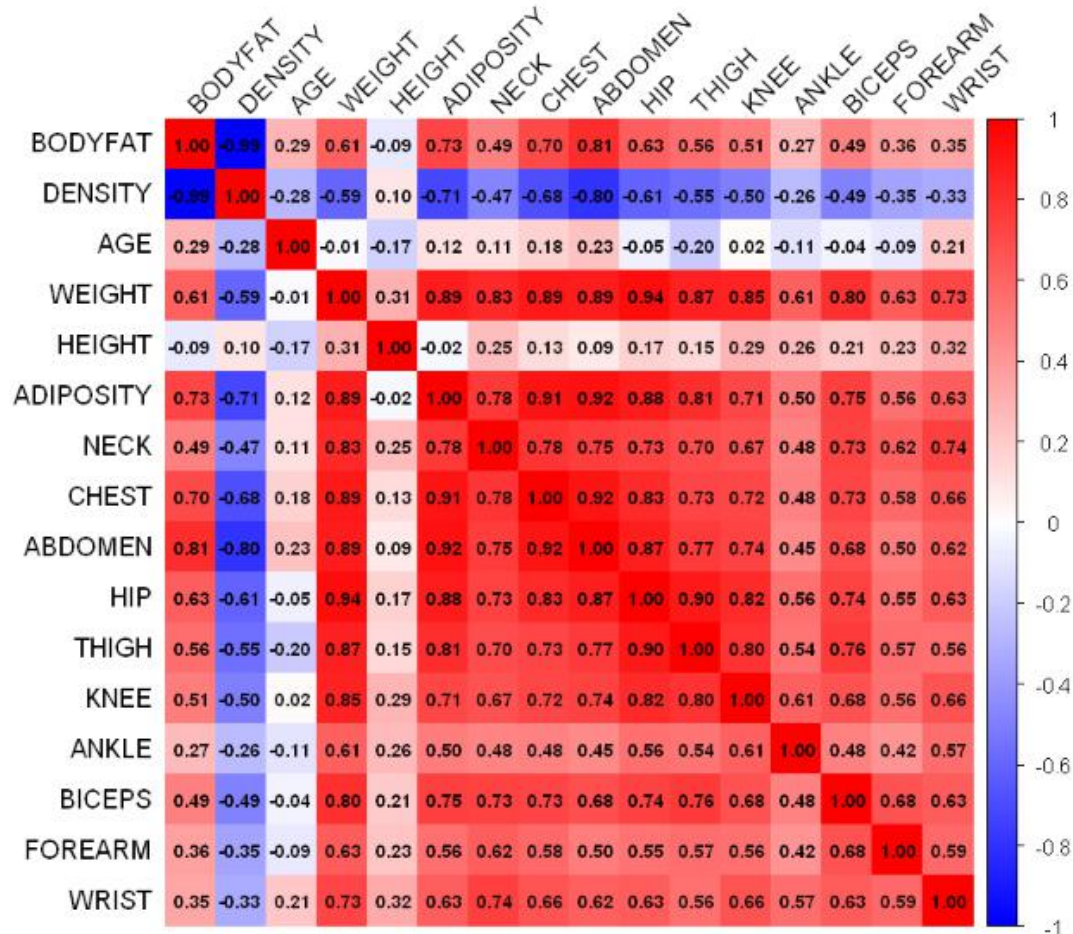
## Group 4

# Data Cleaning

- First, We check the corr. and lof scores.



| Individual (IDNO) | LOF Scores(>2) |
|---|---|
| 39 | 5.616425 |
| 42 | 3.042846 |

- We impute some individuals.

| Individual (IDNO) | Abnormal column | Original Obs. | Imputed Obs. | Reasons | Imputation Method |
|---|---|---|---|---|---|
| 31 | ANKLE | 33.9(cm) | 23.8(cm) | ANKLE's histogram distribution anomaly, Leverage and Cook's distance check | KNN |
| 86 | ANKLE | 33.7(cm) | 22.5(cm) | ANKLE's histogram distribution anomaly, Leverage and Cook's distance check | KNN |
| 221 | WEIGHT | 153.25(lbs) | 173.20(inches) | Mismatch between BMI, HEIGHT, and WEIGHT | BMI and HEIGHT calculate |
| 42 | HEIGHT | 29.5(inches) | 69.43(inches) | Mismatch between BMI, HEIGHT, and WEIGHT | BMI and WEIGHT calculate |

| Individual (IDNO) | Abnormal column | Original Obs. | Imputed Obs. | Reasons | Imputation Method |
|---|---|---|---|---|---|
| 172 | BODYFAT DENSITY | 1.9(%) 1.0983(gm/cm3) | 6.99(%) 1.083(gm/cm3) | BODYFAT and DENSITY calculated using 'Siri's equation' mismatch, Abnormal BODYFAT | Regression by other columns |
| 182 | BODYFAT DENSITY | 0(%) 1.1089(gm/cm3) | 4.53(%) 1.089(gm/cm3) | BODYFAT and DENSITY calculated using 'Siri's equation' mismatch, Abnormal BODYFAT | Regression by other columns |
| 48 | BODYFAT | 6.4(%) | 14.14(%) | BODYFAT and DENSITY calculated using 'Siri's equation' mismatch | DENSITY calculate |
| 76 | BODYFAT | 18.3(%) | 14.09(%) | BODYFAT and DENSITY calculated using 'Siri's equation' mismatch | DENSITY calculate |
| 96 | DENSITY | 1.0991(gm/cm3) | 1.059(gm/cm3) | BODYFAT and DENSITY calculated using 'Siri's equation' mismatch | BODYFAT calculate |

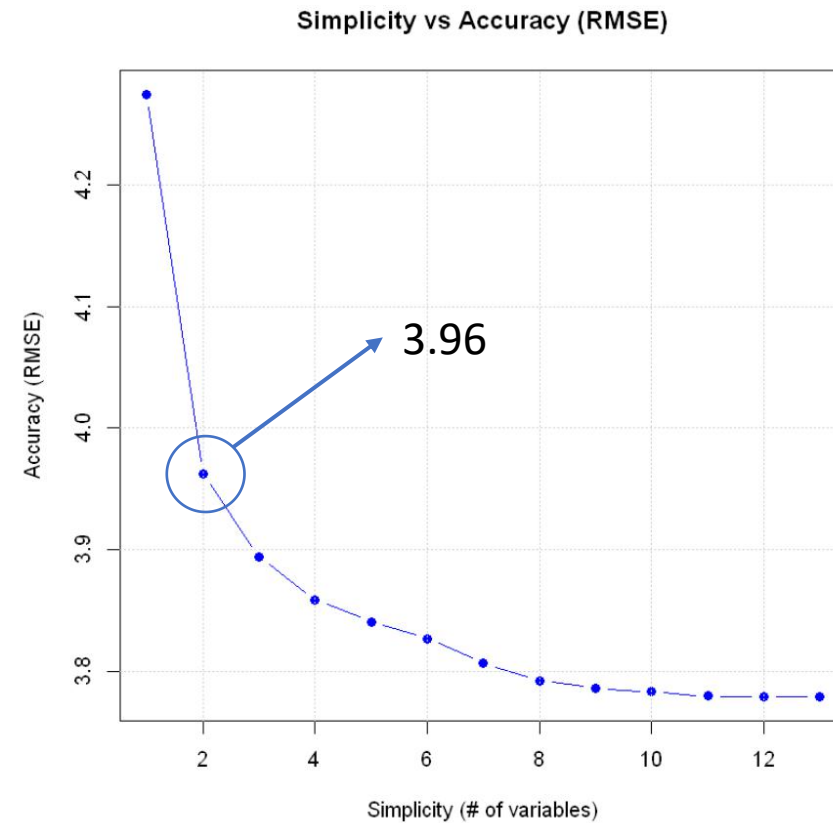- We deleted **five individuals** (39,54,163,175,216.) due to Leverage and Cook's distance check.
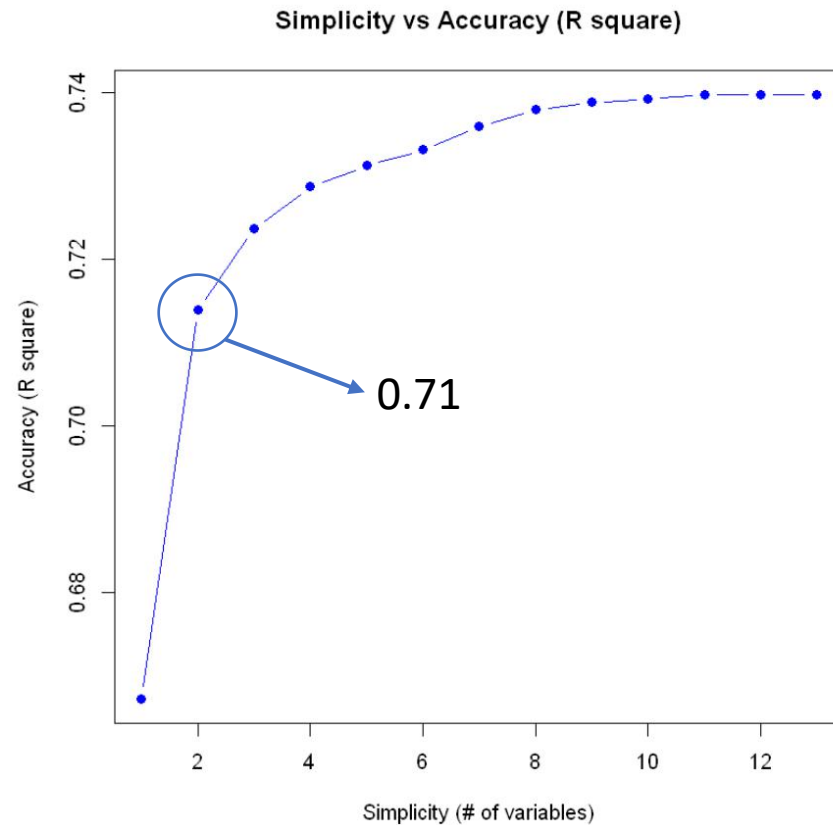
| Individual (IDNO) | Leverage | Cook's distance |
|---|---|---|
| 39 | 0.466 | 0.409 |
| 54 | 0.172 | 0.020 |
| 163 | 0.167 | 0.018 |
| 175 | 0.273 | 0.025 |
| 216 | 0.134 | 0.030 |

Final Cleaned Data: **n=247** (from n=252) with p = 13 predictors

- Remove IDNO
- Remove DENSITY
- Remove BMI – highly collinear
- Predictors: Age (years),Weight (lbs),Height (inches), Adioposity (bmi), Neck circumference (cm), Chest circumference (cm), Abdomen 2 circumference (cm), Hip circumference (cm), Thigh circumference (cm), Knee circumference (cm), Ankle circumference (cm), Biceps (extended) circumference (cm), Forearm circumference (cm), Wrist circumference (cm)

# Final Model

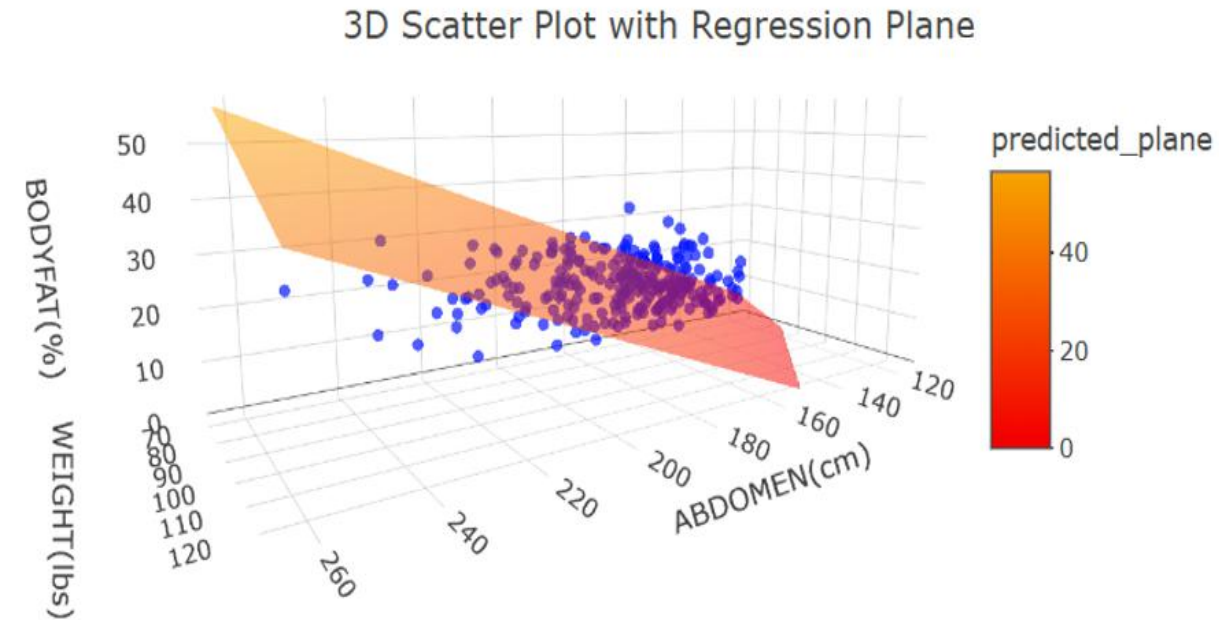BODYFAT = - 41.547 - 0.124*WEIGHT (lbs) + 0.894*ABDOMEN (cm)

# Model Diagnostics

Consider the assumptions of linear regression:

| Assumption | Test | criteria |
|---|---|---|
| Normality of residuals | Shapiro-Wilk test / Q-Q plot | p-value = 0.240 > 0.05 |
| Independence of residuals | Durbin-Watson test | p-value = 0.134 > 0.05 |
| Homoscedasticity | Breusch-Pagan test / residual plot | p-value = 0.496 > 0.05 |
| Multicollinearity | Variance Inflation Factor (VIF) | VIF = 4.27 < 10 |

# Statistical Analysis

| Predictors | P-value |
|---|---|
| WEIGHT | 1.27e-09 |
| ABDOMEN | <2e-16 |
| Overall Model(under F-test) | <2.2e-16 |
| | |
| Residual Standard Error | 3.987 |
| | |
| R^2 | 0.714 |



3D Scatter Plot with Regression Plane

# Robustness and Sensitivity

- **Robustness**

| | BODYFAT | AGE | WEIGHT | HEIGHT | NECK | CHEST | ABDOMEN | HIP | THIGH | KNEE | ANKLE | BICEPS | FOREARM | WRIST |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | \<dbl\> | \<int\> | \<dbl\> | \<dbl\> | \<dbl\> | \<dbl\> | \<dbl\> | \<dbl\> | \<dbl\> | \<dbl\> | \<dbl\> | \<dbl\> | \<dbl\> | \<dbl\> |
| 39 | 33.8 | 46 | 363.15 | 72.25 | 51.2 | 136.2 | 148.1 | 147.7 | 87.3 | 49.1 | 29.6 | 45.0 | 29.0 | 21.4 |

| | fit | lwr | upr |
|---|---|---|---|
| 39 | 45.96312 | 42.44227 | 49.48398 |

USC Units:

$BFP = 86.010 \times \log_{10}(abdomen-neck) - 70.041 \times \log_{10}(height) + 36.76$

bodyfat39= 42.58902

- **Sensitivity**

We introduced Gaussian noise with a "2.5% " standard error:
$$\widetilde{variable} \sim N(variable, 0.025 * \text{mean}(variable))$$
Under this $\widetilde{dataset}$ , the $\widetilde{rmse}$ is 4.22, only less 0.3% predicted bodyfat increase.

# Strengths and Weaknesses

## Final Model:

Body Fat ~ -41.547 - 0.124* WEIGHT(lbs) + 0.894*ABDOMEN(cm)

- **Strengths**
  - A linear model with only two variables
  - Explains ABDOMEN of variation in body fat
  - Satisfies all the fundamental assumptions of linear regression

- **Weaknesses**
  - Less intuitively when interprete the coef. of WEIGHT

# Conclusion

**Fight with mess data.**

1. Correlation
2. Histogram
3. LOF
4. Relationship between data
5. KNN
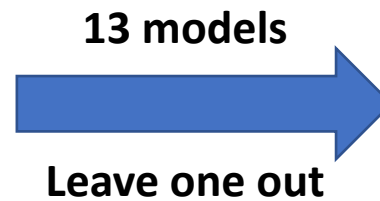6. Setting range
7. Regression
8. Outlier detection

**Checked one by one**

Is this data correct?
How it can be wrong?
Why is value is so strange?
Which one is wrong among these 3?
How to impute it?
Why can't we delete it?
......

# Conclusion

**Compared RMSE in 8 different models**

1. Forward model
2. Stepwise selection
3. Best subset selection
4. LASSO
5. Ridge
6. Elastic Net
7. Random forest
8. XGBoost

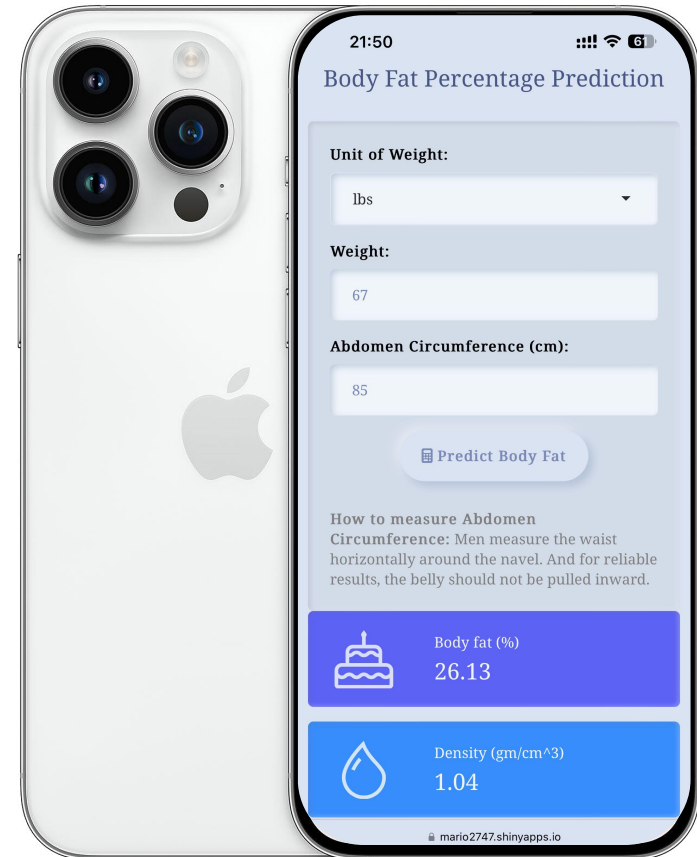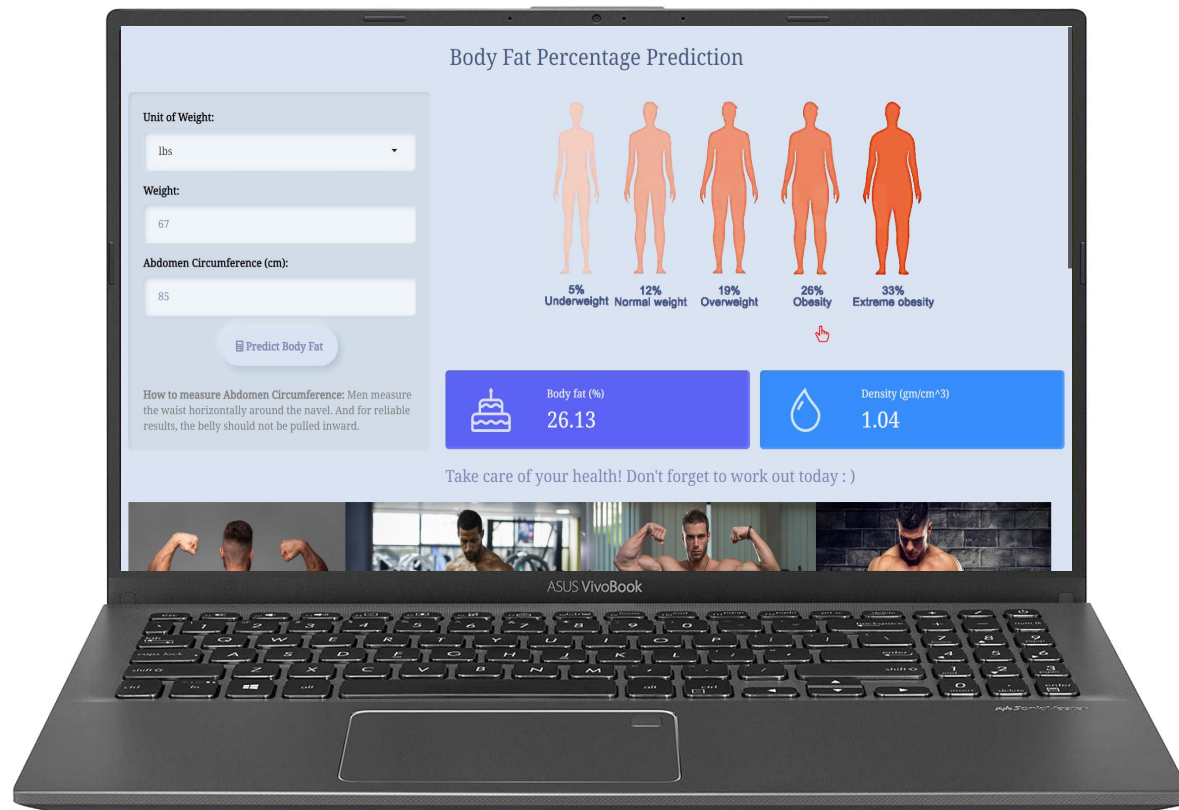**13 models**

**Leave one out**

| Abdomen | Weight |

| Knee | Forearm |

# Conclusion

**Measure your body fat at any time (<span style="color:darkred">with our model</span>) !**

# Body Fat Measurement Model

**Group 4**