

Yen-Chun Chen  
 chenyenc@oregonstate.edu

Q1.

$$\begin{cases} a_1x + b_1y = c_1 \dots \textcircled{1} \\ a_2x + b_2y = c_2 \dots \textcircled{2} \end{cases}$$

(a)

$$\textcircled{2} \times \frac{a_1}{a_2} \Rightarrow a_1x + \frac{a_1b_2}{a_2}y = \frac{a_1c_2}{a_2} \dots \textcircled{3}$$

$$\begin{aligned} \textcircled{3} - \textcircled{1} &\Rightarrow \frac{a_1b_2}{a_2}y - b_1y = \frac{a_1c_2}{a_2} - c_1 \\ &\Rightarrow (a_1b_2 - a_2b_1)y = a_1c_2 - a_2c_1 \end{aligned}$$

Given that  $a_1b_2 - a_2b_1 \neq 0 \Rightarrow y = \frac{a_1c_2 - a_2c_1}{a_1b_2 - a_2b_1} \#$

$$\textcircled{2} \times \frac{b_1}{b_2} \Rightarrow \frac{a_2b_1}{b_2}x + b_1y = \frac{b_1}{b_2}c_2 \dots \textcircled{4}$$

$$\textcircled{4} - \textcircled{1} \Rightarrow \frac{a_2b_1}{b_2}x - a_1x = \frac{b_1}{b_2}c_2 - c_1$$

$$\Rightarrow (a_2b_1 - b_2a_1)x = b_1c_2 - c_1b_2$$

$$\Rightarrow (a_2b_1 - b_2a_1)x = c_1b_2 - c_2b_1$$

Given that  $a_1b_2 - a_2b_1 \neq 0 \Rightarrow x = \frac{c_1b_2 - b_1c_2}{a_1b_2 - a_2b_1} \#$

$$(b) \quad y = \beta_1 \cdot g(x) + \beta_2 \cdot f(x) \xrightarrow{\text{Prediction}} \hat{y}_i = \hat{\beta}_1 \cdot g_i + \hat{\beta}_2 \cdot f_i$$

$$L(\beta_1, \beta_2) = \sum_{i=1}^n (y_i - \beta_1 g_i - \beta_2 f_i)^2$$

To minimize  $\beta_1, \beta_2$ , take derivative with respect to  $\beta_1, \beta_2$  equal to 0

$$\frac{\partial L}{\partial \beta_1} = \cancel{\chi} \sum_{i=1}^n (y_i - \beta_1 g_i - \beta_2 f_i) \cdot (-g_i) = 0$$

$$\rightarrow \sum_{i=1}^n \beta_1 g_i^2 + \sum_{i=1}^n \beta_2 f_i \cdot g_i = \sum_{i=1}^n y_i g_i$$

$$\frac{\partial L}{\partial \beta_2} = \cancel{\chi} \sum_{i=1}^n (y_i - \beta_1 g_i - \beta_2 f_i) \cdot (-f_i) = 0$$

$$\rightarrow \sum_{i=1}^n \beta_1 g_i f_i + \sum_{i=1}^n \beta_2 f_i^2 = \sum_{i=1}^n y_i f_i$$

$$\Rightarrow \begin{cases} (\sum g_i^2) \beta_1 + (\sum f_i g_i) \beta_2 = \sum y_i g_i & \dots (5) \\ (\sum g_i f_i) \beta_1 + (\sum f_i^2) \beta_2 = \sum y_i f_i & \dots (6) \end{cases}$$

According to (a)

$$\begin{cases} a_1 x + b_1 y = c_1 \\ a_2 x + b_2 y = c_2 \end{cases} \Rightarrow \begin{cases} x = \frac{c_1 b_2 - b_1 c_2}{a_1 b_2 - a_2 b_1} \\ y = \frac{a_1 c_2 - c_1 a_2}{a_1 b_2 - a_2 b_1} \end{cases}$$

$$\left\{ \begin{array}{l} \textcircled{5} \\ \textcircled{6} \end{array} \right. \rightarrow \left\{ \begin{array}{l} \hat{\beta}_1 = \frac{\left( \sum_{i=1}^n y_i g_i \right) \left( \sum_{i=1}^n f_i^2 \right) - \left( \sum_{i=1}^n y_i f_i \right) \left( \sum_{i=1}^n f_i g_i \right)}{\left( \sum_{i=1}^n g_i^2 \right) \left( \sum_{i=1}^n f_i^2 \right) - \left( \sum_{i=1}^n f_i g_i \right)^2} \quad (3) \\ \hat{\beta}_2 = \frac{\left( \sum_{i=1}^n y_i f_i \right) \left( \sum_{i=1}^n g_i^2 \right) - \left( \sum_{i=1}^n y_i g_i \right) \left( \sum_{i=1}^n f_i g_i \right)}{\left( \sum_{i=1}^n g_i^2 \right) \left( \sum_{i=1}^n f_i^2 \right) - \left( \sum_{i=1}^n f_i g_i \right)^2} \end{array} \right. \quad (4)$$

#

(C)

$$y = \beta f(x)$$

$$L(\beta) = \sum_{i=1}^n (y - \beta f_i)^2$$

$$0 = \frac{\partial L}{\partial \beta} = \not \cancel{\sum_{i=1}^n (y_i - \beta \cdot f_i)} (-f_i)$$

$$\rightarrow \sum_{i=1}^n \beta \cdot f_i^2 = \sum_{i=1}^n y_i \cdot f_i$$

$$\rightarrow (\sum f_i^2) \beta = \sum y_i f_i$$

$$\rightarrow \hat{\beta} = \frac{\sum_{i=1}^n y_i f_i}{\sum_{i=1}^n f_i^2}$$

#

$$(d) \quad \hat{\beta}_2 = \frac{\left(\sum_{i=1}^n y_i f_i\right)\left(\sum_{i=1}^n g_i^2\right) - \left(\sum_{i=1}^n y_i g_i\right)\left(\sum_{i=1}^n f_i g_i\right)}{\left(\sum_{i=1}^n g_i^2\right) \left(\sum_{i=1}^n f_i^2\right) - \left(\sum_{i=1}^n f_i g_i\right)^2}$$

Given that  $g_i = 0$ , the denominator will be 0 then  $\hat{\beta}_2$  will be undefined. Therefore, it's not able to derive (4) by setting  $g_i = 0$ .

$$\begin{aligned} (e) \quad E[\hat{\beta}] &= E\left[\frac{\sum_{i=1}^n y_i f_i}{\sum_{i=1}^n f_i^2}\right] = E\left[\frac{\sum_{i=1}^n f_i (\beta f_i + \varepsilon_i)}{\sum_{i=1}^n f_i^2}\right] \\ &= E\left[\frac{\cancel{\beta} \sum_{i=1}^n f_i^2}{\sum_{i=1}^n f_i^2}\right] + E\left[\frac{\sum_{i=1}^n \varepsilon_i}{\sum_{i=1}^n f_i^2}\right] \\ &= \beta \frac{n \cdot \bar{f}_i^2}{n \cdot \bar{f}_i^2} + \frac{\sum_{i=1}^n E[\varepsilon_i]}{n \cdot \bar{f}_i^2} \end{aligned}$$

Given  $M_e = E[\varepsilon_i] = 0$

$$\rightarrow E[\hat{\beta}] = \beta \Rightarrow \hat{\beta} \text{ is unbiased } \#$$

$$(f) \quad y_i = \beta x_i^2 + \varepsilon_i \quad , \quad \varepsilon_i \sim N(0, \sigma^2)$$

$$\text{Var}(y_i) = \text{Var}(\beta x_i^2 + \varepsilon_i)$$

$$= \text{Var}(\beta x_i^2) + \text{Var}(\varepsilon_i)$$

$$= 0 + \text{Var}(\varepsilon_i) = \sigma^2$$

$$\text{Var}(\hat{\beta}) = \text{Var}\left(\frac{\sum_{i=1}^n y_i f_i}{\sum_{i=1}^n f_i^2}\right) = \text{Var}\left(\frac{\sum (\beta x_i^2 + \varepsilon_i)(x_i^2)}{\sum x_i^4}\right)$$

$$= \text{Var}\left(\frac{\sum \beta x_i^4}{\sum x_i^4}\right) + \text{Var}\left(\frac{\sum x_i^2 \varepsilon_i}{\sum x_i^4}\right)$$

$$= 0 + \sum \frac{\sum (x_i^2)^2 \text{Var}(\varepsilon_i)}{(\sum x_i^4)^2}$$

$$= \frac{\sum x_i^4}{(\sum x_i^4)^2} \sigma^2 = \frac{\sigma^2}{\sum_{i=1}^n x_i^4} \quad \#$$

$$(g) \quad y = \beta x^2$$

$$\text{From (f)} : \text{var}(\hat{\beta}) = \frac{\sigma^2}{\sum_{i=1}^n x_i^4}$$

When giving larger  $x$ ,  $\sum_{i=1}^n x_i^4$  goes larger,

which decreases the value of  $\text{var}(\hat{\beta})$ .

Therefore, Yes, Larger  $x$  can reduce the uncertainty  
in evaluating  $\hat{\beta}$  #

## Q2

April 17, 2024

The goal of this question is to calculate body fat using Brozek's equation. In this problem, the first column **brozek** is the response variable, and all other columns are treated as the features (we may not use all the features, please only use the ones stated in each question). To start, read the data file fat.csv in the homework folder (containing 252 samples) and split it into two sets. Set 1 includes the first 200 rows of the data (do not count the row associated with the feature/response names), and set 2, which includes the last 52 rows of the data. Name the first set **train** and the second set **test**.

```
[1]: import numpy as np
import pandas as pd
from matplotlib.pyplot import subplots
import statsmodels.api as sm
from ISLP import load_data
from ISLP.models import (ModelSpec as MS,
                         summarize)
```

```
[2]: Fat = pd.read_csv("fat.csv")
train,test = np.split(Fat,[int(200)])
y = train['brozek']
X = pd.DataFrame({'intercept': np.ones(train.shape[0]),
                  'siri': train['siri'],
                  'density': train['density'],
                  'age': train['age'],
                  'weight': train['weight'],
                  'height': train['height'],
                  'adipos': train['adipos'],
                  'free': train['free'],
                  'neck': train['neck'],
                  'chest': train['chest'],
                  'abdom': train['abdom'],
                  'hip': train['hip'],
                  'thigh': train['thigh'],
                  'knee': train['knee'],
                  'ankle': train['ankle'],
                  'biceps': train['biceps'],
                  'forearm': train['forearm'],
                  'wrist': train['wrist']})
```

(a) As a first modeling attempt, consider a linear model using all the 17 features, that is  $\text{brozek} = \beta_0 + \beta_1 \text{siri} + \dots + \beta_{17} \text{wrist}$  report the fitted parameters, the 95% confidence interval for each estimated parameter and the p-values. What is the R2 value, and based on that how good do you see the model fitting the training data?

```
[3]: firstModel = sm.OLS(y, X)
firstFittedModel = firstModel.fit()
firstFittedModel.summary()
```

[3]:

<b>Dep. Variable:</b>	brozek	<b>R-squared:</b>	0.999			
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.999			
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	1.859e+04			
<b>Date:</b>	Wed, 17 Apr 2024	<b>Prob (F-statistic):</b>	7.25e-285			
<b>Time:</b>	20:53:11	<b>Log-Likelihood:</b>	61.275			
<b>No. Observations:</b>	200	<b>AIC:</b>	-86.55			
<b>Df Residuals:</b>	182	<b>BIC:</b>	-27.18			
<b>Df Model:</b>	17					
<b>Covariance Type:</b>	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
<b>intercept</b>	11.7984	4.641	2.542	0.012	2.641	20.956
<b>siri</b>	0.8845	0.013	66.439	0.000	0.858	0.911
<b>density</b>	-9.5175	4.172	-2.281	0.024	-17.750	-1.285
<b>age</b>	-0.0007	0.002	-0.417	0.677	-0.004	0.003
<b>weight</b>	0.0116	0.005	2.429	0.016	0.002	0.021
<b>height</b>	0.0004	0.005	0.075	0.940	-0.009	0.010
<b>adipos</b>	-0.0213	0.016	-1.361	0.175	-0.052	0.010
<b>free</b>	-0.0134	0.006	-2.324	0.021	-0.025	-0.002
<b>neck</b>	-0.0037	0.011	-0.325	0.746	-0.026	0.019
<b>chest</b>	0.0034	0.005	0.631	0.529	-0.007	0.014
<b>abdom</b>	0.0005	0.005	0.090	0.929	-0.010	0.011
<b>hip</b>	-0.0037	0.007	-0.496	0.620	-0.018	0.011
<b>thigh</b>	0.0199	0.008	2.579	0.011	0.005	0.035
<b>knee</b>	-0.0305	0.013	-2.386	0.018	-0.056	-0.005
<b>ankle</b>	0.0037	0.010	0.359	0.720	-0.017	0.024
<b>biceps</b>	-0.0159	0.009	-1.871	0.063	-0.033	0.001
<b>forearm</b>	0.0196	0.010	1.891	0.060	-0.001	0.040
<b>wrist</b>	0.0340	0.027	1.246	0.214	-0.020	0.088
<b>Omnibus:</b>	146.785	<b>Durbin-Watson:</b>	1.875			
<b>Prob(Omnibus):</b>	0.000	<b>Jarque-Bera (JB):</b>	8218.452			
<b>Skew:</b>	2.048	<b>Prob(JB):</b>	0.00			
<b>Kurtosis:</b>	34.136	<b>Cond. No.</b>	1.48e+05			

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

[2] The condition number is large, 1.48e+05. This might indicate that there are strong multicollinearity or other numerical problems.

(b) Based on  $\alpha = 0.05$  and the calculated p-values, 10 features seem problematic.

feature	p-value
age	0.677
height	0.940
adipos	0.175
neck	0.746
chest	0.529
abdom	0.929
hip	0.620
ankle	0.720
biceps	0.063
wrist	0.214

```
[4]: X_test = pd.DataFrame({'intercept': np.ones(test.shape[0]),
                           'siri': test['siri'],
                           'density': test['density'],
                           'age': test['age'],
                           'weight': test['weight'],
                           'height': test['height'],
                           'adipos': test['adipos'],
                           'free': test['free'],
                           'neck': test['neck'],
                           'chest': test['chest'],
                           'abdom': test['abdom'],
                           'hip': test['hip'],
                           'thigh': test['thigh'],
                           'knee': test['knee'],
                           'ankle': test['ankle'],
                           'biceps': test['biceps'],
                           'forearm': test['forearm'],
                           'wrist': test['wrist']})
```

(c) Calculate the prediction error  $e_1$  using your test file.

```
[5]: # the actual value from test dataset
y_actual = test['brozek'].values

# the value predicted by the first model
y_pred_first = firstFittedModel.predict(X_test)

y_actual = np.array(y_actual)
y_pred_first = np.array(y_pred_first)

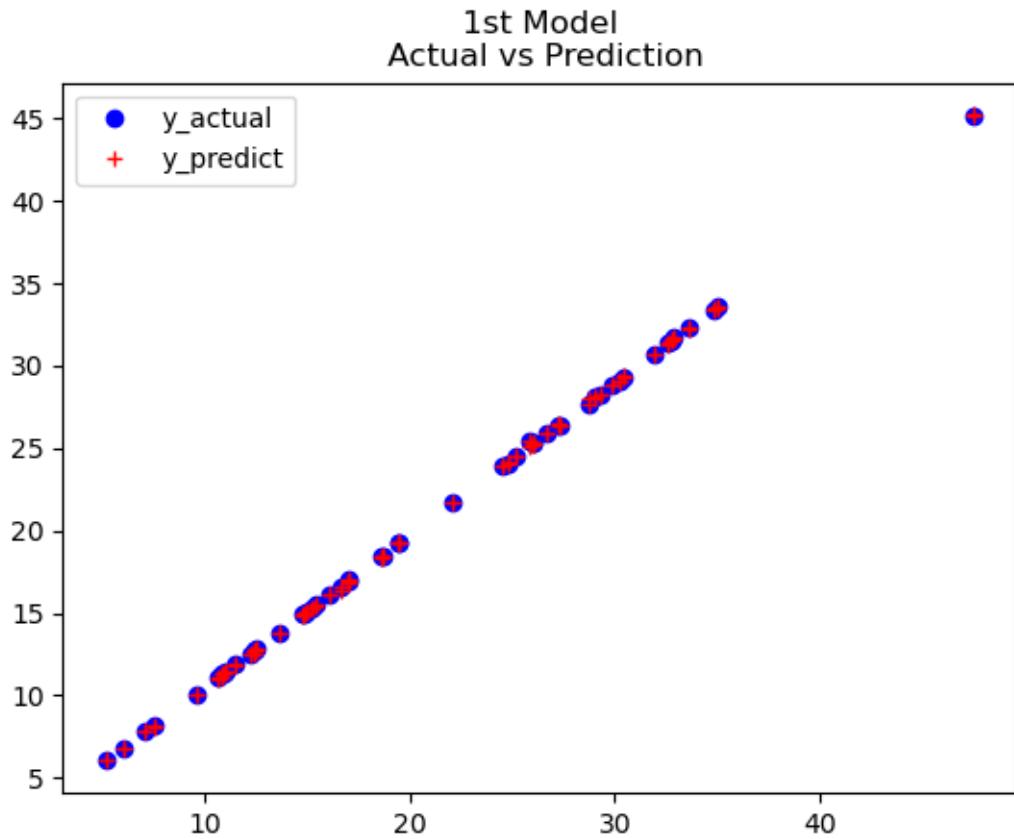
e_1 = np.sqrt(sum(np.square(y_actual - y_pred_first)))
print('Prediction Error e_1 = ', e_1)
```

Prediction Error e\_1 = 0.753742369836337

```
[6]: import matplotlib.pyplot as plt

fig, ax = plt.subplots()
ax.plot(test['siri'], y_actual, "bo", label="y_actual")
ax.plot(np.hstack((test['siri'], X_test['siri'])), np.hstack((y_actual, y_pred_first)), "r+", label="y_predict")
plt.title('1st Model\n Actual vs Prediction')
ax.legend(loc="best")
```

[6]: <matplotlib.legend.Legend at 0x149983c90>



(d) Now consider a second model which uses the features indicated in part (a), with the only difference that density is replaced by inverse density. In other words, your model has the the following parametric form:  

$$\text{brozek} = \beta_0 + \beta_1 \text{siri} + \frac{\beta_2}{\text{density}} + \beta_3 \text{age} + \dots + \beta_{17} \text{wrist}$$
Repeat the steps in part (a) and report the values.

```
[7]: X_second = pd.DataFrame({'intercept': np.ones(train.shape[0]),
                             'siri': train['siri'],
                             'density': 1/train['density'],
                             'age': train['age'],
```

```

'weight': train['weight'],
'height': train['height'],
'adipos': train['adipos'],
'free': train['free'],
'neck': train['neck'],
'chest': train['chest'],
'abdom': train['abdom'],
'hip': train['hip'],
'thigh': train['thigh'],
'knee': train['knee'],
'ankle': train['ankle'],
'biceps': train['biceps'],
'forearm': train['forearm'],
'wrist': train['wrist']})
secondFittedModel = sm.OLS(y, X_second).fit()
secondFittedModel.summary()

```

[7]:

Dep. Variable:	brozek	R-squared:	0.999			
Model:	OLS	Adj. R-squared:	0.999			
Method:	Least Squares	F-statistic:	1.852e+04			
Date:	Wed, 17 Apr 2024	Prob (F-statistic):	1.00e-284			
Time:	20:53:12	Log-Likelihood:	60.917			
No. Observations:	200	AIC:	-85.83			
Df Residuals:	182	BIC:	-26.46			
Df Model:	17					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
intercept	-8.0827	4.536	-1.782	0.076	-17.032	0.867
siri	0.8865	0.013	67.497	0.000	0.861	0.912
density	10.3645	4.867	2.129	0.035	0.761	19.968
age	-0.0007	0.002	-0.423	0.673	-0.004	0.003
weight	0.0109	0.005	2.287	0.023	0.001	0.020
height	0.0005	0.005	0.094	0.925	-0.009	0.010
adipos	-0.0206	0.016	-1.320	0.189	-0.051	0.010
free	-0.0125	0.006	-2.161	0.032	-0.024	-0.001
neck	-0.0038	0.011	-0.332	0.740	-0.026	0.019
chest	0.0031	0.005	0.583	0.561	-0.008	0.014
abdom	0.0005	0.005	0.098	0.922	-0.010	0.011
hip	-0.0039	0.007	-0.529	0.598	-0.019	0.011
thigh	0.0199	0.008	2.580	0.011	0.005	0.035
knee	-0.0305	0.013	-2.386	0.018	-0.056	-0.005
ankle	0.0036	0.010	0.352	0.725	-0.017	0.024
biceps	-0.0158	0.009	-1.850	0.066	-0.033	0.001
forearm	0.0193	0.010	1.858	0.065	-0.001	0.040
wrist	0.0339	0.027	1.239	0.217	-0.020	0.088

<b>Omnibus:</b>	144.182	<b>Durbin-Watson:</b>	1.875
<b>Prob(Omnibus):</b>	0.000	<b>Jarque-Bera (JB):</b>	8277.916
<b>Skew:</b>	1.977	<b>Prob(JB):</b>	0.00
<b>Kurtosis:</b>	34.268	<b>Cond. No.</b>	1.58e+05

Notes:

[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

[2] The condition number is large, 1.58e+05. This might indicate that there are strong multicollinearity or other numerical problems.

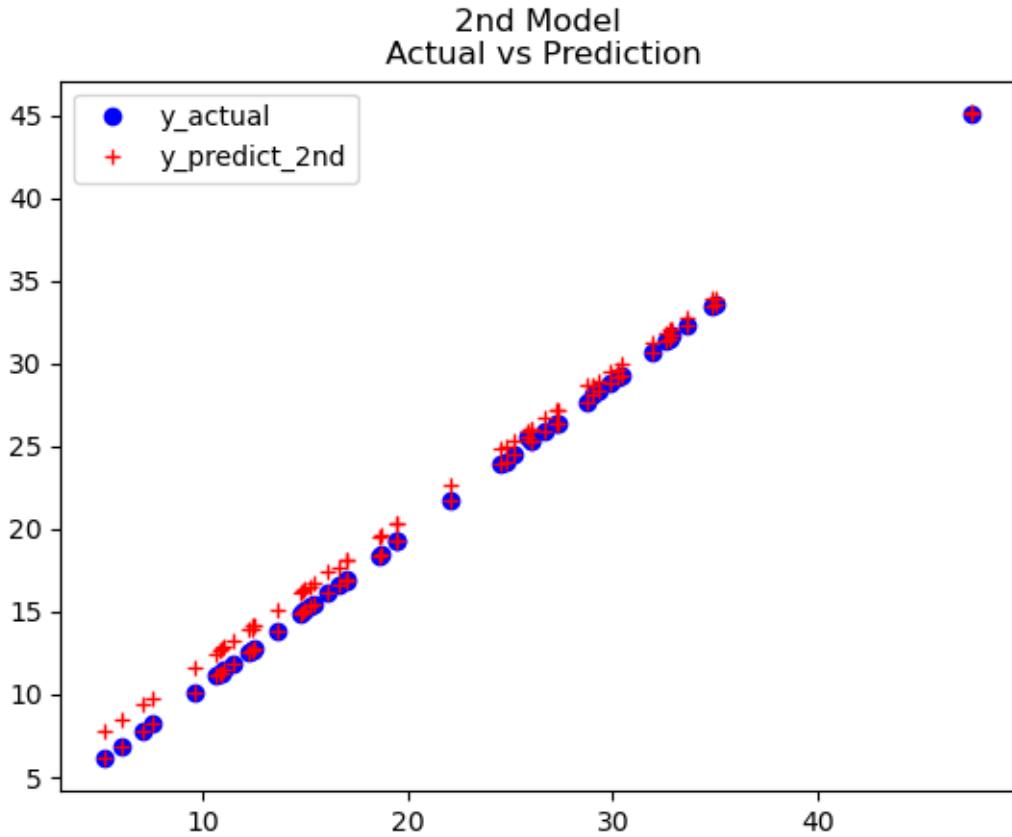
(e) Use a similar formulation in part (c) and calculate the test error (we will call this test error  $e_2$  which corresponds to our second model).

```
[8]: # the value predicted by the second model
y_pred_second = np.array(secondFittedModel.predict(X_test))
e_2 = np.sqrt(sum(np.square(y_actual - y_pred_second)))
print('Second Prediction Error e_2 =', e_2)
```

Second Prediction Error e\_2 = 7.766138616196848

```
[9]: fig, ax = plt.subplots()
ax.plot(test['siri'], y_actual, "bo", label="y_actual")
ax.plot(np.hstack((test['siri'], X_test['siri'])), np.hstack((y_actual, y_pred_second)), "r+", label="y_predict_2nd")
plt.title('2nd Model\n Actual vs Prediction')
ax.legend(loc="best")
```

[9]: <matplotlib.legend.Legend at 0x149a0d410>



(f) Now consider a third model which only uses siri and density as the features, but follows a para-metric formulation as: brozek =  $\beta_0 + \beta_1 \text{siri} + \beta_2 \text{siri}^2 + \frac{\beta_3}{\text{density}} + \beta_4$  density Repeat the steps in part (a) and report the values.

```
[10]: X_third = pd.DataFrame({'intercept': np.ones(train.shape[0]),
                            'siri': train['siri'],
                            'siri_square': train['siri'] ** 2,
                            'inv_density': 1/train['density'],
                            'density': train['density']})

thirdFittedModel = sm.OLS(y, X_third).fit()
thirdFittedModel.summary()
```

[10]:

<b>Dep. Variable:</b>	brozek	<b>R-squared:</b>	0.999			
<b>Model:</b>	OLS	<b>Adj. R-squared:</b>	0.999			
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	7.648e+04			
<b>Date:</b>	Wed, 17 Apr 2024	<b>Prob (F-statistic):</b>	2.51e-310			
<b>Time:</b>	20:53:12	<b>Log-Likelihood:</b>	51.125			
<b>No. Observations:</b>	200	<b>AIC:</b>	-92.25			
<b>Df Residuals:</b>	195	<b>BIC:</b>	-75.76			
<b>Df Model:</b>	4					
<b>Covariance Type:</b>	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
intercept	-1298.6590	567.083	-2.290	0.023	-2417.062	-180.256
siri	0.9618	0.027	35.651	0.000	0.909	1.015
siri_square	-0.0030	0.001	-2.275	0.024	-0.006	-0.000
inv_density	706.0142	305.636	2.310	0.022	103.237	1308.791
density	598.1825	262.950	2.275	0.024	79.592	1116.773
<b>Omnibus:</b>	138.737	<b>Durbin-Watson:</b>	1.955			
<b>Prob(Omnibus):</b>	0.000	<b>Jarque-Bera (JB):</b>	12413.875			
<b>Skew:</b>	1.696	<b>Prob(JB):</b>	0.00			
<b>Kurtosis:</b>	41.447	<b>Cond. No.</b>	2.68e+07			

Notes:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The condition number is large, 2.68e+07. This might indicate that there are strong multicollinearity or other numerical problems.

(g) Based on  $\alpha = 0.05$  and the calculated p-values, which features seem problematic?

All the p-values of the features are less than 0.05

feature	p-value
intercept	0.023
siri	0.00
$siri^2$	0.024
$\frac{1}{density}$	0.022
density	0.024

(h) Repeat part (c) for this model and call the error  $e_3$ .

```
[11]: X_test_third = pd.DataFrame({'intercept': np.ones(test.shape[0]),
                                    'siri': test['siri'],
                                    'siri_square': test['siri'] ** 2,
                                    'inv_density': 1/test['density'],
                                    'density': test['density']})

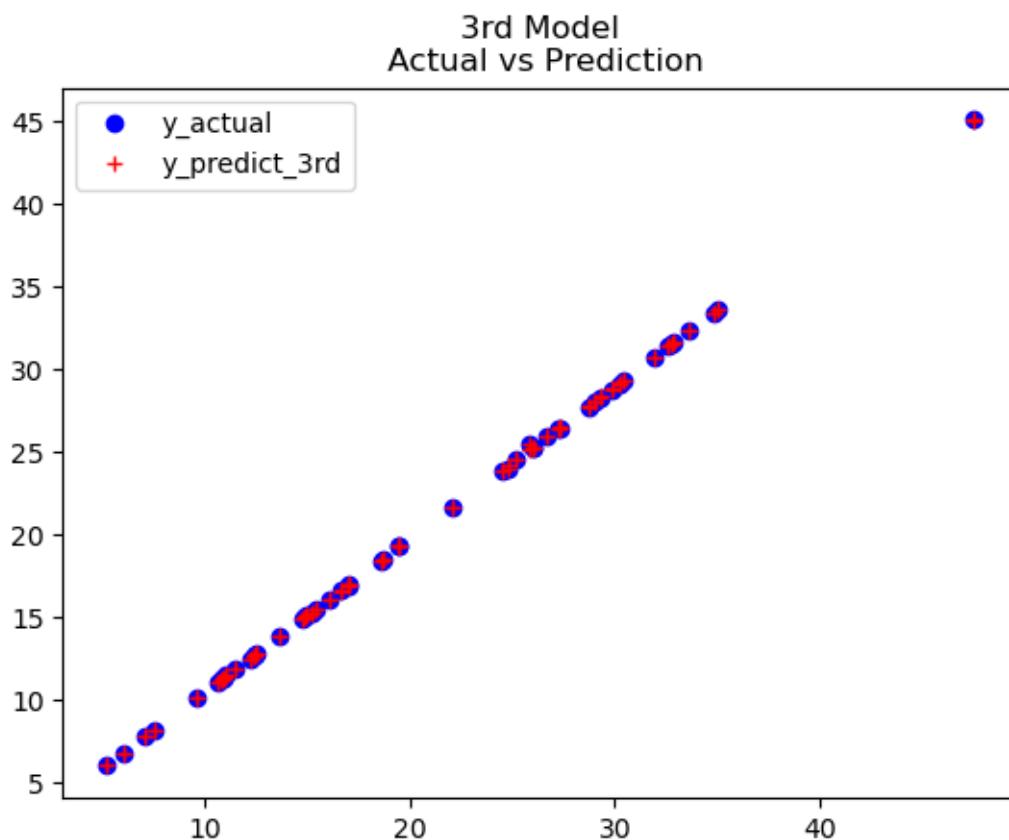
# the value predicted by the third model
y_pred_third = np.array(thirdFittedModel.predict(X_test_third))
```

```
e_3 = np.sqrt(sum(np.square(y_actual - y_pred_third)))
print('Third Prediction Error e_3 =', e_3)
```

Third Prediction Error e\_3 = 0.563361670359713

```
[12]: fig, ax = plt.subplots()
ax.plot(test['siri'], y_actual, "bo", label="y_actual")
ax.plot(np.hstack((test['siri'], X_test_third['siri'])), np.hstack((y_actual, y_pred_third)), "r+", label="y_predict_3rd")
plt.title('3rd Model\n Actual vs Prediction')
ax.legend(loc="best")
```

[12]: <matplotlib.legend.Legend at 0x1486bba50>



- (i) Based on the values  $e_1$ ,  $e_2$  and  $e_3$ , and the model formulations, which model would you pick and why (state two reasons)?

I will pick the third model. (1) It has the lowest prediction error  $0.563361670359713$  so it has the highest prediction accuracy.(2) The R-squared values of the three models are the same but the thrid model has the lowest p-values  $2.51e-310$ .

## Q3

April 17, 2024

In the previous question we observed some redundancy in the features. We would like to try some feature selection heuristic in this question. Consider the same dataset as question 2 (fat.csv), where **brozek** is the response variable and the other 17 columns are the model features. Follow these steps below.

- Form an extended version of the dataset, by appending two more columns. One column corresponding to  $siri^2$  and one column corresponding to  $\frac{1}{density}$ . Your extended dataset should now have 20 columns, where the first column is brozek and used as the response variable, 17 columns identical to the original fat.csv data set, and columns 19 and 20 with the values  $siri^2$  and  $\frac{1}{density}$ , respectively. We will refer to this dataset as the *extended dataset*.
- In a similar way as question 2, split the extended dataset into two sets. Set 1 includes the first 200 rows of the data (do not count the row associated with the feature/response names), and set 2, which includes the last 52 rows of the data. Name the first set **train** and the second set **test**.

```
[34]: import numpy as np
import pandas as pd
from matplotlib.pyplot import subplots
import statsmodels.api as sm
from ISLP import load_data
from ISLP.models import (ModelSpec as MS,
                         summarize)
```

```
[35]: Fat = pd.read_csv("fat.csv")
train,test = np.split(Fat,[int(200)])
y = train['brozek']
X_train_full = pd.DataFrame({'intercept': np.ones(train.shape[0]),
                             'siri': train['siri'],
                             'density': train['density'],
                             'age': train['age'],
                             'weight': train['weight'],
                             'height': train['height'],
                             'adipos': train['adipos'],
                             'free': train['free'],
                             'neck': train['neck'],
                             'chest': train['chest'],
                             'abdom': train['abdom'],
                             'hip': train['hip'],
                             'thigh': train['thigh'],
```

```

'knee':           train['knee'],
'ankle':          train['ankle'],
'biceps':         train['biceps'],
'forearm':        train['forearm'],
'wrists':         train['wrists'],
'siri_squared':   train['siri'] ** 2,
'inv_density':    1/train['density'])

```

(a) Use the training data to fit a model of the following form  $\text{brozek} = \beta_0 + \beta_1 \text{siri} + \dots + \beta_{17} \text{wrists} + \beta_{18} \text{siri}^2 + \beta_{19} \frac{1}{\text{density}}$  report the fitted parameters, the 95% confidence interval for each estimated parameter and the p-values. What is the R2 value?

[36]: fullFittedModel = sm.OLS(y, X\_train\_full).fit()  
fullFittedModel.summary()

[36]:

Dep. Variable:	brozek	R-squared:	0.999			
Model:	OLS	Adj. R-squared:	0.999			
Method:	Least Squares	F-statistic:	1.703e+04			
Date:	Wed, 17 Apr 2024	Prob (F-statistic):	6.09e-282			
Time:	20:53:25	Log-Likelihood:	64.717			
No. Observations:	200	AIC:	-89.43			
Df Residuals:	180	BIC:	-23.47			
Df Model:	19					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
intercept	-950.8057	580.610	-1.638	0.103	-2096.483	194.871
siri	0.9305	0.029	32.063	0.000	0.873	0.988
density	436.7337	269.081	1.623	0.106	-94.225	967.693
age	-0.0007	0.002	-0.414	0.680	-0.004	0.003
weight	0.0162	0.006	2.809	0.006	0.005	0.028
height	-0.0005	0.005	-0.096	0.924	-0.010	0.009
adipos	-0.0235	0.016	-1.505	0.134	-0.054	0.007
free	-0.0204	0.007	-2.780	0.006	-0.035	-0.006
neck	-0.0018	0.011	-0.163	0.870	-0.024	0.020
chest	0.0056	0.005	1.026	0.306	-0.005	0.016
abdom	-0.0006	0.005	-0.109	0.913	-0.011	0.010
hip	0.0008	0.008	0.106	0.916	-0.014	0.016
thigh	0.0174	0.008	2.264	0.025	0.002	0.033
knee	-0.0290	0.013	-2.296	0.023	-0.054	-0.004
ankle	0.0061	0.010	0.591	0.555	-0.014	0.026
biceps	-0.0169	0.008	-2.011	0.046	-0.034	-0.000
forearm	0.0219	0.010	2.110	0.036	0.001	0.042
wrists	0.0343	0.027	1.270	0.206	-0.019	0.088
siri_squared	-0.0026	0.001	-1.956	0.052	-0.005	2.26e-05
inv_density	518.6225	312.993	1.657	0.099	-98.985	1136.230

<b>Omnibus:</b>	98.592	<b>Durbin-Watson:</b>	1.907
<b>Prob(Omnibus):</b>	0.000	<b>Jarque-Bera (JB):</b>	6766.061
<b>Skew:</b>	0.902	<b>Prob(JB):</b>	0.00
<b>Kurtosis:</b>	31.437	<b>Cond. No.</b>	3.19e+07

Notes:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The condition number is large, 3.19e+07. This might indicate that there are strong multicollinearity or other numerical problems.

```
[37]: X_test_full = pd.DataFrame({'intercept': np.ones(test.shape[0]),
                                'siri': test['siri'],
                                'density': test['density'],
                                'age': test['age'],
                                'weight': test['weight'],
                                'height': test['height'],
                                'adipos': test['adipos'],
                                'free': test['free'],
                                'neck': test['neck'],
                                'chest': test['chest'],
                                'abdom': test['abdom'],
                                'hip': test['hip'],
                                'thigh': test['thigh'],
                                'knee': test['knee'],
                                'ankle': test['ankle'],
                                'biceps': test['biceps'],
                                'forearm': test['forearm'],
                                'wrist': test['wrist'],
                                'siri_squared': test['siri'] ** 2,
                                'inv_density': 1/test['density']})
```

(b) Use the test data to calculate the test error (similar to the formulation in part (c) of the previous question), and call it  $e_{full}$ .

```
[38]: # Extract the problematic features that has p-value over 0.03
def get_prediction_error(model, X_pred):
    # the actual value from test dataset
    y_actual = np.array(test['brozek'].values)
    # the value predicted by the full-features model
    y_prediction = np.array(model.predict(X_pred))
    return np.sqrt(sum(np.square(y_actual - y_prediction)))
```

```
[39]: # Calculate the prediction error of full-featured model
print('Prediction Error e_full =', get_prediction_error(fullFittedModel, ↵X_test_full))
```

Prediction Error e\_full = 0.8565466791992765

(c) Let's run a heuristic scheme to perform feature selection (the method is called backward selection and described on page 79 of your textbook, also on the slides). Start with the full model (the model containing all 19 features of the extended dataset) and drop the feature with the highest p-value (or the second largest if the largest p-value is for the intercept), then redo the modeling and drop the next feature with the highest p-value, and continue dropping until all p-values are small and you are left with a set of important features. Implement this approach and stop when all p-values are below 0.03. Which features are selected as the most important ones when your code stops?

```
[40]: # Extract the problematic features that has p-value greater than 0.03, among
      ↪all the p-values.

def get_problematic_pValues(model):
    pv = model.pvalues
    return pv[pv > 0.03]
```

```
[41]: # Extract the problematic features
pv_problematic = get_problematic_pValues(fullFittedModel)
# print(pv_problematic, "\n")

# Training set for iteration
X_train_new = X_train_full

# Run the loop if there is any p-value greater than 0.03
while (pv_problematic > 0.03).any():

    # Find the most problematic feature
    fea_maxP = pv_problematic.idxmax()
    # print(fea_maxP, ": ", pv_problematic[fea_maxP])

    # Drop that feature and get new training data
    X_train_new = X_train_new.drop(fea_maxP, axis = 1)

    # Retrain the model with new data
    newFittedModel = sm.OLS(y, X_train_new).fit()

    # Get new p-values from new model
    pv_problematic = get_problematic_pValues(newFittedModel)

newFittedModel.summary()
```

[41]:

<b>Dep. Variable:</b>	brozek	<b>R-squared (uncentered):</b>	1.000			
<b>Model:</b>	OLS	<b>Adj. R-squared (uncentered):</b>	1.000			
<b>Method:</b>	Least Squares	<b>F-statistic:</b>	5.569e+05			
<b>Date:</b>	Wed, 17 Apr 2024	<b>Prob (F-statistic):</b>	0.00			
<b>Time:</b>	20:53:25	<b>Log-Likelihood:</b>	51.369			
<b>No. Observations:</b>	200	<b>AIC:</b>	-94.74			
<b>Df Residuals:</b>	196	<b>BIC:</b>	-81.55			
<b>Df Model:</b>	4					
<b>Covariance Type:</b>	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
siri	0.9222	0.002	474.600	0.000	0.918	0.926
thigh	0.0142	0.005	3.101	0.002	0.005	0.023
knee	-0.0262	0.010	-2.596	0.010	-0.046	-0.006
inv_density	1.5192	0.267	5.698	0.000	0.993	2.045
<b>Omnibus:</b>	187.007	<b>Durbin-Watson:</b>	1.919			
<b>Prob(Omnibus):</b>	0.000	<b>Jarque-Bera (JB):</b>	13874.811			
<b>Skew:</b>	2.978	<b>Prob(JB):</b>	0.00			
<b>Kurtosis:</b>	43.367	<b>Cond. No.</b>	1.47e+03			

Notes:

- [1]  $R^2$  is computed without centering (uncentered) since the model does not contain a constant.
- [2] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [3] The condition number is large, 1.47e+03. This might indicate that there are strong multicollinearity or other numerical problems.

*siri, thigh, knee,  $\frac{1}{\text{density}}$*  are the most important features after backward selection.

(d) Apply the model developed in part (c) to the test data and call the error  $e_{sel}$ .

```
[42]: # Feature variables from test set
X_test_sel = pd.DataFrame({'siri': test['siri'], 'thigh': test['thigh'], 'knee': test['knee'],
                           'inv_density': 1/test['density']})

# Calculate
print('Prediction Error e_sel =', get_prediction_error(newFittedModel, X_test_sel))
```

Prediction Error e\_sel = 0.6670935385386173

(e) Compare  $e_{full}$  and  $e_{sel}$ . Does the feature selection scheme seem to reduce overfitting?

Yes, it does reduce overfitting. By removing less important features, backward selection helps focus the model on the most relevant predictors. This reduces the complexity of the model and decrease the chance of fitting too close to training data.

(f) Compare  $e_{sel}$  with  $e_3$  from part (h) of question 2. In terms of the test accuracy does your feature selection scheme seem to find the best model?

$e_{sel}$  ( $0.6670935385386173$ ) is greater than  $e_3$  ( $0.563361670359713$ ), therefore, the feature selection scheme seems not to get the best model in this case. The reason might be some useful features that seemed not important were removed too early. This can lead to a loss of valuable information that could be useful in conjunction with other features.