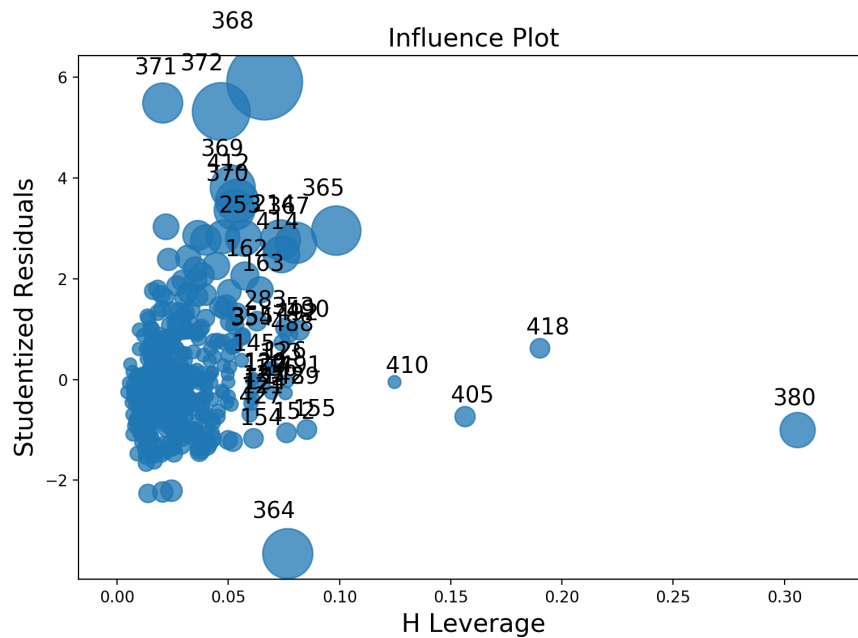***CS 498 HW 6*** Yu Che Wang/ yuchecw2

1. All the points (row number) I removed (indexed on the original dataset) as outlier points by looking at the influence plot:
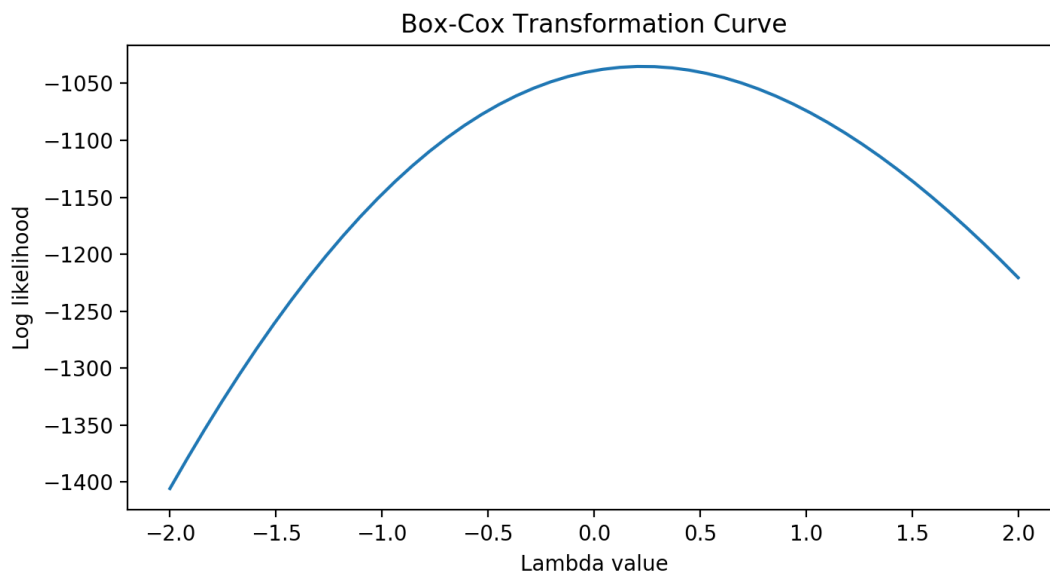   364, 365, 368, 370, 371, 372, 380, 405, 410, 418.
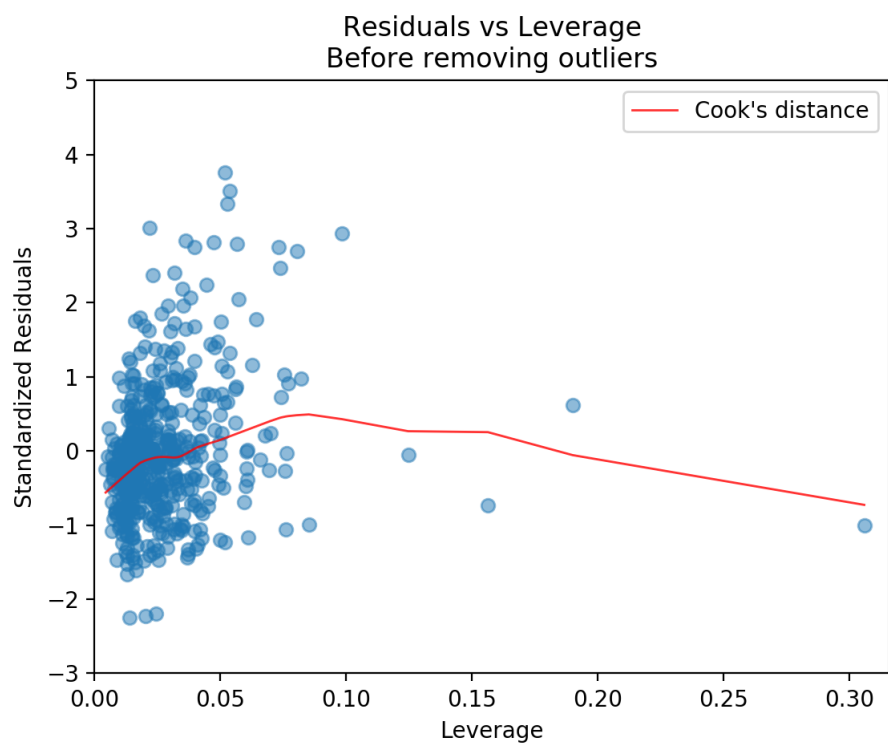   I remove those points with high leverage, standardized residuals, and cook distance. The radius of the data point represents its cook distance.



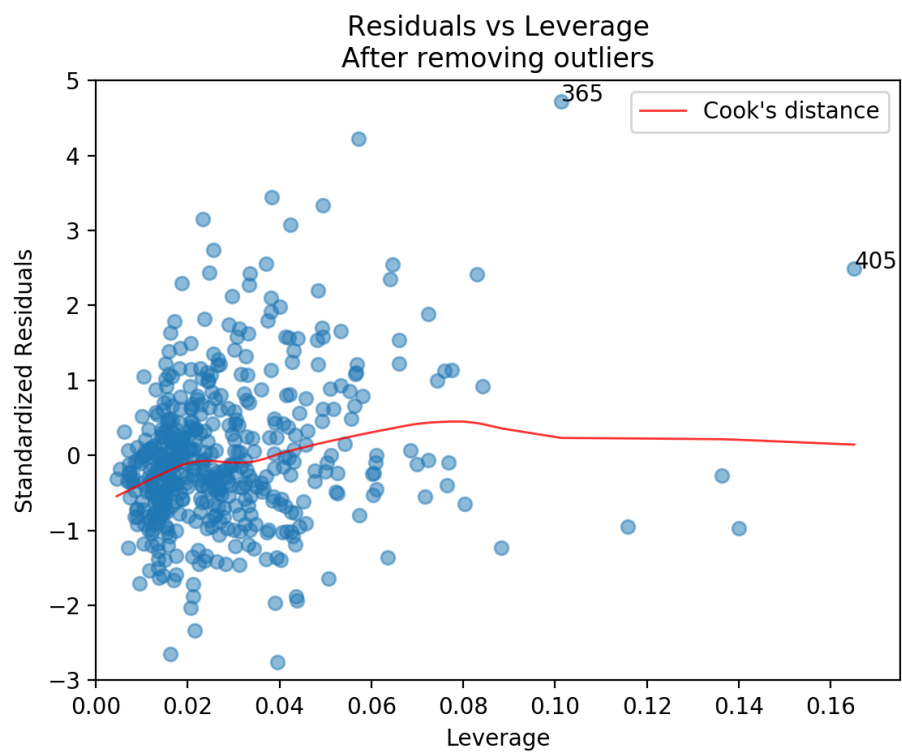2. By using scipy.stats.boxcox, we obtain the best lambda value = 0.219.
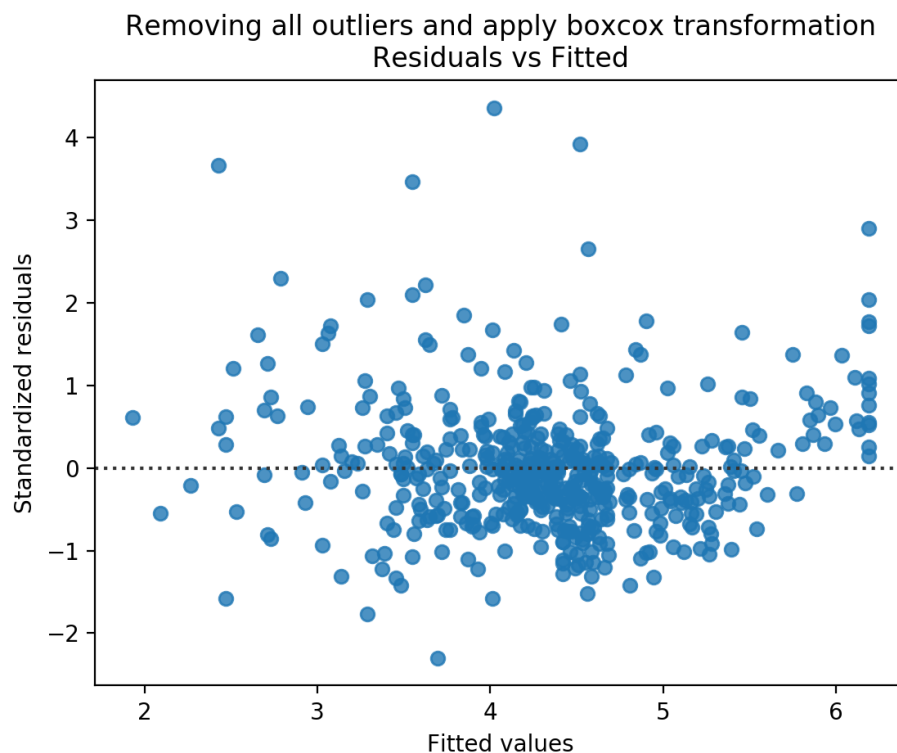
3. Diagnostic plot before removing outliers

**Residuals vs Leverage**
**Before removing outliers**



Diagnostic plot after removing outliers

**Residuals vs Leverage**
**After removing outliers**

4. Raw data



Without any transformation
Residuals vs Fitted

5. Removing all outliers and transforming the dependent variable



Removing all outliers and apply boxcox transformation
Residuals vs Fitted

6. In the first plot, there are some outliers, data that are 4-6 standardized residuals away from the mean. Most of the data points in the second plot are less than two standardized residuals away from the mean. In addition, the data points are evenly distributed across the dotted line.

7. Fitted house price vs True house price



After removing outliers and boxcox transformation

The data points lie on a line, which means that after removing outliers and applying box-cox transformation, the residuals are nearly zero.

8. Code fragments

Linear regression

```
15   # Linear regression
16   X = data[:,:13]
17   y = data[:,13].reshape(-1,1)
18   # Method 1
19   lm = sm.OLS(y, sm.add_constant(X)).fit()
20   # Method 2
21   reg = LinearRegression().fit(X, y)
22   r_squared = reg.score(X,y)
```

Box-cox transformation

```
132   # Boxcox transformation
133   y_boxcox, maxlog = stats.boxcox(y_remove)
```

Box-cox transformation with different values

```
96    n_interval = 50
97    lmbda_list = np.linspace(-2, 2, n_interval)
98    llf_list = np.zeros(n_interval)
99    for i in range(lmbda_list.shape[0]):
100       llf_list[i] = stats.boxcox_llf(lmbda_list[i], y_remove)
```

To apply box-cox transformation with a specific value lmbda, one can use

$$y\_boxcox = stats.boxcox(y, lmbda=lmbda)$$

## 9. Entire code

```python
import numpy as np
import numpy.linalg as la
import pandas as pd
import seaborn as sns
import scipy.stats as stats
import matplotlib.pyplot as plt
from sklearn.linear_model import LinearRegression
from statsmodels.stats.outliers_influence import OLSInfluence
import statsmodels.api as sm
import statsmodels.formula.api as smf
# Load the data
data = np.loadtxt('data.txt')
df = pd.DataFrame(data)

# Linear regression
X = data[:,:13]
y = data[:,13].reshape(-1,1)
# Method 1
lm = sm.OLS(y, sm.add_constant(X)).fit()
# Method 2
reg = LinearRegression().fit(X, y)
r_squared = reg.score(X,y)

# Generate influence plot
fig, ax = plt.subplots(figsize=(9,6))
fig = sm.graphics.influence_plot(lm, alpha=0.001, ax=ax, criterion='cooks')
fig.show()

####################
# Calculate leverage
####################
beta_hat = np.dot(la.inv(np.dot(X.T, X)), np.dot(X.T, y))
hat_matrix = X@la.inv(np.dot(X.T, X))@(X.T)
# plt.bar(list(range(hat_matrix.shape[0])),np.diag(hat_matrix))
# plt.boxplot(np.diag(hat_matrix))
leverage = np.diag(hat_matrix)
# Remove leverage outliers and record their indices
outliers_num_leverage= []
for i in range(leverage.shape[0]):
    if (leverage[i] > np.percentile(leverage, 75)+1.5*stats.iqr(leverage)):
        outliers_num_leverage.append(i)
leverage_remove_outliers = leverage[leverage <= np.percentile(leverage, 75)+1.5*stats.iqr(leverage)]

# Calculate residuals and mean square error
y_predicted = X@beta_hat # Fitted value
e = y - X@beta_hat
N = y.shape[0]
mean_squared_error = np.dot(e.T, e)/N
# r_squred = np.var(X@beta_hat)/np.var(y)

##########################################
# Calculate the cook distance for each data
##########################################
cook_distance = np.zeros((y.shape[0],1))
for i in range(y.shape[0]):
    X_remove = np.delete(X, i, 0)
    y_remove = np.delete(y, i, 0)
    beta_i = np.dot(la.inv(np.dot(X_remove.T, X_remove)), np.dot(X_remove.T, y_remove))
    y_p = X_remove@la.inv(np.dot(X_remove.T, X_remove))@(X_remove.T)@y_remove
    cook_distance[i] = np.dot((y_p-X_remove@beta_i).T, y_p-X_remove@beta_i)/(N*mean_squared_error)
# plt.bar(list(range(y.shape[0])), cook_distance)
# plt.boxplot(cook_distance)
outliers_num_cook_distance = []
for i in range(cook_distance.shape[0]):
    if (cook_distance[i] > np.percentile(cook_distance, 75)+1.5*stats.iqr(cook_distance)):
        outliers_num_cook_distance.append(i)
cook_distance_remove_outliers = cook_distance[cook_distance <= np.percentile(cook_distance, 75)+1.5*stats.iqr(cook_distance)]

################################
# Calculate standardized residuals
################################
s = np.zeros((y.shape[0], 1))
for i in range(y.shape[0]):
    s[i] = e[i]/np.sqrt(mean_squared_error*(1-leverage[i]))
# plt.bar(list(range(y.shape[0])), np.abs(s))
# plt.boxplot(s)
outliers_num_s = []
for i in range(s.shape[0]):
    if (s[i] > np.percentile(s, 75)+1.5*stats.iqr(s)):
        outliers_num_s.append(i)
```

```python
82      # Remove outliers
83      by_eye = True
84      if by_eye:
85          outliers_num = [364, 365, 368, 370, 371, 372, 380, 405, 410, 418]
86      else:
87          outliers_num = np.union1d(outliers_num_leverage, outliers_num_cook_distance)
88          outliers_num = np.union1d(outliers_num, outliers_num_s)
89
90      lm_remove = sm.OLS(y_remove, sm.add_constant(X_remove)).fit()
91      X_remove = np.delete(X, outliers_num, 0)
92      y_remove = np.delete(y, outliers_num, 0)
93      s_remove = np.delete(s, outliers_num, 0)
94
95      # Page1
96      n_interval = 50
97      lmbda_list = np.linspace(-2, 2, n_interval)
98      llf_list = np.zeros(n_interval)
99      for i in range(lmbda_list.shape[0]):
100         llf_list[i] = stats.boxcox_llf(lmbda_list[i], y_remove)
101     fig, ax = plt.subplots(figsize=(8,4))
102     plt.plot(lmbda_list, llf_list)
103     plt.xlabel('Lambda value')
104     plt.ylabel('Log likelihood')
105     plt.title('Box-Cox Transformation Curve')
106     plt.show()
107
108     # Page 2
109     # Generate diagnostic plots before removing outliers
110     before = True # Change this variable to False to generate diagnostic plot after removing outliers
111     if before:
112         model_residuals = lm.resid
113         model_norm_residuals = lm.get_influence().resid_studentized_internal
114         model_leverage = lm.get_influence().hat_matrix_diag
115         model_cooks = lm.get_influence().cooks_distance[0]
116
116
117         diagnostic_plot = plt.figure()
118         plt.scatter(model_leverage, model_norm_residuals, alpha=0.5)
119         sns.regplot(model_leverage, model_norm_residuals, scatter=False, ci=False, lowess=True,
120                     line_kws={'color':'red', 'lw':1, 'alpha':0.8}, label='Cook\'s distance')
121         diagnostic_plot.axes[0].set_xlim(0, max(model_leverage)+0.01)
122         diagnostic_plot.axes[0].set_ylim(-3,5)
123         plt.title('Residuals vs Leverage\nBefore removing outliers')
124         plt.xlabel('Leverage')
125         plt.ylabel('Standardized Residuals')
126         plt.legend()
127
128         leverage_top_3 = np.flip(np.argsort(model_cooks), 0)[:3]
129         for i in leverage_top_3:
130             diagnostic_plot.axes[0].annotate(i, xy=(model_leverage[i], model_norm_residuals[i]))
131         plt.show()
132
133     # Boxcox transformation
134     y_boxcox, maxlog = stats.boxcox(y_remove)
135
136     # Influence plots after removing outliers and boxcox transformation
137     lm_remove = sm.OLS(y_boxcox, sm.add_constant(X_remove)).fit()
138     fig, ax = plt.subplots(figsize=(9,6))
139     fig = sm.graphics.influence_plot(lm_remove, alpha=0.001, ax=ax, criterion='cooks')
140     fig.show()
141
142     # Page 3
143     # Standardized residuals vs Fitted values without any transformation
144     sns.residplot(y_predicted, s)
145     plt.xlabel('Fitted values')
146     plt.ylabel('Standardized residuals')
147     plt.title('Without any transformation\nResiduals vs Fitted')
148     plt.show()
149     # Standardized residuals vs Fitted values removing all outliers and boxcox transformation
150     beta_hat_prime = np.dot(la.inv(np.dot(X_remove.T, X_remove)), np.dot(X_remove.T, y_remove))
151     y_predicted_prime = X_remove@beta_hat_prime # Fitted house price
152     sns.residplot(y_boxcox, np.delete(s, outliers_num))
153     plt.xlabel('Fitted values')
154     plt.ylabel('Standardized residuals')
155     plt.title('Removing all outliers and apply boxcox transformation\nResiduals vs Fitted')
156     plt.show()
157
158     # Page 4
159     # Fitted house price vs True house price
160     plt.scatter(y_predicted_prime, y_remove)
161     plt.plot(y_predicted_prime, LinearRegression().fit(y_predicted_prime, y_remove).predict(y_predicted_prime), 'r')
162     plt.xlabel('Fitted house price')
163     plt.ylabel('True hous price')
164     plt.title('After removing outliers and boxcox transformation')
165     plt.show()
166
```