

# Problem Set 2: The Determinants of Life Expectancy of the Poor

Harvard University  
Spring 2023

Instructor: Gregory Brulich, Ph.D.

Name: Benjamin Zeisberg

- Posted on: 01/30/2023
- Due at: 11:59pm on 02/07/2023

## Suggested Imports

```
In [4]: import pandas as pd
import numpy as np
import matplotlib
import matplotlib.pyplot as plt
import statsmodels.formula.api as sm
import statsmodels.api as smf
from stargazer, stargazer import Stargazer
```

## Background

The Health Inequality Project uses 1.4 billion observations on income from tax records covering the U.S. population from 1999–2014 to construct income-mortality gradients for each geographic region in the United States. The resulting datasets are publicly available at [healthinequality.org](https://healthinequality.org).

The map below shows life expectancy at age 40 for men in the bottom quartile of the income distribution for each commuting zone in the United States.

In this problem set, you will use these data to quantify the determinants of life expectancy for these low income men. The extract of the data set, `healthinequality.dta`, is described below and posted on the course website.

Source: The Health Inequality Project (Chetty, Stepner, Abraham, Lin, Scuderi, Turner, Bergeron, and Cutler 2016)

## Data Description

File: `healthinequality.dta`

The data consist of  $n = 590$  U.S. commuting zones with populations larger than 25,000 in 2000. Commuting zones are geographical aggregations of counties that are similar to metro areas but cover the entire U.S., including rural areas.

For more details on the construction of the variables included in this data set, please see Chetty, Stepner, Abraham, Lin, Scuderi, Turner, Bergeron, and Cutler (2016), which is posted on the course website.

Variable	Definition	Units	Mean
<code>cz</code>	Commuting Zone ID	n/a	n/a
<code>czname</code>	Commuting Zone Name	n/a	n/a
<code>stateabbrv</code>	2-letter state name (U.S. postal code)	n/a	n/a
<code>fips</code>	State FIPS code	n/a	n/a
<code>life_exp</code>	Male life expectancy at age 40 for the bottom quartile of the national income distribution (race adjusted)	Years	76.41
<code>cur_smoke</code>	Fraction of CZ that currently smokes in the bottom quartile of the national income distribution	Decimal, range 0 to 1	0.2792
<code>bmi_obese</code>	Fraction of CZ that is obese in the bottom quartile of the national income distribution	Decimal, range 0 to 1	0.3037
<code>exercise</code>	Fraction of CZ that exercised in the past 30 days in the bottom quartile of the national income distribution	Decimal, range 0 to 1	0.6047

## Data Load

```
In [16]: # Read dataset into a pandas dataframe
health = pd.read_stata("healthinequality.dta")

# Display first 5 rows of data
health.head()
```

```
Out [16]:
```

	<code>cz</code>	<code>czname</code>	<code>fips</code>	<code>stateabbrv</code>	<code>cur_smoke</code>	<code>bmi_obese</code>	<code>exercise</code>	<code>life_exp</code>
0	100	Johnson City	47	TN	0.351208	0.202475	0.546503	75.968956
1	200	Morristown	47	TN	0.419753	0.302632	0.506173	75.419853
2	301	Middlesborough	47	TN	0.364103	0.320856	0.475309	76.098907
3	302	Knoxville	47	TN	0.324125	0.259366	0.581132	75.089340
4	401	Winston-Salem	37	NC	0.323423	0.285605	0.606607	75.910576

## Instructions

Please submit your Problem Set on Canvas. Your submission should include two files:

1. This notebook as a `.ipynb` file with your code and answers to questions
2. A `.pdf` version of this notebook. TODO: Provide general instructions on converting `.ipynb` to `pdf`

## Questions

Note: Short answers should be very succinct. Show your work and intuition clearly: credit is given for explanations and not just having the correct answer

### 1

Use the starter script files to help you get started on this question. The  $R^2$  regression diagnostic statistic measures how much of the variance in the dependent variable can be explained linearly by the covariates in the regression. It equals the ratio between the explained sum of squares and the total sum of squares in a regression:

$$R^2 = \frac{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2}$$

In a simple bivariate regression, it equals the square of the correlation coefficient between the dependent variable and the single independent variable.

1. Estimate a regression of `life_exp` on `cur_smoke`. Explain in words what the coefficient on `cur_smoke` means.
2. What is the  $R^2$  of this regression? Does the  $R^2$  tell us whether the regression does a good job of fitting the data? If not, what does it tell us?
3. Now generate a random variable that is independent of both `life_exp` and `cur_smoke`. Run a regression of `life_exp` on an intercept, `cur_smoke`, and the random variable that you generated. What happens to the  $R^2$ ?
4. Now generate a total of 588 random variables that are independent of each other, `life_exp` and `cur_smoke`. Regress `life_exp` on an intercept, `cur_smoke`, and the 588 random variables that you generated. What happens to the  $R^2$ ? Discuss briefly.
5. Explain (d) using the following simple non-trivial example: with  $N = 2$  observations, what happens if we run a regression with an intercept and one explanatory variable?
6. Now generate a total of 589 random variables. Regress `life_exp` on an intercept, `cur_smoke`, and the 589 independent random variable that you generated. What happens?
7. Use the residual regression formula to explain (f).

```
In [17]: # Your Code Here
# 1. Regression of life_exp on cur_smoke
mod_1 = sm.ols("life_exp ~ cur_smoke", data=health)
res_1 = mod_1.fit(cov_type="HC2")

print(res_1.summary(slim=True))
```

```
OLS Regression Results
=====
Dep. Variable:      life_exp      R-squared:      0.121
Model:              OLS          Adj. R-squared:  0.119
No. Observations:   590          F-statistic:   43.65
Covariance Type:    HC2          Prob (F-statistic): 8.81e-11
=====
                        coef      std err          z      Pr>|z|      [0.025      0.975]
-----
Intercept           78.4493         0.329      238.384         0.000        77.802        79.096
cur_smoke           -7.4469         1.127        -6.601         0.000        -9.658        -5.237
=====
```

Notes:

[1] Standard Errors are heteroscedasticity robust (HC2)

### Question 1.

We estimate the average effect of the fraction of people in a commuting zone that currently smokes in the bottom quartile of the national income distribution on life expectancy to be -7.446 years. Our coefficient `cur_smoke` can be explained as the fraction of people in commuting zone "CZ" that currently smokes in the bottom quartile of the income distribution.

### Question 2.

What is the  $R^2$  of this regression? The  $R^2$  of this regression is 0.121. Does the  $R^2$  tell us whether the regression does a good job of fitting the data? The  $R^2$  does generally not tell us if the regression does a good job of fitting the data, we can see that by plotting the estimated regression line vs the data to see if our estimate is a good or bad. If not, what does it tell us? Generally the  $R^2$  tells us how much of the variance in our dependent variable is explained by our model.

### Question 3.

```
In [22]: # Find number of rows
def random_variables(x):
    N = health.shape[0]
    # Create column names like rand_5, rand_6, ...
    random_column_names = [f"rand_{i}" for i in range(x)]

    # Create a new dataframe with just random columns
    random_df = pd.DataFrame(
        np.random.random(size=(N, x)),
        columns=random_column_names
    )

    # Join old and new dataframes
    new_df = pd.concat([health, random_df], axis=1)
    return new_df

frame = random_variables(1)
```

```
# 1. Regression of life_exp on cur_smoke, one random variable
mod_2 = sm.ols("life_exp ~ cur_smoke + rand_0", data=frame)
res_2 = mod_2.fit(cov_type="HC2")

print(res_2.summary(slim=True))
```

```
OLS Regression Results
=====
Dep. Variable:      life_exp      R-squared:      0.121
Model:              OLS          Adj. R-squared:  0.118
No. Observations:   590          F-statistic:   21.82
Covariance Type:    HC2          Prob (F-statistic): 7.21e-10
=====
                        coef      std err          z      Pr>|z|      [0.025      0.975]
-----
Intercept           78.4493         0.330      237.631         0.000        77.802        79.096
cur_smoke           -7.4469         1.128        -6.601         0.000        -9.658        -5.236
rand_0              0.0900         0.173         0.519         0.604        -0.250         0.430
=====
```

Notes:

[1] Standard Errors are heteroscedasticity robust (HC2)

### What happens to $R^2$ ?

Nothing.  $R^2$  is a measure of how much of the variance in the dependent variable is explained by our model, since we are simply adding one random variable we would expect that there be not much change to the  $R^2$  besides a minor increase. An increase because we now have a random variable with random data, that could explain some of the variance in the dependent variable.

### Question 4.

Now generate a total of 588 random variables. What happens to  $R^2$ ? we can simply call the helper function as defined before to create 588 random variables, and then run our regression.

Answer: We find that our  $R^2$  now is 1.00 meaning that our model explains the variance in our dependent variable. Why? Essentially, we have created enough random variables with random data, that some are correlated with our dependent variable, and so we explained all the variance in our dependent variable with our model.

```
In [26]: frame = random_variables(588)
columns = frame.columns
separator = "+"
string = separator.join(columns[8:1])
mod_3 = sm.ols(f"life_exp ~ cur_smoke + {string}", data=frame)
res_3 = mod_3.fit(cov_type="HC2")
print(res_3.summary())
```



RuntimeWarning: invalid value encountered in divide

```
self.k = np.divide(self.residuals, self.k_constant, self.df_resid)
```

RuntimeWarning: invalid value encountered in double\_scalars

```
return 1 - (np.divide(self.residuals, self.k_constant, self.df_resid))
```

OLS Regression Results

Dep. Variable:	life_exp	R-squared:	0.5	Adj. R-squared:	1.000	
Method: <td>Least Squares<td>F-statistic:<td>3.11e+12</td><td></td><td></td></td></td>	Least Squares <td>F-statistic:<td>3.11e+12</td><td></td><td></td></td>	F-statistic: <td>3.11e+12</td> <td></td> <td></td>	3.11e+12			
Date: <td>Mon, 07 Feb 2022<td>Log-Likelihood:<td></td><td></td><td></td></td></td>	Mon, 07 Feb 2022 <td>Log-Likelihood:<td></td><td></td><td></td></td>	Log-Likelihood: <td></td> <td></td> <td></td>				
No. Observations: <td>589<td></td><td></td><td></td><td></td></td>	589 <td></td> <td></td> <td></td> <td></td>					
DF Residuals: <td>589<td>0 BIC:<td></td><td></td><td></td></td></td>	589 <td>0 BIC:<td></td><td></td><td></td></td>	0 BIC: <td></td> <td></td> <td></td>				
Model: <td></td> <td></td> <td></td> <td></td> <td></td>						
	coef	std err	z	P> z	[0.025	0.975]
Intercept	158.9371	2.57e-05	5.87e+06	0.000	158.937	158.937
cur_smoke	8.8591	nan	nan	nan	nan	nan
rand_0	-7.5158	nan	nan	nan	nan	nan
rand_1	-2.4623	nan	nan	nan	nan	nan
rand_2	7.7262	nan	nan	nan	nan	nan
rand_3	-2.9558	nan	nan	nan	nan	nan
rand_4	-8.2557	nan	nan	nan	nan	nan
rand_5	-8.6829	nan	nan	nan	nan	nan
rand_6	-4.1799	2e-06	-2.09e+06	0.000	-4.180	-4.180
rand_7	-4.7574	nan	nan	nan	nan	nan
rand_8	-1.0872	nan	nan	nan	nan	nan
rand_9	-3.0424	nan	nan	nan	nan	nan
rand_10	-9.4608	nan	nan	nan	nan	nan
rand_11	4.7152	nan	nan	nan	nan	nan
rand_12	-2.6850	nan	nan	nan	nan	nan
rand_13	-2.0554	nan	nan	nan	nan	nan
rand_14	-1.3815	nan	nan	nan	nan	nan
rand_15	5.6105	nan	nan	nan	nan	nan
rand_16	2.2385	1.68e-06	1.33e+06	0.000	2.238	2.238
rand_17	-1.2594	nan	nan	nan	nan	nan
rand_18	-3.3249	2.96e-06	1.14e+05	0.000	0.339	0.339
rand_19	3.1783	nan	nan	nan	nan	nan
rand_20	10.8034	nan	nan	nan	nan	nan
rand_21	-2.4623	nan	nan	nan	nan	nan
rand_22	-9.4769	nan	nan	nan	nan	nan
rand_23	8.2919	nan	nan	nan	nan	nan
rand_24	6.2936	nan	nan	nan	nan	nan
rand_25	4.8665	2.26e-06	1.8e+06	0.000	4.866	4.866
rand_26	2.5067	nan	nan	nan	nan	nan
rand_27	15.9171	nan	nan	nan	nan	nan
rand_28	-2.4623	nan	nan	nan	nan	nan
rand_29	8.6591	nan	nan	nan	nan	nan
rand_30	-1.3123	nan	nan	nan	nan	nan
rand_31	7.5416	nan	nan	nan	nan	nan
rand_32	8.8483	nan	nan	nan	nan	nan
rand_33	0.8987	nan	nan	nan	nan	nan
rand_34	-8.8484	nan	nan	nan	nan	nan
rand_35	-4.1815	nan	nan	nan	nan	nan
rand_36	3.9498	nan	nan	nan	nan	nan
rand_37	15.5916	nan	nan	nan	nan	nan
rand_38	3.2353	nan	nan	nan	nan	nan
rand_39	-3.3249	nan	nan	nan	nan	nan
rand_40	4.0077	nan	nan	nan	nan	nan
rand_41	-2.4492	nan	nan	nan	nan	nan
rand_42	-0.3659	nan	nan	nan	nan	nan
rand_43	0.7587	8.35e-07	9.09e+05	0.000	0.759	0.759
rand_44	-8.4681	nan	nan	nan	nan	nan
rand_45	-2.6467	nan	nan	nan	nan	nan
rand_46	-6.9437	nan	nan	nan	nan	nan
rand_47	1.1080	nan	nan	nan	nan	nan
rand_48	8.1812	nan	nan	nan	nan	nan
rand_49	4.7752	nan	nan	nan	nan	nan
rand_50	3.8653	1.92e-06	2.01e+06	0.000	3.865	3.865
rand_51	9.7533	nan	nan	nan	nan	nan
rand_52	4.8665	3.89e-06	1.23e+05	0.000	-5.185	-5.185
rand_53	-8.6758	4.07e-06	2.13e+06	0.000	-8.676	-8.676
rand_54	1.6034	nan	nan	nan	nan	nan
rand_55	5.8986	nan	nan	nan	nan	nan
rand_56	3.3064	nan	nan	nan	nan	nan
rand_57	0.4567	1.27e-07	3.59e+06	0.000	0.457	0.457
rand_58	-5.8987	nan	nan	nan	nan	nan
rand_59	6.5360	nan	nan	nan	nan	nan
rand_60	-18.6387	nan	nan	nan	nan	nan
rand_61	0.3653	nan	nan	nan	nan	nan
rand_62	4.1298	nan	nan	nan	nan	nan
rand_63	4.2224	nan	nan	nan	nan	nan
rand_64	-7.8058	nan	nan	nan	nan	nan
rand_65	-7.8058	nan	nan	nan	nan	nan
rand_66	-2.6208	nan	nan	nan	nan	nan
rand_67	-4.4356	nan	nan	nan	nan	nan
rand_68	4.2368	3.69e-06	1.15e+06	0.000	4.237	4.237
rand_69	-1.3815	nan	nan	nan	nan	nan
rand_70	-8.2181	nan	nan	nan	nan	nan
rand_71	-13.8221	nan	nan	nan	nan	nan
rand_72	6.2936	nan	nan	nan	nan	nan
rand_73	6.2936	nan	nan	nan	nan	nan
rand_74	-2.9414	nan	nan	nan	nan	nan
rand_75	8.8756	4.34e-06	1.17e+06	0.000	5.876	5.876
rand_76	-1.3815	2.33e-07	5.07e+05	0.000	-1.179	-1.179
rand_77	2.9042	nan	nan	nan	nan	nan
rand_78	9.5188	nan	nan	nan	nan	nan
rand_79	14.5562	3.81e-06	3.82e+06	0.000	14.556	14.556
rand_80	3.3815	nan	nan	nan	nan	nan
rand_81	5.6643	nan	nan	nan	nan	nan
rand_82	-8.5171	nan	nan	nan	nan	nan
rand_83	-2.4623	nan	nan	nan	nan	nan
rand_84	-1.6888	nan	nan	nan	nan	nan
rand_85	-10.8568	nan	nan	nan	nan	nan
rand_86	-7.4983	nan	nan	nan	nan	nan
rand_87	9.5342	nan	nan	nan	nan	nan
rand_88	11.5941	nan	nan	nan	nan	nan
rand_89	13.1167	nan	nan	nan	nan	nan
rand_90	7.2594	nan	nan	nan	nan	nan
rand_91	3.0043	nan	nan	nan	nan	nan
rand_92	5.2759	nan	nan	nan	nan	nan
rand_93	-5.3685	nan	nan	nan	nan	nan
rand_94	-8.4373	nan	nan	nan	nan	nan
rand_95	0.5469	nan	nan	nan	nan	nan
rand_96	-2.2591	nan	nan	nan	nan	nan
rand_97	-5.3685	nan	nan	nan	nan	nan
rand_98	-4.9554	nan	nan	nan	nan	nan
rand_99	-18.2926	nan	nan	nan	nan	nan
rand_100	8.8932	nan	nan	nan	nan	nan
rand_101	2.5281	nan	nan	nan	nan	nan
rand_102	11.5892	nan	nan	nan	nan	nan
rand_103	1.0551	nan	nan	nan	nan	nan
rand_104	4.0538	nan	nan	nan	nan	nan
rand_105	28.5129	nan	nan	nan	nan	nan
rand_106	11.4659	nan	nan	nan	nan	nan
rand_107	13.3549	nan	nan	nan	nan	nan
rand_108	-0.1599	nan	nan	nan	nan	nan
rand_109	0.5272	nan	nan	nan	nan	nan
rand_110	-2.4623	2.85e-06	1.75e+06	0.000	4.906	4.906
rand_111	-3.2747	nan	nan	nan	nan	nan
rand_112	4.4333	nan	nan	nan	nan	nan
rand_113	-8.5688	nan	nan	nan	nan	nan
rand_114	-6.5925	nan	nan	nan	nan	nan
rand_115	2.2482	nan	nan	nan	nan	nan
rand_116	4.9792	nan	nan	nan	nan	nan
rand_117	-1.3815	3.18e-06	3.12e+05	0.000	0.000	0.000
rand_118	-1.0330	nan	nan	nan	nan	nan
rand_119	-6.9555	nan	nan	nan	nan	nan
rand_120	-1.7322	nan	nan	nan	nan	nan
rand_121	-1.2171	nan	nan	nan	nan	nan
rand_122	6.5143	nan	nan	nan	nan	nan
rand_123	7.2896	nan	nan	nan	nan	nan
rand_124	-2.9558	nan	nan	nan	nan	nan
rand_125	8.4542	nan	nan	nan	nan	nan
rand_126	2.5180	nan	nan	nan	nan	nan
rand_127	6.2187	nan	nan	nan	nan	nan
rand_128	2.2706	nan	nan	nan	nan	nan
rand_129	-3.6944	nan	nan	nan	nan	nan
rand_130	-2.9248	nan	nan	nan	nan	nan
rand_131	-6.1728	nan	nan	nan	nan	nan
rand_132	-3.2587	1.64e-06	1.99e+06	0.000	-3.251	-3.251
rand_133	1.9553	nan	nan	nan	nan	nan
rand_134	3.8932	nan	nan	nan	nan	nan
rand_135	-6.8715	nan	nan	nan	nan	nan
rand_136	-5.9142	nan	nan	nan	nan	nan
rand_137	-1.5856	nan	nan	nan	nan	nan
rand_138	4.3985	nan	nan	nan	nan	nan
rand_139	-14.4252	nan	nan	nan	nan	nan
rand_140	-9.7862	nan	nan	nan	nan	nan
rand_141	-4.2281	nan	nan	nan	nan	nan
rand_142	8.7684	nan	nan	nan	nan	nan
rand_143	-11.8498	nan	nan	nan	nan	nan
rand_144	-7.0588	nan	nan	nan	nan	nan
rand_145	5.8352	nan	nan	nan	nan	nan
rand_146	3.6383	nan	nan	nan	nan	nan
rand_147	-1.5948	nan	nan	nan	nan	nan
rand_148	-9.2811	nan	nan	nan	nan	nan
rand_149	-6.0557	nan	nan	nan	nan	nan
rand_150	-9.7156	nan	nan	nan	nan	nan
rand_151	-2.7736	nan	nan	nan	nan	nan
rand_152	-8.8651	nan	nan	nan	nan	nan
rand_153	-0.7678	nan	nan	nan	nan	nan
rand_154	-3.3153	nan	nan	nan	nan	nan
rand_155	-5.5225	nan	nan	nan	nan	nan
rand_156	-3.2139	2.55e-06	-1.26e+06	0.000	-3.214	-3.214
rand_157	7.6446	nan	nan	nan	nan	nan
rand_158	-5.5225	3.03e-06	-1.82e+06	0.000	-5.525	-5.525
rand_159	-6.2186	nan	nan	nan	nan	nan
rand_160	-5.5189	nan	nan	nan	nan	nan
rand_161	-1.0358	nan	nan	nan	nan	nan
rand_162	7.8989	nan	nan	nan	nan	nan
rand_163	-19.1564	nan	nan	nan	nan	nan
rand_164	1.7499	nan	nan	nan	nan	nan
rand_165	-2.9558	nan	nan	nan	nan	nan
rand_166	4.1292	1.29e-06	3.2e+06	0.000	4.129	4.129
rand_167	-0.1135	nan	nan	nan	nan	nan
rand_168	7.3116	nan	nan	nan	nan	nan
rand_169	-4.9854	nan	nan	nan	nan	nan
rand_170	-17.7557	nan	nan	nan	nan	nan
rand_171	-7.8039	nan	nan	nan	nan	nan
rand_172	-6.4525	nan	nan	nan	nan	nan
rand_173	9.2341	nan	nan	nan	nan	nan
rand_174	8.1612	nan	nan	nan	nan	nan
rand_175	-0.6638	2.2e-07	-3.01e+06	0.000	-0.663	-0.663
rand_176	1.6328	nan	nan	nan	nan	nan
rand_177	3.0872	nan	nan	nan	nan	nan
rand_178	-3.3783	nan	nan	nan	nan	nan
rand_179	-2.9558	nan	nan	nan	nan	nan
rand_180	-11.4945	nan	nan	nan	nan	nan
rand_181	-0.2258	nan	nan	nan	nan	nan
rand_182	-2.9245	nan	nan	nan	nan	nan
rand_183	-20.2887	nan	nan	nan	nan	nan
rand_184	-0.8340	nan	nan	nan	nan	nan
rand_185	-6.3853	9.89e-07	-6.38e+06	0.000	-6.385	-6.385
rand_186	1.5562	1.56e-06	1.25e+05	0.000	0.154	0.154
rand_187	1.5185	nan	nan	nan	nan	nan
rand_188	-11.0145	nan	nan	nan	nan	nan
rand_189	-16.4232	nan	nan	nan	nan	nan
rand_190	13.5880	nan	nan	nan	nan	nan
rand_191	-10.3184	nan	nan	nan	nan	nan
rand_192	14.4101	nan	nan	nan	nan	nan
rand_193	-20.1471	nan	nan	nan	nan	nan
rand_194	-3.0611	nan	nan	nan	nan	nan
rand_195	-3.3883	nan	nan	nan	nan	nan
rand_196	-8.4028	nan	nan	nan	nan	nan
rand_197	6.5838	nan	nan	nan	nan	nan
rand_198	4.4129	nan	nan	nan	nan	nan
rand_199	7.4496	nan	nan	nan	nan	nan
rand_200	-2.9027	1.15e-06	6.80e+06	0.000	0.793	0.793
rand_201	6.5159	nan	nan	nan	nan	nan
rand_202	-3.3183	1.45e-06	-2.29e+06	0.000	-3.318	-3.318
rand_203	3.4258	nan	nan	nan	nan	nan
rand_204	6.8206	nan	nan	nan	nan	nan
rand_205</						



RuntimeWarning: divide by zero encountered in divide

RuntimeWarning: invalid value encountered in double\_scalars

```
return 1 - np.divide(self.nodes - self.f_constant, self.df_resid)
```

Dep. Variable: life\_exp R-squared: 0.808

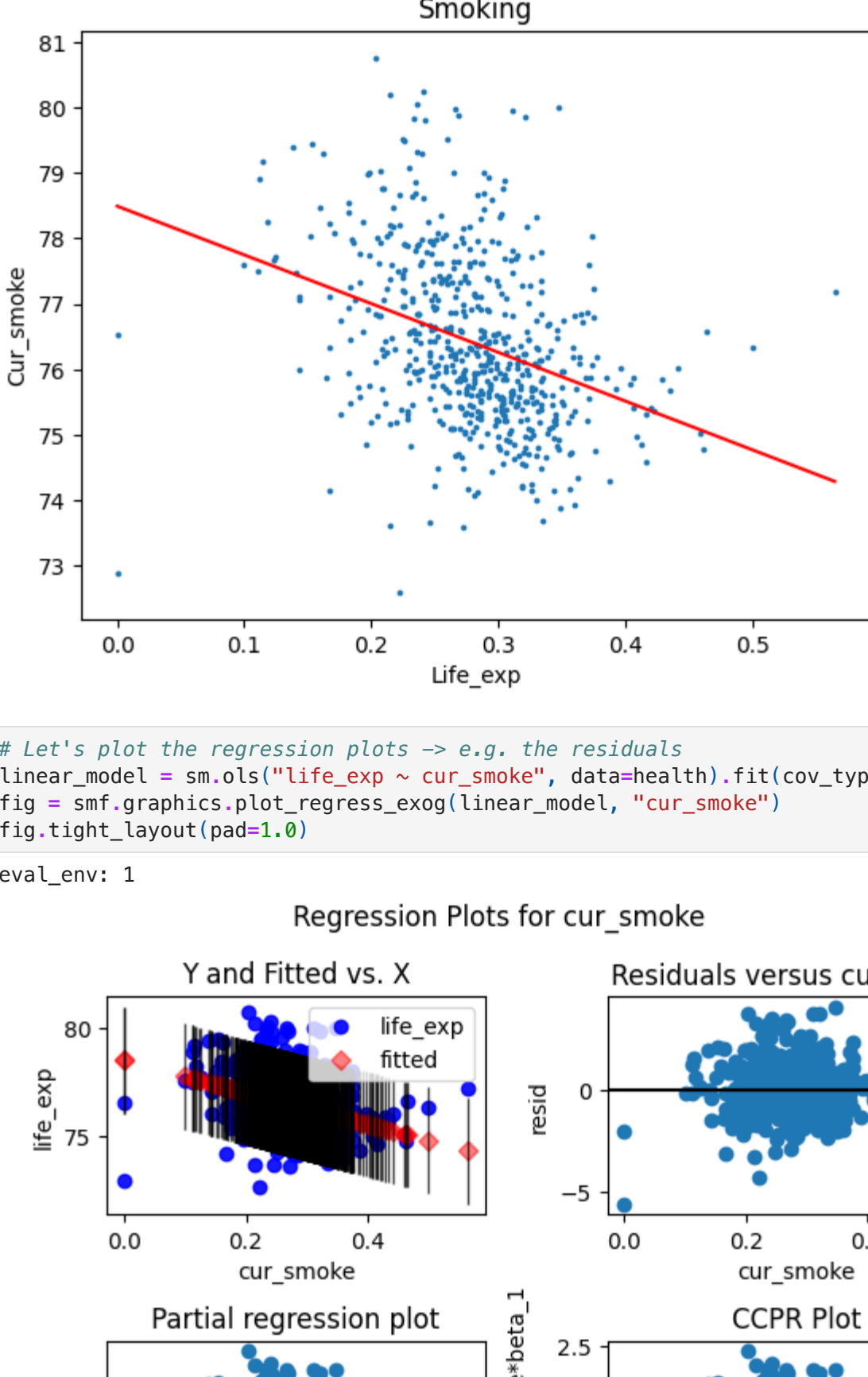
No. Observations: 590 F-statistic: 2.057e+15

Covariance Type: HC2 Prob (F-statistic): nan

	coef	std err	z	P> z	[0.025	0.975]
Intercept	167.6481	4.96e-05	3.38e+06	0.000	167.648	167.648
cur_smoke	-0.0029	1.38e-06	-2.75e+06	0.000	-12.769	-12.769
bmi_obese	0.7086	4.97e-07	1.41e+06	0.000	0.701	0.701
cur_smoke	0.2638	9.04e-06	1.28e+06	0.000	0.624	0.624
bmi_obese	-0.2873	1.72e-06	-1.67e+06	0.000	-0.287	-0.287
cur_smoke	-2.6919	1.72e-06	-1.55e+06	0.000	-2.662	-2.662
bmi_obese	-7.2949	nan	nan	nan	nan	nan
cur_smoke	5.5086	3.22e-06	1.71e+06	0.000	5.501	5.501
bmi_obese	-4.2873	1.7e-06	-2.47e+06	0.000	-4.287	-4.287
cur_smoke	11.0727	1.72e-06	6.40e+06	0.000	11.073	11.073
bmi_obese	5.9978	1.95e-06	3.07e+06	0.000	5.998	5.998
cur_smoke	-4.2873	1.93e-06	-2.27e+06	0.000	-4.287	-4.287
bmi_obese	11.2880	3.67e-06	3.07e+06	0.000	11.288	11.288
cur_smoke	-24.8881	8.16e-06	-3.05e+06	0.000	-24.881	-24.881
bmi_obese	-0.0571	6.59e-06	-8.64e+05	0.000	-20.880	-20.880
cur_smoke	12.1327	5.36e-06	3.04e+06	0.000	12.133	12.133
bmi_obese	-11.7737	5.17e-06	-2.28e+06	0.000	-11.774	-11.774
cur_smoke	-6.2015	2.39e-06	-2.6e+06	0.000	-6.201	-6.201
bmi_obese	-4.4716	nan	nan	nan	nan	nan
cur_smoke	-6.2288	3.27e-06	-1.9e+06	0.000	-6.229	-6.229
bmi_obese	4.1066	nan	nan	nan	nan	nan
cur_smoke	-7.6184	1.98e-06	-3.82e+06	0.000	-7.618	-7.618
bmi_obese	5.2012	1.68e-06	3.14e+06	0.000	5.201	5.201
cur_smoke	4.4125	nan	nan	nan	nan	nan
bmi_obese	-10.1534	4.47e-06	-2.29e+06	0.000	-10.213	-10.213
cur_smoke	-4.7452	1.29e-06	-3.74e+06	0.000	-4.746	-4.746
bmi_obese	-17.2384	5.6e-06	-3.08e+06	0.000	-17.238	-17.238
cur_smoke	12.2849	1.26e-06	6.92e+06	0.000	12.285	12.285
bmi_obese	-14.7287	1.74e-06	-8.02e+05	0.000	-14.73	-14.73
cur_smoke	-7.6564	1.59e-06	-4.71e+06	0.000	-7.656	-7.656
bmi_obese	-7.7155	3.48e-07	-2.22e+07	0.000	-7.716	-7.716
cur_smoke	-9.4257	nan	nan	nan	nan	nan
bmi_obese	-3.0529	8.24e-07	-4.32e+06	0.000	-3.053	-3.053
cur_smoke	11.2042	2.52e-06	4.45e+06	0.000	11.204	11.204
bmi_obese	6.3083	3.06e-06	2.06e+06	0.000	6.308	6.308
cur_smoke	-4.9723	1.67e-06	-2.95e+06	0.000	-4.973	-4.973
bmi_obese	-15.6719	7.28e-06	-2.15e+06	0.000	-15.672	-15.672
cur_smoke	-19.3179	7.05e-06	-2.74e+06	0.000	-19.318	-19.318
bmi_obese	-9.9076	1.21e-06	-4.93e+06	0.000	-9.908	-9.908
cur_smoke	14.7287	1.63e-06	3.1e+06	0.000	14.73	14.73
bmi_obese	-1.5977	1.26e-06	-1.27e+06	0.000	-1.598	-1.598
cur_smoke	19.1078	3.12e-06	3.52e+06	0.000	19.1075	19.1075
bmi_obese	8.0457	1.98e-06	4.07e+06	0.000	8.046	8.046
cur_smoke	-1.4375	1.28e-06	-1.13e+06	0.000	-1.437	-1.437
bmi_obese	11.3142	7.95e-07	1.42e+07	0.000	11.314	11.314
cur_smoke	4.2158	1.37e-06	2.13e+06	0.000	4.216	4.216
bmi_obese	-2.7599	1.81e-06	-2.78e+06	0.000	-2.759	-2.759
cur_smoke	-15.5259	5.39e-06	-2.83e+06	0.000	-15.527	-15.527
bmi_obese	25.8002	8.43e-06	2.97e+06	0.000	25.802	25.802
cur_smoke	14.6538	4.99e-06	2.93e+06	0.000	14.653	14.653
bmi_obese	-4.8995	nan	nan	nan	nan	nan
cur_smoke	-4.7452	1.29e-06	-3.74e+06	0.000	-4.746	-4.746
bmi_obese	-23.9493	8.53e-06	-2.81e+06	0.000	-23.949	-23.949
cur_smoke	-7.9999	1.6e-06	-4.99e+06	0.000	-8.000	-8.000
bmi_obese	-15.9477	6.48e-06	-3.08e+06	0.000	-15.948	-15.948
cur_smoke	9.6618	3.64e-06	2.65e+06	0.000	9.661	9.661
bmi_obese	3.2888	nan	nan	nan	nan	nan
cur_smoke	9.0454	3.07e-06	2.95e+06	0.000	9.045	9.045
bmi_obese	-2.4161	1.86e-06	-1.09e+06	0.000	-2.416	-2.416
cur_smoke	11.8093	3.57e-06	3.34e+06	0.000	11.809	11.809
bmi_obese	-4.9453	nan	nan	nan	nan	nan
cur_smoke	-4.7292	1.44e-06	-3.29e+06	0.000	-4.729	-4.729
bmi_obese	8.0757	nan	nan	nan	nan	nan
cur_smoke	-1.8816	1.67e-06	-1.47e+06	0.000	-1.882	-1.882
bmi_obese	-3.7749	2.59e-06	-1.05e+06	0.000	-3.775	-3.775
cur_smoke	-5.5781	2.91e-06	-1.92e+06	0.000	-5.578	-5.578
bmi_obese	-4.0853	1.12e-06	-4.01e+06	0.000	-4.086	-4.086
cur_smoke	-12.1012	6.06e-06	-2.72e+06	0.000	-12.101	-12.101
bmi_obese	-12.0112	nan	nan	nan	nan	nan
cur_smoke	1.0568	5.72e-07	1.85e+06	0.000	1.057	1.057
bmi_obese	-2.1117	2.49e-06	-1.02e+06	0.000	-2.112	-2.112
cur_smoke	-1.2658	nan	nan	nan	nan	nan
bmi_obese	-7.6108	2.67e-06	-2.84e+06	0.000	-7.610	-7.610
cur_smoke	-1.1171	1.37e-06	-1.11e+06	0.000	-1.117	-1.117
bmi_obese	-1.2286	1.71e-06	-1.78e+06	0.000	-1.229	-1.229
cur_smoke	-1.2239	1.45e-06	-8.43e+05	0.000	-1.224	-1.224
bmi_obese	15.9722	5.88e-06	2.72e+06	0.000	15.972	15.972
cur_smoke	-4.6326	1.69e-06	-2.73e+06	0.000	-4.634	-4.634
bmi_obese	0.4238	2.1e-06	2.02e+05	0.000	0.423	0.423
cur_smoke	-5.4143	5.32e-07	-1.02e+07	0.000	-5.414	-5.414
bmi_obese	1.8816	1.67e-06	-1.47e+06	0.000	1.882	1.882
cur_smoke	0.4945	1.88e-06	2.63e+05	0.000	0.495	0.495
bmi_obese	12.6629	5.11e-06	2.48e+06	0.000	12.663	12.663
cur_smoke	-19.9282	4.53e-06	-2.44e+06	0.000	-19.928	-19.928
bmi_obese	17.1216	5.25e-06	3.26e+06	0.000	17.122	17.122
cur_smoke	-3.6781	6.59e-06	-2.18e+06	0.000	-3.678	-3.678
bmi_obese	14.7327	2.83e-06	7.25e+06	0.000	14.733	14.733
cur_smoke	8.0457	1.98e-06	4.07e+06	0.000	8.046	8.046
bmi_obese	-10.3581	4.33e-06	-2.53e+06	0.000	-10.358	-10.358
cur_smoke	5.1431	2.12e-06	2.42e+06	0.000	5.143	5.143
bmi_obese	-11.4129	5.86e-06	-2.59e+06	0.000	-11.413	-11.413
cur_smoke	3.7143	1.21e-06	2.13e+06	0.000	3.714	3.714
bmi_obese	4.3198	2.47e-06	1.75e+06	0.000	4.320	4.320
cur_smoke	6.3083	3.06e-06	2.06e+06	0.000	6.308	6.308
bmi_obese	-1.5977	1.26e-06	-1.27e+06	0.000	-1.598	-1.598
cur_smoke	19.1078	3.12e-06	3.52e+06	0.000	19.1075	19.1075
bmi_obese	8.0457	1.98e-06	4.07e+06	0.000	8.046	8.046
cur_smoke	-1.4375	1.28e-06	-1.13e+06	0.000	-1.437	-1.437
bmi_obese	11.3142	7.95e-07	1.42e+07	0.000	11.314	11.314
cur_smoke	4.2158	1.37e-06	2.13e+06	0.000	4.216	4.216
bmi_obese	-2.7599	1.81e-06	-2.78e+06	0.000	-2.759	-2.759
cur_smoke	-15.5259	5.39e-06	-2.83e+06	0.000	-15.527	-15.527
bmi_obese	25.8002	8.43e-06	2.97e+06	0.000	25.802	25.802
cur_smoke	14.6538	4.99e-06	2.93e+06	0.000	14.653	14.653
bmi_obese	-4.8995	nan	nan	nan	nan	nan
cur_smoke	-4.7452	1.29e-06	-3.74e+06	0.000	-4.746	-4.746
bmi_obese	-23.9493	8.53e-06	-2.81e+06	0.000	-23.949	-23.949
cur_smoke	-7.9999	1.6e-06	-4.99e+06	0.000	-8.000	-8.000
bmi_obese	-15.9477	6.48e-06	-3.08e+06	0.000	-15.948	-15.948
cur_smoke	9.6618	3.64e-06	2.65e+06	0.000	9.661	9.661
bmi_obese	3.2888	nan	nan	nan	nan	nan
cur_smoke	9.0454	3.07e-06	2.95e+06	0.000	9.045	9.045
bmi_obese	-2.4161	1.86e-06	-1.09e+06	0.000	-2.416	-2.416
cur_smoke	11.8093	3.57e-06	3.34e+06	0.000	11.809	11.809
bmi_obese	-4.9453	nan	nan	nan	nan	nan
cur_smoke	-4.7292	1.44e-06	-3.29e+06	0.000	-4.729	-4.729
bmi_obese	8.0757	nan	nan	nan	nan	nan
cur_smoke	-1.8816	1.67e-06	-1.47e+06	0.000	-1.882	-1.882
bmi_obese	-3.7749	2.59e-06	-1.05e+06	0.000	-3.775	-3.775
cur_smoke	-5.5781	2.91e-06	-1.92e+06	0.000	-5.578	-5.578
bmi_obese	-4.0853	1.12e-06	-4.01e+06	0.000	-4.086	-4.086
cur_smoke	-12.1012	6.06e-06	-2.72e+06	0.000	-12.101	-12.101
bmi_obese	-12.0112	nan	nan	nan	nan	nan
cur_smoke	1.0568	5.72e-07	1.85e+06	0.000	1.057	1.057
bmi_obese	-2.1117	2.49e-06	-1.02e+06	0.000	-2.112	-2.112
cur_smoke	-1.2658	nan	nan	nan	nan	nan
bmi_obese	-7.6108	2.67e-06	-2.84e+06	0.000	-7.610	-7.610
cur_smoke	-1.1171	1.37e-06	-1.11e+06	0.000	-1.117	-1.117
bmi_obese	-1.2286	1.71e-06	-1.78e+06	0.000	-1.229	-1.229
cur_smoke	-1.2239	1.45e-06	-8.43e+05	0.000	-1.224	-1.224
bmi_obese	15.9722	5.88e-06	2.72e+06	0.000	15.972	15.972
cur_smoke	-4.6326	1.69e-06	-2.73e+06	0.000	-4.634	-4.634
bmi_obese	0.4238	2.1e-06	2.02e+05	0.000	0.423	0.423
cur_smoke	-5.4143	5.32e-07	-1.02e+07	0.000	-5.414	-5.414
bmi_obese	1.8816	1.67e-06	-1.47e+06	0.000	1.882	1.882
cur_smoke	0.4945	1.88e-06	2.63e+05	0.000	0.495	0.495
bmi_obese	12.6629	5.11e-06	2.48e+06	0.000	12.663	12.663
cur_smoke	-19.9282	4.53e-06	-2.44e+06	0.000	-19.928	-19.928
bmi_obese	17.1216	5.25e-06	3.26e+06	0.000	17.122	17.122
cur_smoke	-3.6781	6.59e-06	-2.18e+06	0.000	-3.678	-3.678
bmi_obese	14.7327	2.83e-06	7.25e+06	0.000	14.733	14.733
cur_smoke	8.0457	1.98e-06	4.07e+06	0.000	8.046	8.046
bmi_obese	-10.3581	4.33e-06	-2.53e+06	0.000	-10.358	-10.358
cur_smoke	5.1431	2.12e-06	2.42e+06	0.000	5.143	5.143
bmi_obese	-11.4129	5.86e-06	-2.59e+06	0.000	-11.413	-11.413
cur_smoke	3.7143	1.21e-06	2.13e+06	0.000	3.714	3.714
bmi_obese	4.3198	2.47e-06	1.75e+06	0.000	4.320	4.320
cur_smoke	6.3083	3.06e-06	2.06e+06	0.000	6.308	6.308
bmi_obese	-1.5977	1.26e-06	-1.27e+06	0.000	-1.598	-1.598
cur_smoke	19.1078	3.12e-06	3.52e+06	0.000	19.1075	19.1075
bmi_obese	8.0457	1.98e-06	4.07e+06	0.000	8.046	8.046
cur_smoke	-1.4375	1.28e-06	-1.13e+06	0.000	-1.437	-1.437
bmi_obese	11.3142	7.95e-07	1.42e+07	0.000	11.314	11.314
cur_smoke	4.2158	1.37e-06	2.13e+06	0.000	4.216	4.216
bmi_obese	-2.7599	1.81e-06	-2.78e+06	0.000	-2.759	-2.759
cur_smoke	-15.5259	5.39e-06	-2.83e+06	0.000	-15.527	-15.527
b						



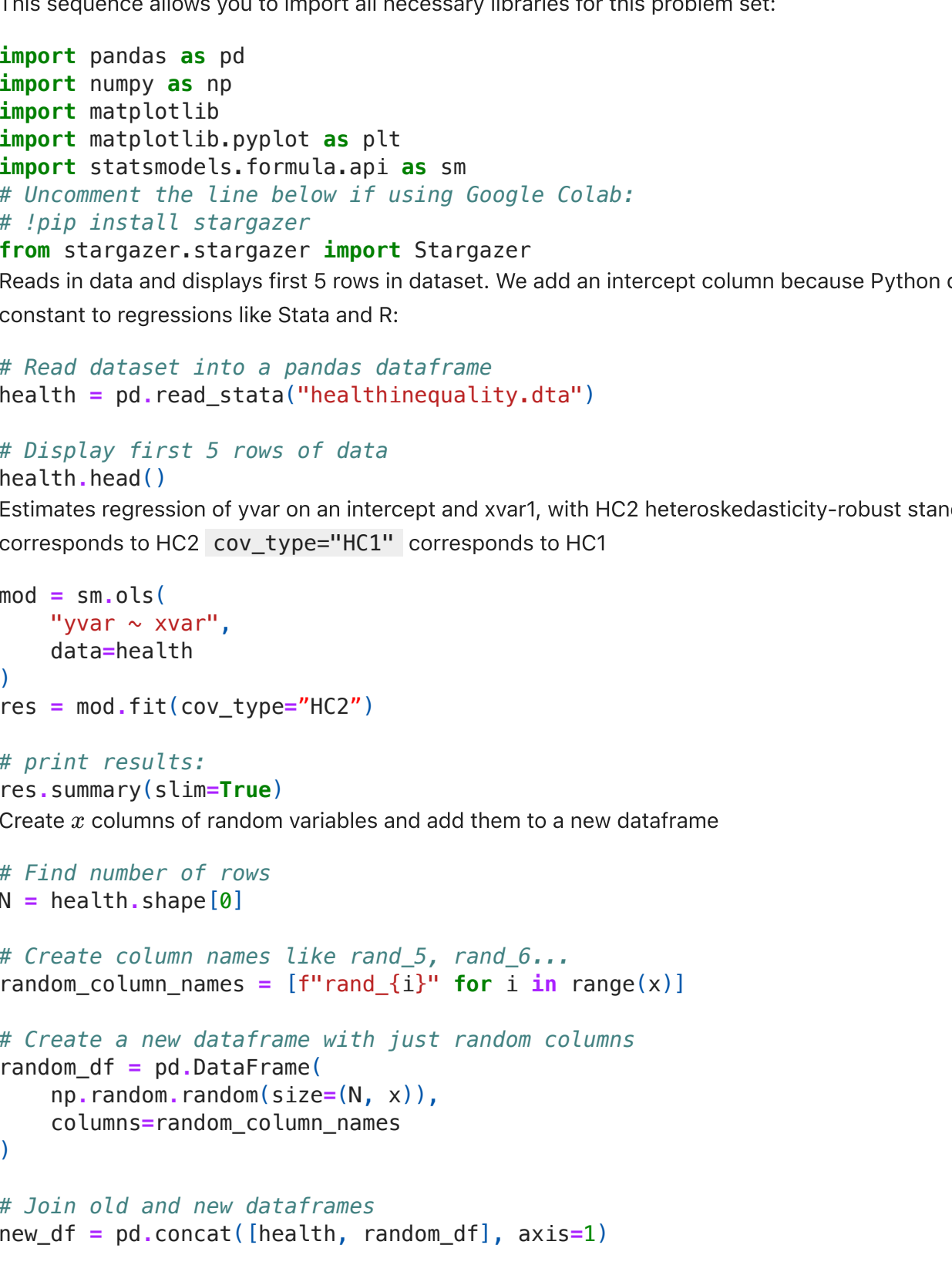
Out[137]: Text(0, 0.5, 'Cur\_smoke')



```
In [32]: # Let's plot the regression plots -> e.g. the residuals
linear_model = sm.ols("life_exp ~ cur_smoke", data=health).fit(cov_type="HC2")
fig = smf.graphics.plot_regress_exog(linear_model, "cur_smoke")
fig.tight_layout(pad=1.0)
```

eval\_env: 1

Regression Plots for cur\_smoke



## Sample Code

This sequence allows you to import all necessary libraries for this problem set:

```
import pandas as pd
import numpy as np
import matplotlib
import matplotlib.pyplot as plt
import statsmodels.formula.api as sm
# Uncomment the line below if using Google Colab:
# !pip install stargazer
from stargazer.stargazer import Stargazer
Reads in data and displays first 5 rows in dataset. We add an intercept column because Python does not automatically add a
constant to regressions like Stata and R:

# Read dataset into a pandas dataframe
health = pd.read_stata("healthinequality.dta")

# Display first 5 rows of data
health.head()
Estimates regression of year on an intercept and xvar1, with HC2 heteroskedasticity-robust standard errors. cov_type="HC2"
corresponds to HC2 cov_type="HC1" corresponds to HC1

mod = sm.ols(
    "yvar ~ xvar",
    data=health
)
res = mod.fit(cov_type="HC2")

# Print results:
res.summary(slim=True)
Create x columns of random variables and add them to a new dataframe

# Find number of rows
N = health.shape[0]

# Create column names like rand_5, rand_6,...
random_column_names = [f"rand_{i}" for i in range(x)]

# Create a new dataframe with just random columns
random_df = pd.DataFrame(
    np.random.random(size=(N, x)),
    columns=random_column_names
)

# Join old and new dataframes
new_df = pd.concat([health, random_df], axis=1)

# Extra hint: Look at the python "join" function to
# create a string out of the list of variable names
Create regression table with custom column labels

# Estimate Regressions:
mod1 = sm.ols(
    "yvar1 ~ xvar1 + xvar2 + xvar3",
    data=health
)
res1 = mod.fit(cov_type="HC2")
mod2 = sm.ols(
    "yvar1 ~ xvar1 + xvar2",
    data=health
)
res2 = mod.fit(cov_type="HC2")
mod3 = sm.ols(
    "yvar2 ~ xvar2 + xvar3",
    data=health
)
res3 = mod.fit(cov_type="HC2")

# Create Table
table = Stargazer(models)

# Label columns
# This list of ls should be the same length as the
# number of columns
table.custom_columns(["yvar1", "yvar1", "yvar2"],
separators=[1, 1, 1])

# Display table
table
Calculate residuals from regression of variable yvar on variables xvar1 and xvar2

# Estimate Regressions:
mod1 = sm.ols(
    "yvar ~ xvar1 + xvar2",
    data=health
)
res1 = mod.fit(cov_type="HC2")

# Find residuals
residuals = res1.resid
Draw a scatter plot of variable x1 against variable y1 and add a line best fit

# Sets up space for graphing
fig, ax = plt.subplots(1, 1)

# Plots scatter plot (size of points is set by s=3)
ax.scatter(x1, y1, s=3)

# Adds a line best fit for the data to the plot
ax.plot(
    np.polynomial.Polynomial.fit(
        x1, y1, 1
    ).linspace(2),
    "r--"
)

# Label axes
ax.set_title("Graph Title")
ax.set_xlabel("X Axis Label")
ax.set_ylabel("Y Axis Label")
```