

Module 2
Lesson 3 - Multiple Regression

WESLEYAN
UNIVERSITY



© Creative Commons, 2015

```
56
57
58
59
60 # center quantitative IVs for regression analysis
61 sub1['numbercigsmoked_c'] = (sub1['numbercigsmoked'] - sub1['numbercigsmoked'].mean())
62
63 print (sub1['numbercigsmoked_c'].mean())
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
```



```
In [24]: print (sub1['numbercigsmoked_c'].mean())
```

```
....:
```

```
-1.8371020199707542e-14
```

```
In [25]:
```

0.0000000000000001837

```
42 # bivariate bar graph
43 seaborn.factorplot(x="MAJORDEPLIFE", y="NDSymptoms", data=sub2, kind="bar", ci=None)
44 plt.xlabel('Major Life Depression')
45 plt.ylabel('Mean Number Nicotine Dependence Symptoms')
46
47
48 # center quantitative IVs for regression analysis
49 sub1['numbercigsmoked_c'] = (sub1['numbercigsmoked'] - sub1['numbercigsmoked'].mean())
50
51 print (sub1['numbercigsmoked_c'].mean())
```

```
53
54 reg2 = smf.ols('NDSymptoms ~ MAJORDEPLIFE + numbercigsmoked_c', data=sub1).fit()
55 print (reg2.summary())
```

```
In [25]: reg2 = smf.ols('NDSymptoms ~ MAJORDEPLIFE + numbercigsmoked_c', data=sub1).fit()
...: print (reg2.summary())
...:
```

OLS Regression Results


```
=====
Dep. Variable:          NDSymptoms    R-squared:                0.132
Model:                  OLS           Adj. R-squared:           0.131
Method:                 Least Squares  F-statistic:              99.87
Date:                   Fri, 23 Oct 2015  Prob (F-statistic):      4.28e-41
Time:                   13:47:53        Log-Likelihood:          -2593.9
No. Observations:      1313           AIC:                     5194.
Df Residuals:          1310           BIC:                     5209.
Df Model:               2
Covariance Type:        nonrobust
=====
```

	coef	std err	t	P> t	[95.0% Conf. Int.]
Intercept	2.1946	0.056	38.940	0.000	2.084 2.305
MAJORDEPLIFE	1.3424	0.109	12.327	0.000	1.129 1.556
numbercigsmoked_c	0.0358	0.006	6.432	0.000	0.025 0.047

```
=====
Omnibus:                70.355    Durbin-Watson:           2.066
Prob(Omnibus):          0.000    Jarque-Bera (JB):        50.154
Skew:                   0.372    Prob(JB):                1.29e-11
Kurtosis:               2.398    Cond. No.                20.4
=====
```

Warnings:

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

Disthymia

Number
Nicotine
Dependence
symptoms

```
reg3 = smf.ols('NDSymptoms ~ DYSLIFE', data=sub1).fit()
In [20]: reg3 = smf.ols('NDSymptoms ~ DYSLIFE', data=sub1).fit()
print (reg3.summary())
```

```
...:
OLS Regression Results
=====
Dep. Variable:      NDSymptoms      R-squared:      0.023
Model:              OLS             Adj. R-squared:  0.022
Method:             Least Squares   F-statistic:    30.35
Date:               Fri, 23 Oct 2015 Prob (F-statistic): 4.34e-08
Time:               13:49:37        Log-Likelihood: -2679.2
No. Observations:   1317            AIC:            5362.
Df Residuals:       1315            BIC:            5373.
Df Model:           1
Covariance Type:    nonrobust
=====
                    coef    std err          t      P>|t|      [95.0% Conf. Int.]
-----
Intercept          2.4785      0.053    46.958      0.000      2.375      2.582
DYSLIFE            1.1378      0.207     5.509      0.000      0.733      1.543
=====
Omnibus:            123.634    Durbin-Watson:      2.079
Prob(Omnibus):      0.000    Jarque-Bera (JB):    63.384
Skew:               0.374    Prob(JB):            1.72e-14
Kurtosis:           2.228    Cond. No.            4.07
=====
```

Warnings:

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

```
In [21]:
```



```
reg4 = smf.ols('NDSymptoms ~ DYSLIFE + MAJORDEPLIFE', data=sub1).fit()
In [27]: reg4 = smf.ols('NDSymptoms ~ DYSLIFE + MAJORDEPLIFE', data=sub1).fit()
print (reg4.summary())
```

OLS Regression Results

```
=====
Dep. Variable:          NDSymptoms      R-squared:                0.106
Model:                  OLS             Adj. R-squared:           0.105
Method:                 Least Squares    F-statistic:              78.23
Date:                  Fri, 23 Oct 2015   Prob (F-statistic):       7.98e-33
Time:                  13:53:26          Log-Likelihood:          -2620.2
No. Observations:      1317             AIC:                     5246.
Df Residuals:          1314             BIC:                     5262.
Df Model:               2
Covariance Type:       nonrobust
=====
```

	coef	std err	t	P> t	[95.0% Conf. Int.]
Intercept	2.1808	0.057	38.152	0.000	2.069 2.293
DYSLIFE	0.3477	0.210	1.656	0.098	-0.064 0.760
MAJORDEPLIFE	1.2993	0.117	11.103	0.000	1.070 1.529

```
=====
Omnibus:                70.772      Durbin-Watson:           2.061
Prob(Omnibus):          0.000      Jarque-Bera (JB):        51.246
Skew:                   0.380      Prob(JB):                7.45e-12
Kurtosis:               2.402      Cond. No.                 4.61
=====
```

Warnings:

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

```
In [28]:
```



```
reg5 = smf.ols('NDSymptoms ~ DYSLIFE + MAJORDEPLIFE + numcigsmoked_c + age_c + SEX', data=sub1).fit()
print (reg4.summary())
```

OLS Regression Results

```
=====
Dep. Variable:      NDSymptoms      R-squared:      0.136
Model:              OLS             Adj. R-squared: 0.133
Method:             Least Squares   F-statistic:    41.08
Date:               Fri, 23 Oct 2015 Prob (F-statistic): 2.39e-39
Time:               13:56:07         Log-Likelihood: -2591.2
No. Observations:   1313            AIC:            5194.
Df Residuals:       1307            BIC:            5226.
Df Model:           5
Covariance Type:    nonrobust
=====
```

	coef	std err	t	P> t	[95.0% Conf. Int.]
Intercept	2.2550	0.156	14.480	0.000	1.949 2.561
DYSLIFE	0.2746	0.209	1.316	0.188	-0.135 0.684
MAJORDEPLIFE	1.2975	0.116	11.161	0.000	1.069 1.526
numcigsmoked_c	0.0353	0.006	6.257	0.000	0.024 0.046
age_c	-0.0400	0.022	-1.806	0.071	-0.083 0.003
SEX	-0.0439	0.099	-0.442	0.658	-0.238 0.151

```
=====
Omnibus:            69.558      Durbin-Watson:      2.075
Prob(Omnibus):      0.000      Jarque-Bera (JB):    48.596
Skew:               0.361      Prob(JB):           2.80e-11
Kurtosis:           2.394      Cond. No.           38.5
=====
```

Warnings:

[1] Standard Errors assume that the covariance matrix of the parameters is correct.

```
reg5 = smf.ols('NDSymptoms ~ DYSLIFE + MAJORDEPLIFE + numcigsmoked_c + age_c + SEX', data=sub1).fit()
print (reg5.summary())
```

OLS Regression Results

```
=====
Dep. Variable:      NDSymptoms      R-squared:      0.136
Model:              OLS             Adj. R-squared: 0.133
Method:             Least Squares   F-statistic:    41.88
Date:               Fri, 23 Oct 2015 Prob (F-statistic): 2.39e-39
Time:               13:56:07         Log-Likelihood: -2591.2
No. Observations:   1313            AIC:            5194.
Df Residuals:       1307            BIC:            5226.
Df Model:           5
Covariance Type:    nonrobust
=====
```

	coef	std err	t	P> t	[95.0% Conf. Int.]
Intercept	2.2550	0.156	14.480	0.000	1.949 2.561
DYSLIFE	0.2746	0.209	1.316	0.188	-0.135 0.684
MAJORDEPLIFE	1.2975	0.116	11.161	<u>0.000</u>	1.069 1.526
numcigsmoked_c	0.0353	0.006	6.257	<u>0.000</u>	0.024 0.046
age_c	-0.0400	0.022	-1.806	0.071	-0.083 0.003
SEX	-0.0439	0.099	-0.442	0.658	-0.238 0.151

```
=====
Omnibus:            89.558      Durbin-Watson:      2.075
Prob(Omnibus):      0.000      Jarque-Bera (JB):    48.596
Skew:               0.361      Prob(JB):            2.80e-11
Kurtosis:           2.394      Cond. No.             10.5
=====
```

Warnings:

Module 2
Lesson 4 - Confidence Intervals

WESLEYAN
UNIVERSITY



© Creative Commons, 2015

OLS Regression Results

```

=====
Dep. Variable:      NDSymptoms      R-squared:      0.136
Model:              OLS             Adj. R-squared:  0.133
Method:             Least Squares   F-statistic:    41.08
Date:               Fri, 23 Oct 2015 Prob (F-statistic): 2.38e-39
Time:               13:56:07         Log-Likelihood: -2591.2
No. Observations:   1313            AIC:            5194.
Df Residuals:       1307            BIC:            5226.
Df Model:           5
Covariance Type:    nonrobust
=====

```

	coef	std err	t	P> t	[95.0% Conf. Int.]	
Intercept	2.2550	0.156	14.400	0.000	1.949	2.561
OYSLIFE	0.2746	0.209	1.316	0.188	-0.135	0.684
MAJORDEPLIFE	1.2975	0.116	11.161	0.000	1.069	1.526
numbercigs smoked_c	0.0353	0.006	6.257	0.000	0.024	0.046
age_c	-0.0400	0.022	-1.806	0.071	-0.083	0.003
SEX	-0.0439	0.099	-0.442	0.658	-0.238	0.151

```

=====
Omnibus:           69.558      Durbin-Watson:      1.979
Prob(Omnibus):     0.000      Jarque-Bera (JB):   48.596
Skew:              0.361      Prob(JB):           2.80e-11
Kurtosis:          2.394      Cond. No.           38.5
=====

```

Warnings:

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


```

...:
                                OLS Regression Results
=====
Dep. Variable:                  NDSymptoms    R-squared:                  0.136
Model:                          OLS          Adj. R-squared:             0.133
Method:                        Least Squares  F-statistic:                41.08
Date:                          Fri, 23 Oct 2015 Prob (F-statistic):         2.39e-39
Time:                           13:56:07      Log-Likelihood:            -2591.2
No. Observations:                1313         AIC:                       5194.
Df Residuals:                    1307         BIC:                       5226.
Df Model:                         5
Covariance Type:                 nonrobust
=====

```

	coef	std err	t	P> t	[95.0% Conf. Int.]
Intercept	2.2550	0.156	14.480	0.000	1.949 2.561
DYSLIFE	0.2746	0.209	1.316	0.188	-0.135 0.684
MAJORDEPLIFE	1.2975	0.116	11.161	0.000	1.069 1.526
numercigsmoked_c	0.0353	0.006	6.257	0.000	0.024 0.046
age_c	-0.0400	0.022	-1.806	0.071	-0.083 0.003
SEX	-0.0439	0.099	-0.442	0.658	-0.238 0.151

```

=====
Omnibus:                        69.558    Durbin-Watson:              2.075
Prob(Omnibus):                  0.000    Jarque-Bera (JB):           48.596
Skew:                           0.361    Prob(JB):                   2.80e-11
Kurtosis:                       2.394    Cond. No.                    38.5
=====

```

Warnings:

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

Module 3
Lesson 3 - Polynomial Regression

WESLEYAN
UNIVERSITY



© Creative Commons, 2015



In [3]:

Console

IPython console

History log

Permissions: RM

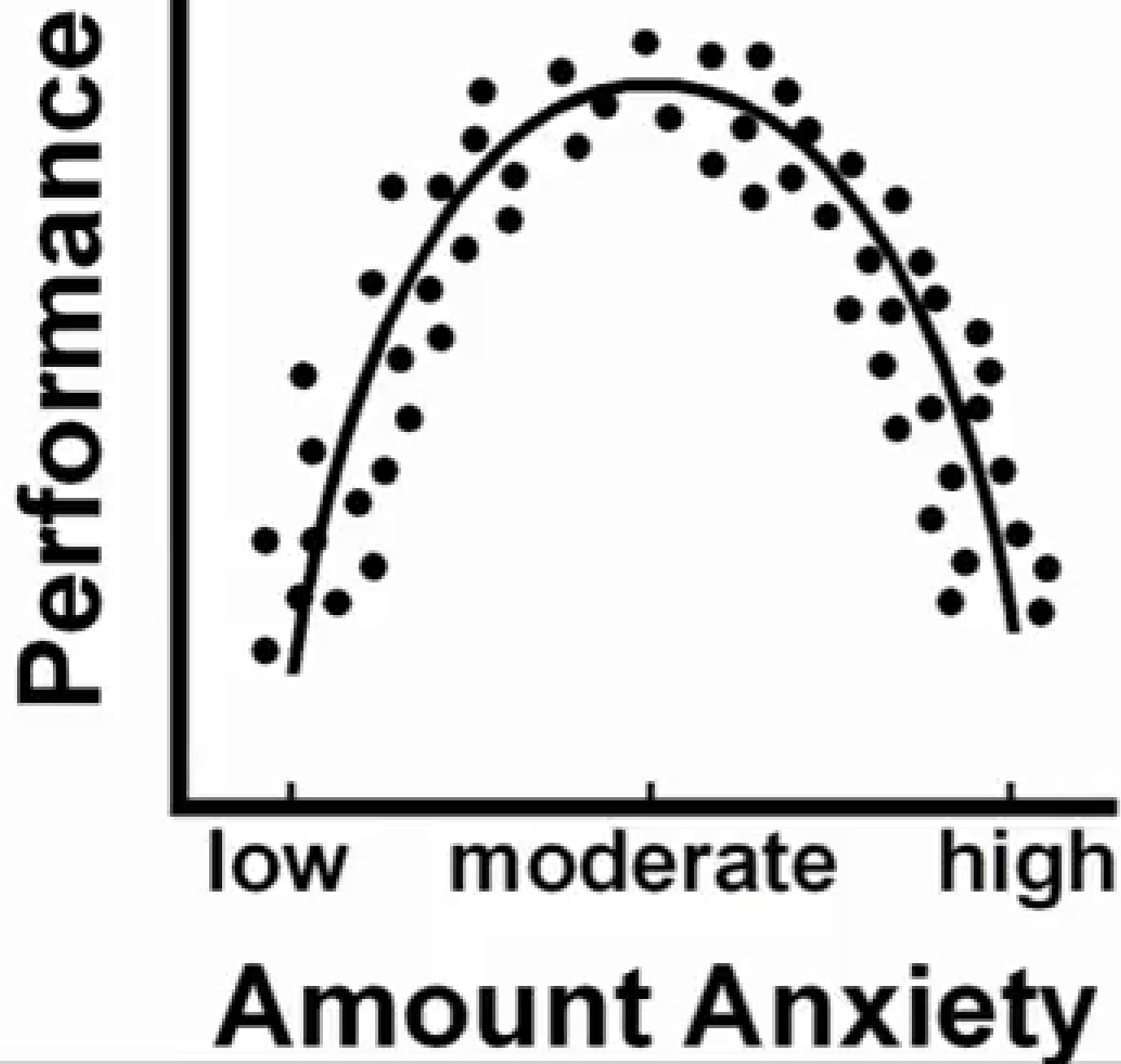
End-of-lines: CRLF

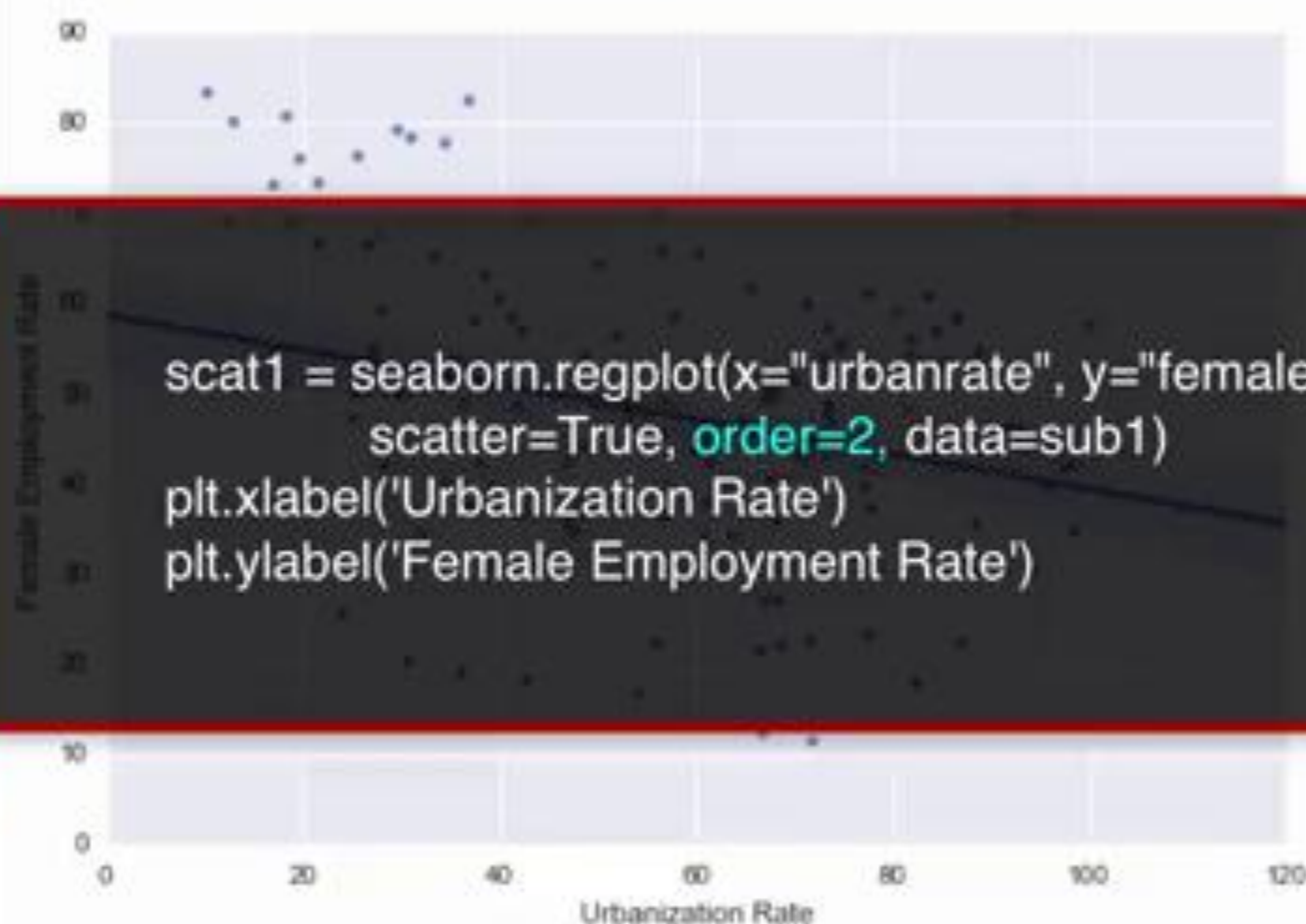
Encoding: UTF-8

Line: 27

Column: 1

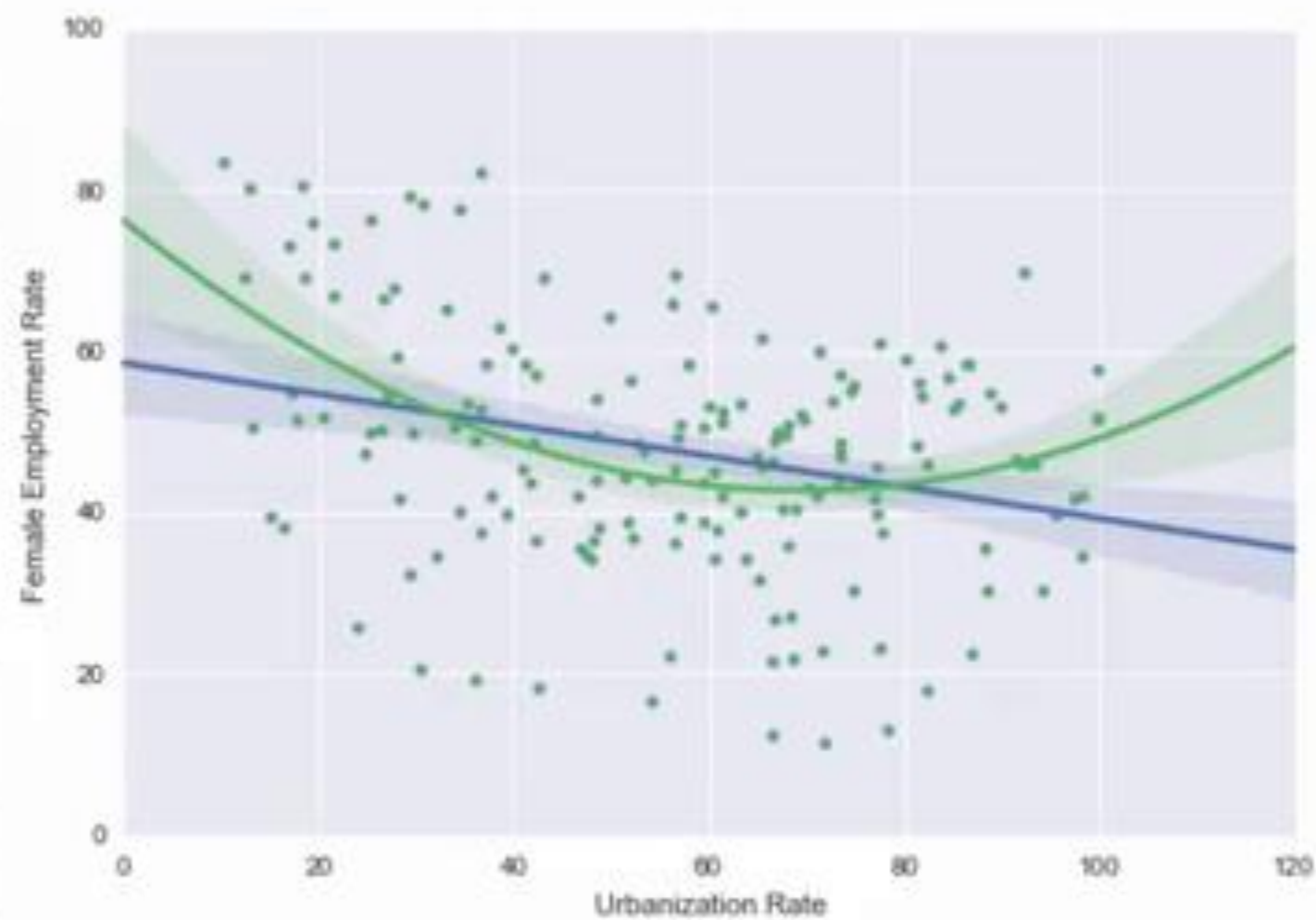
Memory: 83 %





```
scat1 = seaborn.regplot(x="urbanrate", y="femaleemployrate",  
                        scatter=True, order=2, data=sub1)  
plt.xlabel('Urbanization Rate')  
plt.ylabel('Female Employment Rate')
```

In [14]:



In [16]:

```
23 del I
24
25 sub1 = data[['urbanrate', 'femaleemployrate', 'internetuserate']].dropna()
26
27
28 scat1 = seaborn.regplot(x="urbanrate", y="femaleemployrate", scatter=True, data=sub1)
29 plt.xlabel('Urbanization Rate')
30 plt.ylabel('Female Employment Rate')
31
32 scat1 = seaborn.regplot(x="urbanrate", y="femaleemployrate", scatter=True, order=2, data=sub1)
33 plt.xlabel('Urbanization Rate')
34 plt.ylabel('Female Employment Rate')
35
36 # center quantitative IVs for regression analysis
37 sub1['urbanrate_c'] = (sub1['urbanrate'] - sub1['urbanrate'].mean())
38 sub1['internetuserate_c'] = (sub1['internetuserate'] - sub1['internetuserate'].mean())
39
40 reg1 = smf.ols('femaleemployrate ~ urbanrate_c', data=sub1).fit()
41 print (reg1.summary())
42
43
44
45
46
47
48
49
50
51
52
53
54
55
```

```
In [21]: reg1 = smf.ols('femaleemployrate ~ urbanrate_c', data=sub1).fit()  
...: print (reg1.summary())
```

OLS Regression Results

```
=====
```

Dep. Variable:	femaleemployrate	R-squared:	0.092
Model:	OLS	Adj. R-squared:	0.086
Method:	Least Squares	F-statistic:	16.69
Date:	Fri, 23 Oct 2015	Prob (F-statistic):	6.84e-05
Time:	14:41:44	Log-Likelihood:	-678.68
No. Observations:	167	AIC:	1361.
Df Residuals:	165	BIC:	1368.
Df Model:	1		
Covariance Type:	nonrobust		

```
=====
```

	coef	std err	t	P> t	[95.0% Conf. Int.]
Intercept	47.6024	1.096	43.416	0.000	45.438 49.767
urbanrate_c	-0.1927	0.047	-4.086	0.000	-0.286 -0.100

```
=====
```

```
Omnibus:                2.347    Durbin-Watson:           1.868  
Prob(Omnibus):           0.309    Jarque-Bera (JB):         2.409  
Skew:                    -0.269    Prob(JB):                 0.300  
Kurtosis:                2.763    Cond. No.                 23.2  
=====
```

Warnings:

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

```
In [22]:
```



```
In [21]: reg1 = smf.ols('femaleemployrate ~ urbanrate_c', data=sub1).fit()  
...: print (reg1.summary())
```

OLS Regression Results

```
=====
```

Dep. Variable:	femaleemployrate	R-squared:	0.092
Model:	OLS	Adj. R-squared:	0.086
Method:	Least Squares	F-statistic:	16.69
Date:	Fri, 23 Oct 2015	Prob (F-statistic):	6.84e-05
Time:	14:41:44	Log-Likelihood:	-678.68
No. Observations:	167	AIC:	1361.
Df Residuals:	165	BIC:	1368.
Df Model:	1		
Covariance Type:	nonrobust		

```
=====
```

	coef	std err	t	P> t	[95.0% Conf. Int.]
Intercept	47.6024	1.096	43.416	0.000	45.438 49.767
urbanrate_c	-0.1927	0.047	-4.086	0.000	-0.286 -0.100

```
=====
```

```
Omnibus:                2.347    Durbin-Watson:           1.868  
Prob(Omnibus):          0.309    Jarque-Bera (JB):        2.409  
Skew:                   -0.269    Prob(JB):                0.300  
Kurtosis:               2.763    Cond. No.                23.2  
=====
```

Warnings:

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

```
In [22]:
```

```
31 scat1 = seaborn.regplot(x="urbanrate", y="femaleemployrate", scatter=True, order=2, data=sub1)
32 plt.xlabel('Urbanization Rate')
33 plt.ylabel('Female Employment Rate')
34
35
36 # center quantitative IVs for regression analysis
37 sub1['urbanrate_c'] = (sub1['urbanrate'] - sub1['urbanrate'].mean())
38 sub1['internetuserate_c'] = (sub1['internetuserate'] - sub1['internetuserate'].mean())
39
40 reg1 = smf.ols('femaleemployrate ~ urbanrate_c', data=sub1).fit()
41 print (reg1.summary())
42
43
44 # regression model with second order polynomial (quadratic term)
45 reg2 = smf.ols('femaleemployrate ~ urbanrate_c + I(urbanrate_c**2)', data=sub1).fit()
46 print (reg2.summary())
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
```

OLS Regression Results

```

=====
Dep. Variable:      femaleemployrate    R-squared:      0.160
Model:              OLS                 Adj. R-squared:  0.150
Method:             Least Squares       F-statistic:    15.60
Date:               Fri, 23 Oct 2015     Prob (F-statistic): 6.30e-07
Time:               14:45:30            Log-Likelihood: -672.19
No. Observations:   167                 AIC:            1350.
Df Residuals:       164                 BIC:            1360.
Df Model:           2
Covariance Type:    nonrobust
=====

```

	coef	std err	t	P> t	[95.0% Conf. Int.]
Intercept	43.8428	1.478	29.659	0.000	40.924 46.762
urbanrate_c	-0.1751	0.046	-3.827	0.000	-0.266 -0.085
I(urbanrate_c ** 2)	0.0070	0.002	3.641	0.000	0.003 0.011

```

=====
Omnibus:            3.627    Durbin-Watson:      1.898
Prob(Omnibus):      0.163    Jarque-Bera (JB):  3.677
Skew:               -0.351    Prob(JB):          0.159
Kurtosis:           2.811    Cond. No.          1.08e+03
=====

```

Warnings:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The condition number is large, 1.08e+03. This might indicate that there are strong multicollinearity or other numerical problems.

In [24]:

OLS Regression Results

```

=====
Dep. Variable:      femaleemployrate    R-squared:                0.160
Model:              OLS                 Adj. R-squared:           0.150
Method:             Least Squares       F-statistic:              15.60
Date:               Fri, 23 Oct 2015    Prob (F-statistic):       6.30e-07
Time:               14:45:30           Log-Likelihood:           -672.19
No. Observations:   167                AIC:                      1350.
Df Residuals:       164                BIC:                      1360.
Df Model:           2
Covariance Type:    nonrobust
=====

```

	coef	std err	t	P> t	[95.0% Conf. Int.]
Intercept	43.8428	1.478	29.659	0.000	40.924 46.762
urbanrate_c	-0.1751	0.046	-3.827	0.000	-0.266 -0.085
I(urbanrate_c ** 2)	0.0070	0.002	3.641	0.000	0.003 0.011

```

=====
Omnibus:            3.627    Durbin-Watson:           1.898
Prob(Omnibus):      0.163    Jarque-Bera (JB):         3.677
Skew:               -0.351    Prob(JB):                 0.159
Kurtosis:           2.811    Cond. No.                  1.08e+03
=====

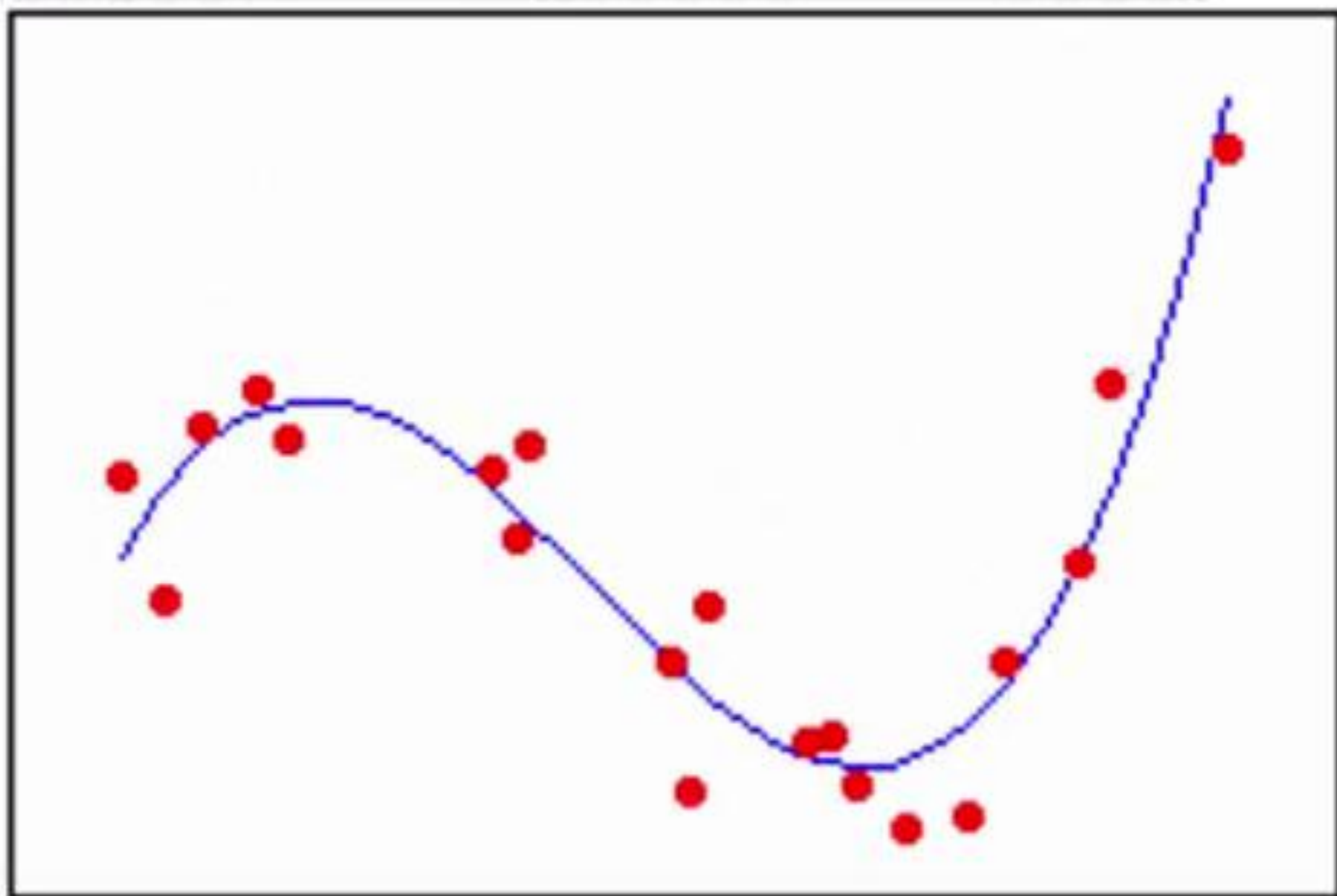
```

Warnings:

- [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [2] The condition number is large, 1.08e+03. This might indicate that there are strong multicollinearity or other numerical problems.

In [24]:

OLS Regression Results



In [24]: |

Permissions: RW

End-of-lines: CRLF

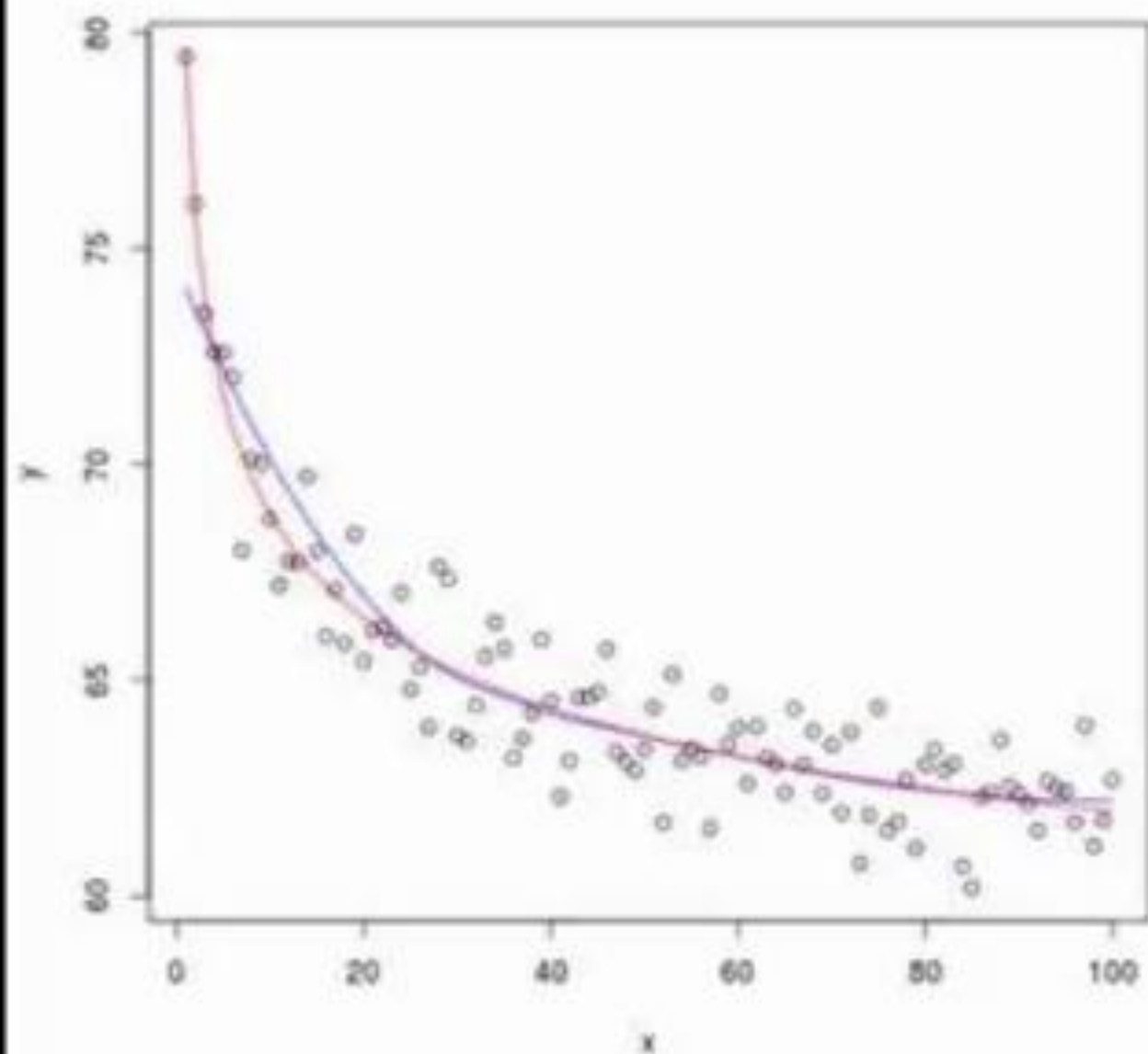
Encoding: UTF-8

Line: 58

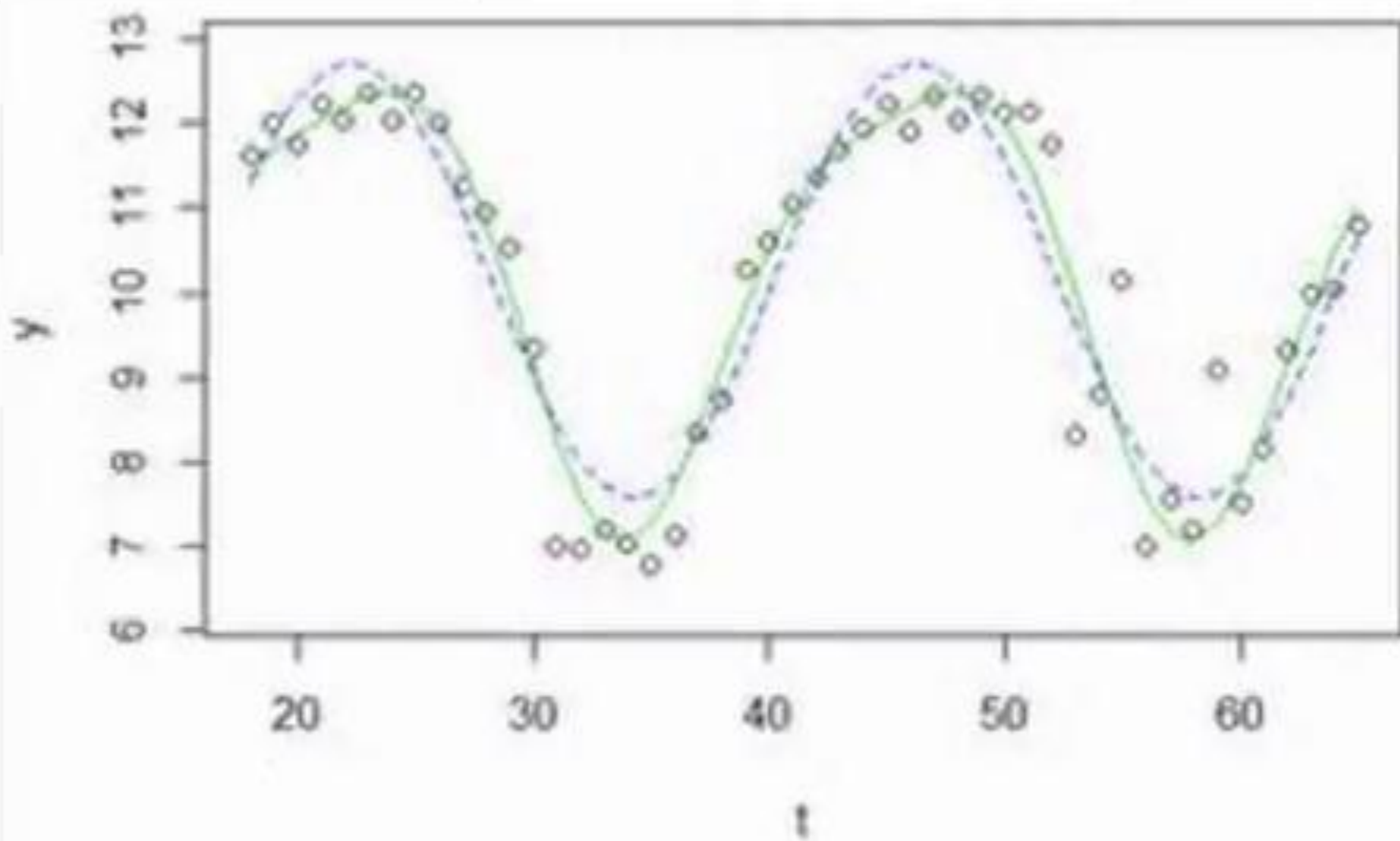
Column: 1

Memory: 22 %

Overfitting



Overfitting



Bias-Variance Tradeoff



Module 3
Lesson 4 - Evaluating Model Fit, part 1

WESLEYAN
UNIVERSITY



© Creative Commons, 2015

Misspecification

```
# Plot linear line
plot1 = scatter.regplot(x='urbanrate', y='femaleemployrate', scatter=True, data=sub1)
plt.xlabel('Urbanization Rate')
plt.ylabel('Female Employment Rate')

# Plot second order polynomial line
plot2 = scatter.regplot(x='urbanrate', y='femaleemployrate', scatter=True, order=2, data=sub1)
plt.xlabel('Urbanization Rate')
plt.ylabel('Female Employment Rate')

# Linear regression coefficients for regression analysis
sub1['urbanrate_s'] = (sub1['urbanrate'] - sub1['urbanrate'].mean())
sub1['interurbanrate_s'] = (sub1['interurbanrate'] + sub1['urbanrate_s'])
sub1[['urbanrate_s', 'interurbanrate_s']].dropna(inplace=True)

reg1 = smf.ols('femaleemployrate ~ urbanrate_s', data=sub1).fit()
print(reg1.summary())

# Regression model with second order polynomial (quadratic) term
reg2 = smf.ols('femaleemployrate ~ urbanrate_s + I(urbanrate_s**2)', data=sub1).fit()
print(reg2.summary())

# Adding interaction and other explanatory variables
reg3 = smf.ols('femaleemployrate ~ urbanrate_s + I(urbanrate_s**2) + interurbanrate_s',
              data=sub1).fit()
print(reg3.summary())
```

Usage

Here you can get help of any object by pressing **Object** in front of it, either in the Editor or the Console.

Help can also be shown automatically after writing a left parenthesis next to an object. You can activate this behavior in **Preferences > Object Inspector**.

New to Spyder? Read our [tutorial](#).

Object Inspector Variable Explorer File Explorer

Python console

Code Editor

OLS Regression Results

Dep. Variable:	femaleemployrate	R-squared:	0.186
Model:	OLS	Adj. R-squared:	0.165
Method:	Least Squares	F-statistic:	11.92
Date:	Fri, 23 Oct 2015	Prob (F-statistic):	4.25e-07
Time:	17:29:39	Log-Likelihood:	-679.17
No. Observations:	167	AIC:	1348.
Df Residuals:	163	BIC:	1361.
Df Model:	3		
Covariance Type:	nonrobust		

	coef	std err	t	P> t	10.0% Conf. Int.	Int.
Intercept	43.9886	1.457	29.991	0.000	41.092	46.885
urbanrate_s	-8.2008	0.602	-4.186	0.000	-9.383	-6.117
I(urbanrate_s ** 2)	0.0067	0.002	3.523	0.001	0.003	0.010
interurbanrate_s	0.1038	0.032	3.000	0.003	0.041	0.166

Omnibus:	2.937	Durbin-Watson:	1.893
Prob(Omnibus):	0.231	Jarque-Bera (JB):	2.000
Skew:	-0.264	Prob(JB):	0.368
Kurtosis:	3.905	cond. No.	1.09e+03

Warnings:

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

Specification: The process of developing a regression model

```
# Import libraries
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np
import pandas as pd
import statsmodels.api as sm
import statsmodels.formula.api as smf

# Load data
data = pd.read_csv('data/urbanrate.csv')

# Visualize data
sns.scatterplot(x='urbanrate', y='femaleemployrate', data=data)
plt.xlabel('Urbanization Rate')
plt.ylabel('Female Employment Rate')

# Create a linear regression model
reg1 = smf.ols('femaleemployrate ~ urbanrate', data=data).fit()
print(reg1.summary())

# Create a polynomial regression model
reg2 = smf.ols('femaleemployrate ~ urbanrate + I(urbanrate**2)', data=data).fit()
print(reg2.summary())

# Add interaction and quadratic terms
reg3 = smf.ols('femaleemployrate ~ urbanrate + I(urbanrate**2) + internetuse_rate', data=data).fit()
print(reg3.summary())
```

Usage

Here you can get help of any object by pressing **Ctrl** or **Cmd** + **?**, either in the Editor or the Console.

Help can also be shown automatically after writing a left parenthesis next to an object. You can activate this behavior in **Preferences > Object Inspector**.

New to Spyder? Read our [tutorial](#).

Object Inspector Variable Explorer File Explorer

Python console

Code Editor

OLS Regression Results

Dep. Variable:	femaleemployrate	R-squared:	0.188
Model:	OLS	Adj. R-squared:	0.165
Method:	Least Squares	F-statistic:	11.92
Date:	Fri, 23 Oct 2015	Prob (F-statistic):	4.25e-07
Time:	17:29:39	Log-Likelihood:	-679.17
No. Observations:	167	AIC:	1348.
Df Residuals:	163	BIC:	1361.
Df Model:	3		
Covariance Type:	nonrobust		

	coef	std err	t	P> t	10% Conf. Int.	Int.
Intercept	43.9886	1.457	29.991	0.000	41.092	46.885
urbanrate_c	-8.2008	0.602	-4.186	0.000	-9.383	-6.117
I(urbanrate_c ** 2)	0.0067	0.002	3.523	0.001	0.003	0.010
internetuse_rate	0.1038	0.032	3.000	0.003	0.040	0.166

Omnibus:	2.937	Durbin-Watson:	1.893
Prob(Omnibus):	0.231	Jarque-Bera (JB):	2.000
Skew:	-0.264	Prob(JB):	0.368
Kurtosis:	3.905	Cond. No.	1.09e+03

Warnings:

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

Console Python console File explorer

```

84 # center quantitative 1s for regression analysis
85 sub1['urbanrate_c'] = (sub1['urbanrate'] - sub1['urbanrate'].mean())
86 sub1['internetuserate_c'] = (sub1['internetuserate'] - sub1['internetuserate'].mean())
87 sub1[['urbanrate_c', 'internetuserate_c']].describe()
88
89 reg1 = smf.ols('femaleemployrate ~ urbanrate_c', data=sub1).fit()
90 print (reg1.summary())
91
92 # regression model with second order polynomial (quadratic term)
93 reg2 = smf.ols('femaleemployrate ~ urbanrate_c + I(urbanrate_c**2)', data=sub1).fit()
94 print (reg2.summary())
95
96 # adding internet use rate explanatory variable
97 reg3 = smf.ols('femaleemployrate ~ urbanrate_c + I(urbanrate_c**2) + internetuserate_c',
98               data=sub1).fit()
99 print (reg3.summary())

```

DFFITS Statistic

Bartlett's Test

Kolmogorov-Smirnoff Test

Durbin-Watson Test

Variance Inflation Factor

```

84 # center quantitative IVs for regression analysis
85 sub1['urbanrate_c'] = (sub1['urbanrate'] - sub1['urbanrate'].mean())
86 sub1['internetuserate_c'] = (sub1['internetuserate'] - sub1['internetuserate'].mean())
87 sub1[['urbanrate_c', 'internetuserate_c']].describe()
88
89 reg1 = smf.ols('femaleemployrate ~ urbanrate_c', data=sub1).fit()
90 print (reg1.summary())
91
92 # regression model with second order polynomial (quadratic term)
93 reg2 = smf.ols('femaleemployrate ~ urbanrate_c + I(urbanrate_c**2)', data=sub1).fit()
94 print (reg2.summary())
95
96 # adding internet use rate explanatory variable
97 reg3 = smf.ols('femaleemployrate ~ urbanrate_c + I(urbanrate_c**2) + internetuserate_c',
98               data=sub1).fit()
99 print (reg3.summary())

```

$$\text{Femaleemployrate} = \beta_0 + \beta_1 (\text{urbanrate}) + \beta_2 (\text{urbanrate})^2 + \beta_3 (\text{internetuserate}) + \epsilon$$

OLS Regression Results

```
=====
Dep. Variable:      femaleemployrate    R-squared:                0.180
Model:              OLS                 Adj. R-squared:           0.165
Method:             Least Squares       F-statistic:              11.92
Date:              Fri, 23 Oct 2015     Prob (F-statistic):      4.25e-07
Time:              17:29:30             Log-Likelihood:          -670.17
No. Observations:   167                 AIC:                     1348.
Df Residuals:       163                 BIC:                     1361.
Df Model:           3
Covariance Type:    nonrobust
=====
```

	coef	std err	t	P> t	[95.0% Conf. Int.]	
Intercept	43.9886	1.467	29.991	0.000	41.092	46.885
urbanrate_c	-0.2600	0.062	-4.186	0.000	-0.383	-0.137
I(urbanrate_c ** 2)	0.0067	0.002	3.523	0.001	0.003	0.010
internetuserate_c	0.1038	0.052	2.000	0.047	0.001	0.206

```
=====
Omnibus:            2.037    Durbin-Watson:           1.893
Prob(Omnibus):      0.361    Jarque-Bera (JB):       2.000
Skew:               -0.264    Prob(JB):               0.368
Kurtosis:           2.905    Cond. No.                1.09e+03
=====
```

Warnings:

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



Console

IPython console

History log


```
97 reg3 = smf.ols('femaleemployrate ~ urbanrate_c + I(urbanrate_c**2) + internetuserate_c',
98               data=sub1).fit()
99 print (reg3.summary())
100
101
102
103 #####
104 # Regression diagnostic plots
105 #####
106
107 #Q-Q plot for normality
108 fig1=sm.qqplot(reg3.resid, line='r')
109
110
111
112
113
114
115
116
117
118
119
120
121
122
```

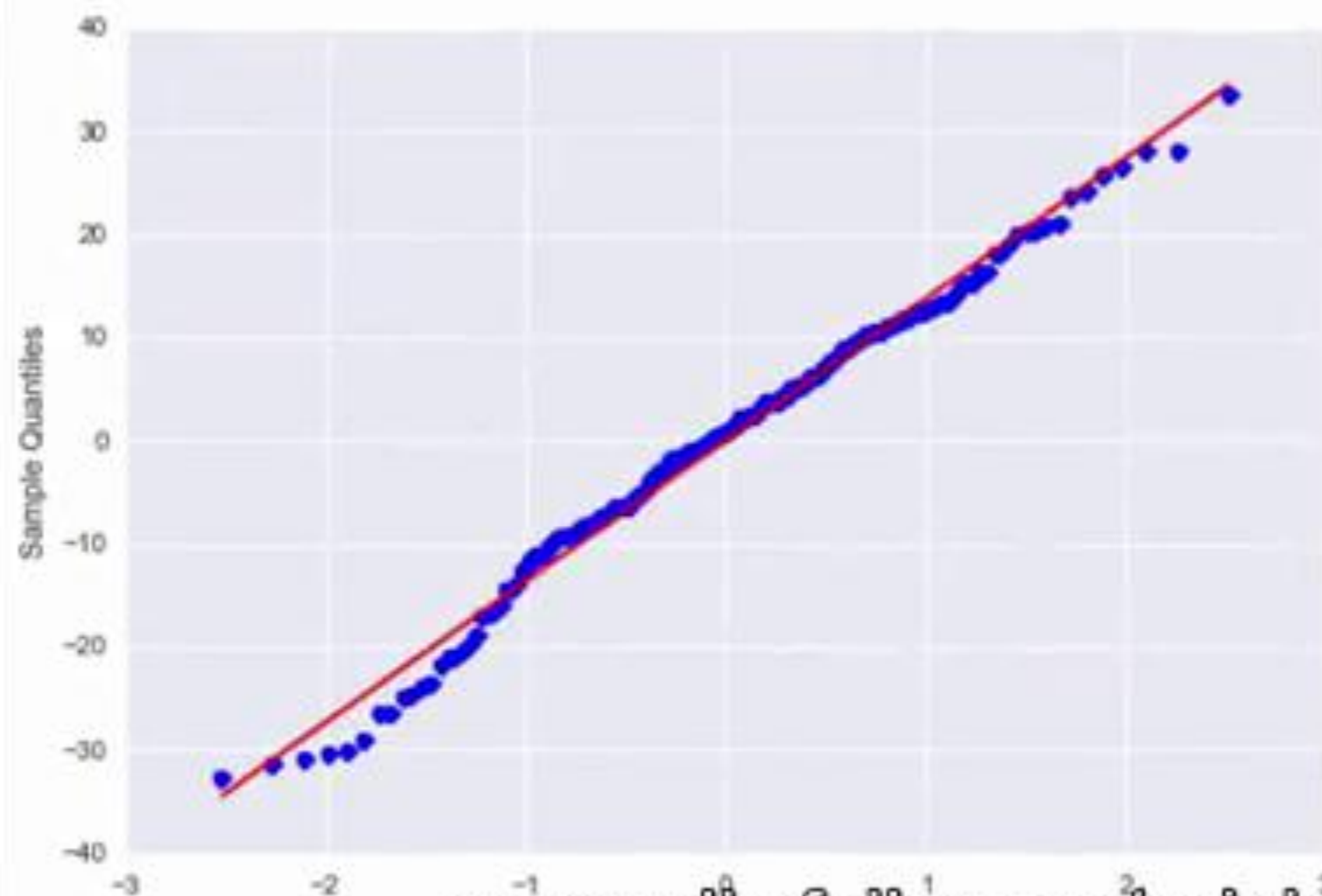
What we're looking for is to see if
the points follow a straight line.

```
97 reg3 = smf.ols('femaleemployrate ~ urbanrate_c + I(urbanrate_c**2) + internetuserate_c',
98               data=sub1).fit()
99 print (reg3.summary())
100
101
102
103 #####
104 # Regression diagnostic plots
105 #####
106
107 #Q-Q plot for normality
108 fig1=sm.qqplot(reg3.resid, line='r')
109
110
111
112
113
114
115
116
117
118
119
120
121
122
```

Meaning that the model estimated residuals are what we would expect

```
97 reg3 = smf.ols('femaleemployrate ~ urbanrate_c + I(urbanrate_c**2) + internetuserate_c',
98               data=sub1).fit()
99 print (reg3.summary())
100
101
102
103 #####
104 # Regression diagnostic plots
105 #####
106
107 #Q-Q plot for normality
108 fig1=sm.qqplot(reg3.resid, line='r')
109
110
111
112
113
114
115
116
117
118
119
120
121
122
```

if the residuals were
normally distributed.



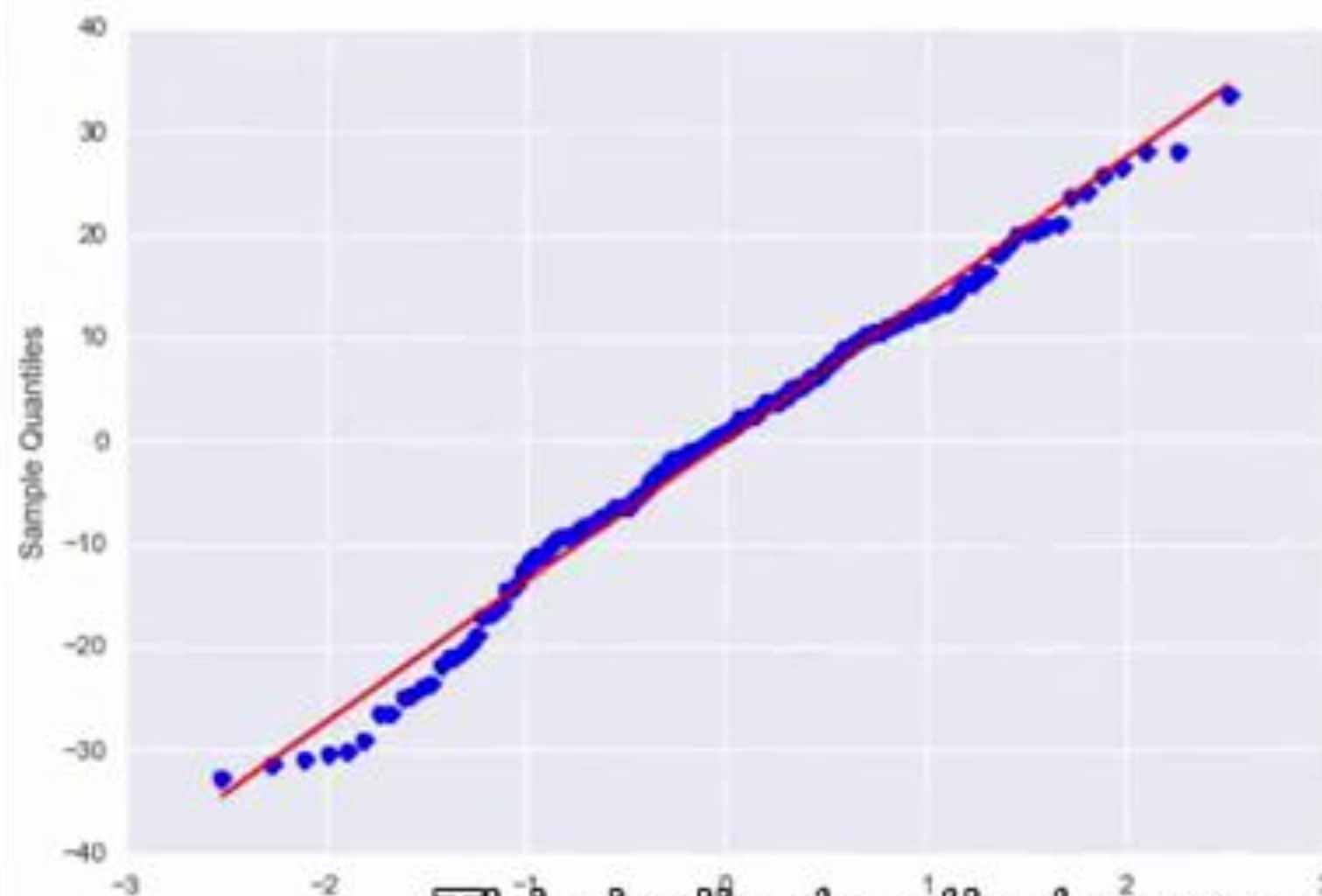
generally follow a straight line, but deviate at the lower and higher quantiles.

In [95]:

Console

IPython console

History log



This indicates that our residuals did not follow perfect normal distribution.

In [95]:

Console

IPython console

History log

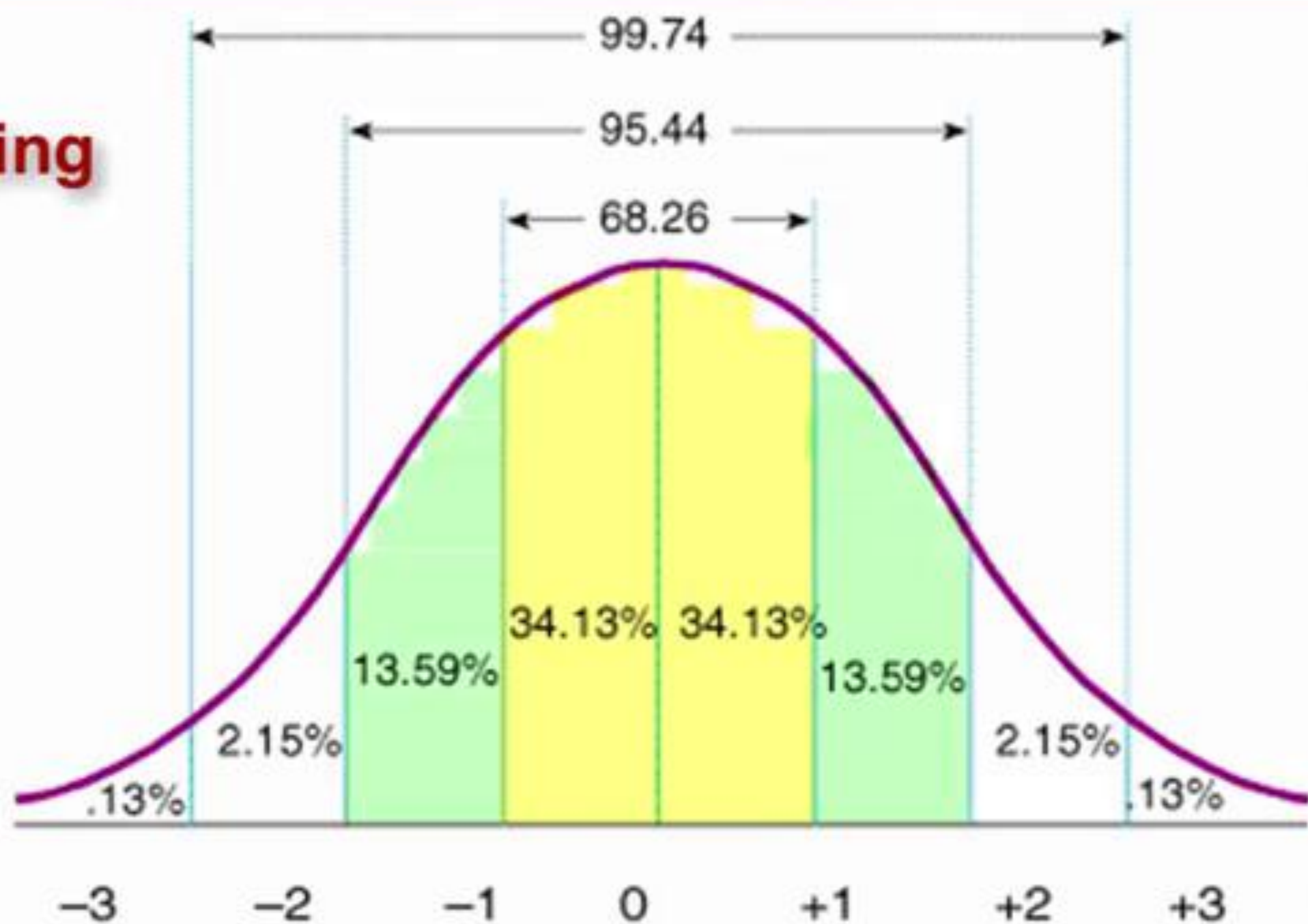
Module 3
Lesson 5 - Evaluating Model Fit, part 2

WESLEYAN
UNIVERSITY



© Creative Commons, 2015

Normalizing



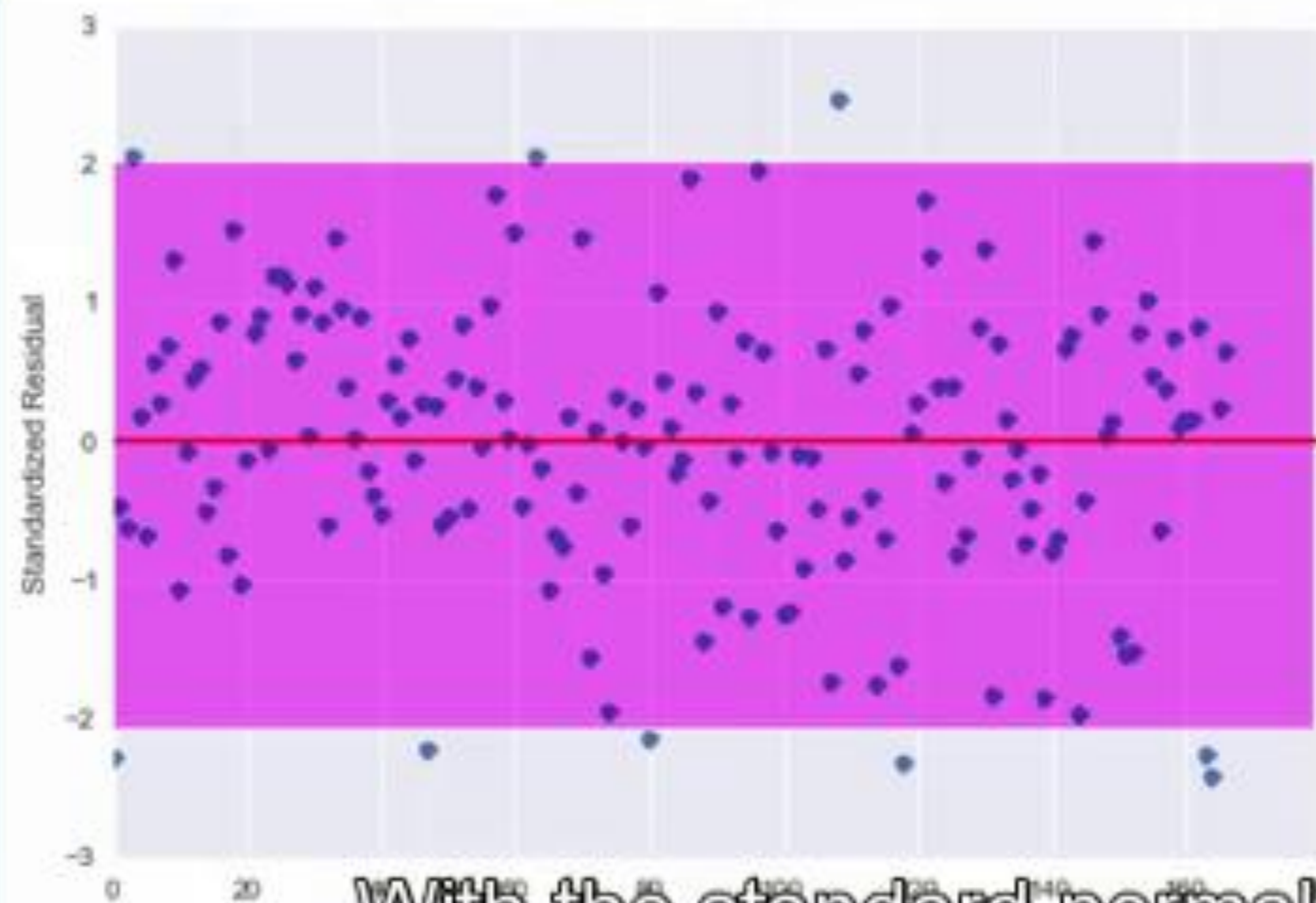
```
100
101
102
103 ##### Regression Diagnostic plots #####
104 # Regression Diagnostic plots
105 #####
106
107 # qq plot for normality
108 fig1=sm.qqplot(reg3.resid, line='r')
109
110
111 stdres=pandas.DataFrame(reg3.resid_pearson)
112 fig2=plt.plot(stdres, 'o', ls='none')
113 l = plt.axhline(y=0, color='r')
114 plt.ylabel('Standardized Residual')
115 plt.xlabel('Observation Number')
116 print(fig2)
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
```

```
14 print(fig1)
15
16
17 # simple plot of residuals
18 stdres=pandas.DataFrame(reg3.resid_pearson)
19 fig2=plt.plot(stdres, 'o', ls='none')
20 l = plt.axhline(y=0, color='r')
21 plt.ylabel('Standardized Residual')
22 plt.xlabel('Observation Number')
23 print(fig2)
```

```
24
25 # additional regression diagnostic plots
26 fig3 = plt.figure(figsize(12,8))
27 fig3 = sm.graphics.plot_regress_exog(reg3, "internetuserate_c", fig=fig3)
```

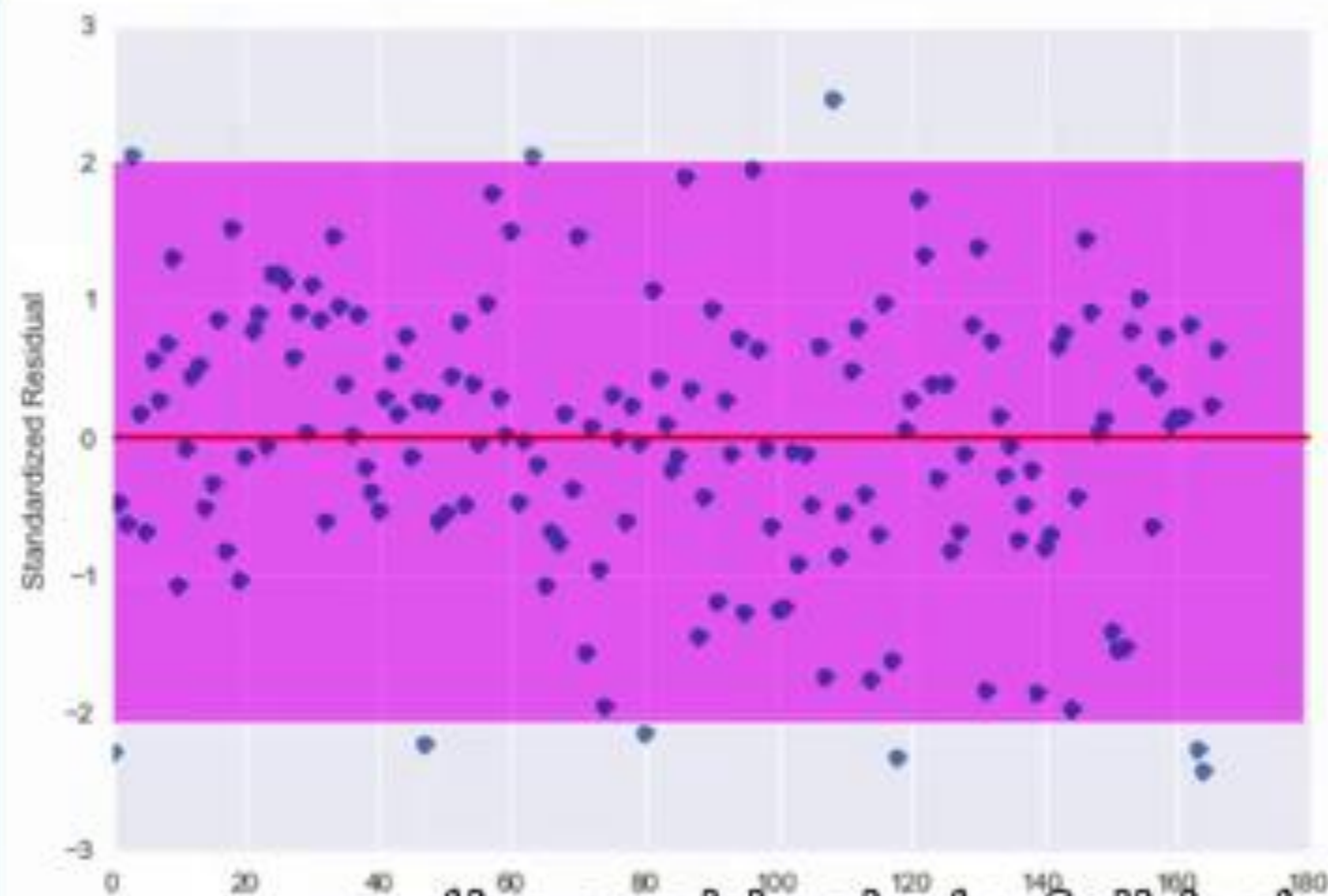


```
[<matplotlib.lines.Line2D object at 0x000000000BD56C18>]
```



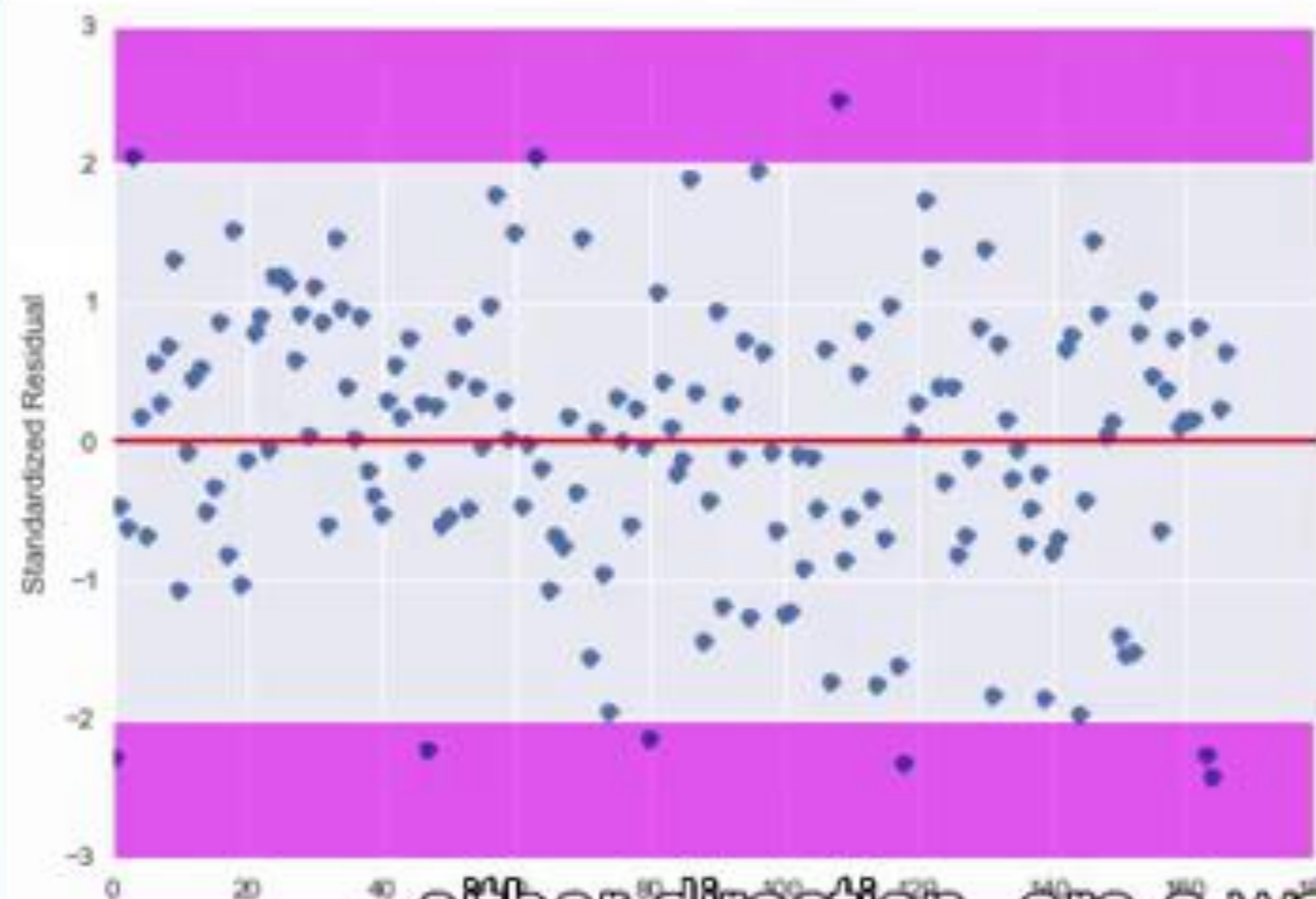
With the standard normal distribution,
we would expect 95% of the values of

```
[<matplotlib.lines.Line2D object at 0x000000000BD56C18>]
```



the residuals to fall between two
standard deviations of the mean.

```
[<matplotlib.lines.Line2D object at 0x000000000BD56C18>]
```

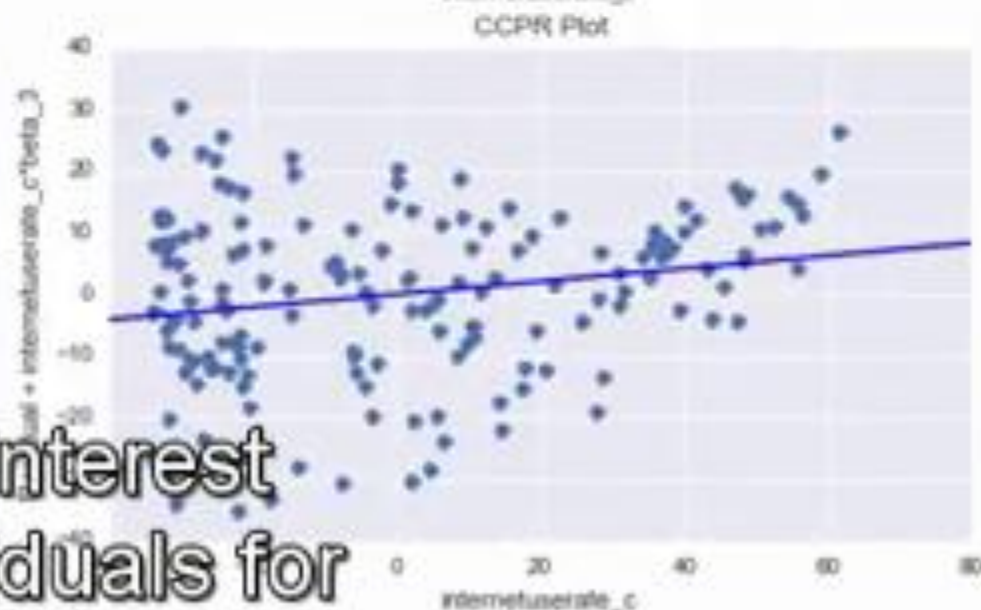
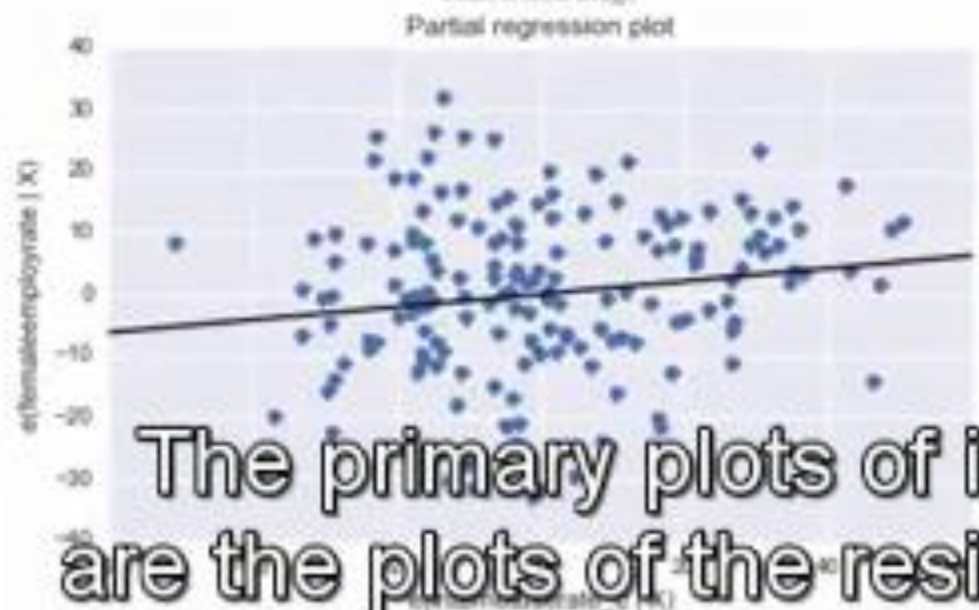
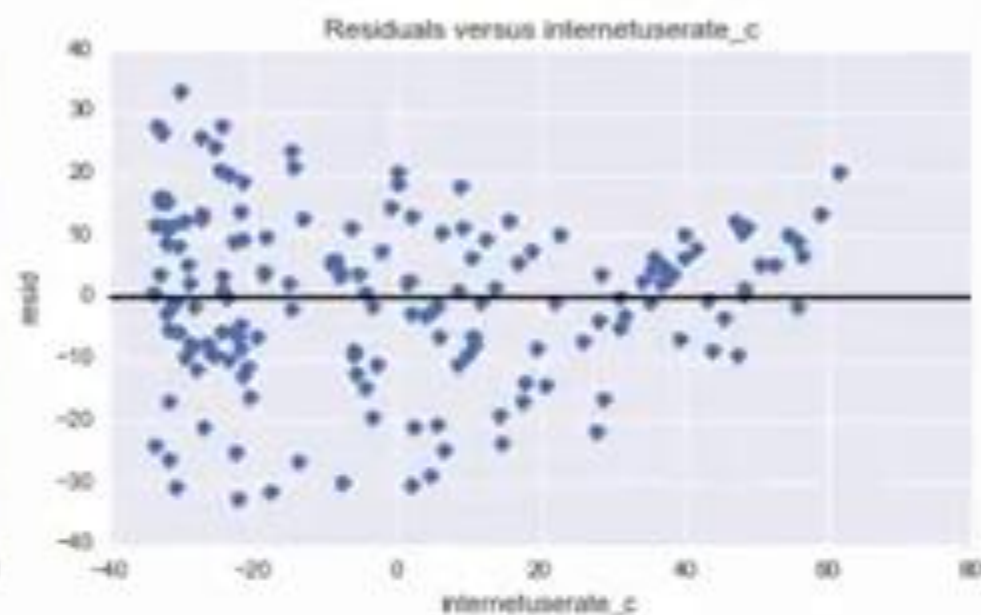
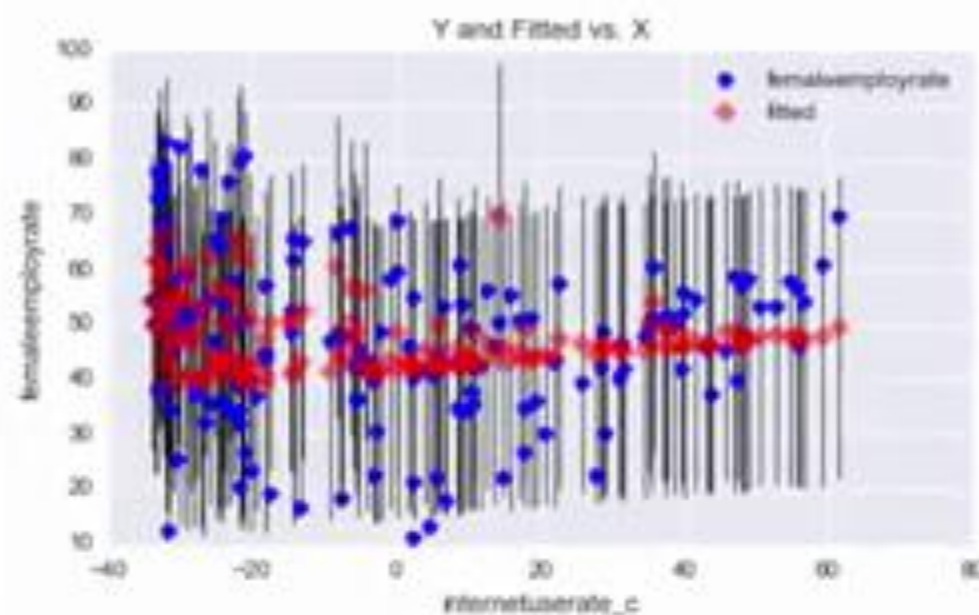


either direction, are a warning sign
that we may have some outliers.

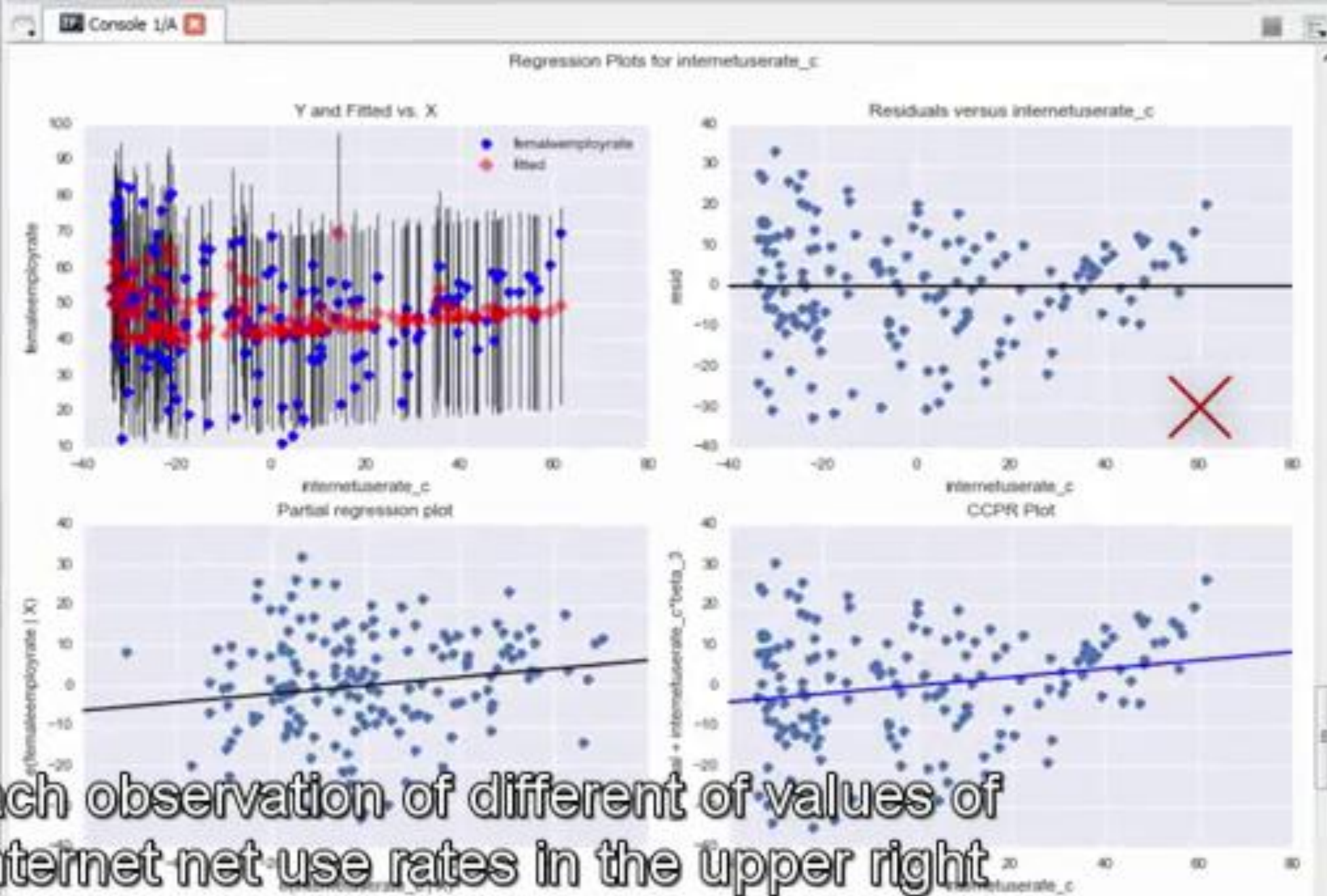
```
14 print(fig1)
15
16
17 # simple plot of residuals
18 stdres=pandas.DataFrame(reg3.resid_pearson)
19 fig2=plt.plot(stdres, 'o', ls='none')
20 l = plt.axhline(y=0, color='r')
21 plt.ylabel('Standardized Residual')
22 plt.xlabel('Observation Number')
23 print(fig2)
```

```
24
25 # additional regression diagnostic plots
26 fig3 = plt.figure(figsize(12,8))
27 fig3 = sm.graphics.plot_regress_exog(reg3, "internetuserate_c", fig=fig3)
```


Regression Plots for internetuserate_c

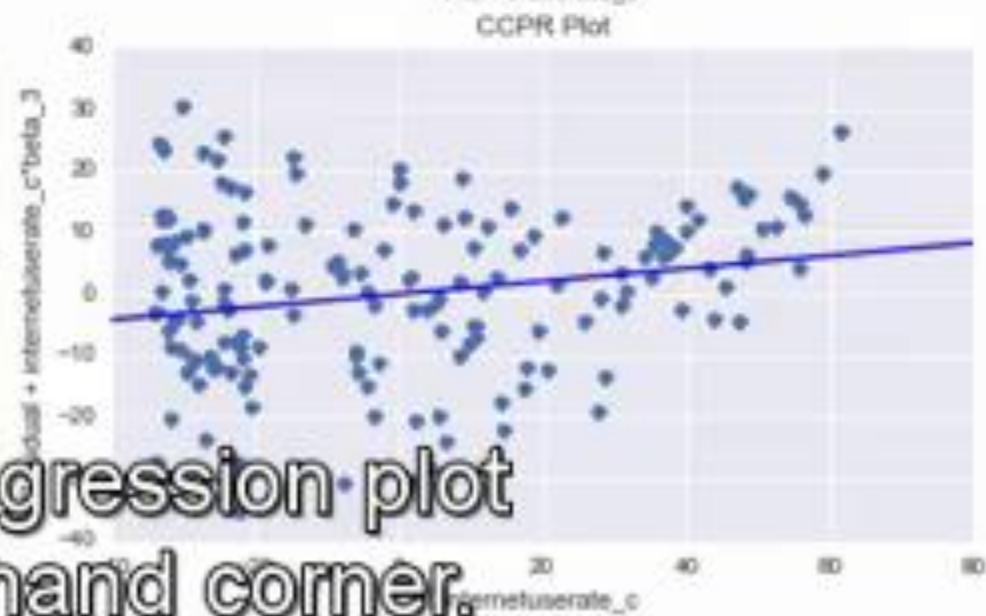
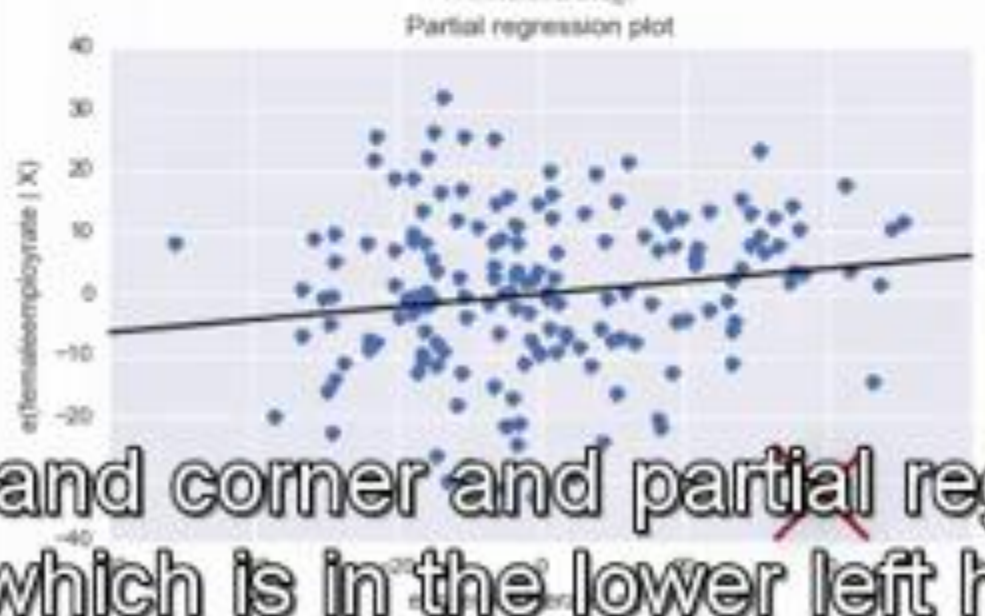
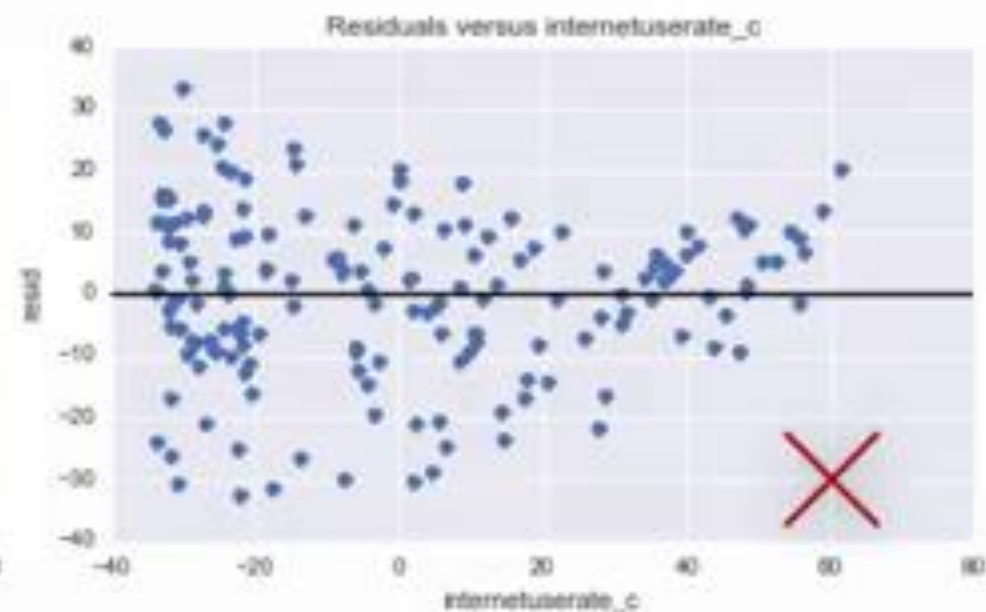
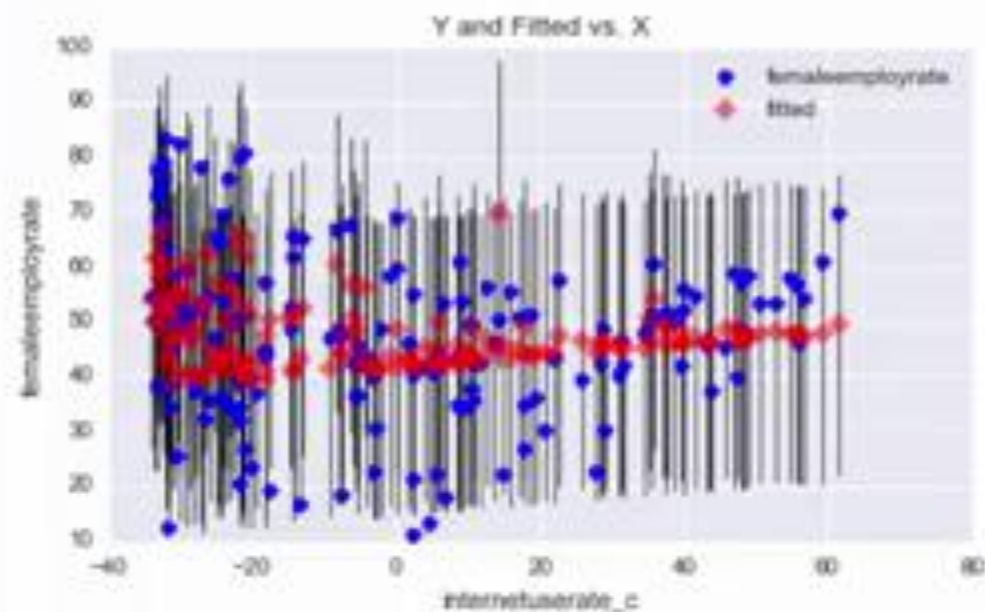


The primary plots of interest are the plots of the residuals for



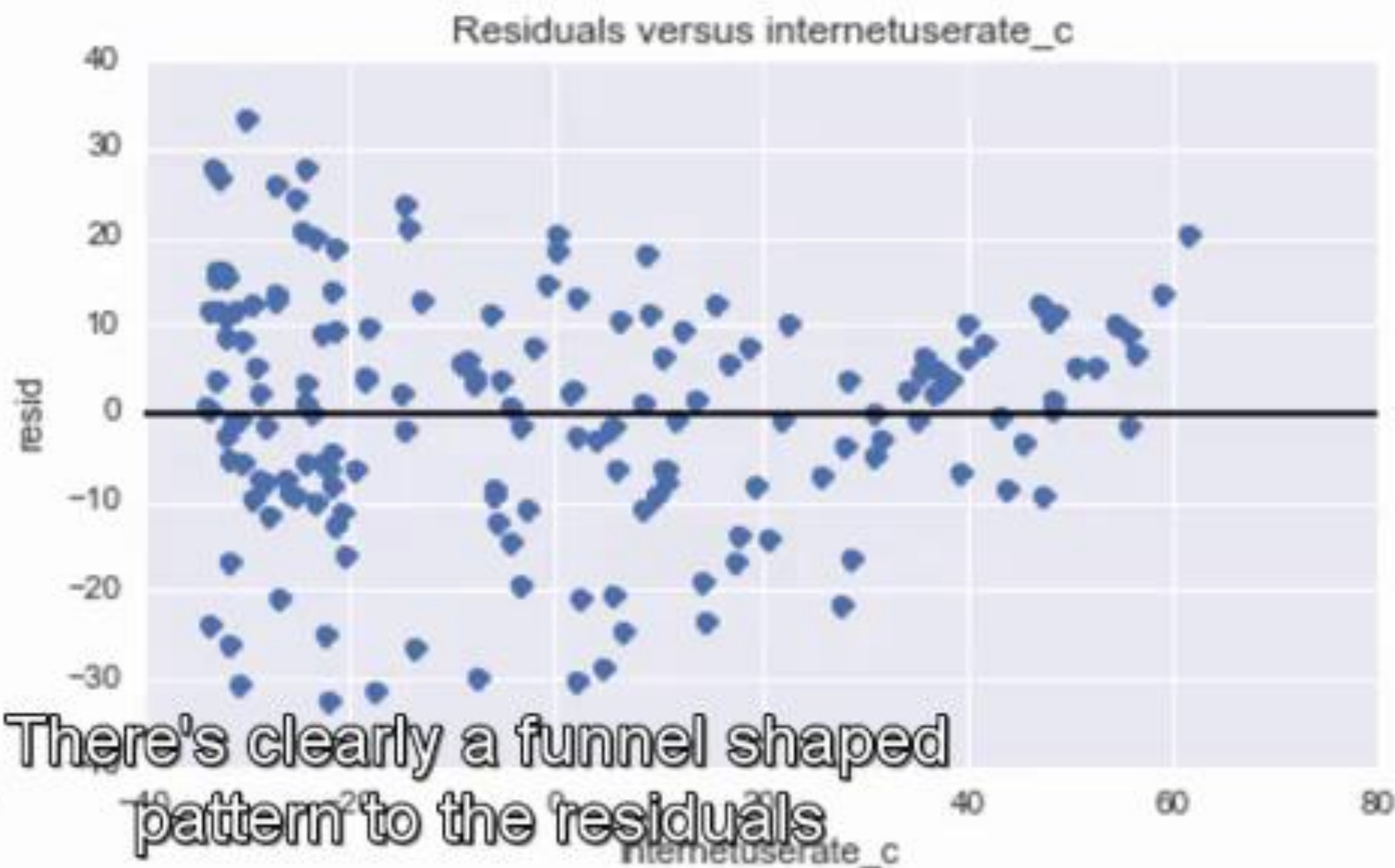
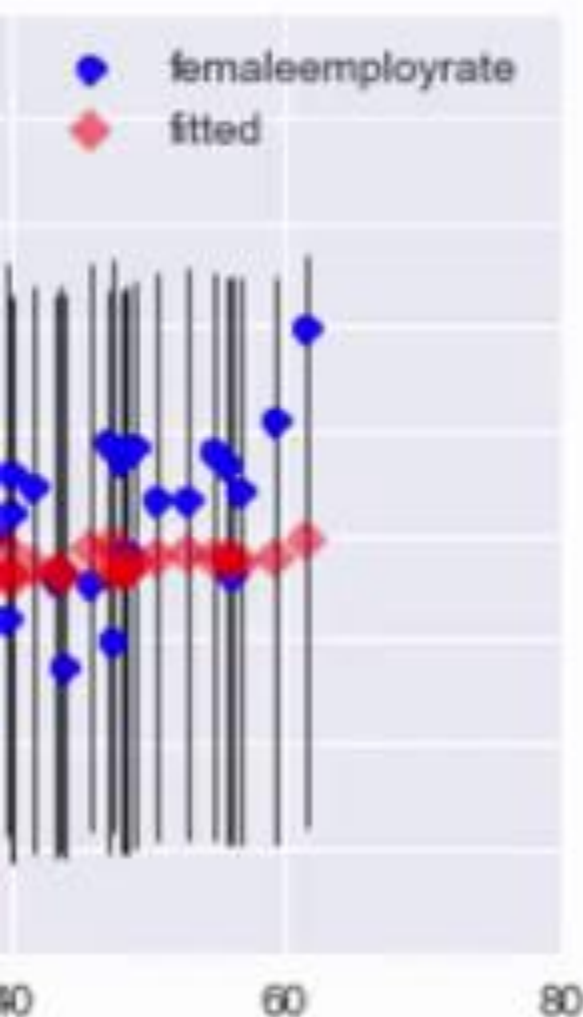
each observation of different of values of
Internet net use rates in the upper right

Regression Plots for internetuserate_c



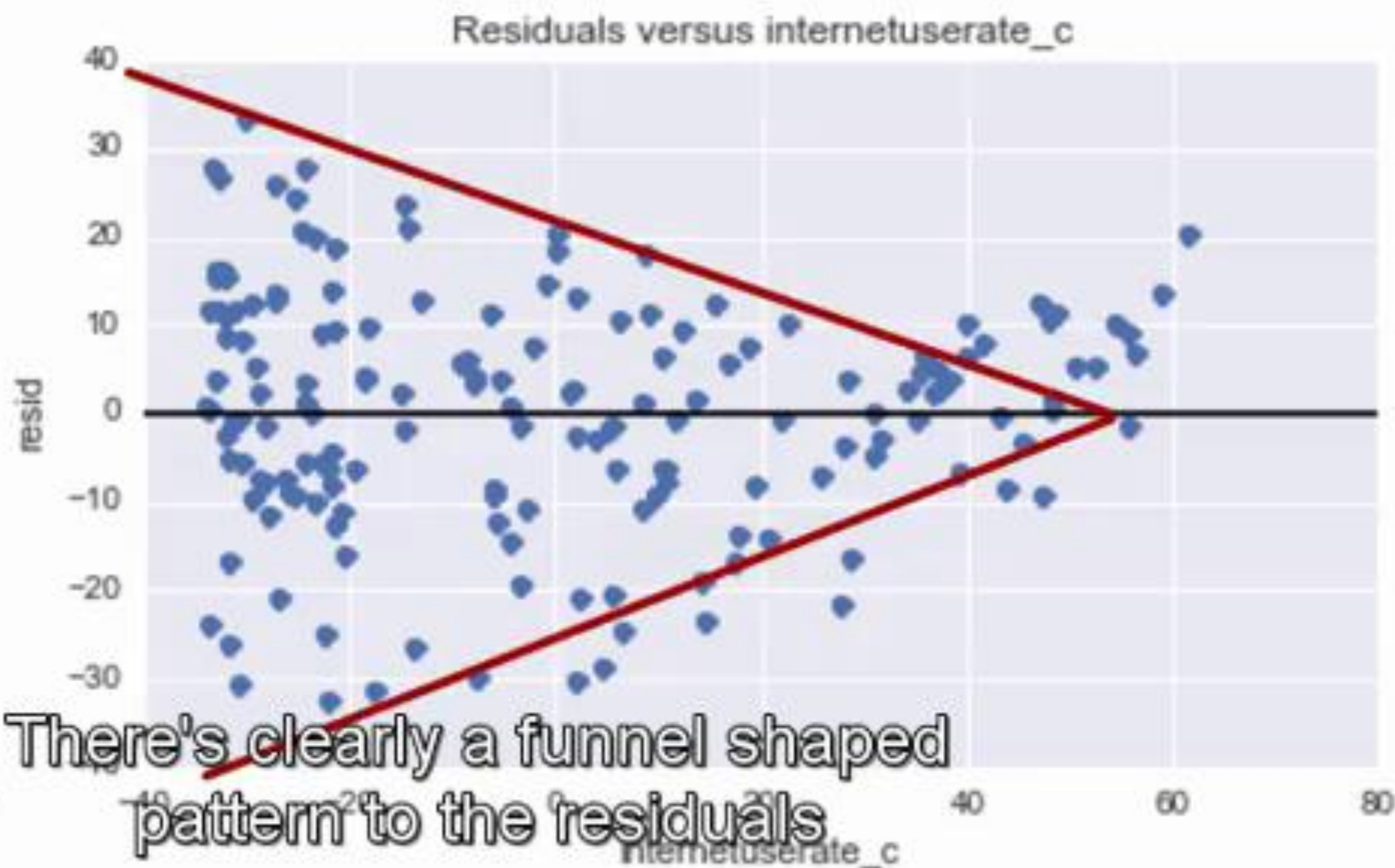
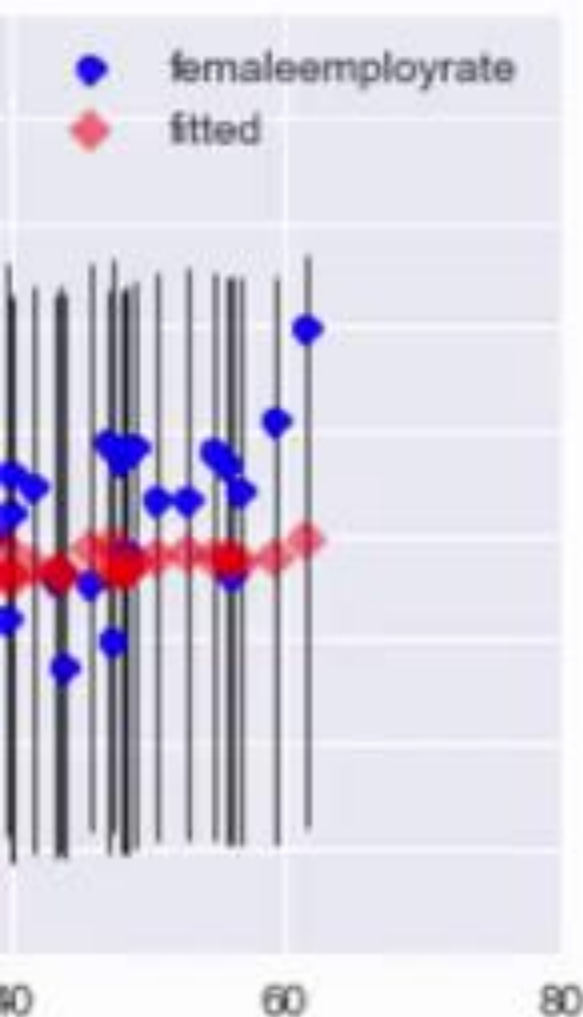
hand corner and partial regression plot
which is in the lower left hand corner.

Regression Plots for internetuserate_c



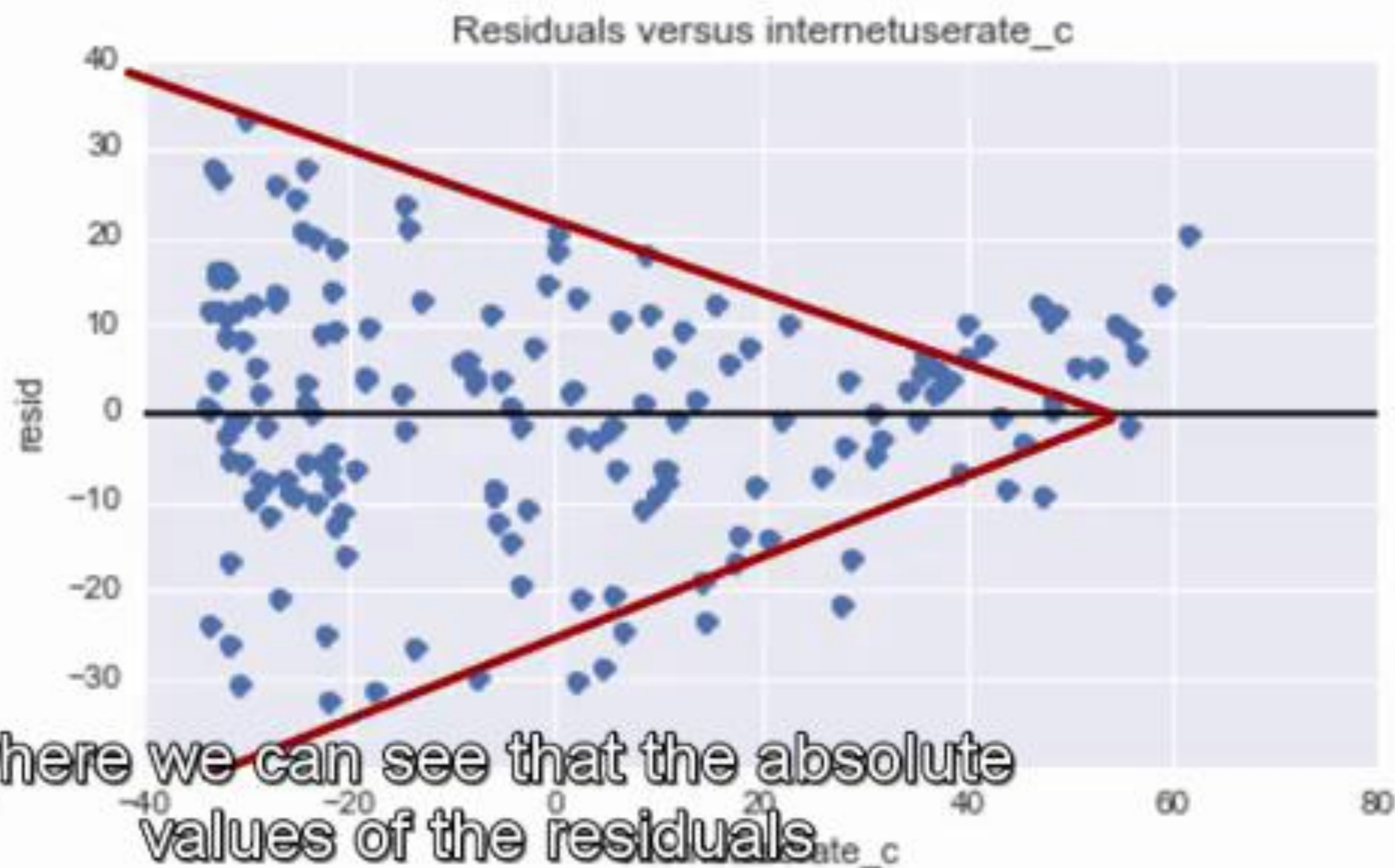
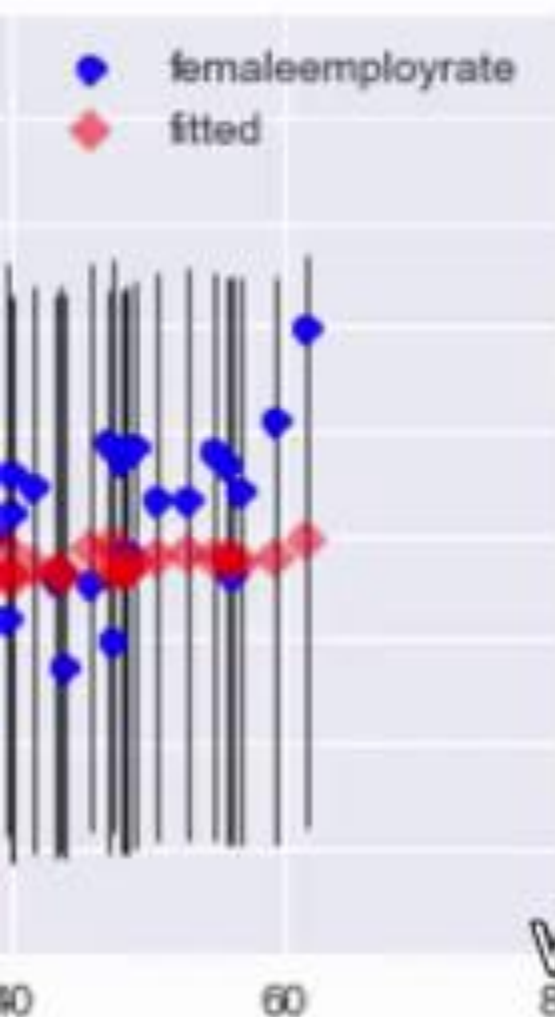
There's clearly a funnel shaped pattern to the residuals

Regression Plots for internetuserate_c

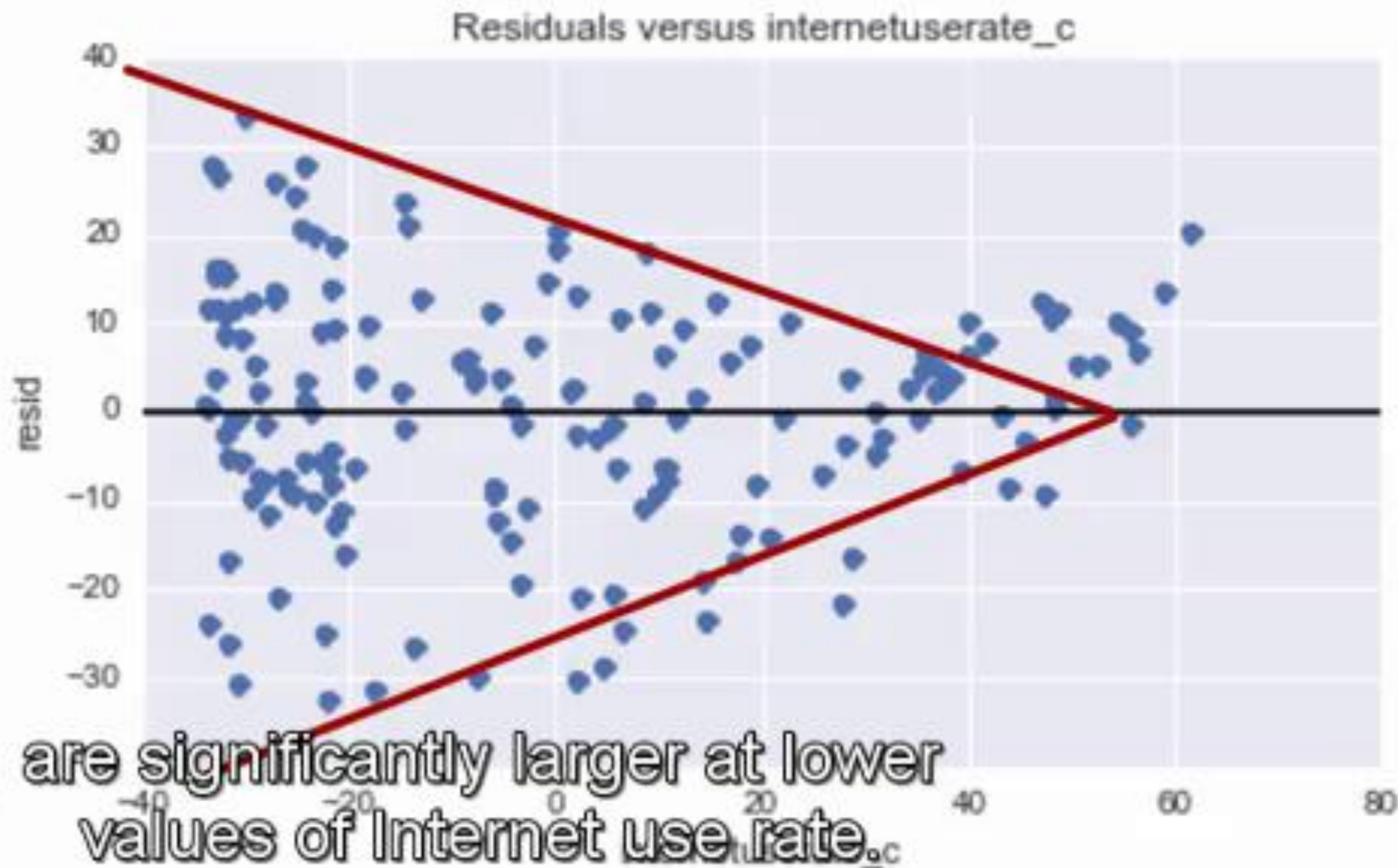
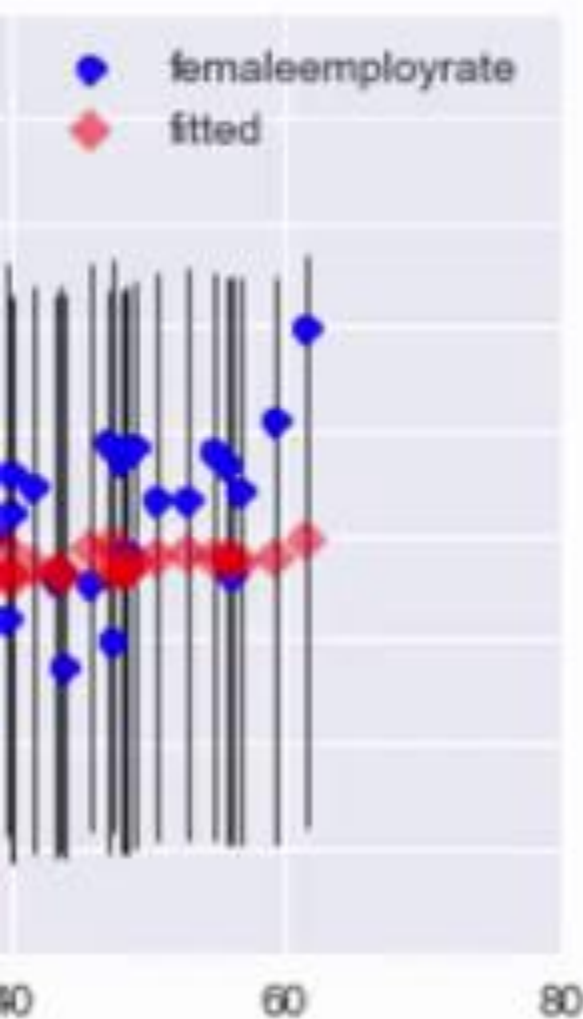


There's clearly a funnel shaped pattern to the residuals

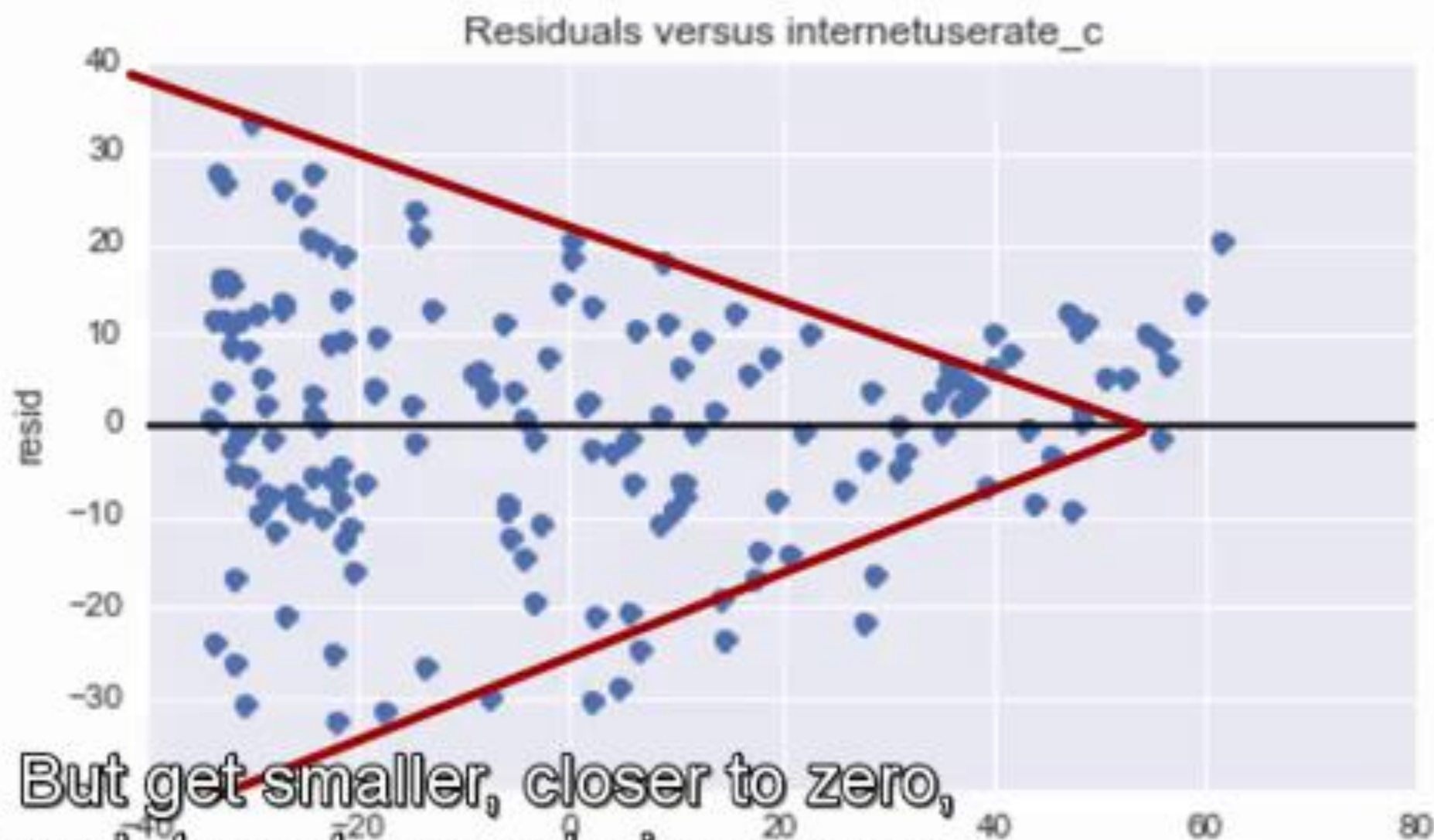
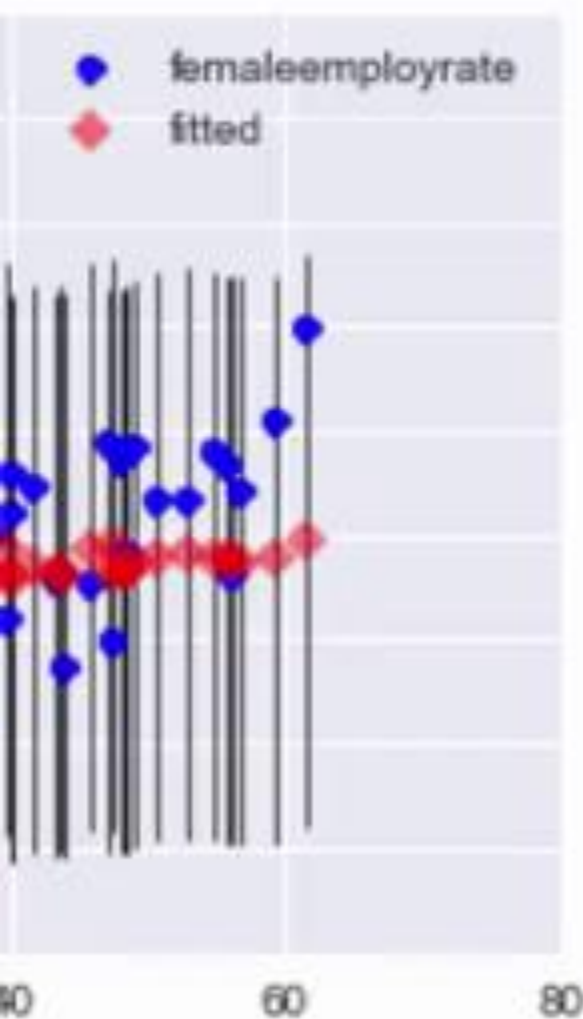
Regression Plots for internetuserate_c



Regression Plots for internetuserate_c

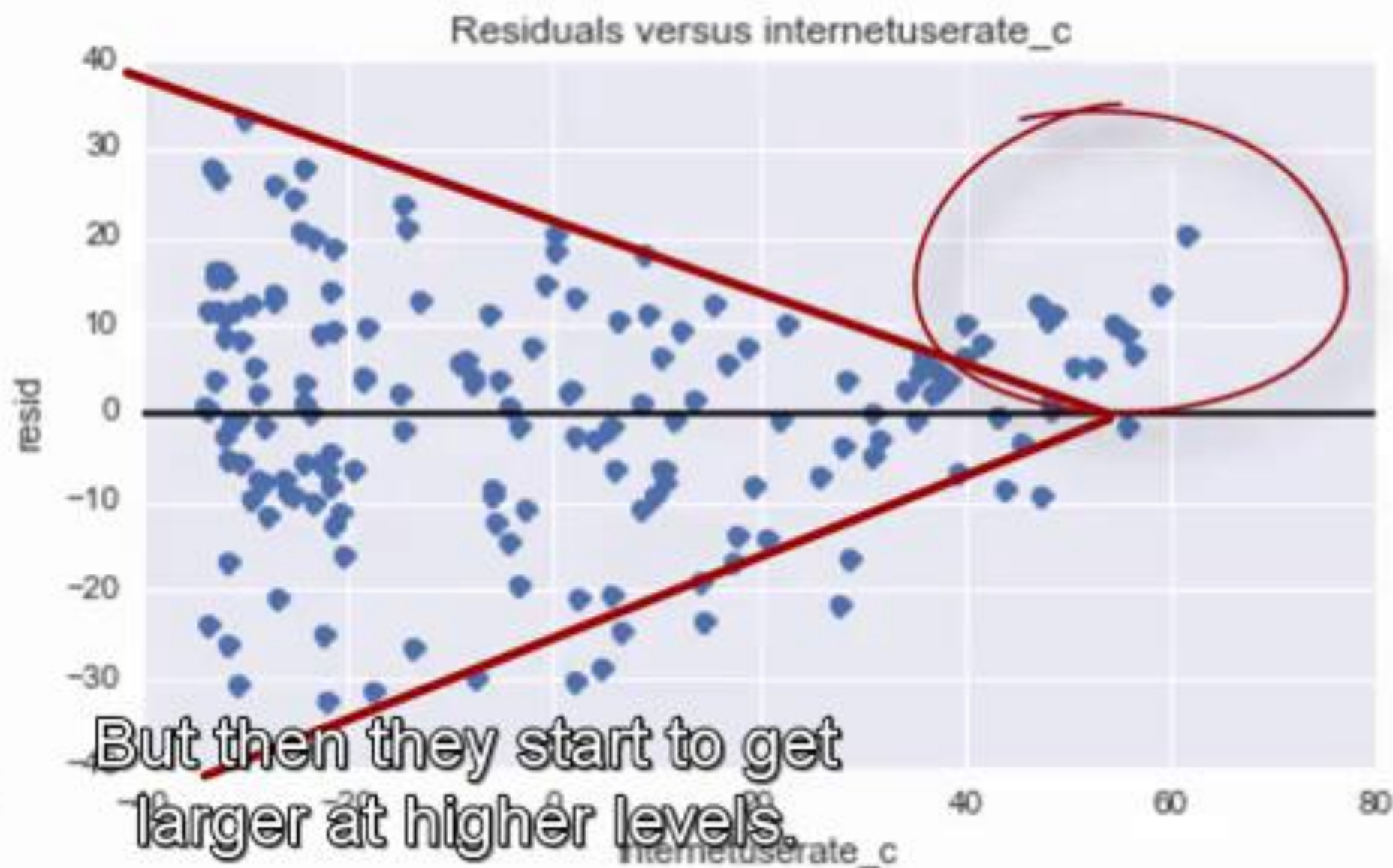
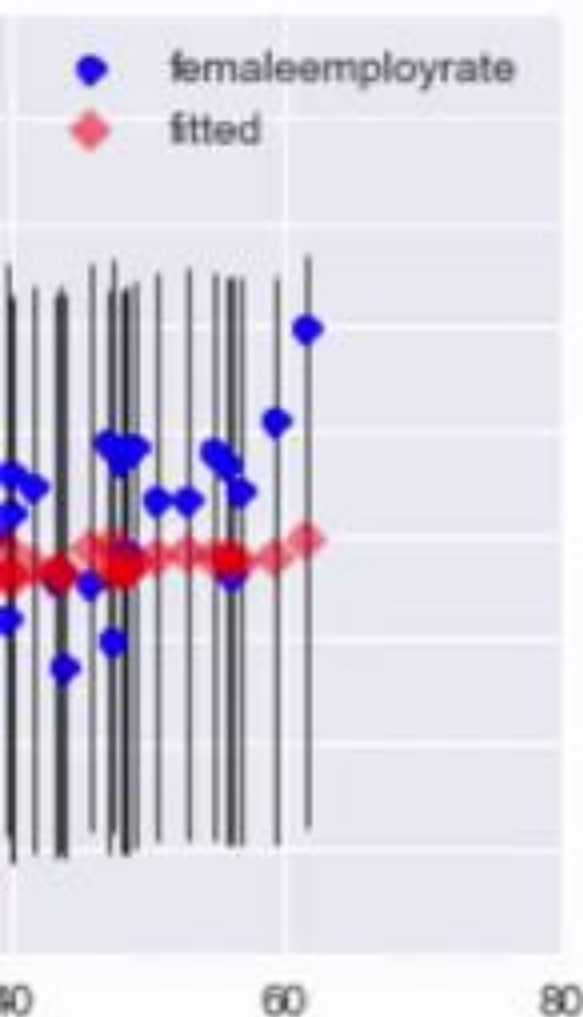


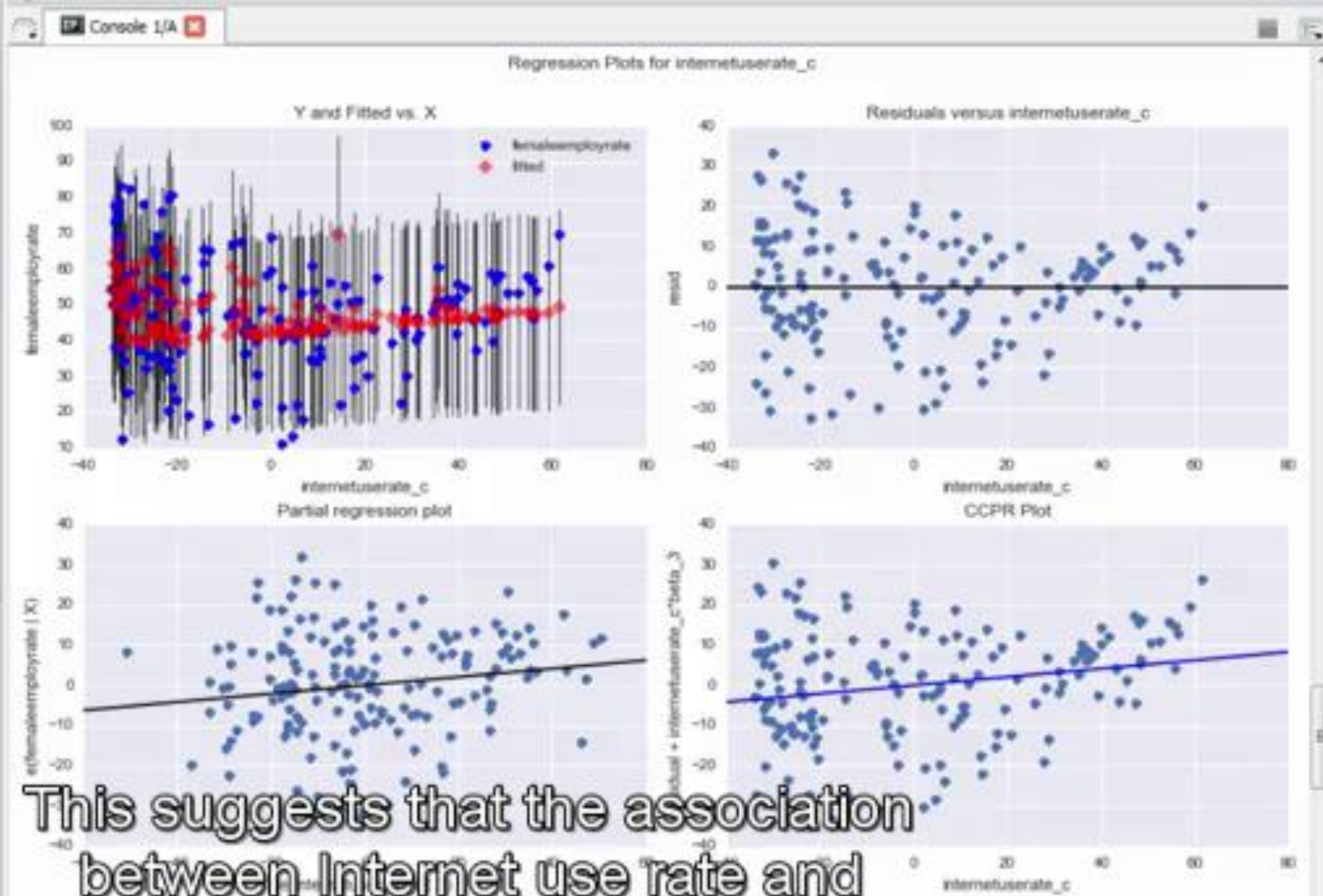
Regression Plots for internetuserate_c

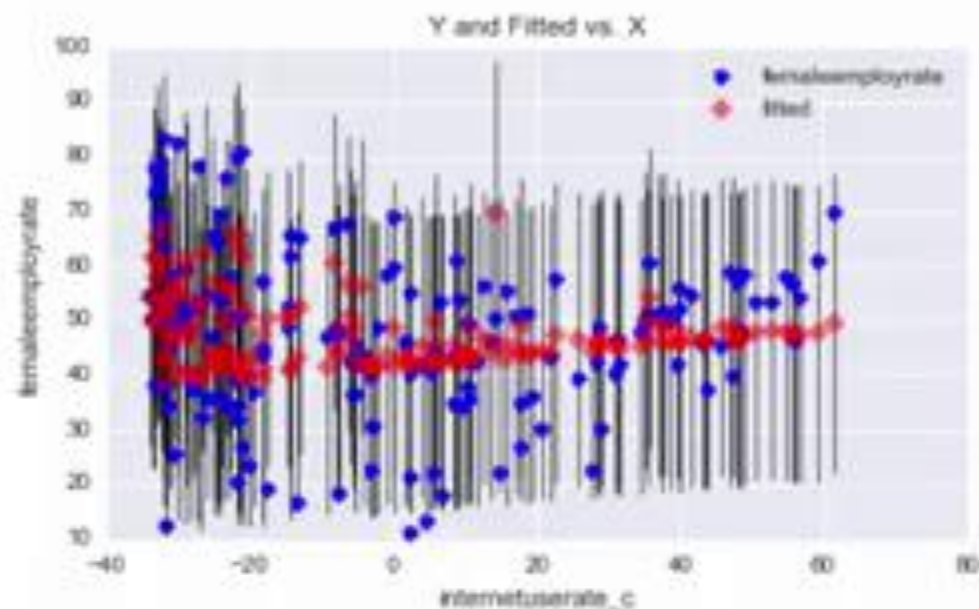


But get smaller, closer to zero,
as Internet use rate increases.

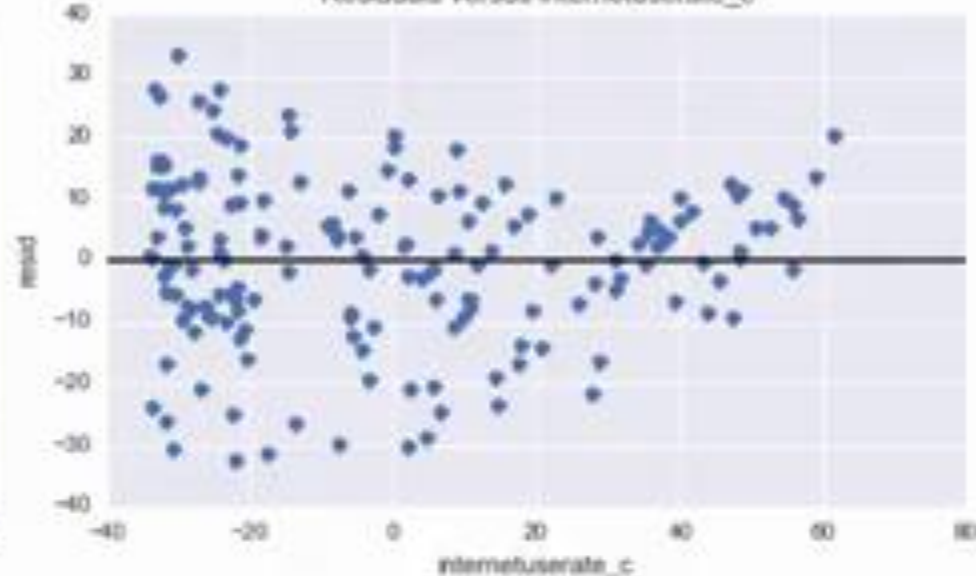
Regression Plots for internetuserate_c







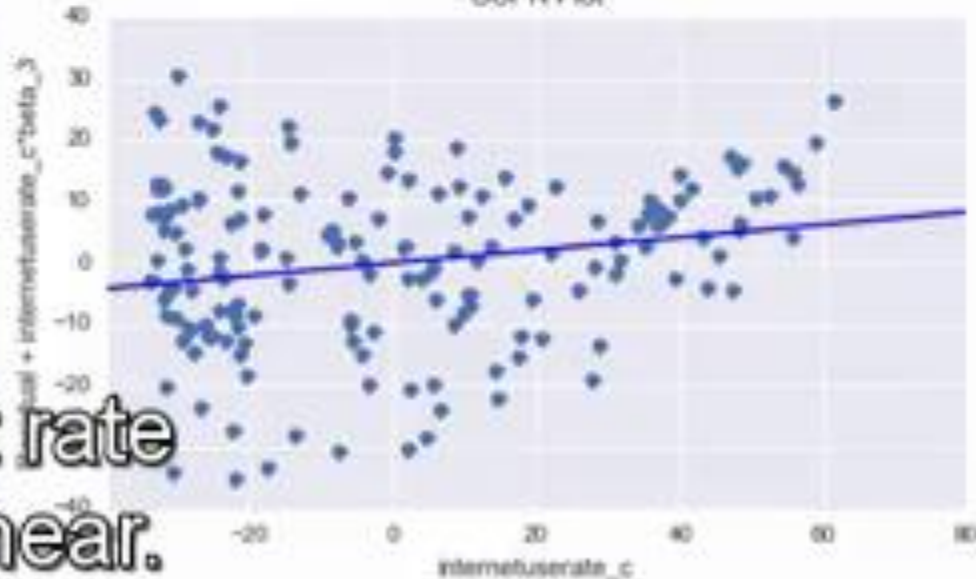
Residuals versus internetuserate_c



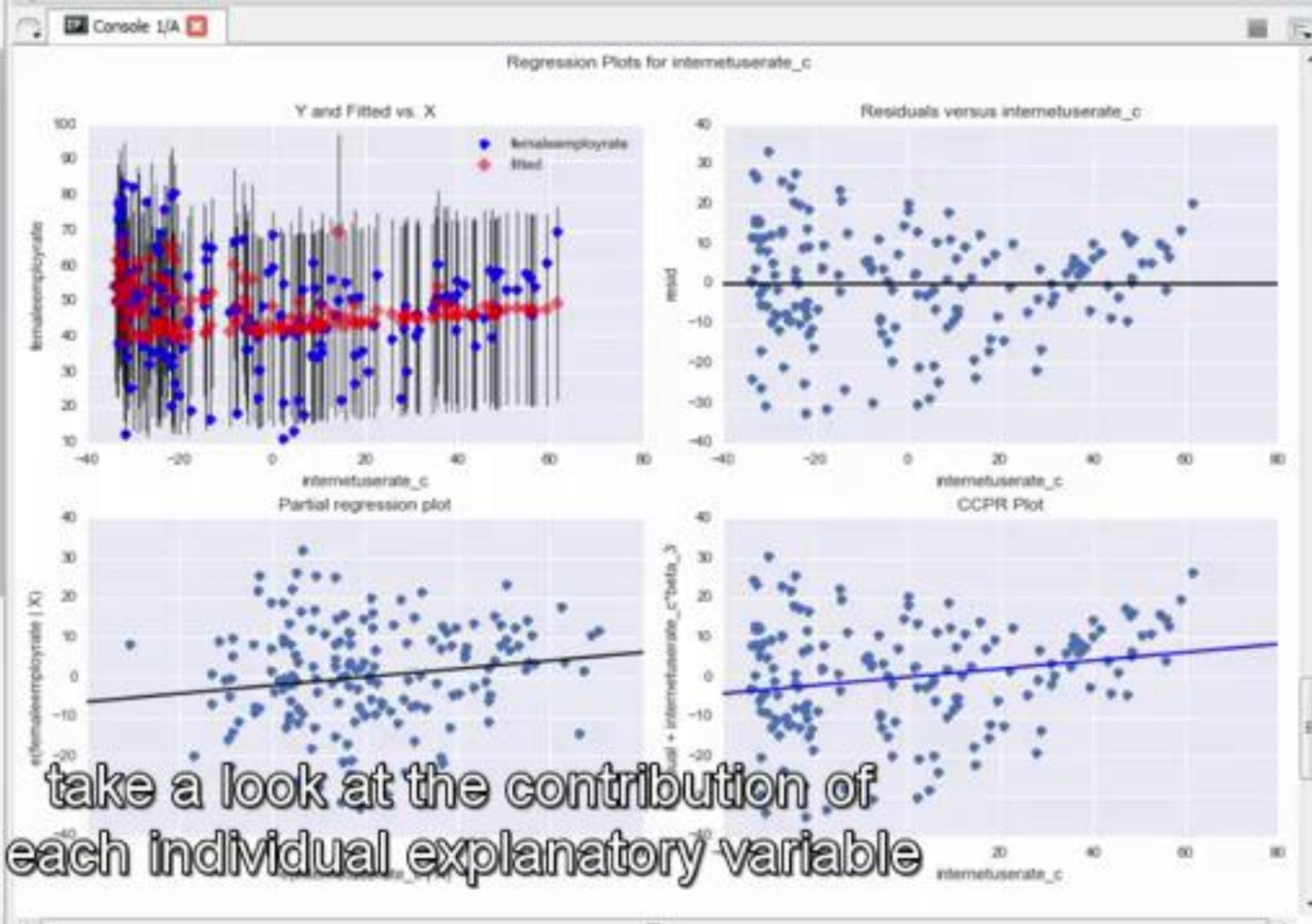
Partial regression plot

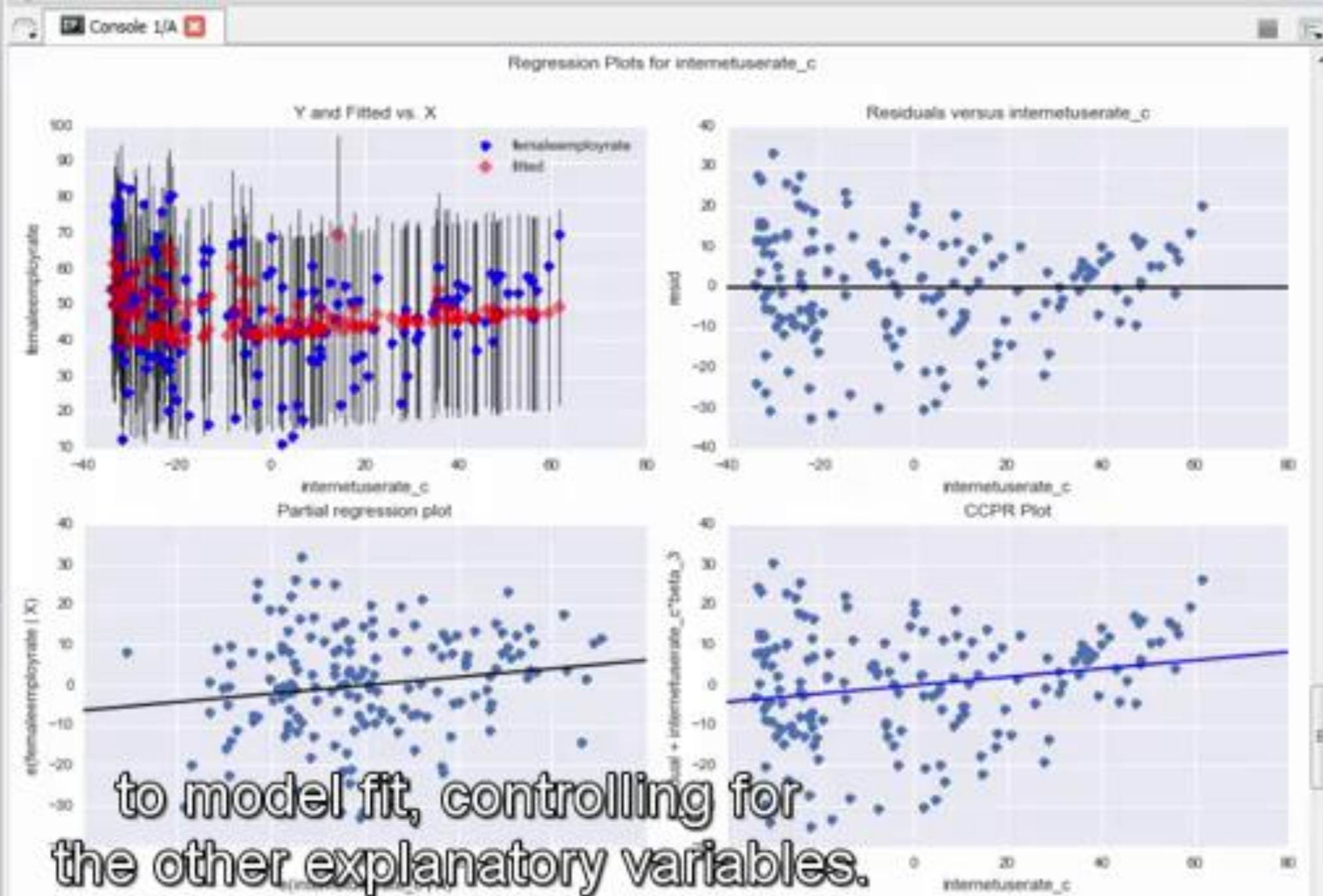


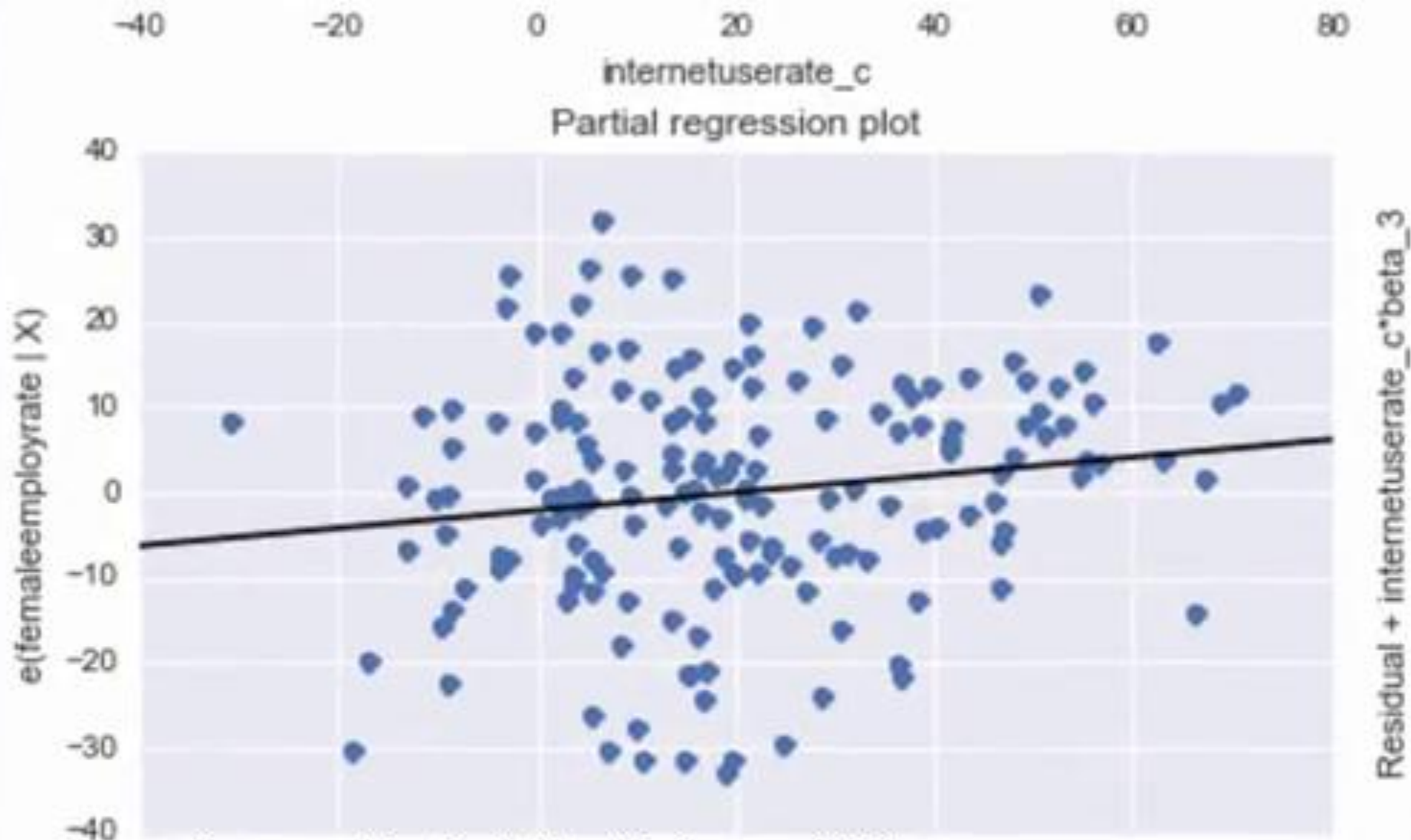
CCPR Plot



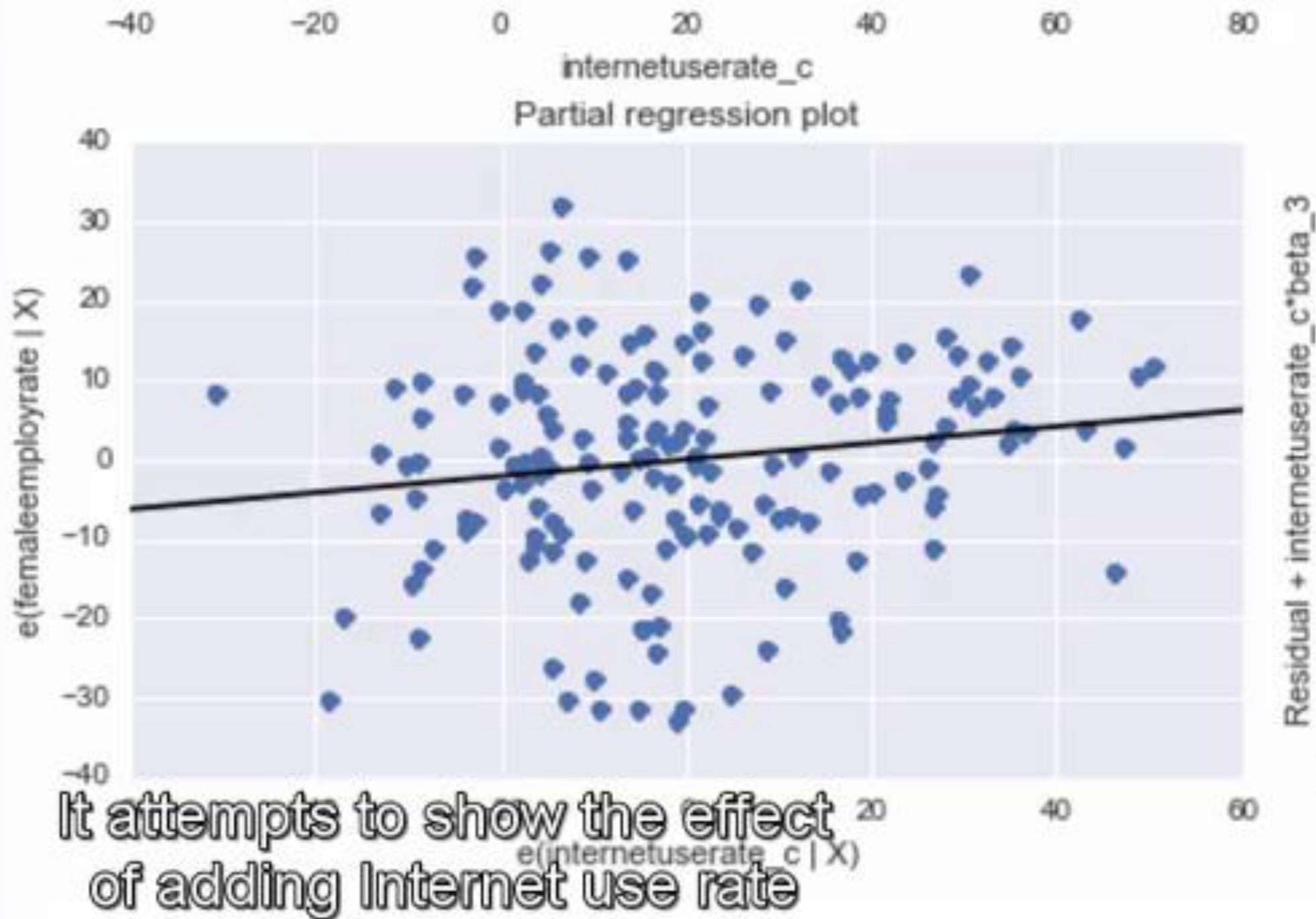
female employment rate
may also be curvilinear.

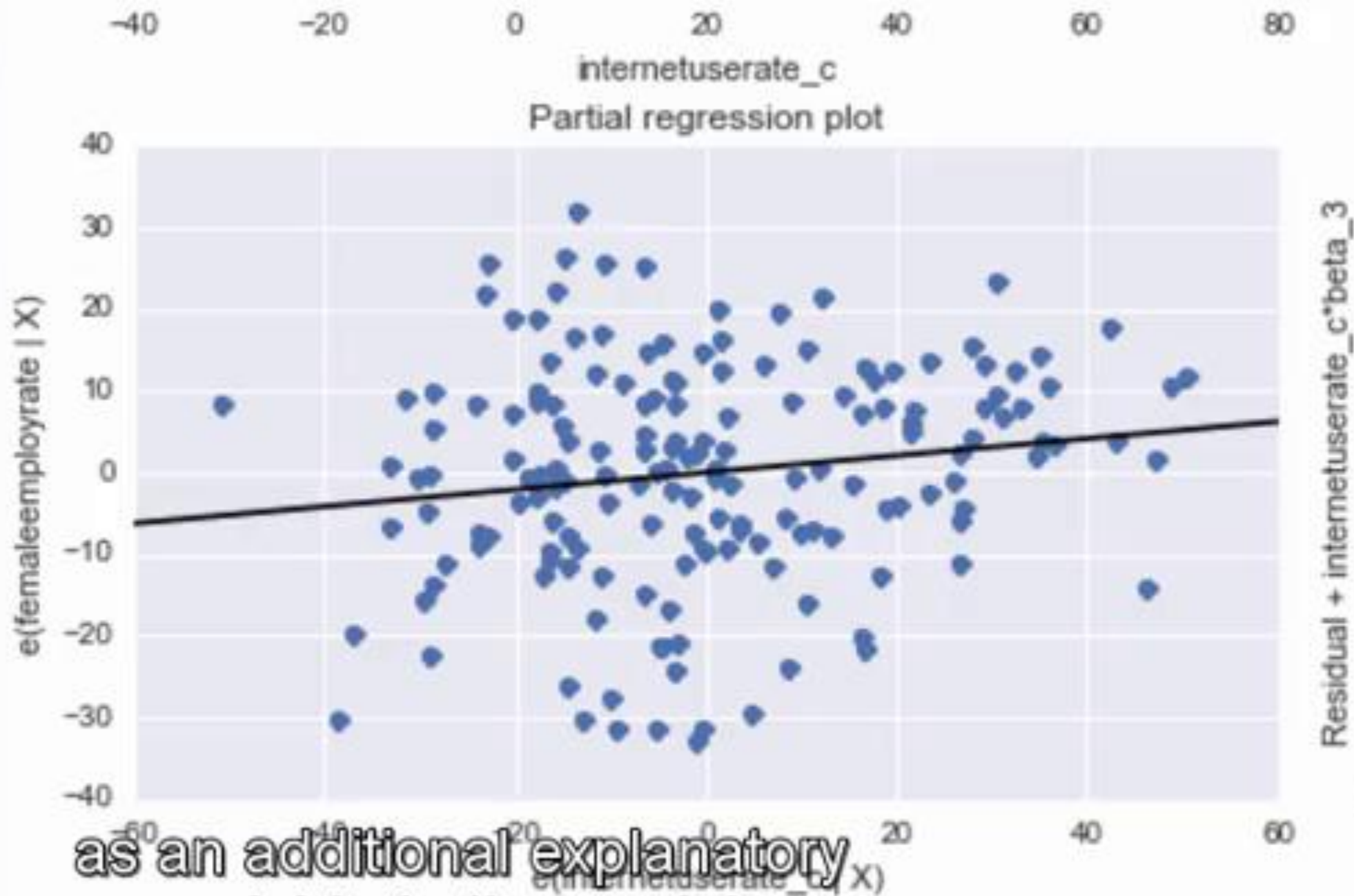


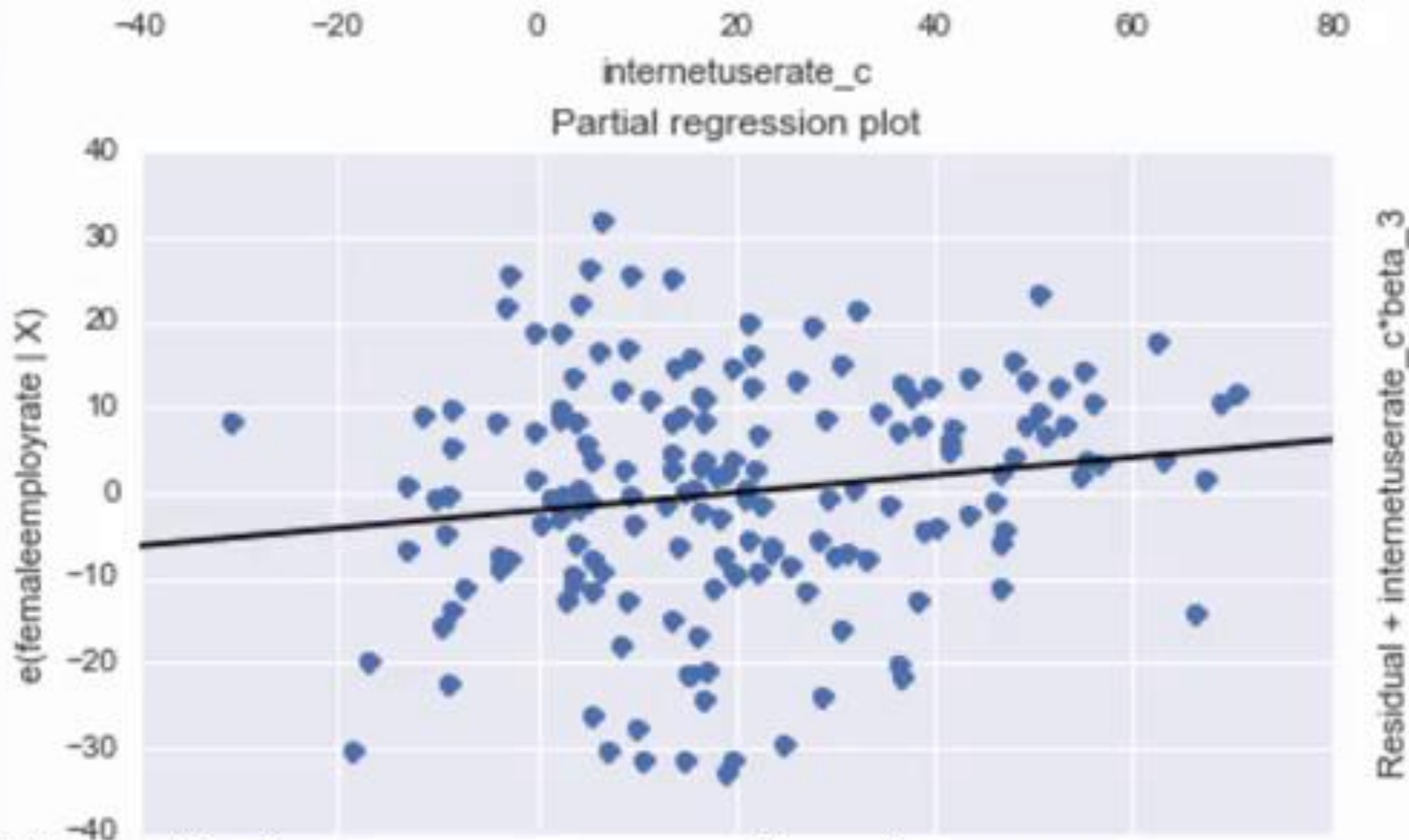




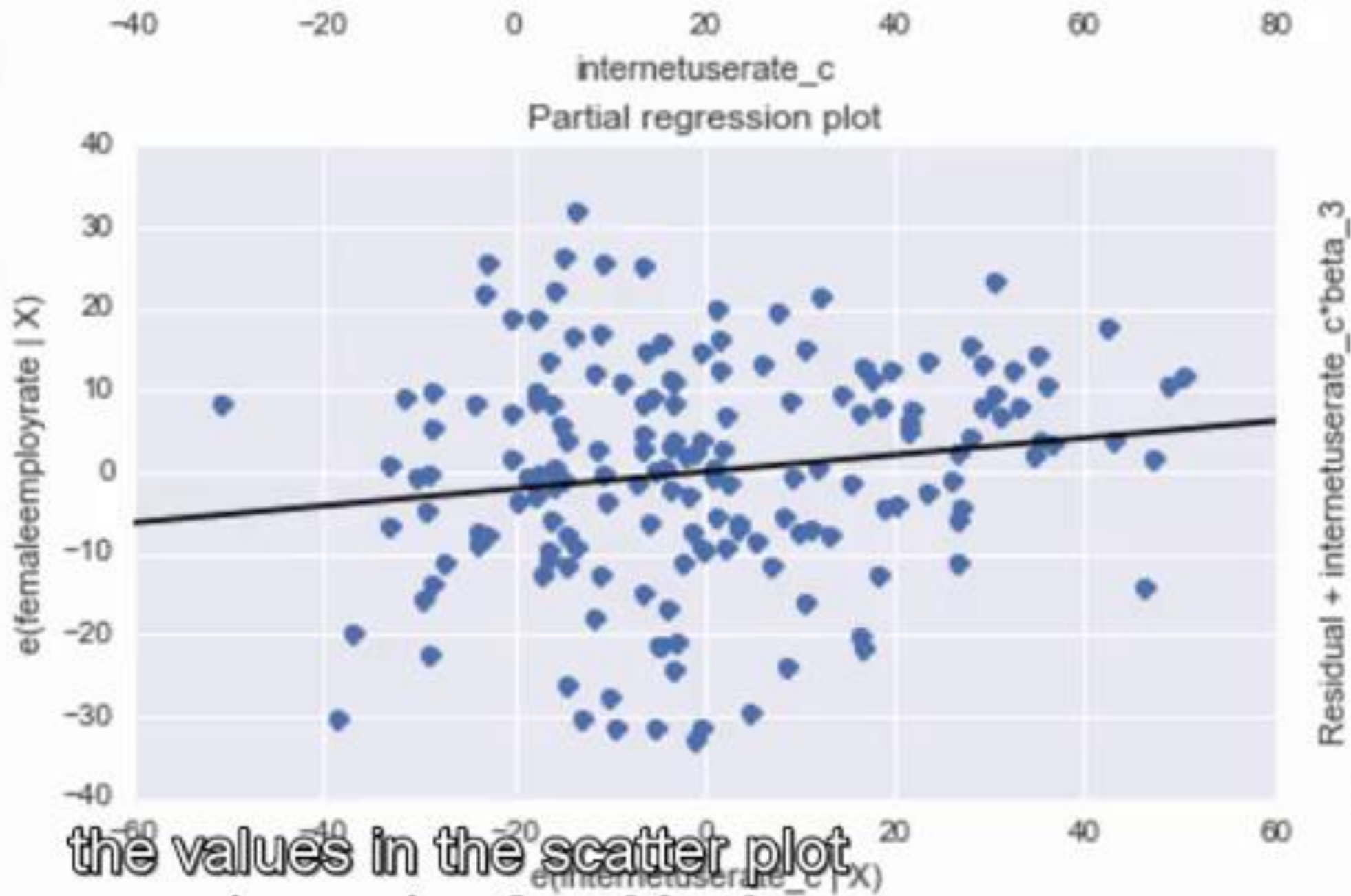
One type of plot that does this,
is the partial regression residual plot.



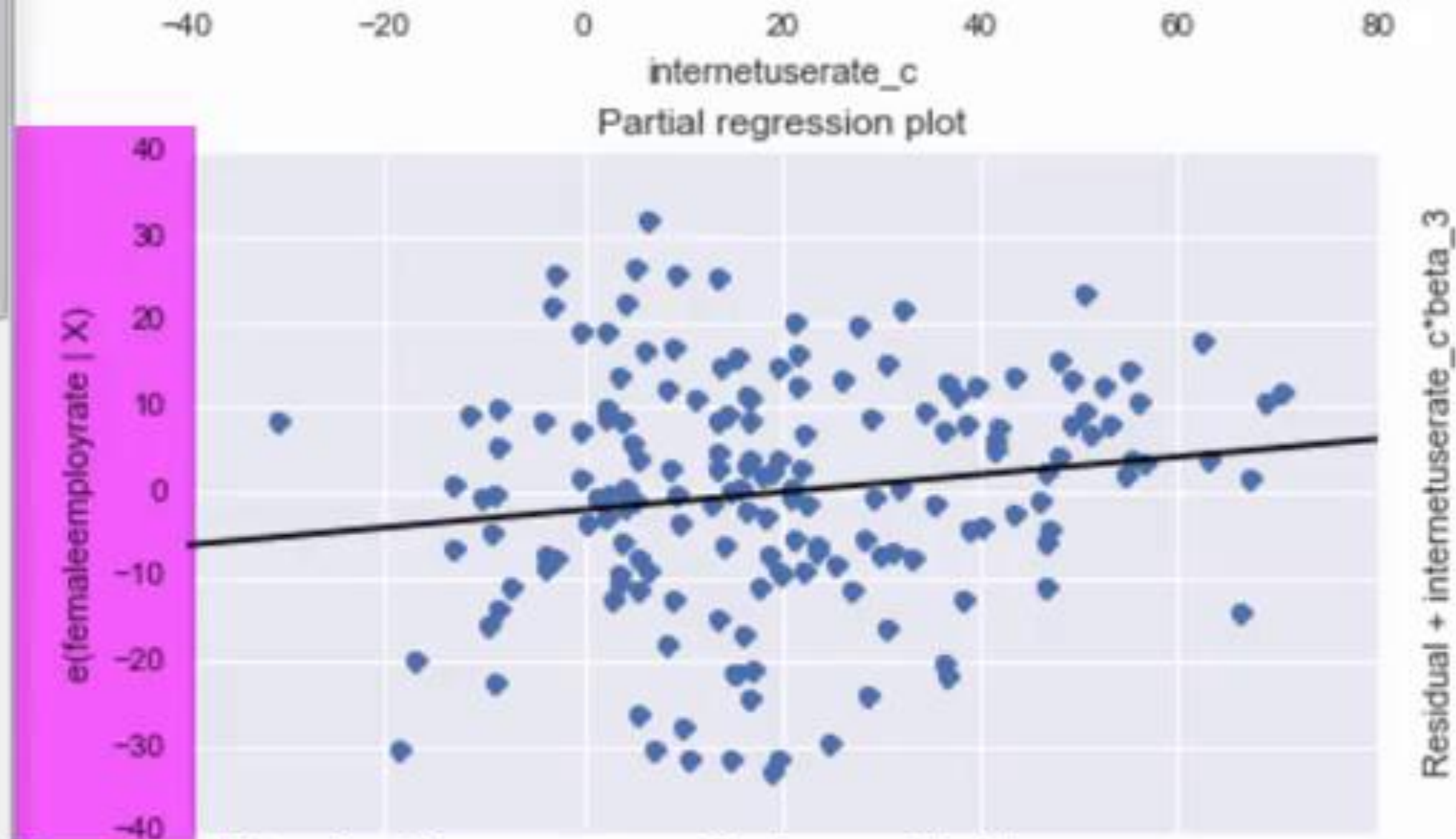




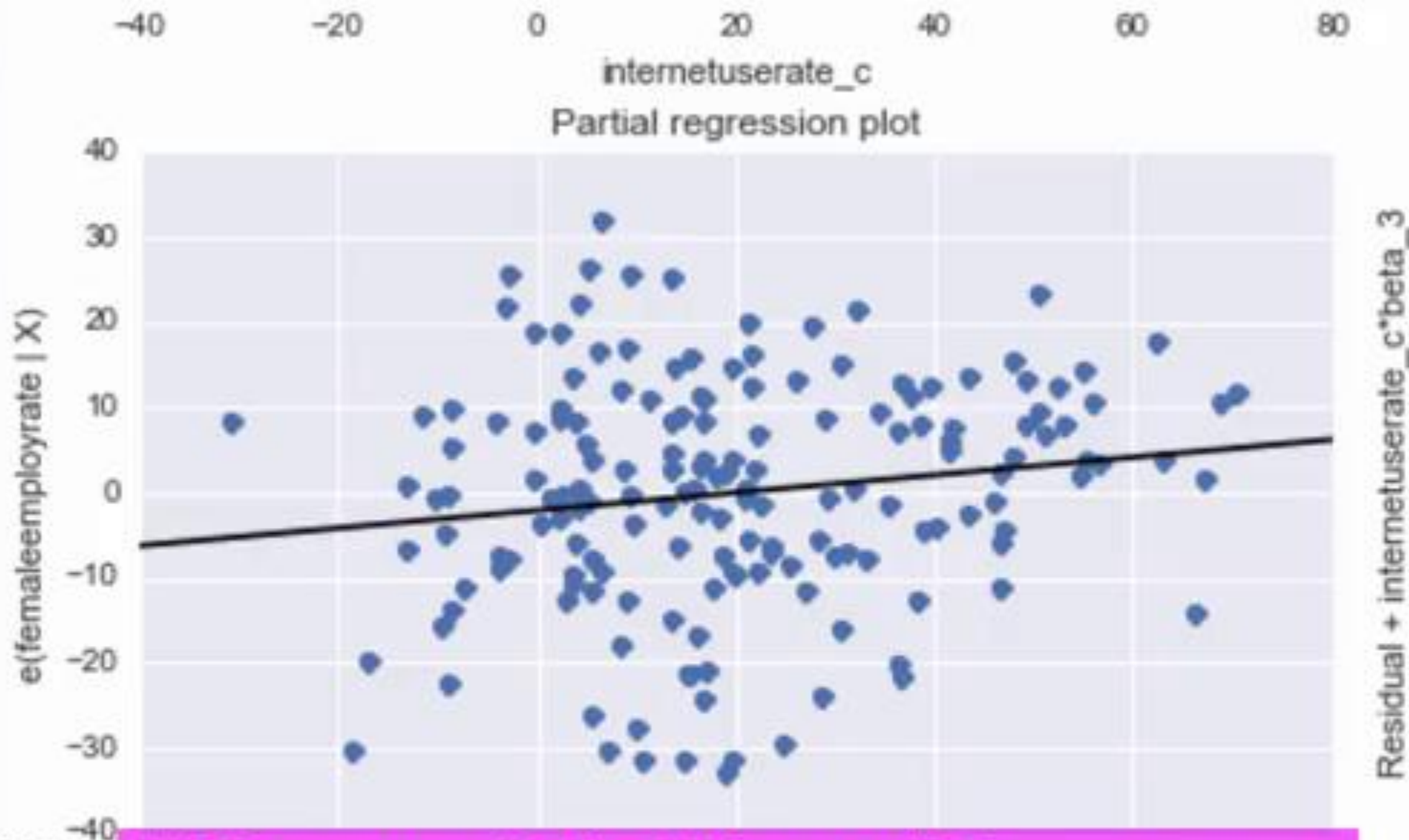
Given that one or more explanatory variables are already in the model.



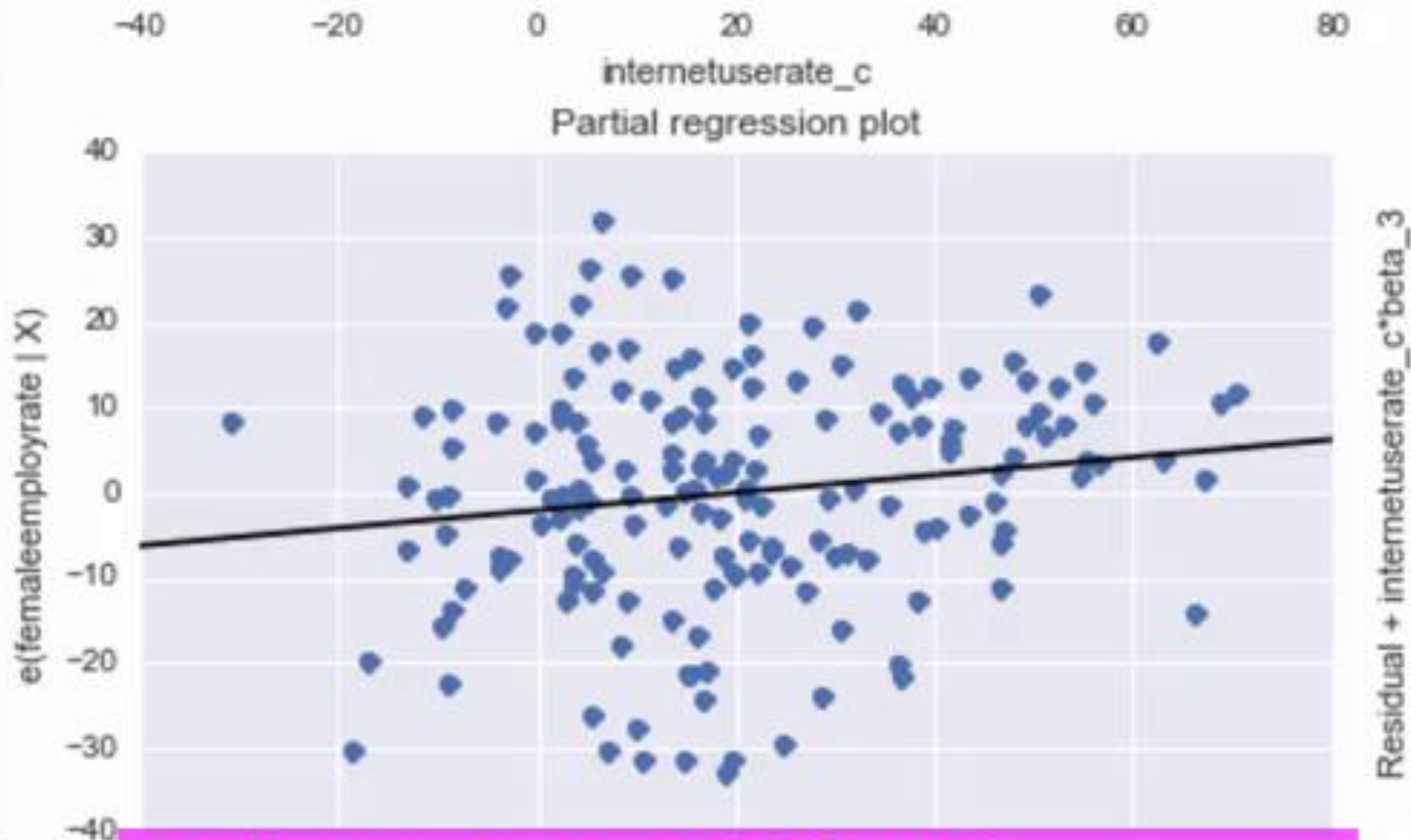
the values in the scatter plot
are two sets of residuals.



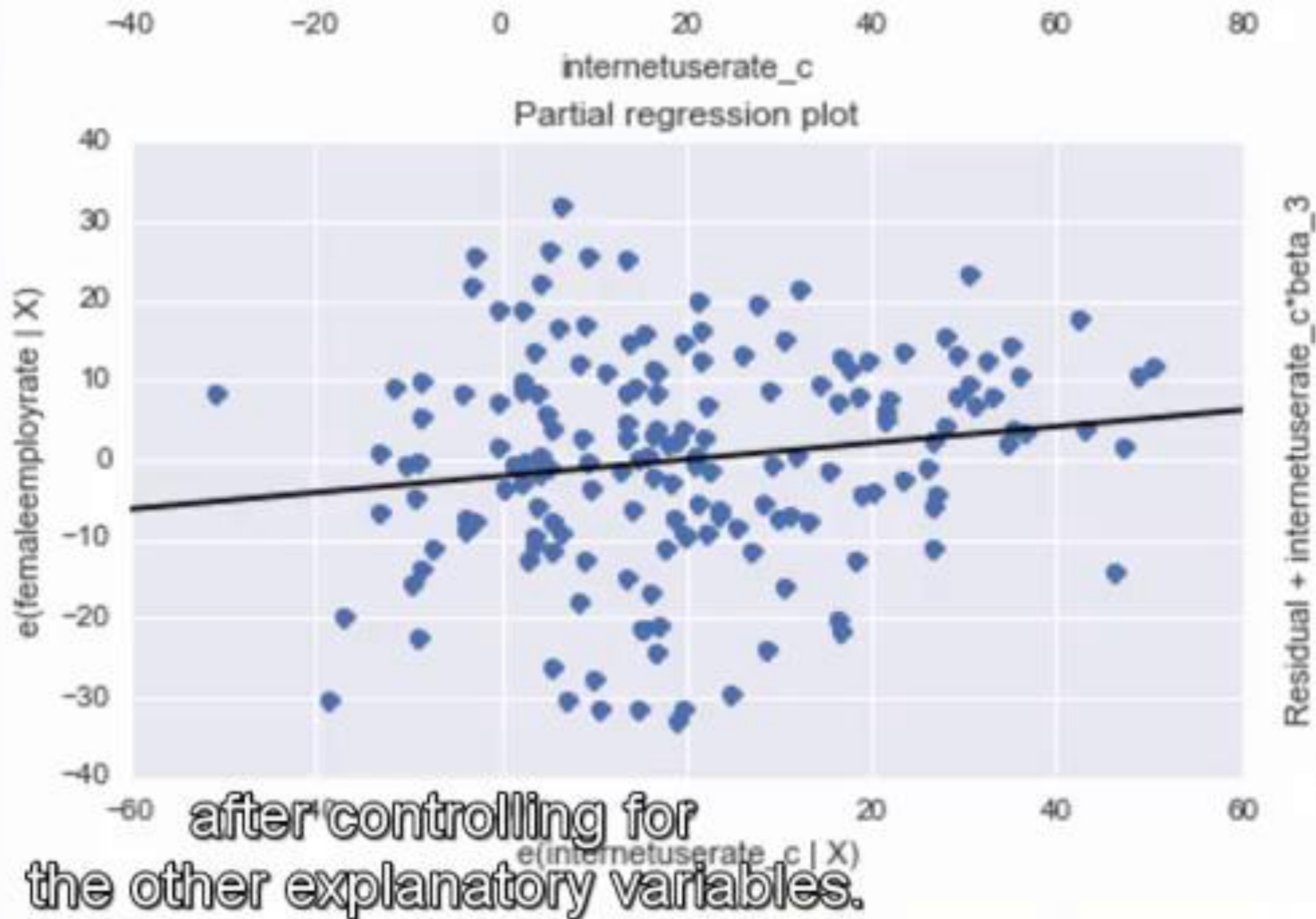
The residuals from a model predicting the female employment rate response from

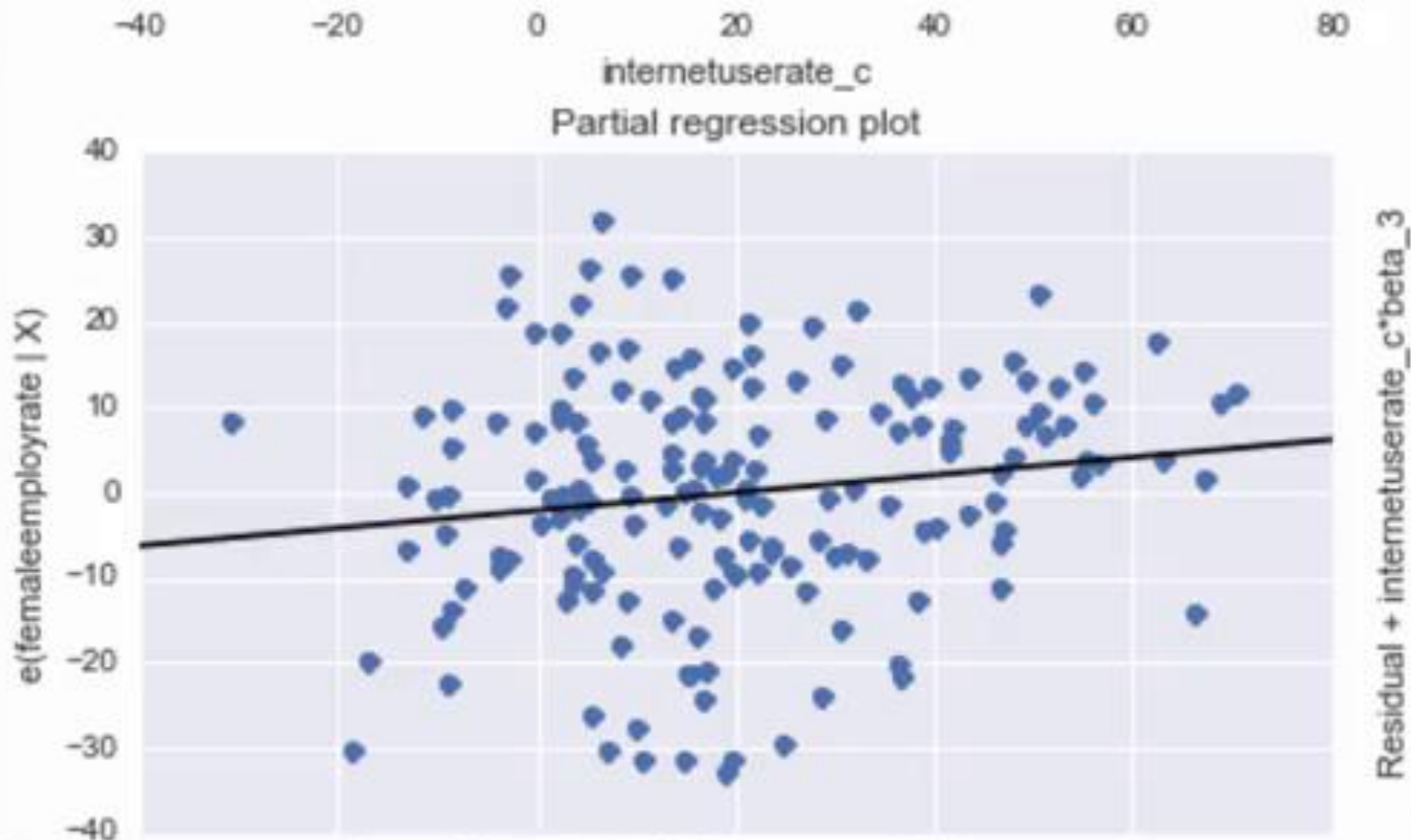


What this means is that the partial regression plot shows the relationship

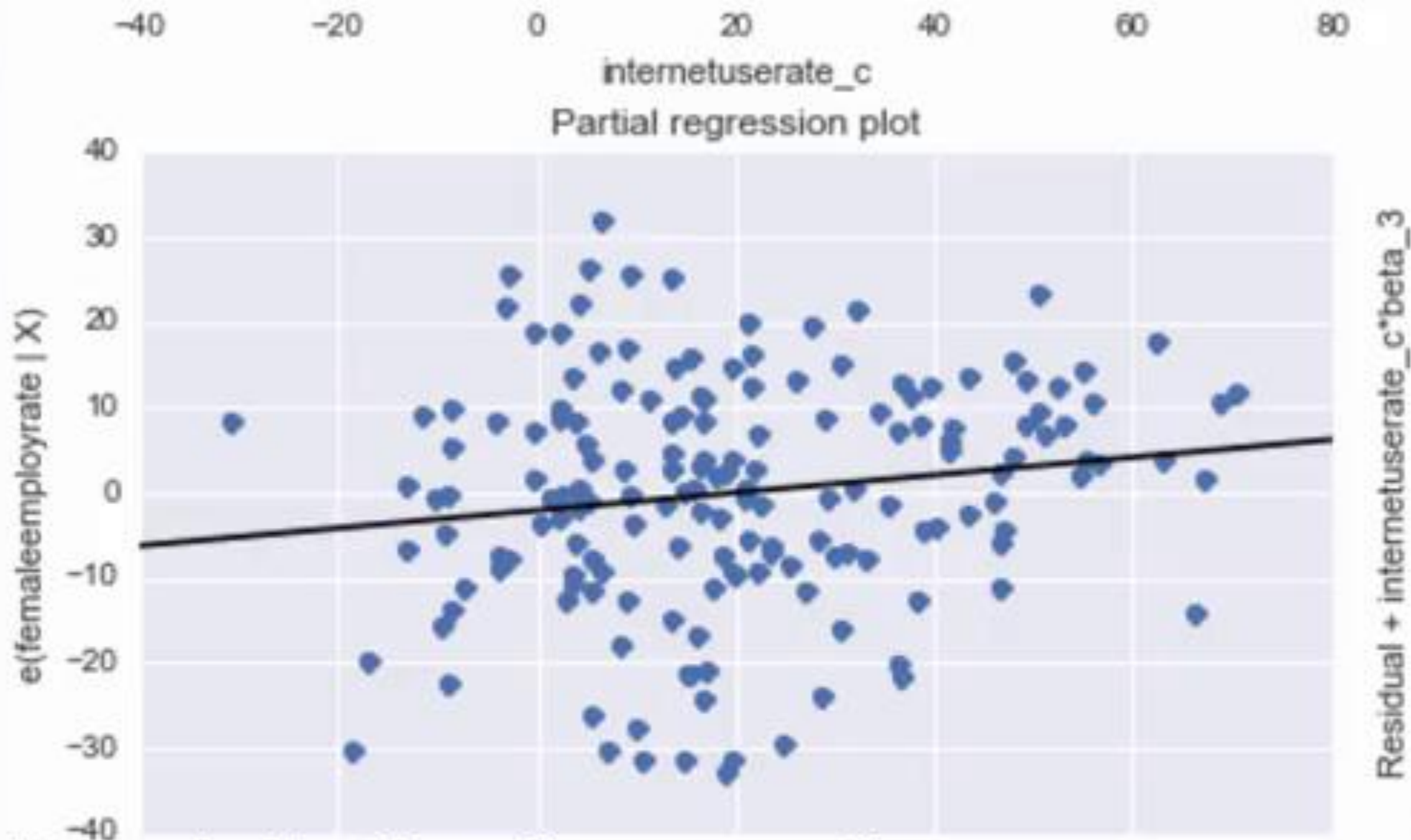


between the response variable and
specific explanatory variable,

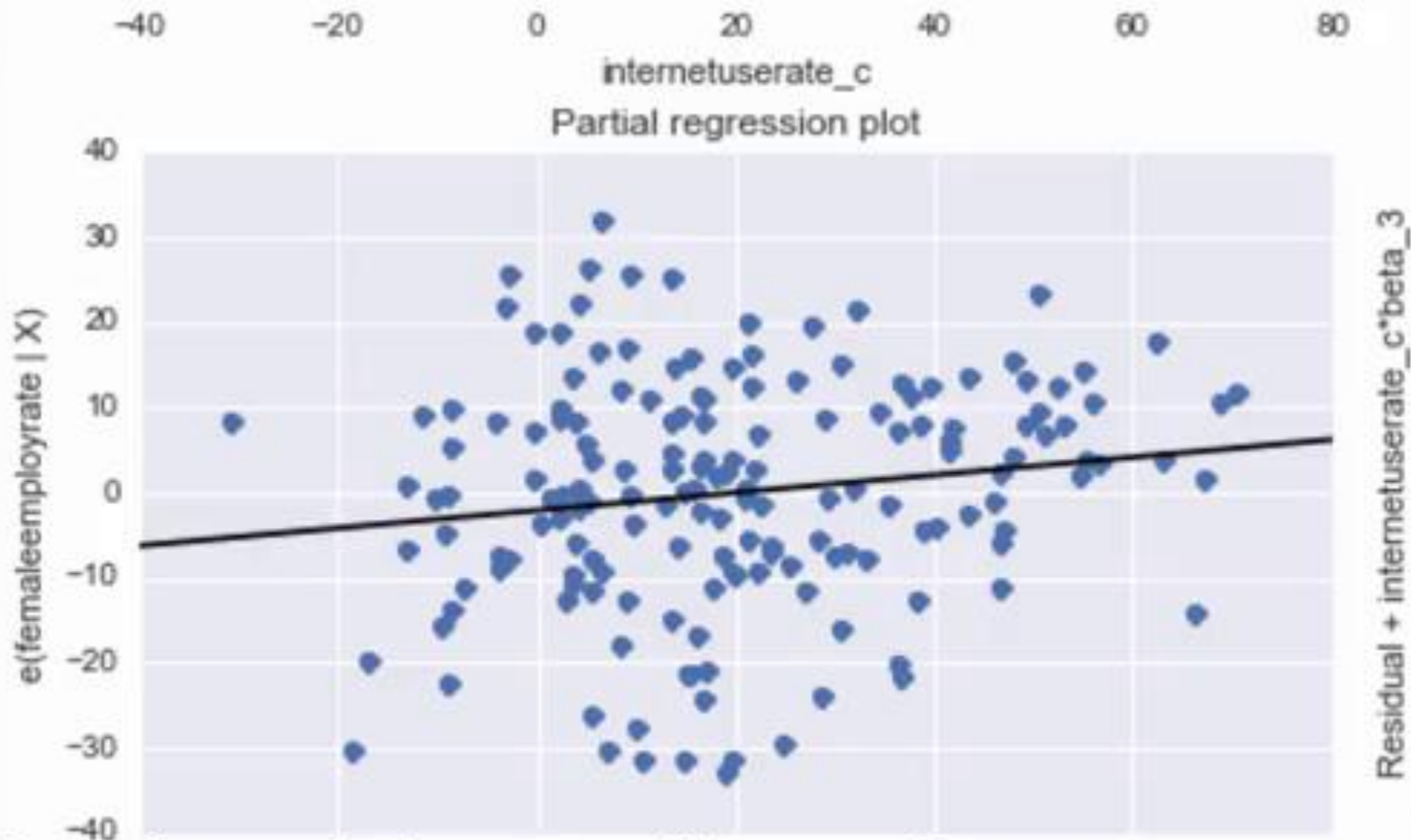




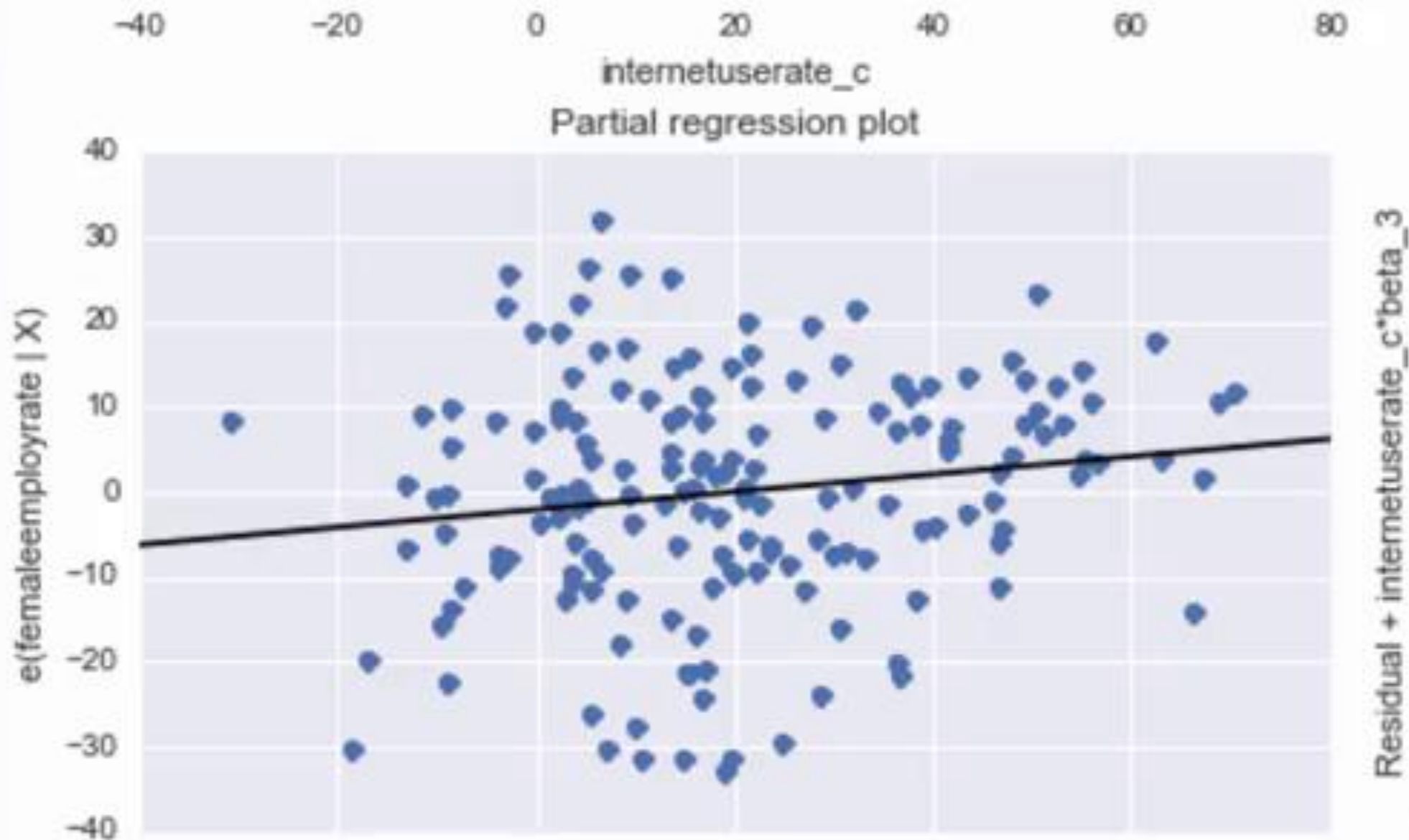
We can examine the plot to see if the Internet use rate residuals show a linear,



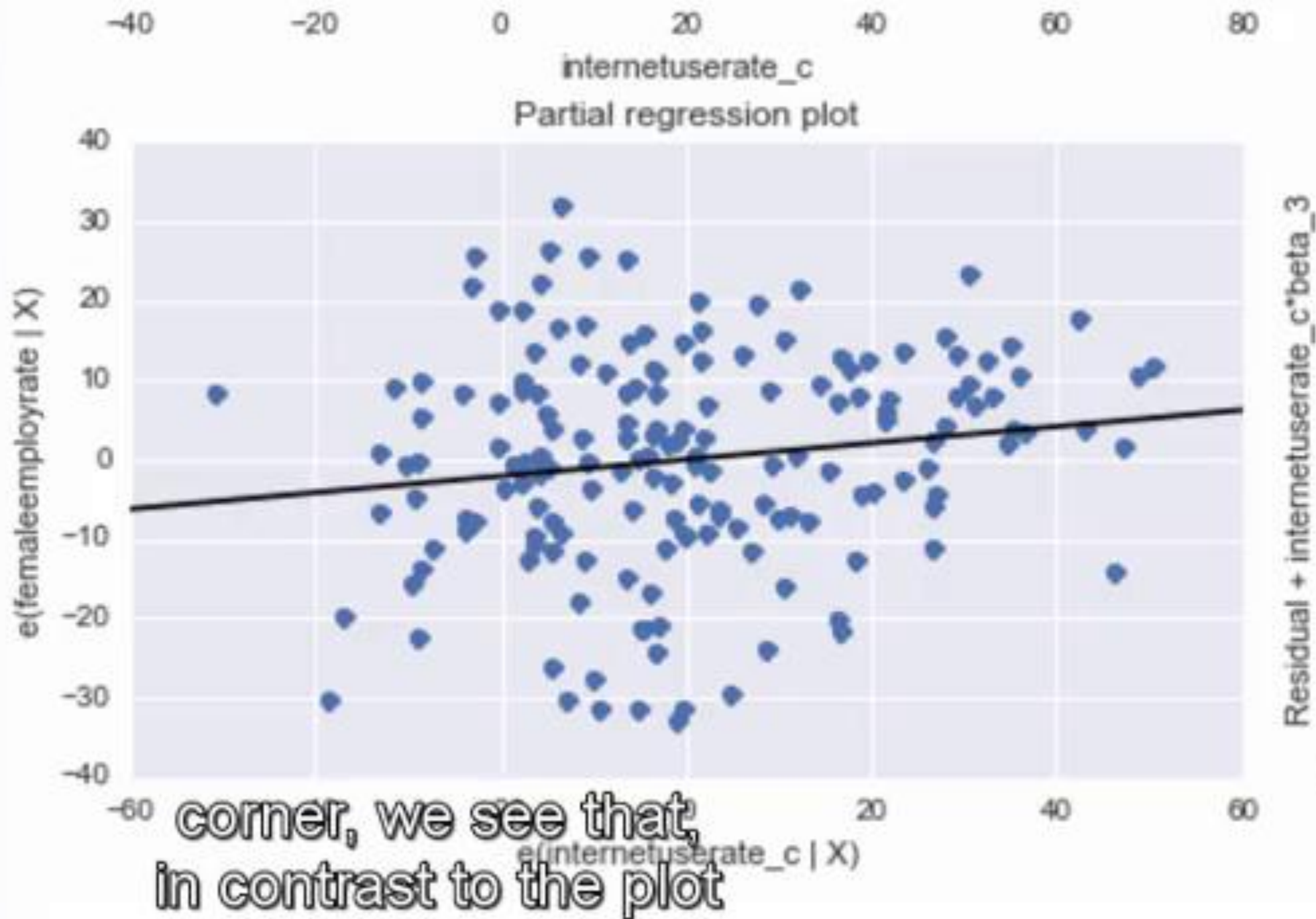
it meets the linearity assumption
in the multiple regression.

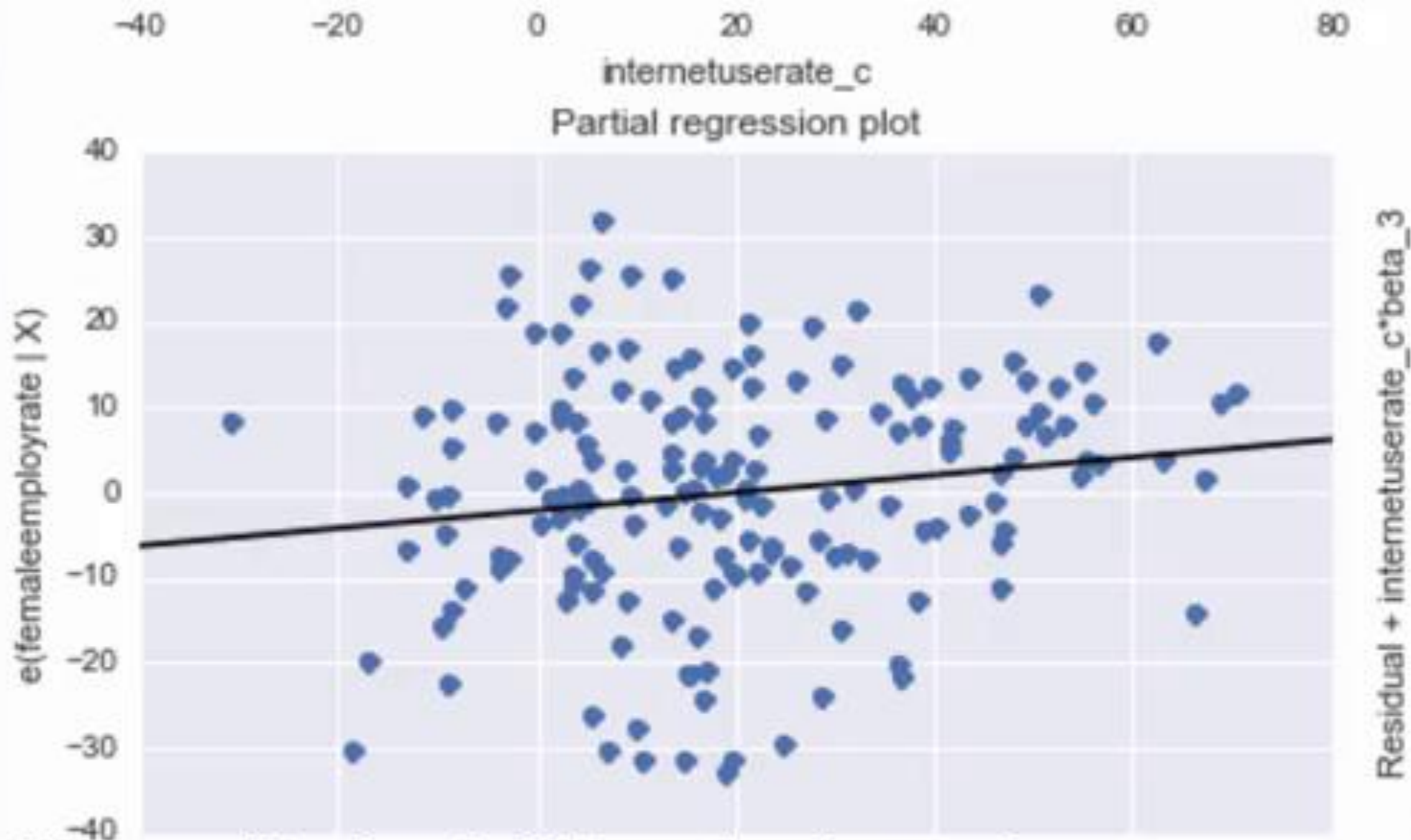


If there⁶ is an obvious non-linear pattern,
this would be additional support for

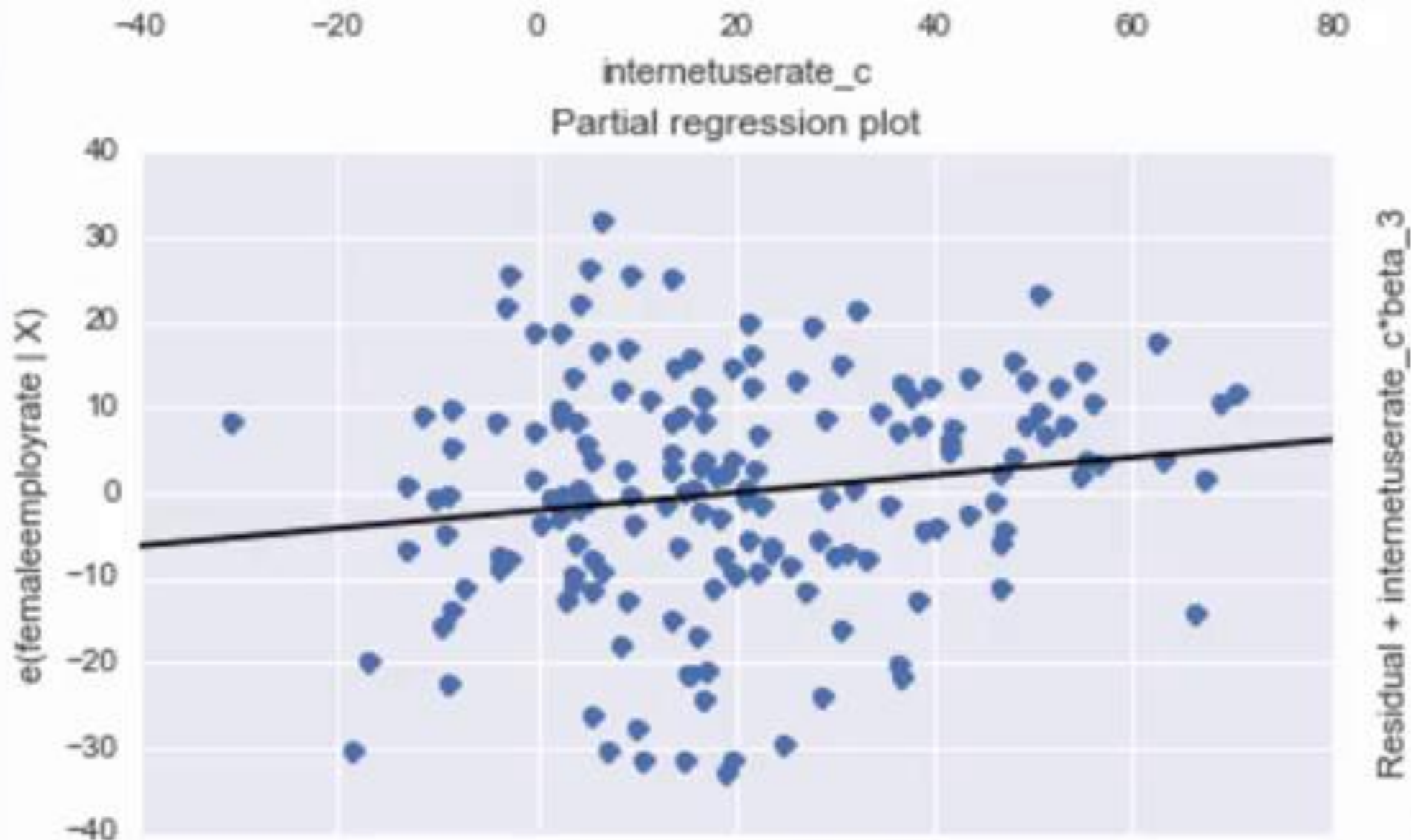


adding a polynomial term for
Internet use rate to the model.

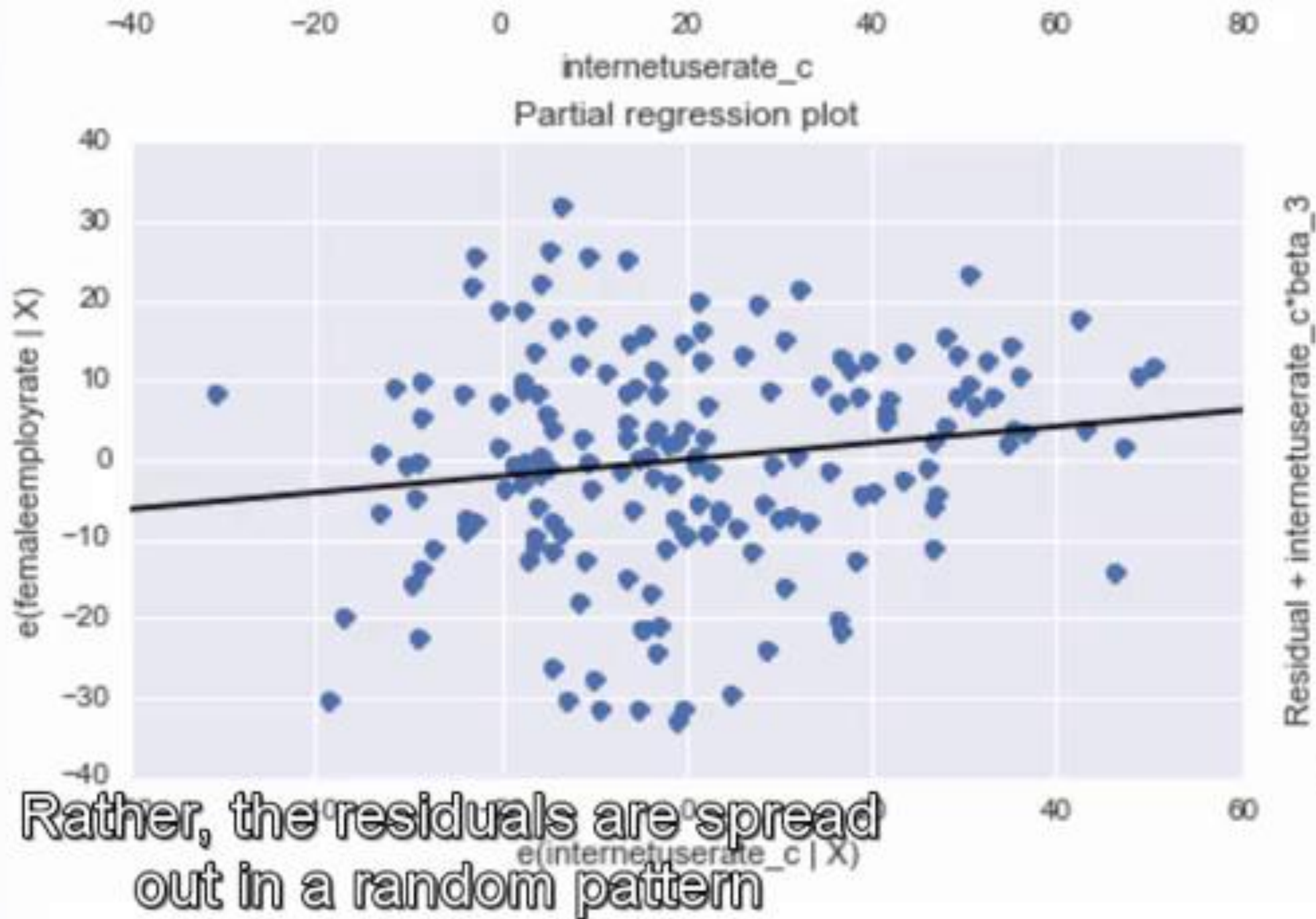


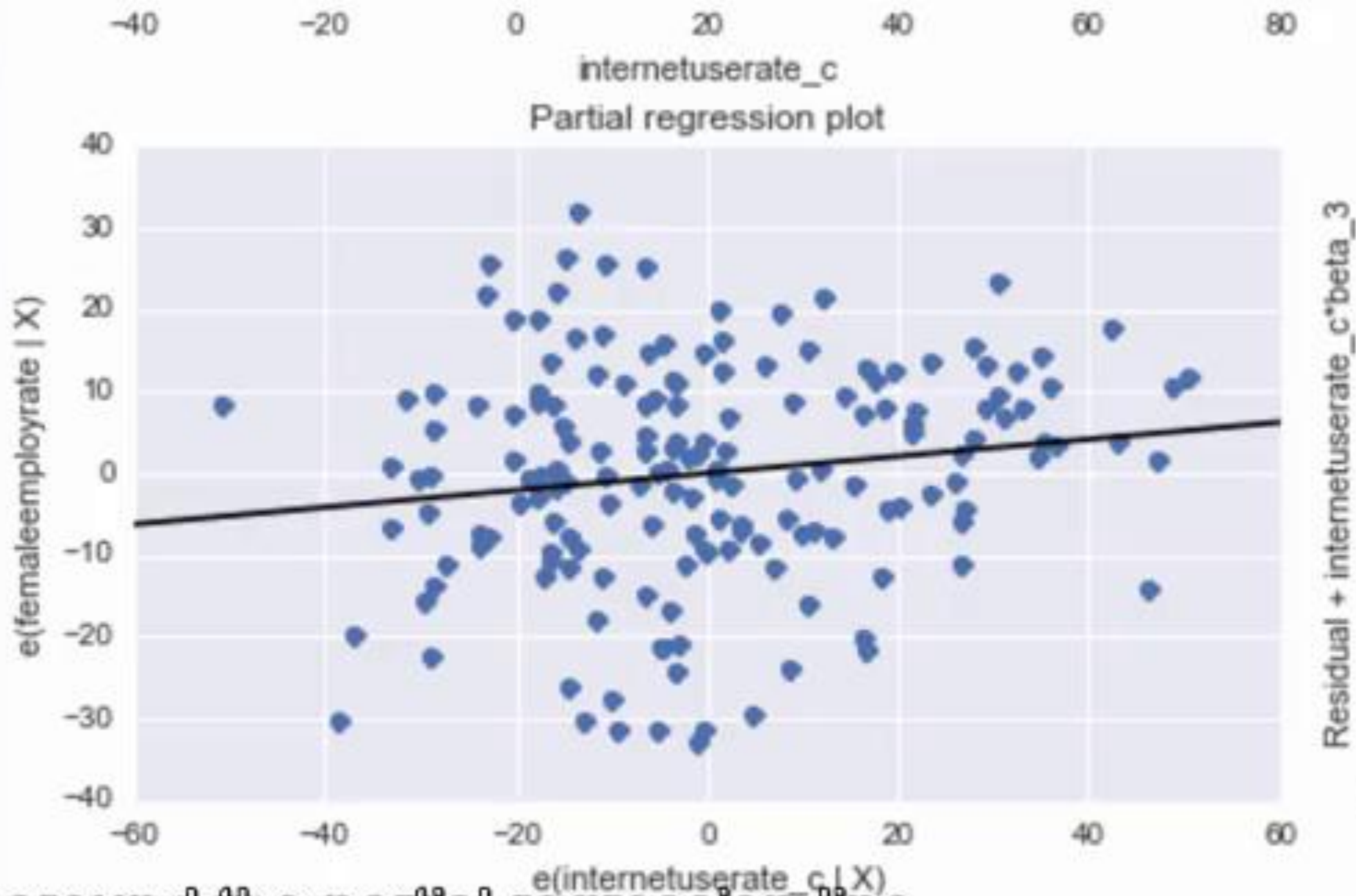


of the residuals at different values of Internet use rate without adjusting for

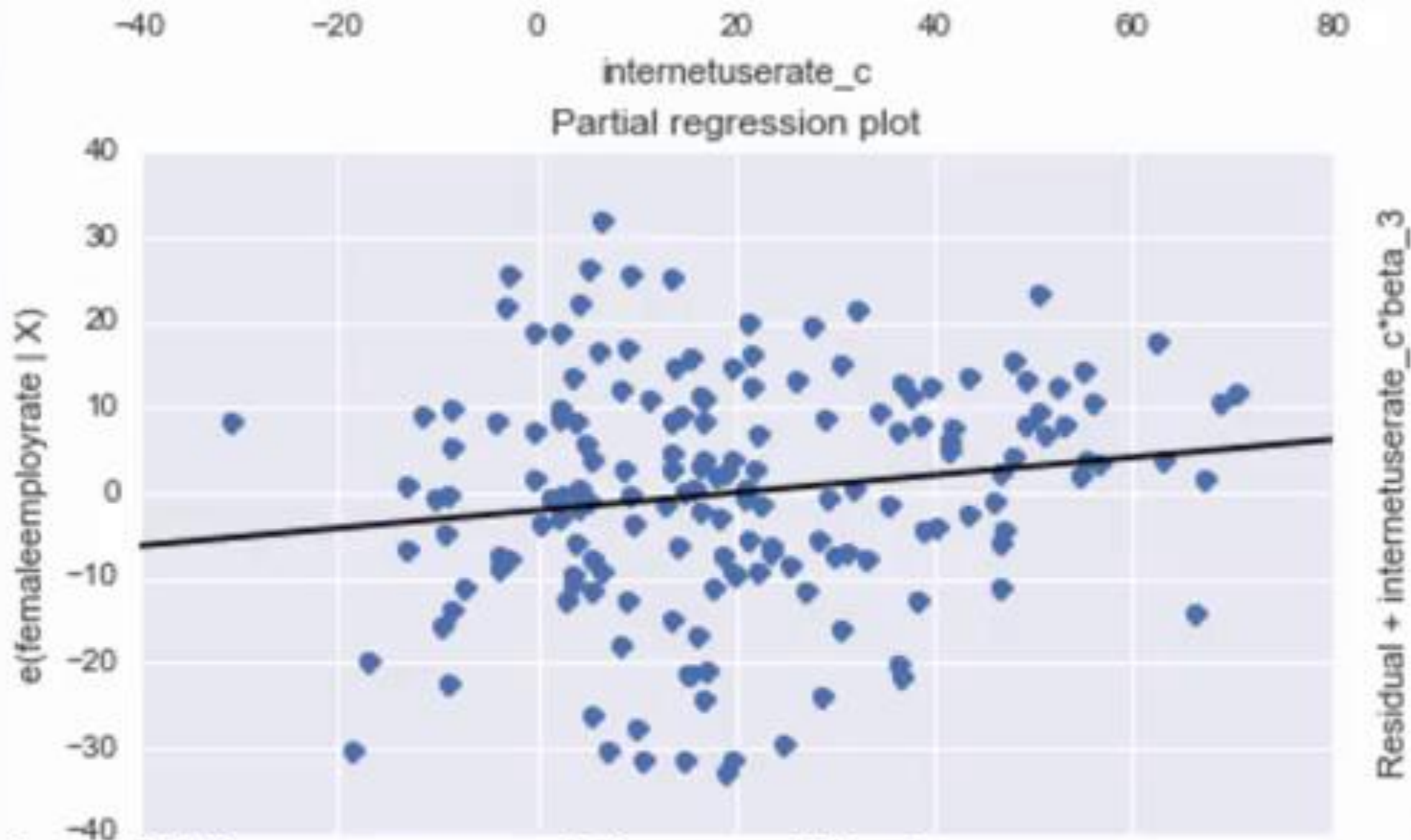


Internet use does not clearly indicate a non-linear association.

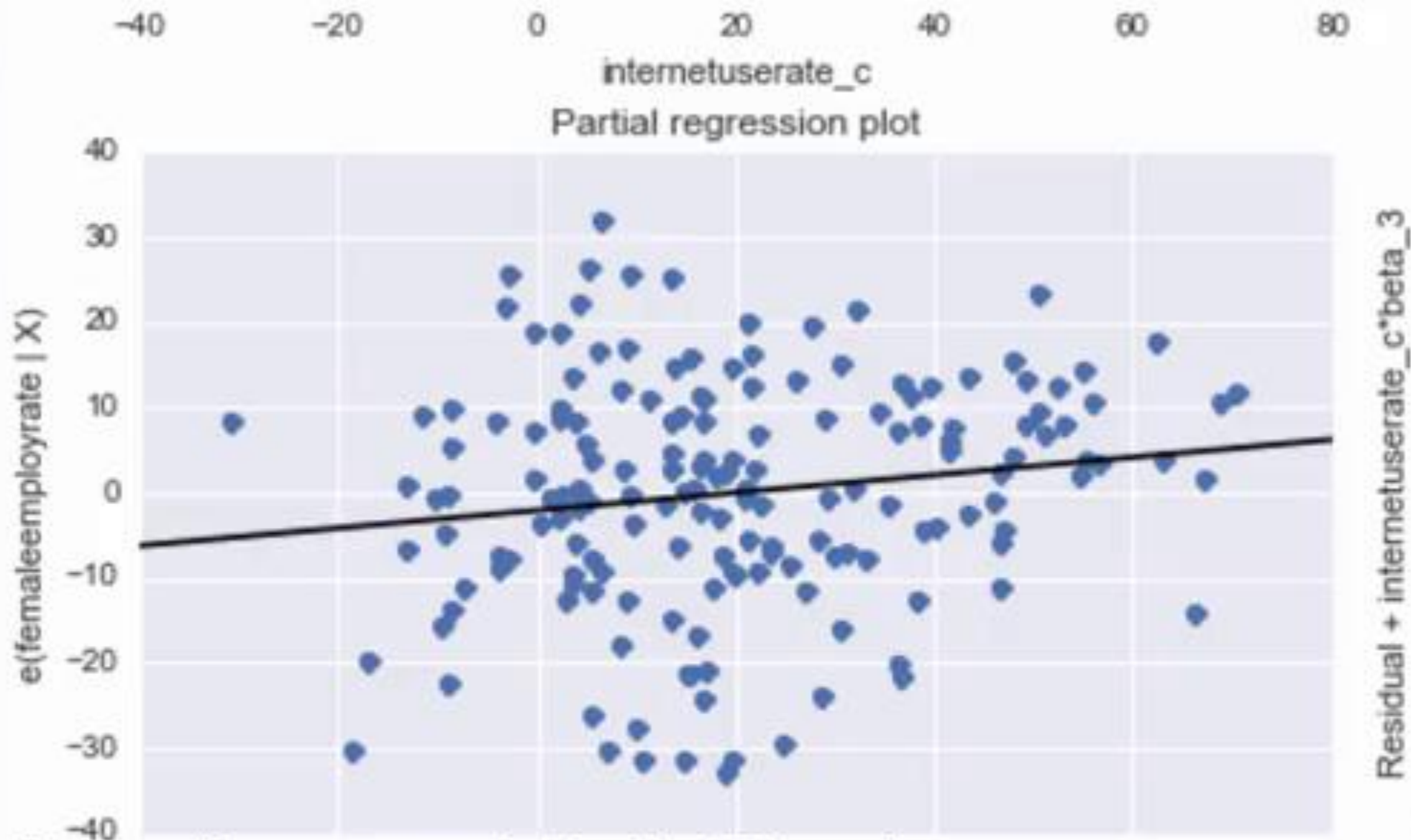




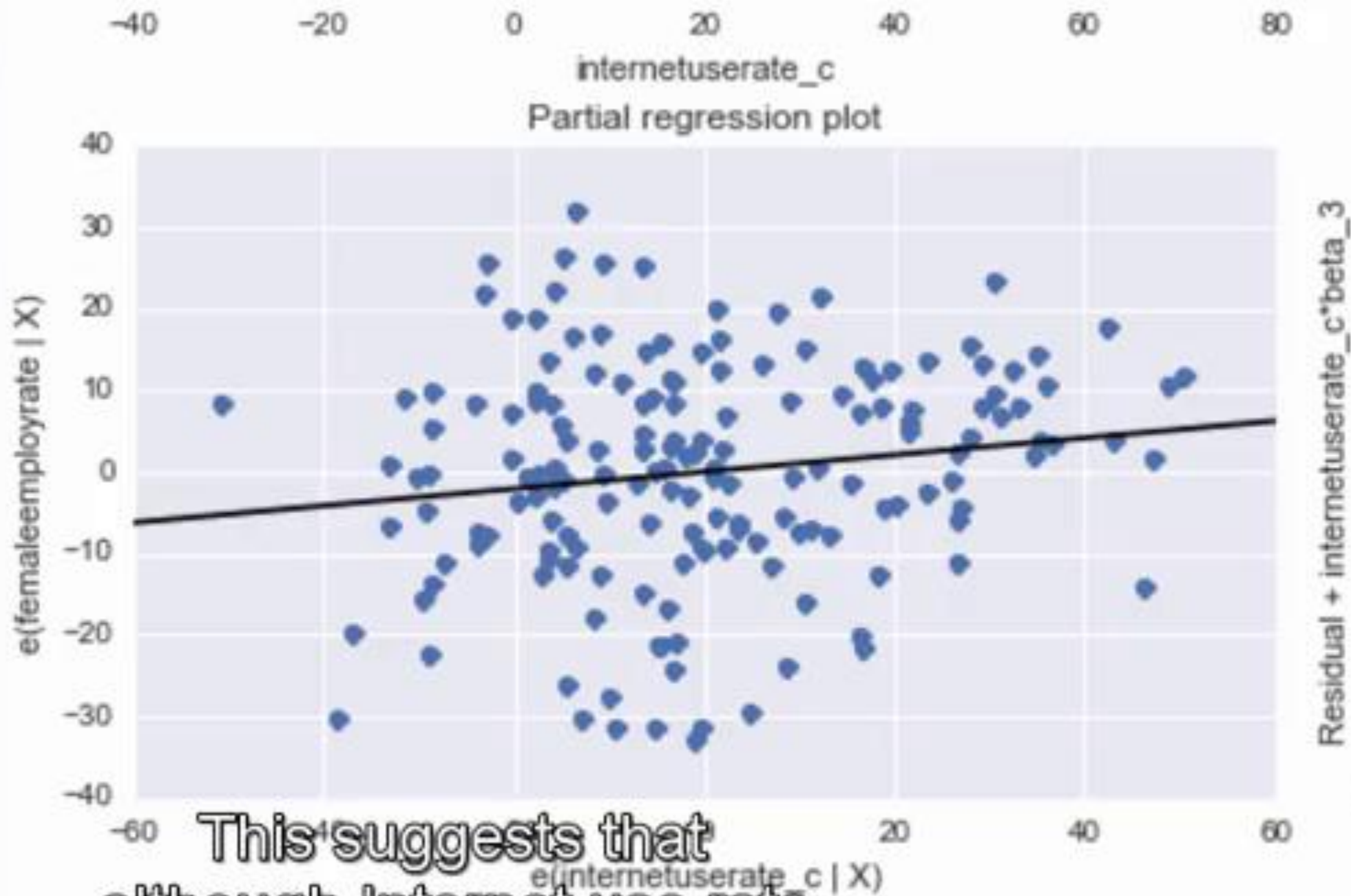
around the partial regression line.



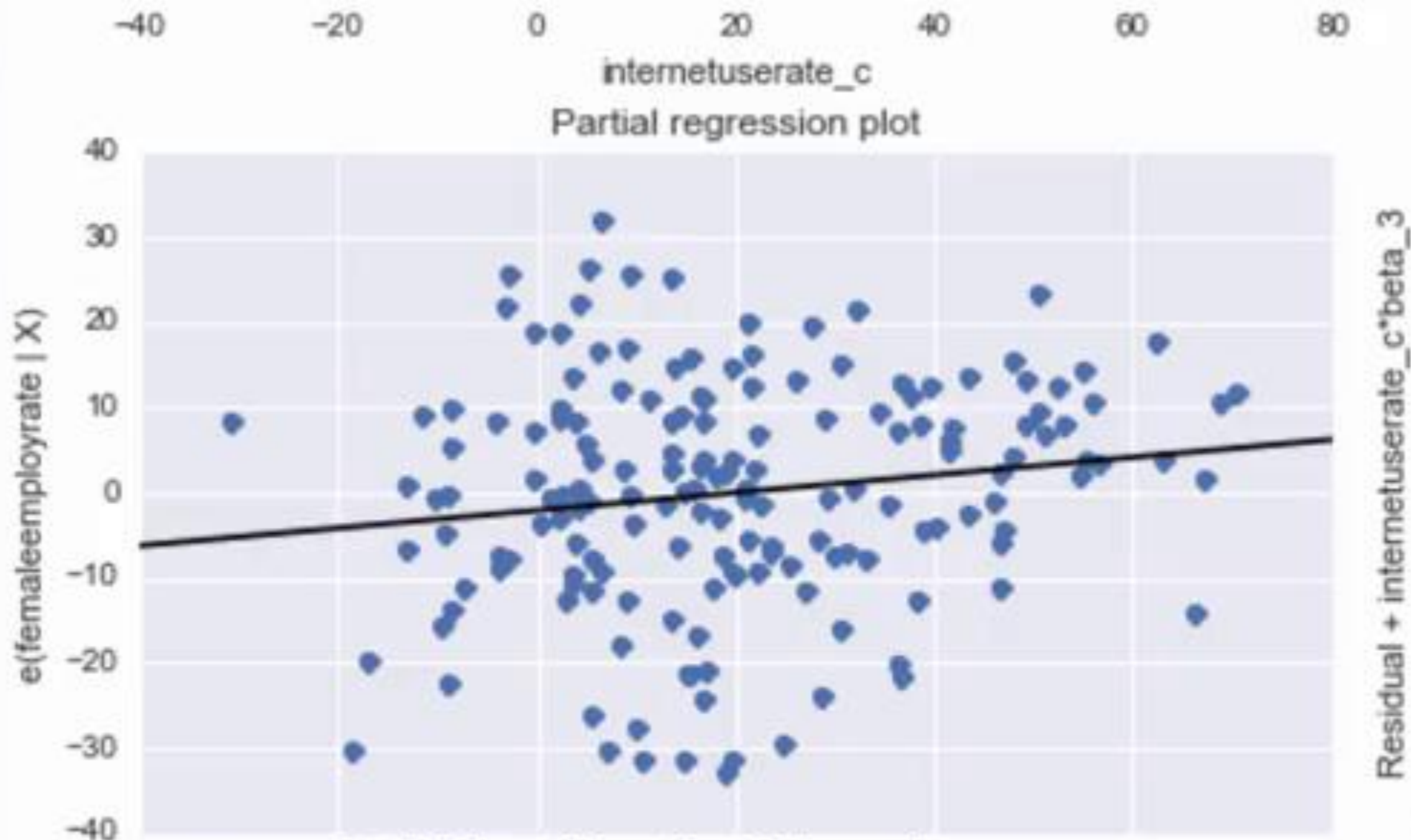
In addition, many of the residuals are pretty far from this line,



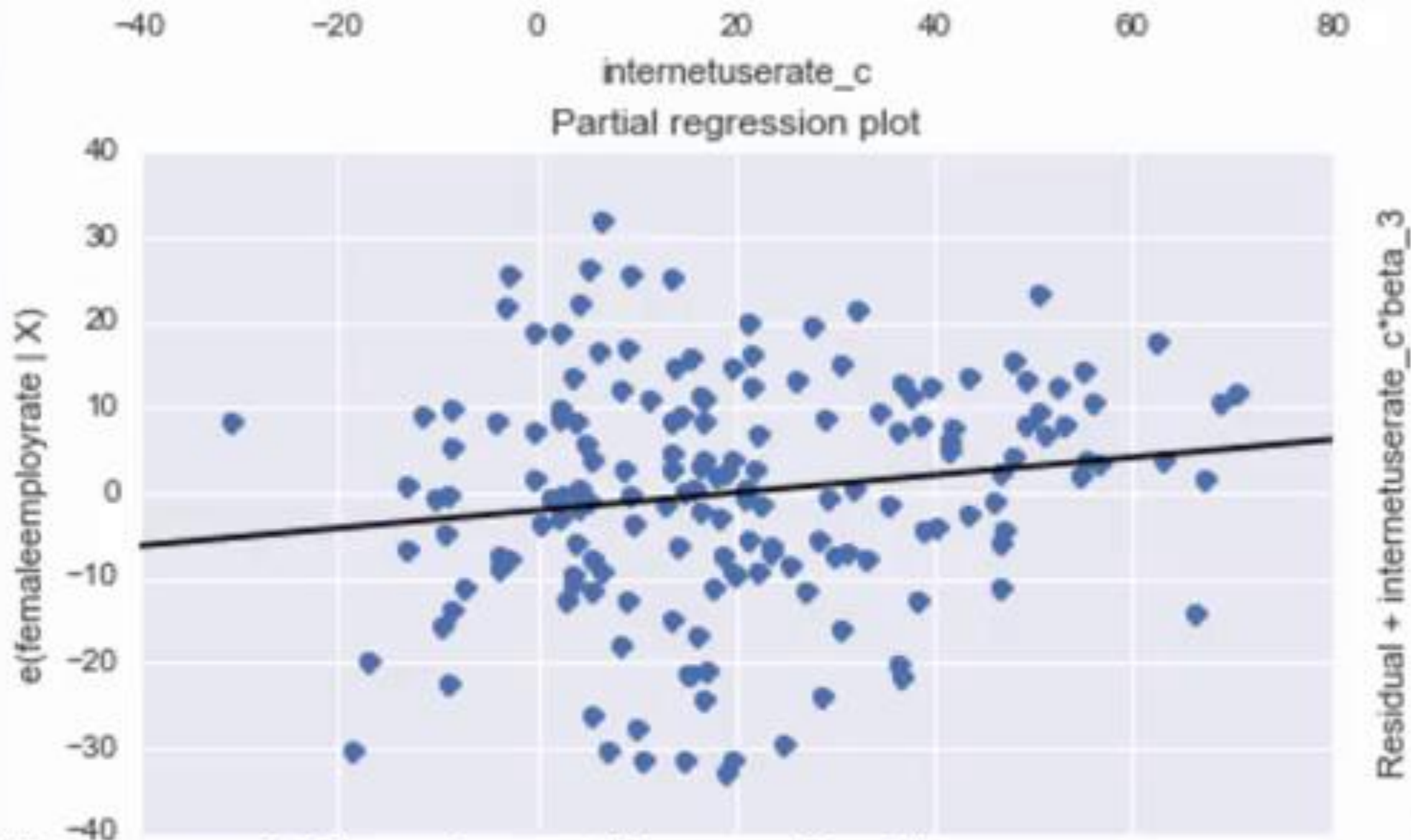
indicating a great deal of female employment rate prediction error.



This suggests that
although Internet use rate



shows a statistically significant
association with female employment rate,



this association is pretty weak after
controlling for urbanization rate.


```
60 l = plt.axhline(y=0, color='r')
61 plt.ylabel('Standardized Residual')
62 plt.xlabel('Observation Number')
63 print(fig2)
64
65
66 # additional regression diagnostic plots
67 fig3 = plt.figure(figsize=(12,8))
68 fig3 = sm.graphics.plot_regress_exog(reg3, "internetuserate_c", fig=fig3)
69
70
71 # leverage plot
72 fig4=sm.graphics.influence_plot(reg3, size=8)
73 print(fig4)
74
75
76
77
78 |
```

Finally, we can examine a leverage plot
to identify observations that have

```
60 l = plt.axhline(y=0, color='r')
61 plt.ylabel('Standardized Residual')
62 plt.xlabel('Observation Number')
63 print(fig2)
64
65
66 # additional regression diagnostic plots
67 fig3 = plt.figure(figsize=(12,8))
68 fig3 = sm.graphics.plot_regress_exog(reg3, "internetuserate_c", fig=fig3)
69
70
71 # leverage plot
72 fig4=sm.graphics.influence_plot(reg3, size=8)
73 print(fig4)
74
75
76
77
78 |
```

an unusually large influence on
the estimation of the predicted value of

```
60 l = plt.axhline(y=0, color='r')
61 plt.ylabel('Standardized Residual')
62 plt.xlabel('Observation Number')
63 print(fig2)
64
65
66 # additional regression diagnostic plots
67 fig3 = plt.figure(figsize=(12,8))
68 fig3 = sm.graphics.plot_regress_exog(reg3, "internetuserate_c", fig=fig3)
69
70
71 # leverage plot
72 fig4=sm.graphics.influence_plot(reg3, size=8)
73 print(fig4)
74
75
76
77
78 |
```

of how much the predicted scores for
the other observations would differ

```
60 l = plt.axhline(y=0, color='r')
61 plt.ylabel('Standardized Residual')
62 plt.xlabel('Observation Number')
63 print(fig2)
64
65
66 # additional regression diagnostic plots
67 fig3 = plt.figure(figsize=(12,8))
68 fig3 = sm.graphics.plot_regress_exog(reg3, "internetuserate_c", fig=fig3)
69
70
71 # leverage plot
72 fig4=sm.graphics.influence_plot(reg3, size=8)
73 print(fig4)
74
75
76
77
78 |
```

if the observations in question
were not included in the analysis.


```
60 l = plt.axhline(y=0, color='r')
61 plt.ylabel('Standardized Residual')
62 plt.xlabel('Observation Number')
63 print(fig2)
64
65
66 # additional regression diagnostic plots
67 fig3 = plt.figure(figsize=(12,8))
68 fig3 = sm.graphics.plot_regress_exog(reg3, "internetuserate_c", fig=fig3)
69
70
71 # leverage plot
72 fig4=sm.graphics.influence_plot(reg3, size=8)
73 print(fig4)
74
75
76
77
78 |
```

The leverage always takes on values between zero and one.

```
60 l = plt.axhline(y=0, color='r')
61 plt.ylabel('Standardized Residual')
62 plt.xlabel('Observation Number')
63 print(fig2)
64
65
66 # additional regression diagnostic plots
67 fig3 = plt.figure(figsize=(12,8))
68 fig3 = sm.graphics.plot_regress_exog(reg3, "internetuserate_c", fig=fig3)
69
70
71 # leverage plot
72 fig4=sm.graphics.influence_plot(reg3, size=8)
73 print(fig4)
74
75
76
77
78 |
```

A point with zero leverage has no effect on the regression model.

```
60 l = plt.axhline(y=0, color='r')
61 plt.ylabel('Standardized Residual')
62 plt.xlabel('Observation Number')
63 print(fig2)
64
65
66 # additional regression diagnostic plots
67 fig3 = plt.figure(figsize=(12,8))
68 fig3 = sm.graphics.plot_regress_exog(reg3, "internetuserate_c", fig=fig3)
69
70
71 # leverage plot
72 fig4=sm.graphics.influence_plot(reg3, size=8)
73 print(fig4)
74
75
76
77
78 |
```

And outliers are observations with residuals greater than 2 or less than -2.

```
60 l = plt.axhline(y=0, color='r')
61 plt.ylabel('Standardized Residual')
62 plt.xlabel('Observation Number')
63 print(fig2)
64
65
66 # additional regression diagnostic plots
67 fig3 = plt.figure(figsize=(12,8))
68 fig3 = sm.graphics.plot_regress_exog(reg3, "internetuserate_c", fig=fig3)
69
70
71 # leverage plot
72 fig4=sm.graphics.influence_plot(reg3, size=8)
73 print(fig4)
74
75
76
77
78 |
```

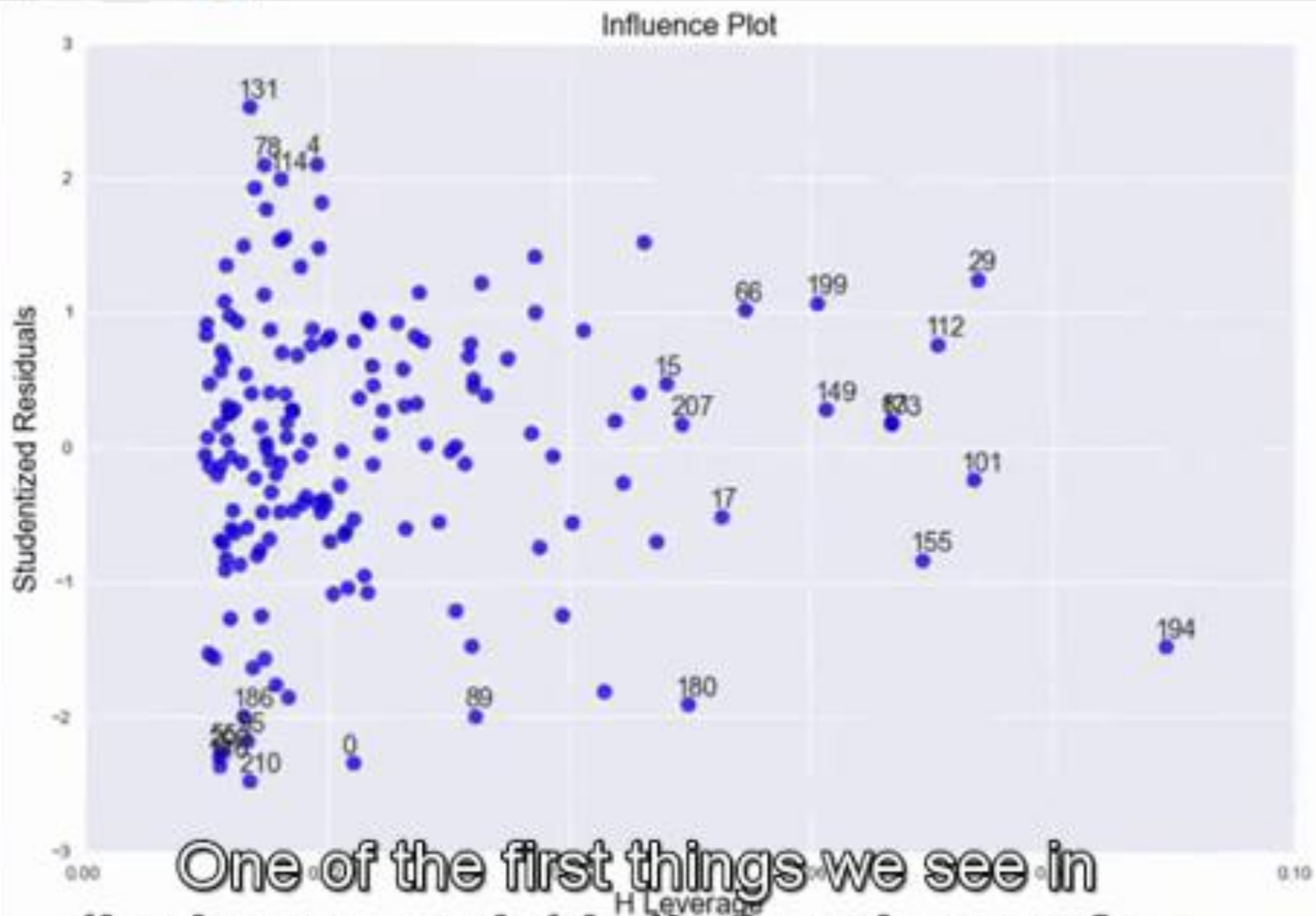
but this time we use
the code `influence_plot`.

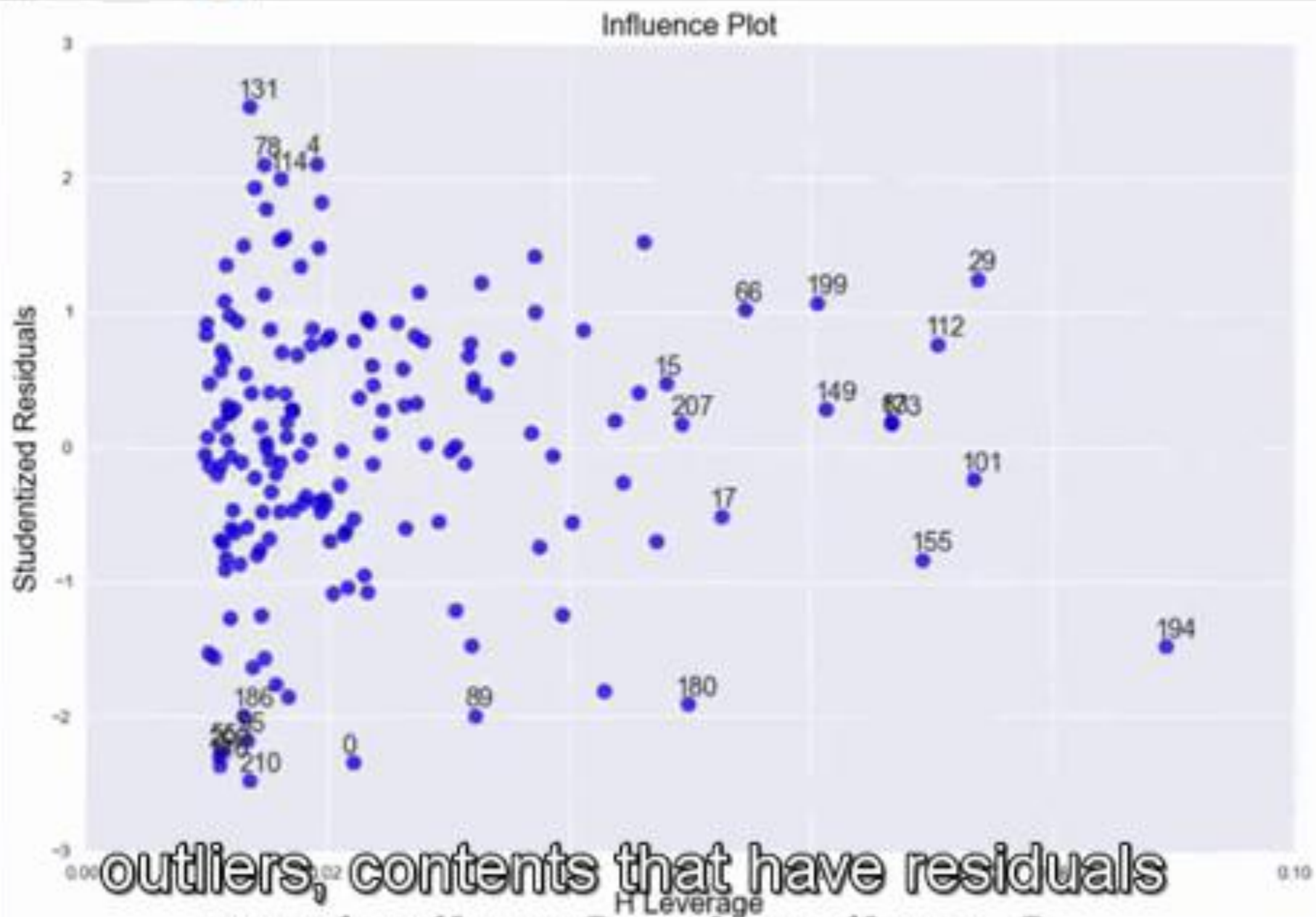

```
60 l = plt.axhline(y=0, color='r')
61 plt.ylabel('Standardized Residual')
62 plt.xlabel('Observation Number')
63 print(fig2)
64
65
66 # additional regression diagnostic plots
67 fig3 = plt.figure(figsize=(12,8))
68 fig3 = sm.graphics.plot_regress_exog(reg3, "internetuserate_c", fig=fig3)
69
70
71 # leverage plot
72 fig4=sm.graphics.influence_plot(reg3, size=8)
73 print(fig4)
74
75
76
77
78 |
```

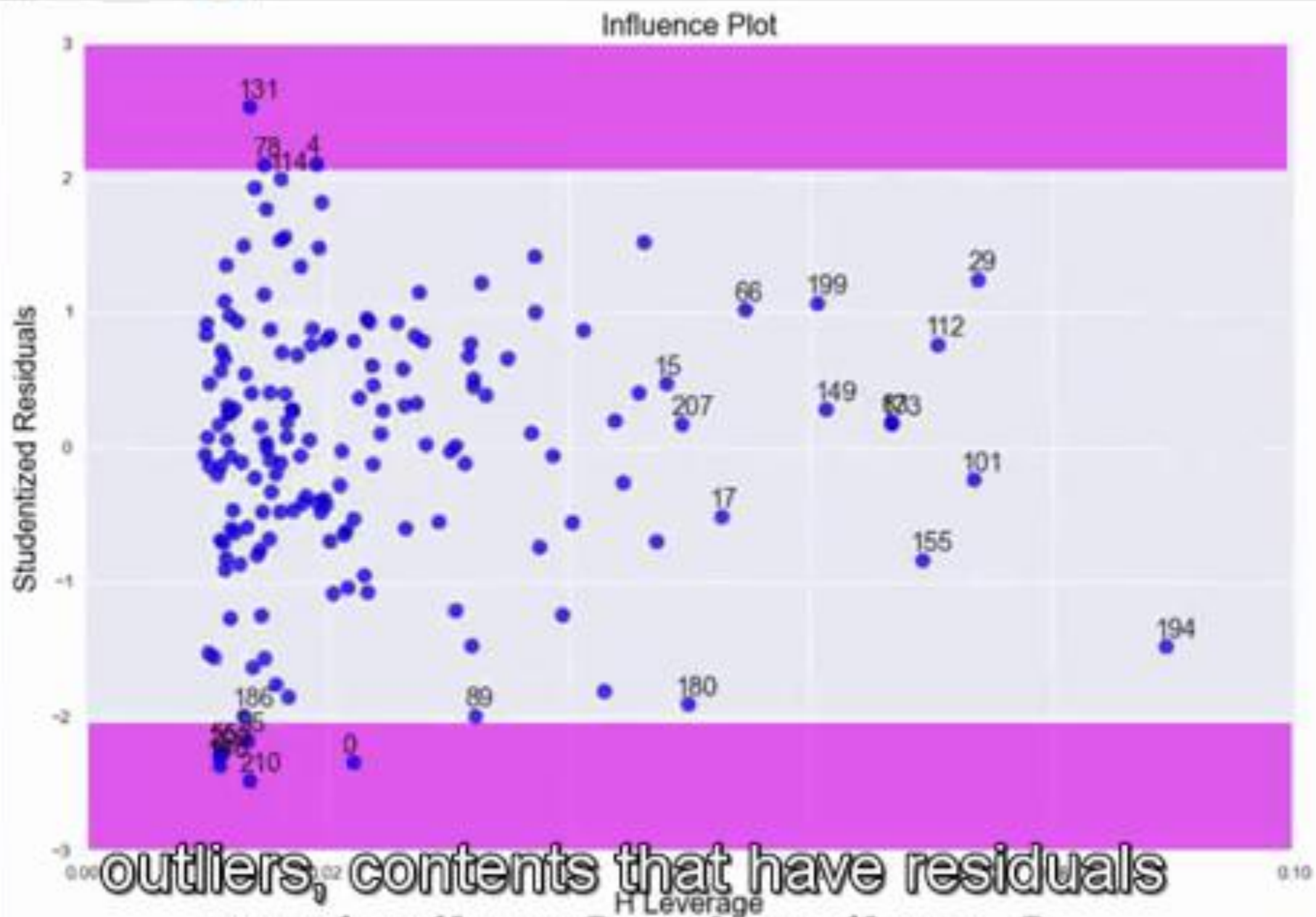
Size=8 is an option to make the points on the plot smaller than the default size so

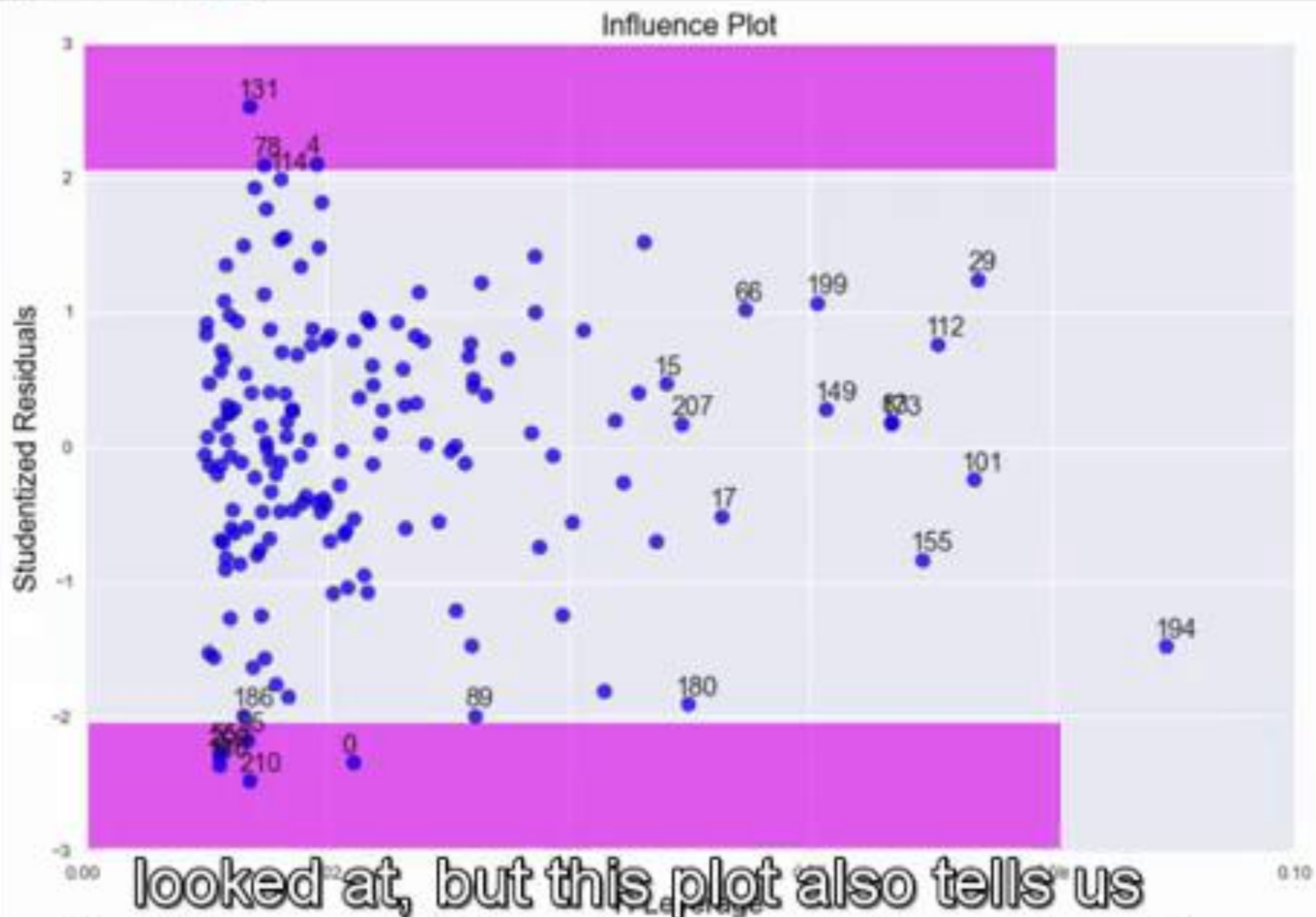
```
60 l = plt.axhline(y=0, color='r')
61 plt.ylabel('Standardized Residual')
62 plt.xlabel('Observation Number')
63 print(fig2)
64
65
66 # additional regression diagnostic plots
67 fig3 = plt.figure(figsize=(12,8))
68 fig3 = sm.graphics.plot_regress_exog(reg3, "internetuserate_c", fig=fig3)
69
70
71 # leverage plot
72 fig4=sm.graphics.influence_plot(reg3, size=8)
73 print(fig4)
74
75
76
77
78 |
```

that they're easier to distinguish.

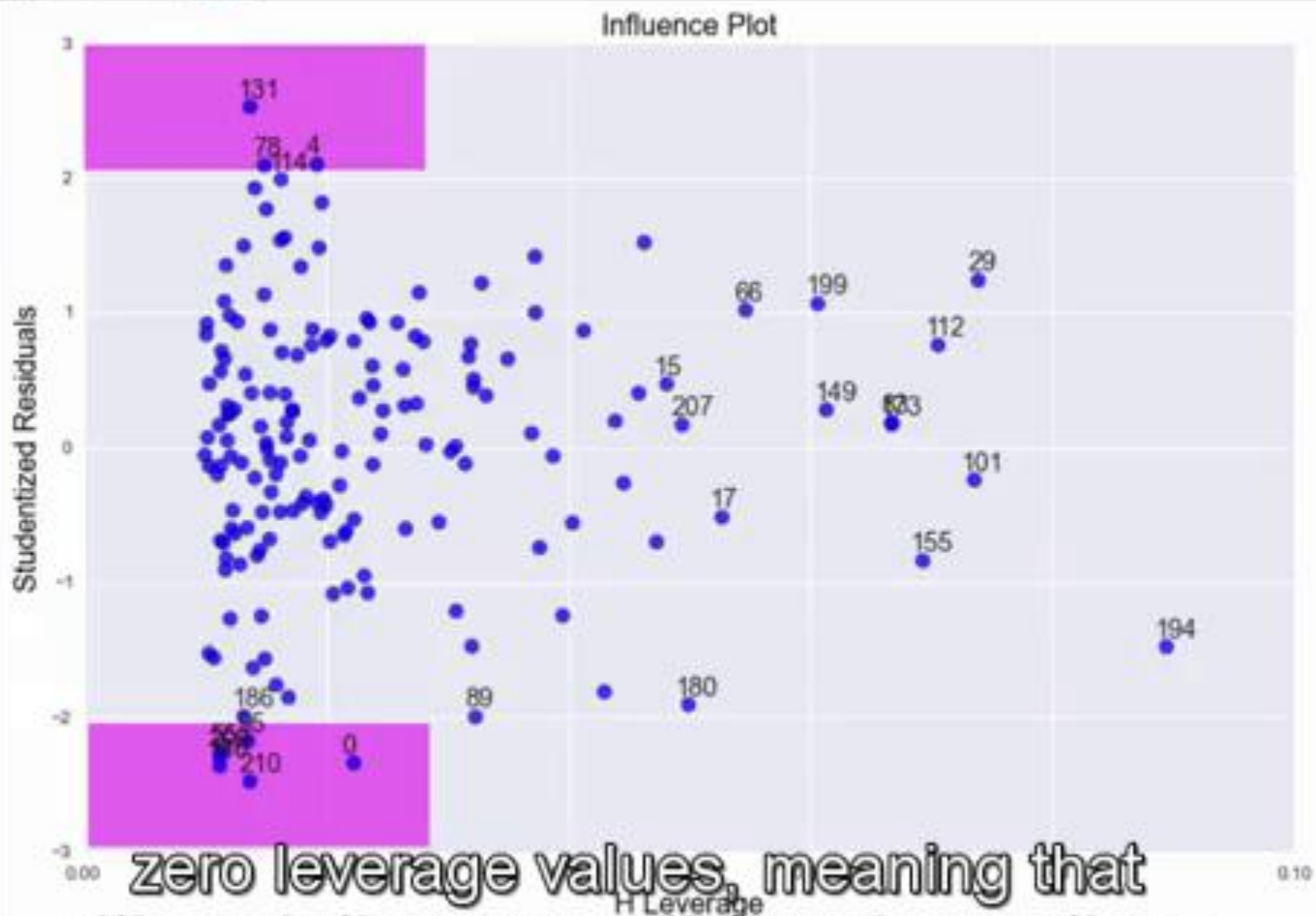




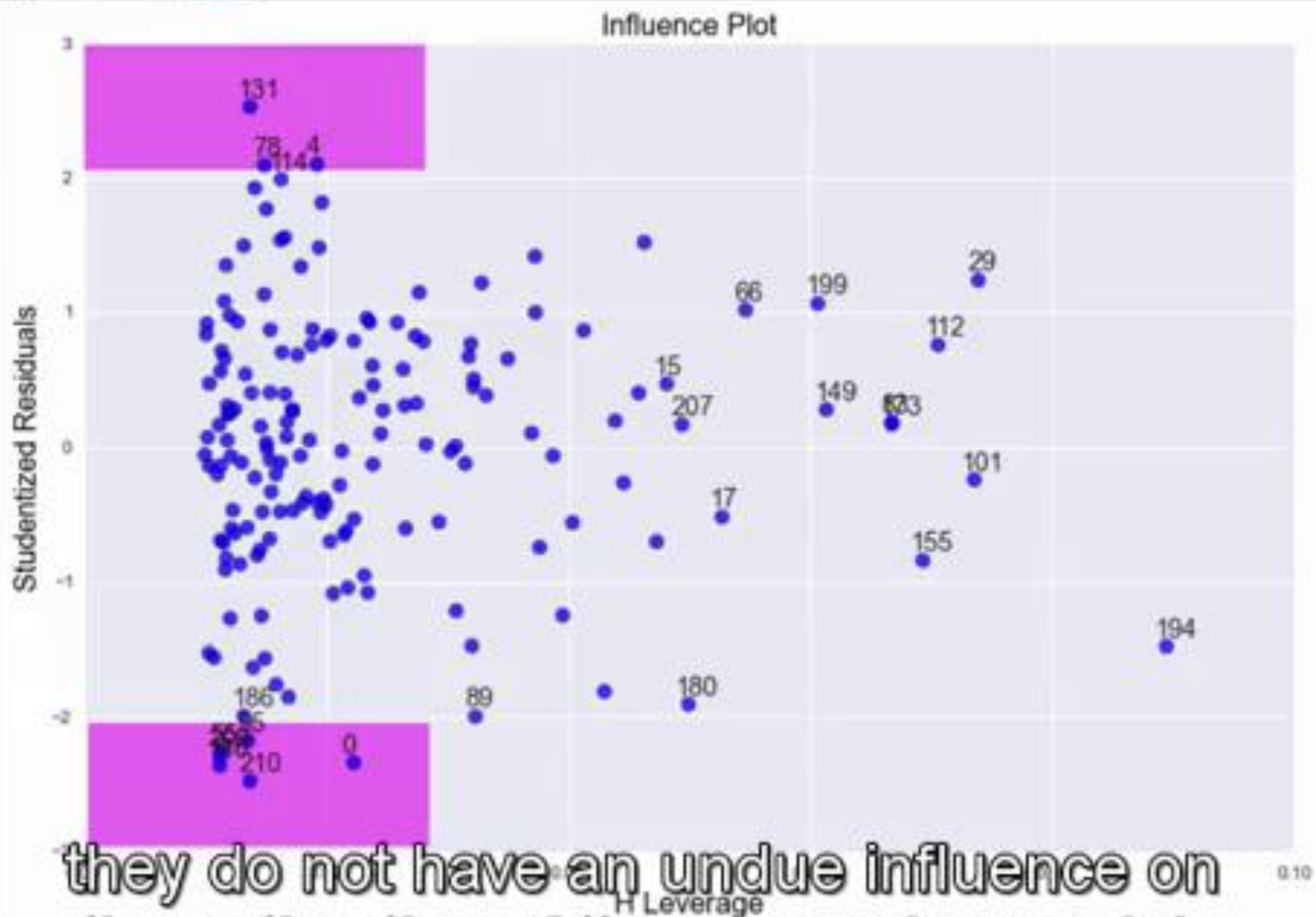




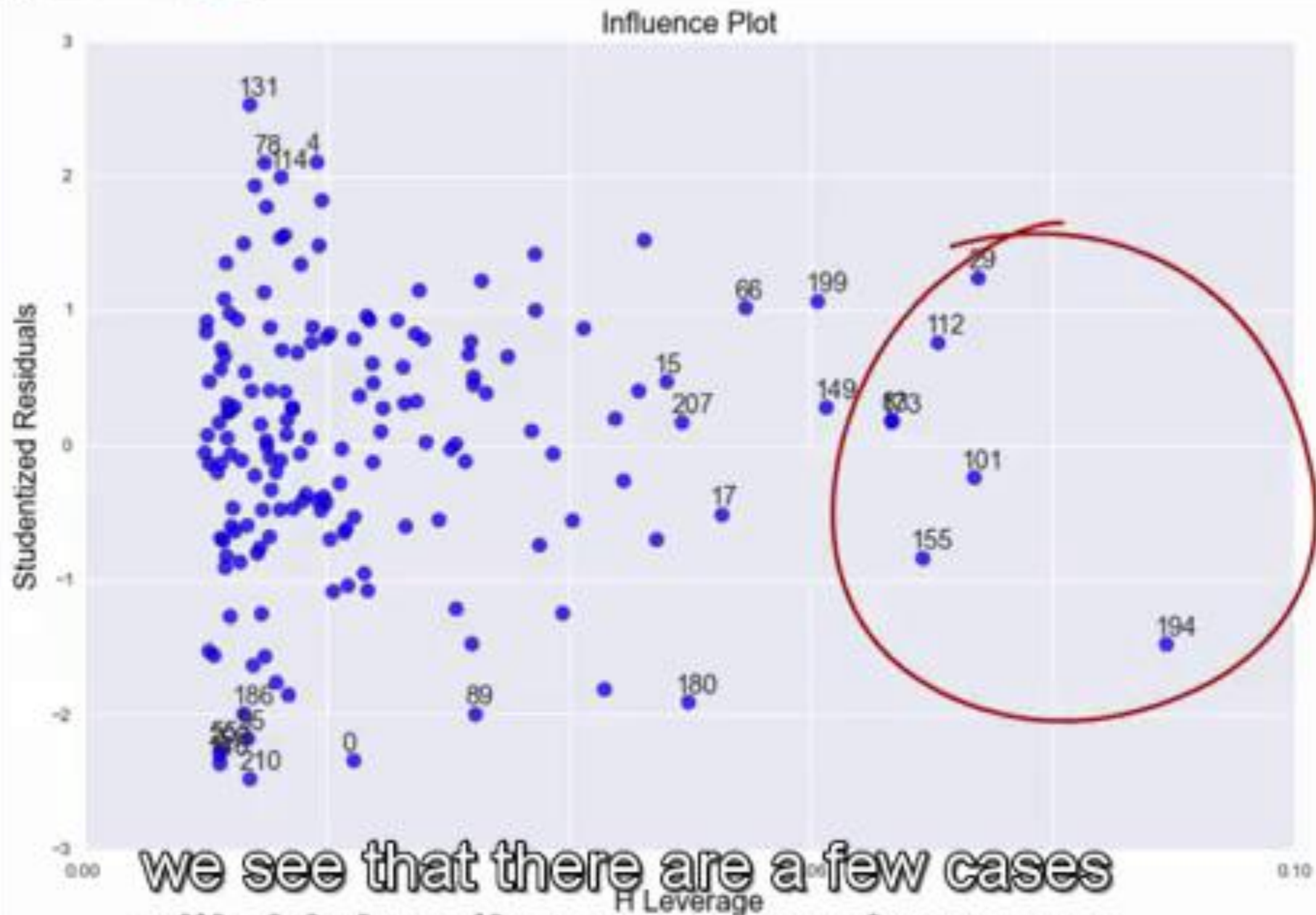
looked at, but this plot also tells us
that these outliers have small or close to

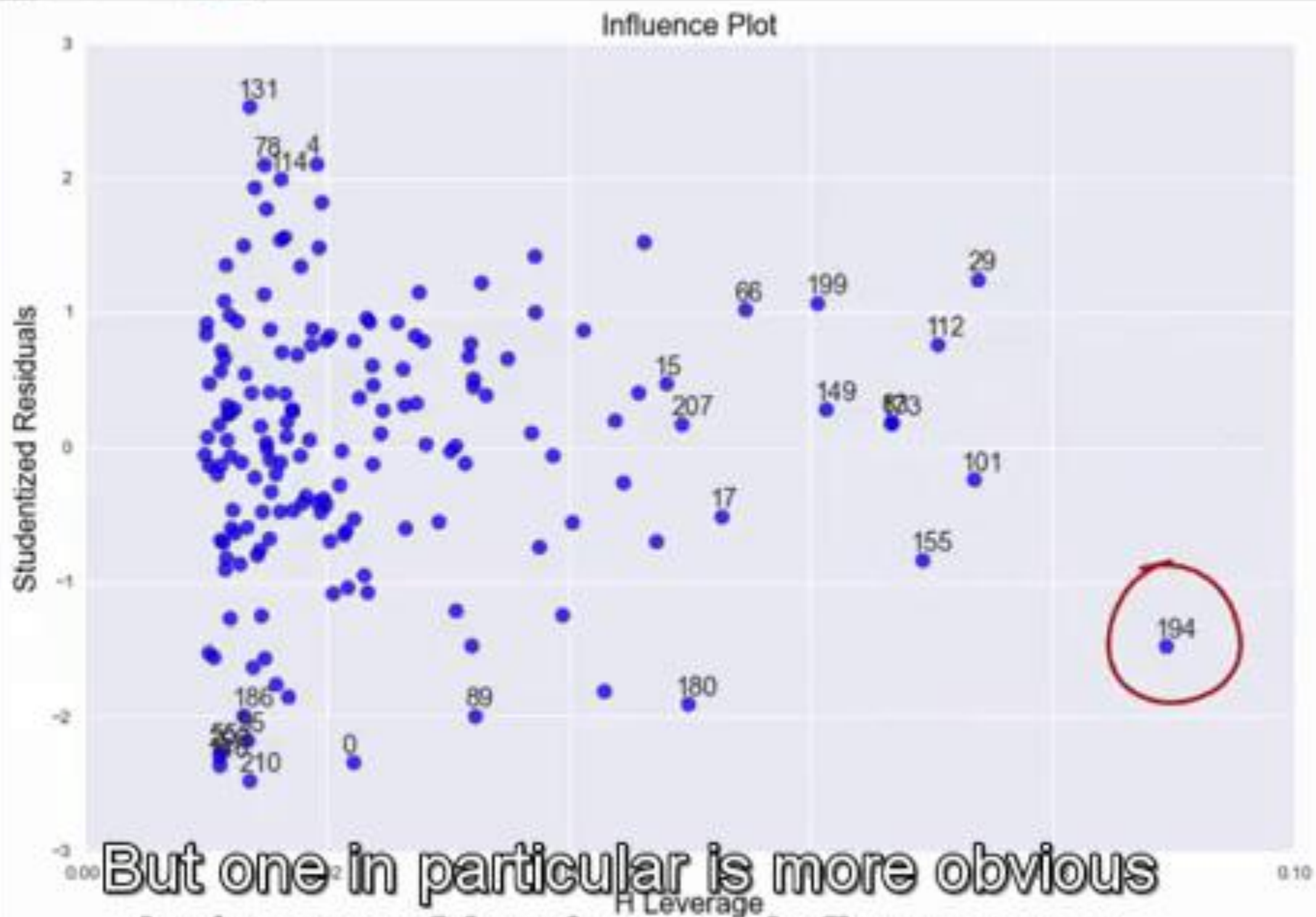


zero leverage values, meaning that
although they are outlying observations,



they do not have an undue influence on
the estimation of the regression model.

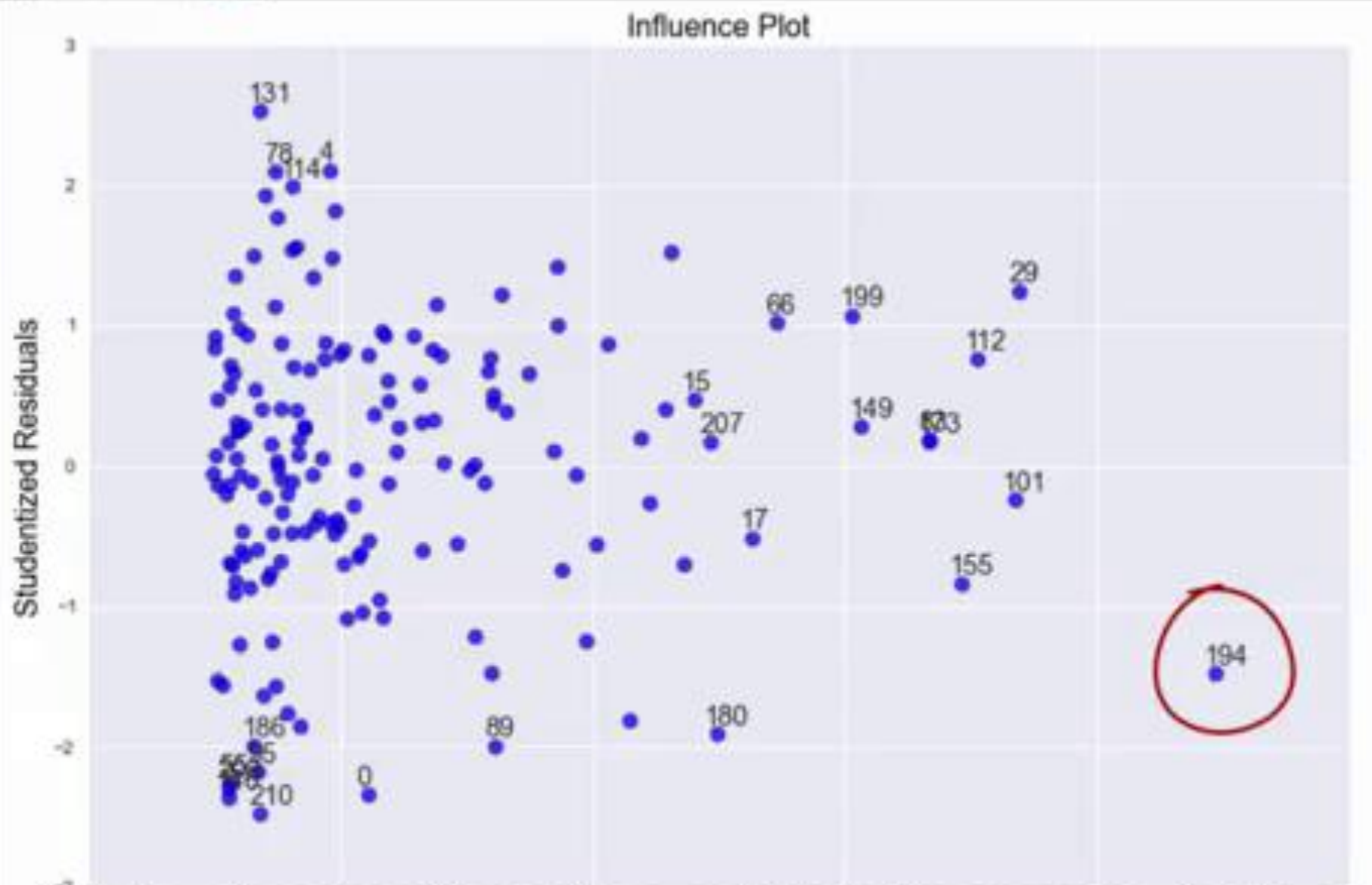




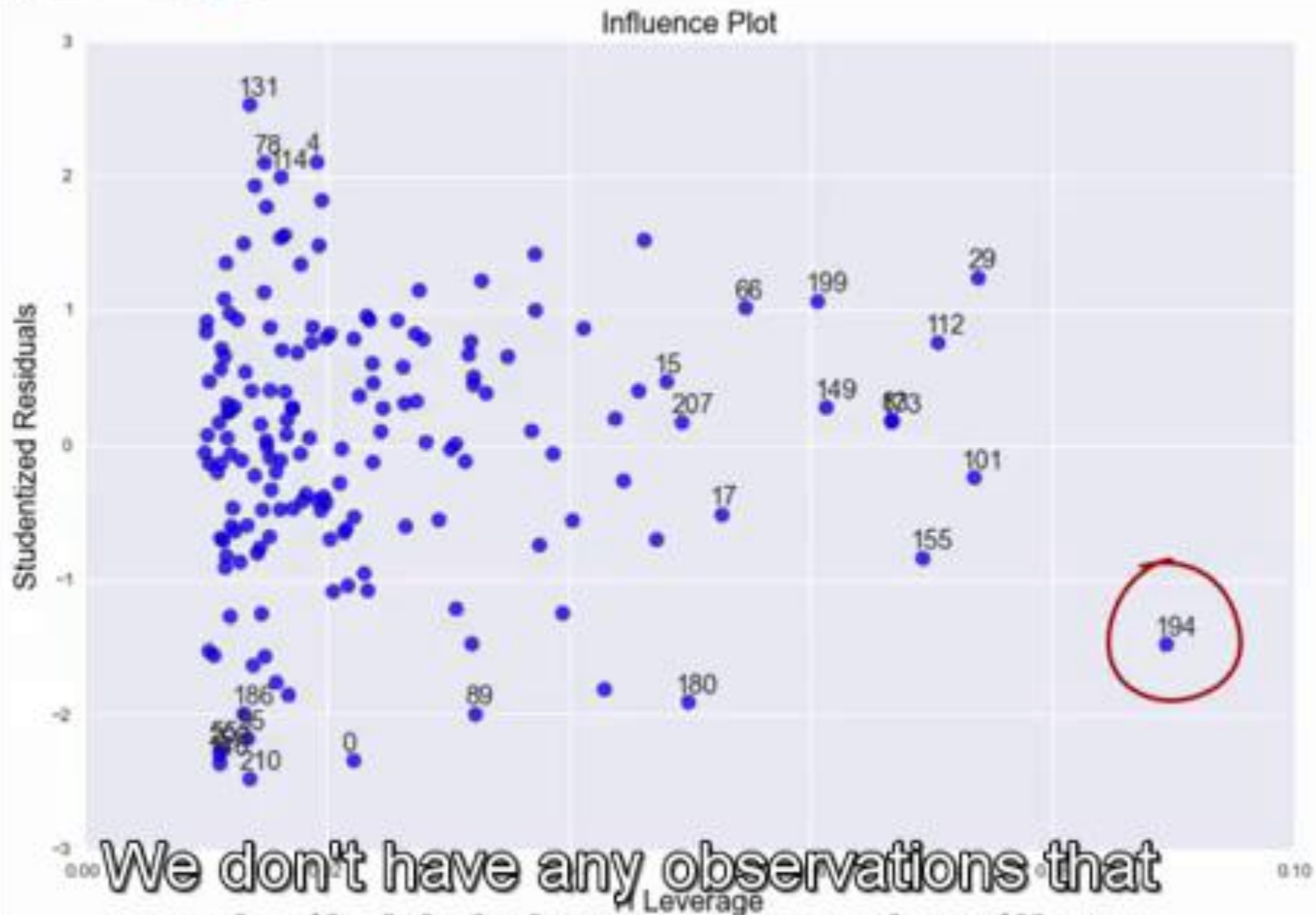
But one in particular is more obvious
in terms of having an influence on



the estimation of the predicted
value of female employment rate.



This observation has a high leverage but
is not an outlier.



We don't have any observations that are both high leverage and outliers.