# Linear Regression
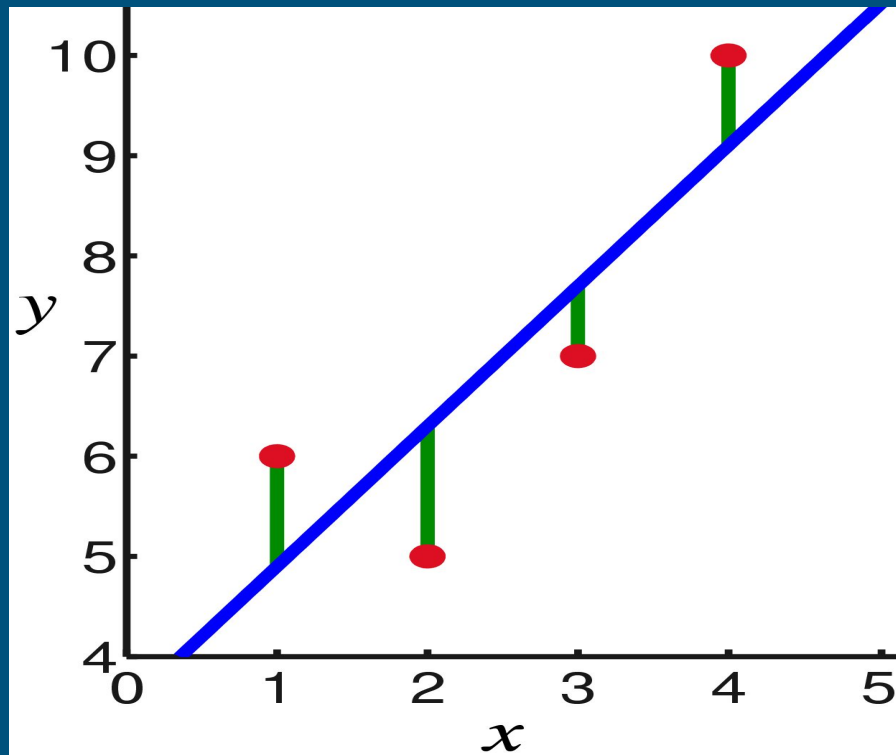
Colin Bennie, Caleb Hamblen, Sissi Shen, Justin Yang

# How to fit a Simple Linear Model

- Goal is to estimate:

  yi = Bo + B1xi + εi

- Want to minimize (yi - yi_hat)

# Simple Linear Model Assumptions

- $e_i$ random so that the following are true:

  - $E(\varepsilon_i) = 0$

  - $var(\varepsilon_i) = \sigma^2$

  - $Cov(\varepsilon_i, \varepsilon_j) = 0$

  - $E_i \sim N(0, \sigma^2)$

- Y needs to be linear to Bo and B1

# Results

- $E(y_i) = B_o + B_1 x_i$ : ***this is the regression function***

- $var(y_i) = \boldsymbol{\sigma}\verb|^|2$

- $Y_i|x_i = x \sim N(B_o + B_1 x_i, \boldsymbol{\sigma}\verb|^|2)$

- $B_o\_hat = y\_bar - B_1\_hat(x\_bar)$

- $B_1\_hat = \sum(x_i - x\_bar)\sum(y_i - y\_bar) / \sum(x_i - x\_bar)\verb|^|2$

- $Y_i\_hat = B_o\_hat + B_1\_hat(x_i)$

# How to implement in Python: Sample Data

```
In [1]:  1  import numpy as np
         2  import pandas as pd
         3  import matplotlib.pyplot as plt
         4  import statsmodels.formula.api as smf
```

```
In [2]:  1  df = pd.read_csv('Automobile_data.csv')
         2  df.head()
```

Out[2]:

| | symboling | normalized-losses | make | fuel-type | aspiration | num-of-doors | body-style | drive-wheels | engine-location | wheel-base | ... | engine-size | fuel-system | bore | stroke | compression-ratio | horsepower |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 3 | ? | alfa-romero | gas | std | two | convertible | rwd | front | 88.6 | ... | 130 | mpfi | 3.47 | 2.68 | 9.0 | 111 |
| 1 | 3 | ? | alfa-romero | gas | std | two | convertible | rwd | front | 88.6 | ... | 130 | mpfi | 3.47 | 2.68 | 9.0 | 111 |
| 2 | 1 | ? | alfa-romero | gas | std | two | hatchback | rwd | front | 94.5 | ... | 152 | mpfi | 2.68 | 3.47 | 9.0 | 154 |
| 3 | 2 | 164 | audi | gas | std | four | sedan | fwd | front | 99.8 | ... | 109 | mpfi | 3.19 | 3.4 | 10.0 | 102 |
| 4 | 2 | 164 | audi | gas | std | four | sedan | 4wd | front | 99.4 | ... | 136 | mpfi | 3.19 | 3.4 | 8.0 | 115 |

5 rows × 26 columns

```
In [3]:  1  df.columns
```

```
Out[3]:  Index(['symboling', 'normalized-losses', 'make', 'fuel-type', 'aspiration',
                'num-of-doors', 'body-style', 'drive-wheels', 'engine-location',
                'wheel-base', 'length', 'width', 'height', 'curb-weight', 'engine-type',
                'num-of-cylinders', 'engine-size', 'fuel-system', 'bore', 'stroke',
                'compression-ratio', 'horsepower', 'peak-rpm', 'city-mpg',
                'highway-mpg', 'price'],
               dtype='object')
```

# How to implement in Python: Model Fitting

```
In [8]:  1  import statsmodels.formula.api as smf
         2
         3  model = smf.ols("price ~ horsepower", data = dfnew).fit()
         4  b = model.params['Intercept']
         5  m = model.params['horsepower']
         6  print(f'intercept={b}, slope={m}')
```
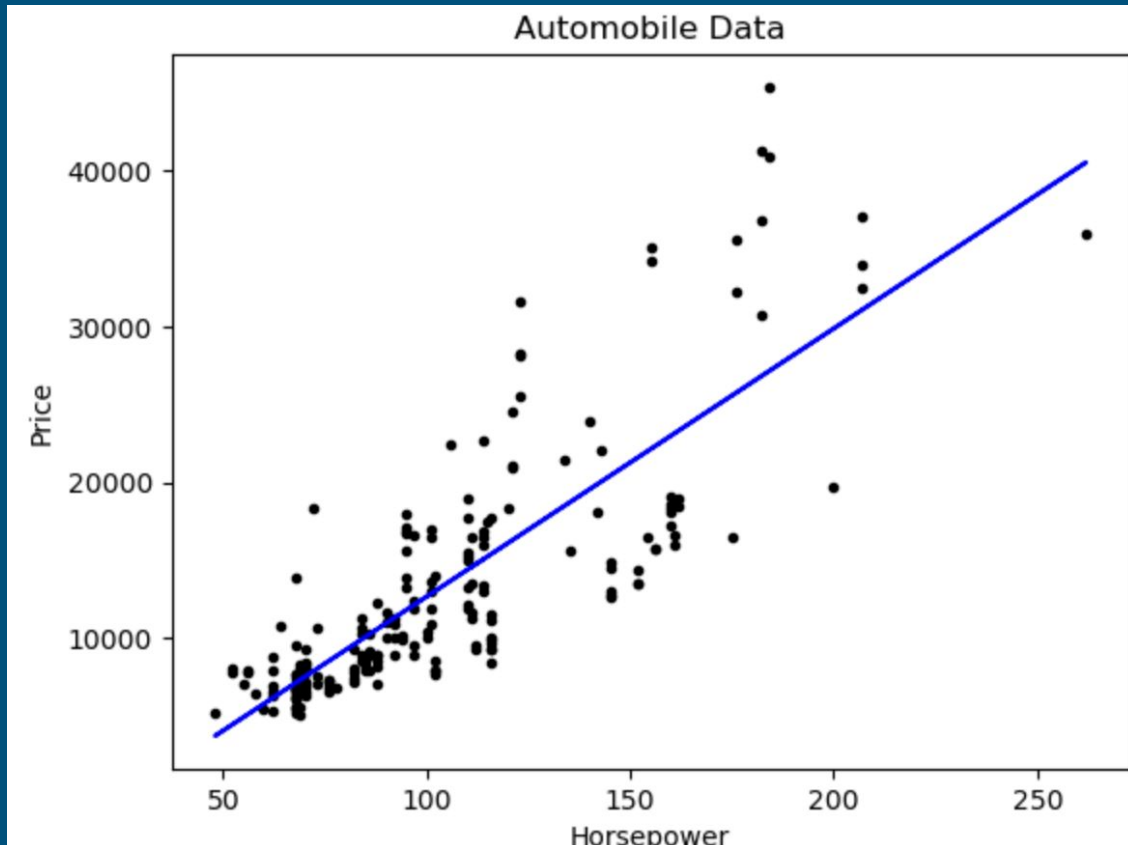
intercept=-4562.174995667492, slope=172.20625117310604
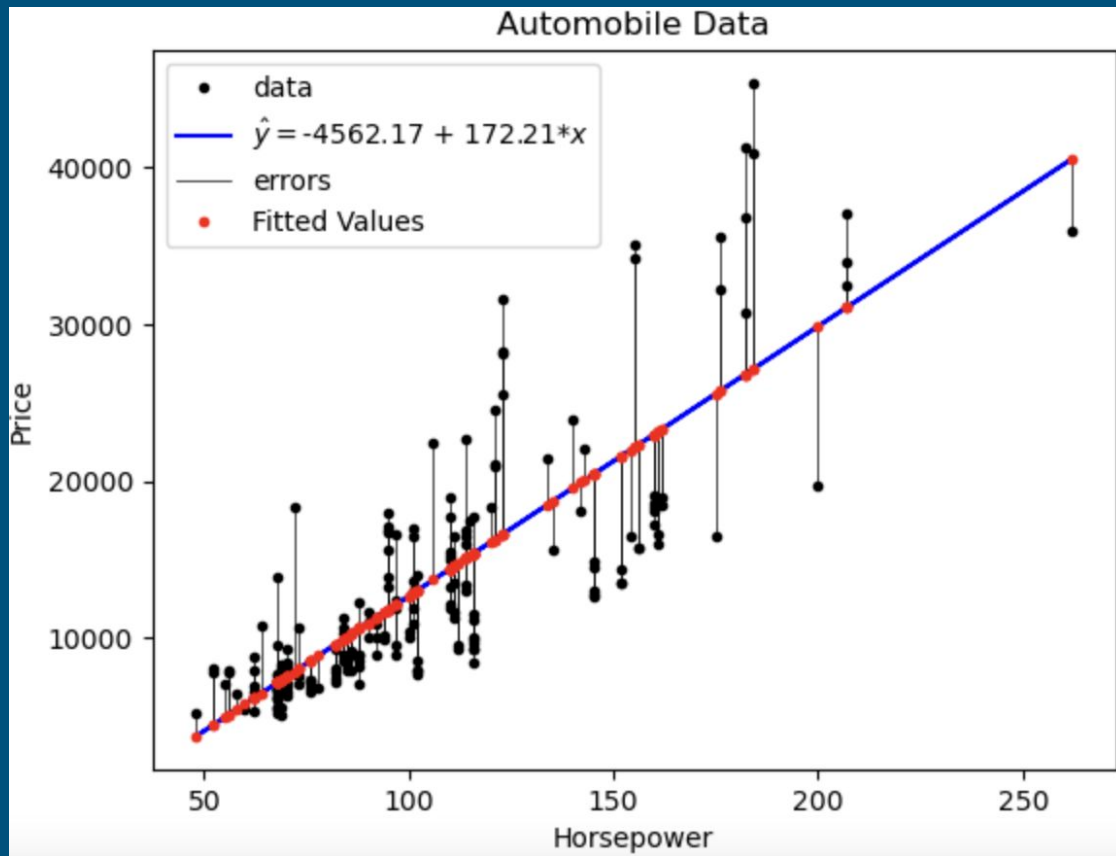
```
In [7]:  1  from sklearn import linear_model
         2
         3  model = linear_model.LinearRegression()
         4  model.fit(dfnew[['horsepower']], dfnew['price'])
         5  b = model.intercept_
         6  w = model.coef_[0]
         7  print(f'intercept={b}, slope={w}')
```

intercept=-4562.1749956674885, slope=172.20625117310607

# How the model looks

# How the model looks: residuals

# Predicting using Model

```python
# Sklearn
X = dfnew[['horsepower']]
y = dfnew['price']
model = linear_model.LinearRegression()
model.fit(X, y)

XX = [[1000]] # Horsepower of a Formula 1 Car
model.predict(XX)[0]
```

Out[11]: 167644.0761774386

```python
# Statsmodels
model = smf.ols("price ~ horsepower", data = dfnew).fit()

XX = np.array([1000])# Horsepower of a Formula 1 Car
prediction = model.params[0] + model.params[1] * XX
prediction[0]
```

Out[12]: 167644.07617743855

# Model Diagnosis



```
In [14]:    1  model.summary()
```

Out[14]:

OLS Regression Results

| | | | |
|---|---|---|---|
| Dep. Variable: | price | R-squared: | 0.657 |
| Model: | OLS | Adj. R-squared: | 0.655 |
| Method: | Least Squares | F-statistic: | 377.3 |
| Date: | Fri, 29 Sep 2023 | Prob (F-statistic): | 1.19e-47 |
| Time: | 16:40:42 | Log-Likelihood: | -1963.3 |
| No. Observations: | 199 | AIC: | 3931. |
| Df Residuals: | 197 | BIC: | 3937. |
| Df Model: | 1 | | |
| Covariance Type: | nonrobust | | |

| | coef | std err | t | P>|t| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| Intercept | -4562.1750 | 974.995 | -4.679 | 0.000 | -6484.943 | -2639.407 |
| horsepower | 172.2063 | 8.866 | 19.424 | 0.000 | 154.722 | 189.690 |

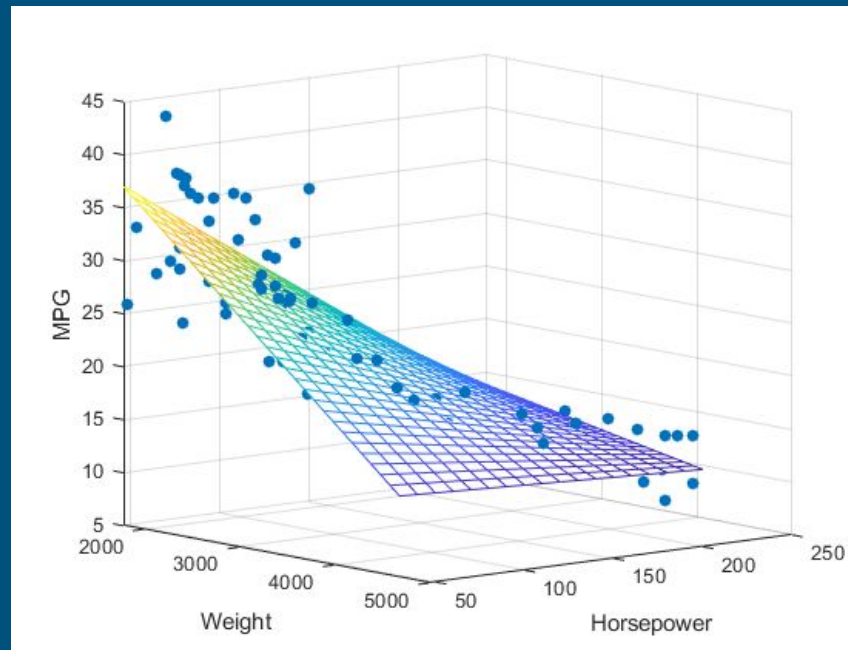| | | | |
|---|---|---|---|
| Omnibus: | 38.494 | Durbin-Watson: | 0.749 |
| Prob(Omnibus): | 0.000 | Jarque-Bera (JB): | 64.496 |
| Skew: | 1.013 | Prob(JB): | 9.89e-15 |
| Kurtosis: | 4.916 | Cond. No. | 323. |

```
In [15]:    1  import statsmodels.api as sm
            2
            3  aov_table = sm.stats.anova_lm(model, typ=1)
            4  aov_table
```

Out[15]:

| | df | sum_sq | mean_sq | F | PR(>F) |
|---|---|---|---|---|---|
| horsepower | 1.0 | 8.280790e+09 | 8.280790e+09 | 377.283543 | 1.189128e-47 |
| Residual | 197.0 | 4.323845e+09 | 2.194845e+07 | NaN | NaN |

# How to fit a Multiple Linear Regression

- Uses two or more independent variables

- Fitting a plane instead of a line

- $Y = b_0 + b_1X_1 + b_2X_2 + ... + b_nX_n + \varepsilon$

- Still want to minimize (Yi - Yi_hat)

# Multiple Linear Model Assumptions

- **Linearity**: The relationship between the dependent and independent variables is linear.

- **Independence and Normality of Errors**: The error terms ($\varepsilon$) are independent of each other and follow a normal distribution.

- **Homoscedasticity**: The variance of the errors is constant across all levels of independent variables.

- **No Multicollinearity**: The independent variables are not highly correlated with each other.

# Results

- Regression Function: $E(Y_i) = b_o + b_1 X_{i1} + b_2 X_{i2} + \ldots + b_n * X_{in}$

- Coefficient Estimates:

  - $Bo\_hat = y\_bar - B1\_hat(x\_bar1) - B2\_hat(x\_bar2) - \ldots - Bn\_hat(x\_barn)$

  - $Bn\_hat = \sum(x_{i\_n} - x\_bar\_n)\sum(y_i - y\_bar) / \sum(x_{i\_n} - x\_bar\_n)^2$

  - $B1\_hat = \sum(x_{i1} - x\_bar1)\sum(y_i - y\_bar) / \sum(x_{i1} - x\_bar1)^2$

- Fitted Values using estimated coefficients:

  - $Y_i\_hat = Bo\_hat + B1\_hat X1_i + B2\_hat X2_i + \ldots + Bn\_hat * Xn_i$

# Multiple Linear Regression: Python

```python
1  # Load in data for regression analysis
2  mlr_example = pd.read_csv('Automobile_data.csv')
3
4  # Identify target variable and potential predictors
5  print(mlr_example.columns)
```
✓ 0.0s                                                                    Python

```
Index(['symboling', 'normalized-losses', 'make', 'fuel-type', 'aspiration',
       'num-of-doors', 'body-style', 'drive-wheels', 'engine-location',
       'wheel-base', 'length', 'width', 'height', 'curb-weight', 'engine-type',
       'num-of-cylinders', 'engine-size', 'fuel-system', 'bore', 'stroke',
       'compression-ratio', 'horsepower', 'peak-rpm', 'city-mpg',
       'highway-mpg', 'price'],
      dtype='object')
```

```python
1  # Create regression model using smf.ols
2  model = smf.ols('price ~ engine_size + horsepower + highway_mpg + fuel_type', data=model_data).fit()
3  model.summary()
```
✓ 0.0s                                                                    Python

# MLR: Model Diagnosis

- $R^2$ and Adj-$R^2$
- Predictors all have their own coefficients and t test results

| | | OLS Regression Results | | | |
|---|---|---|---|---|---|
| Dep. Variable: | price | R-squared: | 0.812 | | |
| Model: | OLS | Adj. R-squared: | 0.808 | | |
| Method: | Least Squares | F-statistic: | 209.7 | | |
| Date: | Sat, 30 Sep 2023 | Prob (F-statistic): | 2.84e-69 | | |
| Time: | 15:26:07 | Log-Likelihood: | -1903.4 | | |
| No. Observations: | 199 | AIC: | 3817. | | |
| Df Residuals: | 194 | BIC: | 3833. | | |
| Df Model: | 4 | | | | |
| Covariance Type: | nonrobust | | | | |

| | coef | std err | t | P>\|t\| | [0.025 | 0.975] |
|---|---|---|---|---|---|---|
| Intercept | 3071.2513 | 3175.980 | 0.967 | 0.335 | -3192.632 | 9335.134 |
| fuel_type[T.gas] | -3882.1745 | 910.199 | -4.265 | 0.000 | -5677.331 | -2087.018 |
| engine_size | 102.6287 | 11.343 | 9.048 | 0.000 | 80.258 | 124.999 |
| horsepower | 58.0126 | 14.923 | 3.888 | 0.000 | 28.581 | 87.444 |
| highway_mpg | -174.3590 | 61.706 | -2.826 | 0.005 | -296.059 | -52.659 |

| | | | | | |
|---|---|---|---|---|---|
| Omnibus: | 11.078 | Durbin-Watson: | 0.841 | | |
| Prob(Omnibus): | 0.004 | Jarque-Bera (JB): | 20.221 | | |
| Skew: | 0.243 | Prob(JB): | 4.07e-05 | | |
| Kurtosis: | 4.484 | Cond. No. | 2.25e+03 | | |

# MLR: Model Diagnosis

- Sequential vs. Partial F Tests

```
1  # Sequential F Test on regression model
2
3  sequential_aov_table = sm.stats.anova_lm(model, typ=1)
4  sequential_aov_table
✓ 0.0s
```

|            | df    | sum_sq        | mean_sq       | F          | PR(>F)       |
|------------|-------|---------------|---------------|------------|--------------|
| fuel_type  | 1.0   | 1.496960e+08  | 1.496960e+08  | 12.267655  | 5.720085e-04 |
| engine_size| 1.0   | 9.503851e+09  | 9.503851e+09  | 778.844709 | 7.616789e-70 |
| horsepower | 1.0   | 4.863750e+08  | 4.863750e+08  | 39.858640  | 1.818820e-09 |
| highway_mpg| 1.0   | 9.742880e+07  | 9.742880e+07  | 7.984332   | 5.212118e-03 |
| Residual   | 194.0 | 2.367285e+09  | 1.220250e+07  | NaN        | NaN          |

Reduced Model: *price ~ fuel_type + engine_size*

Full Model: *price ~ fuel_type + engine_size + horsepower*

```
1  # Partial F Test on regression model
2
3  partial_aov_table = sm.stats.anova_lm(model, typ=2)
4  partial_aov_table
[43]  ✓ 0.0s
```

|            | sum_sq        | df    | F          | PR(>F)       |
|------------|---------------|-------|------------|--------------|
| fuel_type  | 2.219861e+08  | 1.0   | 18.191857  | 3.118291e-05 |
| engine_size| 9.989917e+08  | 1.0   | 81.867806  | 1.530876e-16 |
| horsepower | 1.844180e+08  | 1.0   | 15.113133  | 1.389653e-04 |
| highway_mpg| 9.742880e+07  | 1.0   | 7.984332   | 5.212118e-03 |
| Residual   | 2.367285e+09  | 194.0 | NaN        | NaN          |

Reduced Model: *price ~ fuel_type + engine_size + highway_mpg*

Full Model: *price ~ fuel_type + engine_size + highway_mpg + horsepower*

# MLR: Multicollinearity and VIF Check

```python
1   # additional imports needed to check for multicollinearity
2   from statsmodels.stats.outliers_influence import variance_inflation_factor
3   from patsy import dmatrices
4
5   # Calculate the VIF for each predictor variable included in model
6   vif_y, vif_x = dmatrices('price ~ engine_size + horsepower + highway_mpg + fuel_type',
7                            data=model_data, return_type='dataframe')
8
9   vif = pd.DataFrame()
10  vif['VIF Factor'] = [variance_inflation_factor(vif_x.values, i) for i in range(vif_x.shape[1])]
11  vif['Features'] = vif_x.columns
12  print(vif)
```
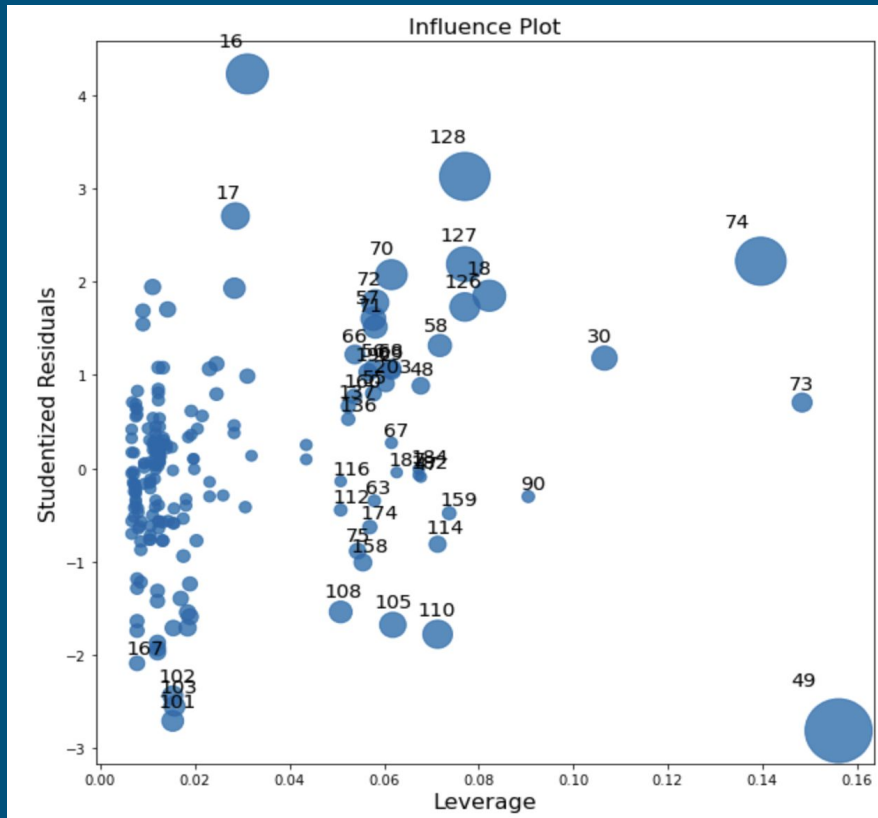
✓ 0.0s

```
   VIF Factor          Features
0  164.497754         Intercept
1    1.221390  fuel_type[T.gas]
2    3.639274       engine_size
3    5.095827        horsepower
4    2.898497       highway_mpg
```

# Link to Github with Code

https://github.com/cwbennie/comms_code_demo23

# MLR: Influential Points

- Leverage
- Standardized Residuals
- Cook's Distance

# How the model looks: Code to plot

```python
In [9]:     1  import matplotlib.pyplot as plt
            2
            3  X = dfnew[['horsepower']]
            4  y = dfnew['price']
            5  model = linear_model.LinearRegression()
            6  model.fit(X, y)
            7  b = model.intercept_
            8  w = model.coef_[0]
            9
           10  x = dfnew.horsepower # we need a 1D array for plotting (and a 2D array for .fit() above)
           11  plt.plot(x, y, '.', color='black', label='data')
           12  plt.title('Automobile Data')
           13  plt.xlabel('Horsepower')
           14  plt.ylabel('Price')
           15
           16  y_hat = model.predict(X) # equivalent to y_hat = w * X[:, 0] + b
           17  plt.plot(x, y_hat, color='red',
           18          label=f'$\\hat{{y}}=${round(b, 2)} + ({round(w, 2)})$x$')
           19
           20  plt.show()
```