

OOP Concept in Machine Learning 2

Simple Linear Regression

- [Gradient Descent Python Implementation](#)
- [Scipy Implementation](#)
- [Scikit-Learn Implementation](#)
- [Statsmodel Implementation](#)
- [Multi-step visual of Gradient Descent](#)
- [Animating the Gradient Descent](#)

Importing needed libraries

```
[ ] 1 import numpy as np
     2 import pandas as pd
     3 import matplotlib.pyplot as plt
     4 import seaborn as sns
     5 %matplotlib inline
     6
     7 from google.colab import files
     8 uploaded = files.upload()
```

Loading our housing dataset

We will load our data on house sales in King County to predict house prices using simple (one input) linear regression

```
[ ] 1 dataset = pd.read_csv('kc_house_data.csv')
```

We want to be able to predict y which is our price variable.

```
[ ] 1 Y = dataset[['price']]
```

```
[ ] 1 X = dataset.drop(['price', 'id', 'date'], axis=1)
```



```
1 X.info()
```



```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 21613 entries, 0 to 21612
Data columns (total 18 columns):
 #   Column              Non-Null Count  Dtype
---  -
 0   bedrooms            21613 non-null  int64
 1   bathrooms            21613 non-null  float64
 2   sqft_living          21613 non-null  int64
 3   sqft_lot             21613 non-null  int64
 4   floors              21613 non-null  float64
 5   waterfront           21613 non-null  int64
 6   view                21613 non-null  int64
 7   condition            21613 non-null  int64
 8   grade               21613 non-null  int64
 9   sqft_above           21613 non-null  int64
10  sqft_basement        21613 non-null  int64
11  yr_built+            21613 non-null  int64
```

```
[ ] 1 #list our columns
    2 columns = X.columns
    3 columns
```

```
Index(['bedrooms', 'bathrooms', 'sqft_living', 'sqft_lot', 'floors',
      'waterfront', 'view', 'condition', 'grade', 'sqft_above',
      'sqft_basement', 'yr_built', 'yr_renovated', 'zipcode', 'lat', 'long',
      'sqft_living15', 'sqft_lot15'],
      dtype='object')
```

```
[ ] 1 #show first 5 records
    2 X.head()
```

	bedrooms	bathrooms	sqft_living	sqft_lot	floors	waterfront	view	condition	grade	sqft_above	sqft_basement	yr_built
0	3	1.00	1180	5650	1.0	0	0	3	7	1180	0	1955
1	3	2.25	2570	7242	2.0	0	0	3	7	2170	400	1951
2	2	1.00	770	10000	1.0	0	0	3	6	770	0	1933

```
1 X.describe()
```

	bedrooms	bathrooms	sqft_living	sqft_lot	floors	waterfront	view
count	21613.000000	21613.000000	21613.000000	2.161300e+04	21613.000000	21613.000000	21613.000000
mean	3.370842	2.114757	2079.899736	1.510697e+04	1.494309	0.007542	0.234303
std	0.930062	0.770163	918.440897	4.142051e+04	0.539989	0.086517	0.766318
min	0.000000	0.000000	290.000000	5.200000e+02	1.000000	0.000000	0.000000
25%	3.000000	1.750000	1427.000000	5.040000e+03	1.000000	0.000000	0.000000
50%	3.000000	2.250000	1910.000000	7.618000e+03	1.500000	0.000000	0.000000
75%	4.000000	2.500000	2550.000000	1.068800e+04	2.000000	0.000000	0.000000
max	33.000000	8.000000	13540.000000	1.651359e+06	3.500000	1.000000	4.000000

```
1 dataset = dataset.drop(['id', 'date'], axis=1)
```

```
1 dataset.corr(method='pearson')
```

	price	bedrooms	bathrooms	sqft_living	sqft_lot	floors	waterfront	view
price	1.000000	0.308350	0.525138	0.702035	0.089661	0.256794	0.266369	0.397293
bedrooms	0.308350	1.000000	0.515884	0.576671	0.031703	0.175429	-0.006582	0.079532
bathrooms	0.525138	0.515884	1.000000	0.754665	0.087740	0.500653	0.063744	0.187737
sqft_living	0.702035	0.576671	0.754665	1.000000	0.172826	0.353949	0.103818	0.284611
sqft_lot	0.089661	0.031703	0.087740	0.172826	1.000000	-0.005201	0.021604	0.074710
floors	0.256794	0.175429	0.500653	0.353949	-0.005201	1.000000	0.023698	0.029444
waterfront	0.266369	-0.006582	0.063744	0.103818	0.021604	0.023698	1.000000	0.401857
view	0.397293	0.079532	0.187737	0.284611	0.074710	0.029444	0.401857	1.000000
condition	0.036362	0.028472	-0.124982	-0.058753	-0.008958	-0.263768	0.016653	0.045990

```
1 plt.subplots(figsize=(10,8))
2 sns.heatmap(dataset.corr())
```

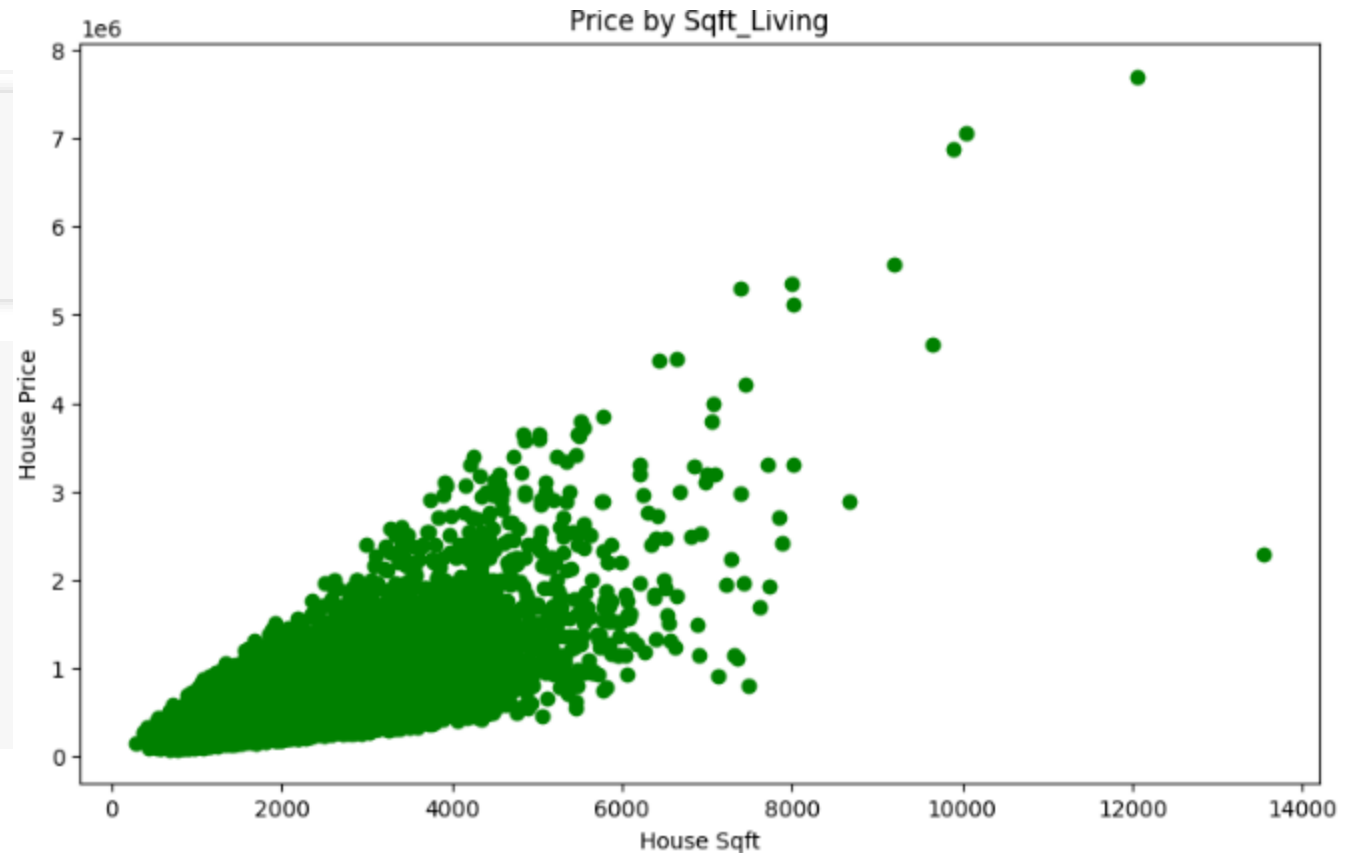


It is **Simple Linear Regression** when we have one dependent variable (feature) and one independent variable. Here we will pick `sqft_living` as our independent variable `x`.

Our goal is to estimate $\hat{y} = x\theta_1 + \theta_0$, where θ_1 is our coefficient and θ_0 is our `y` intercept. To estimate \hat{y} we need to find a function such as $\hat{y} = h(x) = x\theta_1 + \theta_0$

```
1 x = X[['sqft_living']]
2 y = Y
```

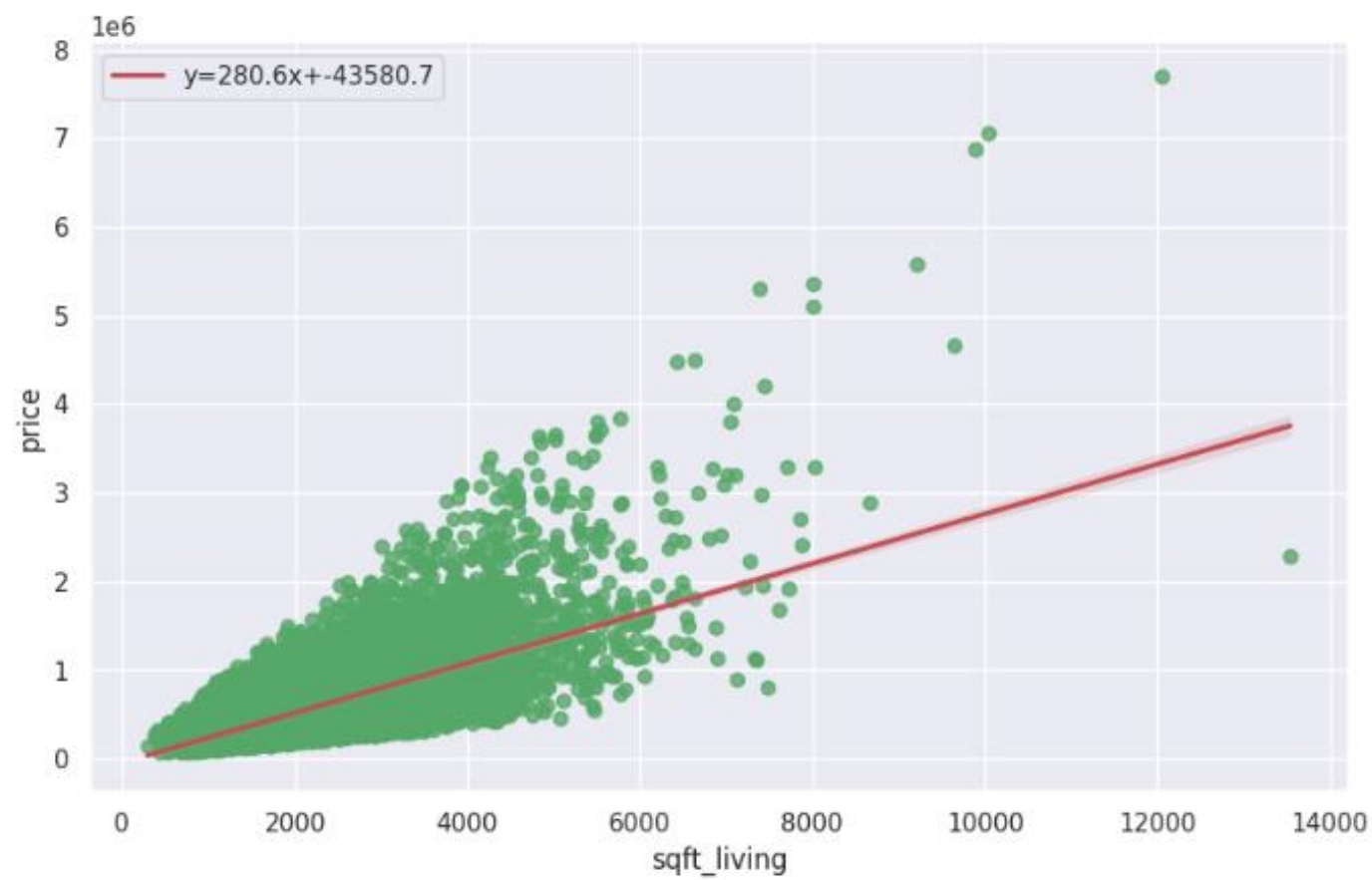
```
[ ] 1 plt.figure(figsize=(10,6))
    2 plt.xlabel('House Sqft')
    3 plt.ylabel('House Price')
    4 plt.title('Price by Sqft_Living')
    5 plt.scatter(x,y, marker='o', color='g')
```



Simple Linear Regression Implementations:

1. Using `seaborn.regplot()` and `scipy.stats`

```
[ ] 1 from scipy import stats
    2 sns.set(color_codes=True)
    3
    4 slope, intercept, r_value, p_value, std_err = stats.linregress(dataset['sqft_living'], dataset['price'])
    5
    6 f = plt.figure(figsize=(10,6))
    7 data = dataset[['price', 'sqft_living']]
    8 ax = sns.regplot(x='sqft_living', y='price', data=data,
    9                 scatter_kws={"color": "g"},
   10                 line_kws={'color': 'r', 'label': "y={0:.1f}x+{1:.1f}".format(slope, intercept)})
   11 ax.legend()
```



```
1 print(slope, intercept)
```

```
280.6235678974483 -43580.74309447408
```

```
1 print(std_err)
```

```
1.9363985519989133
```

2. Manual Method : Gradient Descent Implementation

[Top](#)

Equations

Objective of Linear Regression is to minimize the cost function:

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

where the hypothesis $h_{\theta}(x)$ is given by the linear model:

$$h_{\theta}(x) = \theta^T X = \theta_1 X_1 + \theta_0$$

In batch gradient descent, each iteration performs the update:

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

```
[ ] 1 x = X[['sqft_living']]  
    2 y = Y
```

```
[▶] 1 xg = x.values.reshape(-1,1)  
    2 yg = y.values.reshape(-1,1)  
    3 xg = np.concatenate((np.ones(len(x)).reshape(-1,1), x), axis=1)
```

Implementing the Cost Function $J(\theta)$ in python

```
[▶] 1 def computeCost(x, y, theta):  
    2     m = len(y)  
    3     h_x = x.dot(theta)  
    4     j = np.sum(np.square(h_x - y))*(1/(2*m))  
    5     return j
```

```

▶ 1 def gradientDescent(x, y, theta, alpha, iteration):
  2     print('Running Gradient Descent...')
  3     j_hist = []
  4     m = len(y)
  5     for i in range(iteration):
  6         j_hist.append(computeCost(x, y, theta))
  7         h_x = x.dot(theta)
  8         theta = theta - ((alpha/m) * ((np.dot(x.T, (h_x-y) ))))
  9         #theta[0] = theta[0] - ((alpha/m) * (np.sum((h_x-y))))
10     return theta, j_hist

```

```

[ ] 1 theta = np.zeros((2,1))
    2 iteration = 2000
    3 alpha = 0.001
    4
    5 theta, cost = gradientDescent(xg, yg, theta, alpha, iteration)
    6 print('Theta found by Gradient Descent: slope = {} and intercept {}'.format(theta[1], theta[0]))

```

Running Gradient Descent...

<ipython-input-170-4c6aec17be90>:4: RuntimeWarning: overflow encountered in square

```
j = np.sum(np.square(h_x - y))*(1/(2*m))
```

<ipython-input-171-e81683f4f12d>:8: RuntimeWarning: invalid value encountered in subtract

```
theta = theta - ((alpha/m) * ((np.dot(x.T, (h_x-y) ))))
```

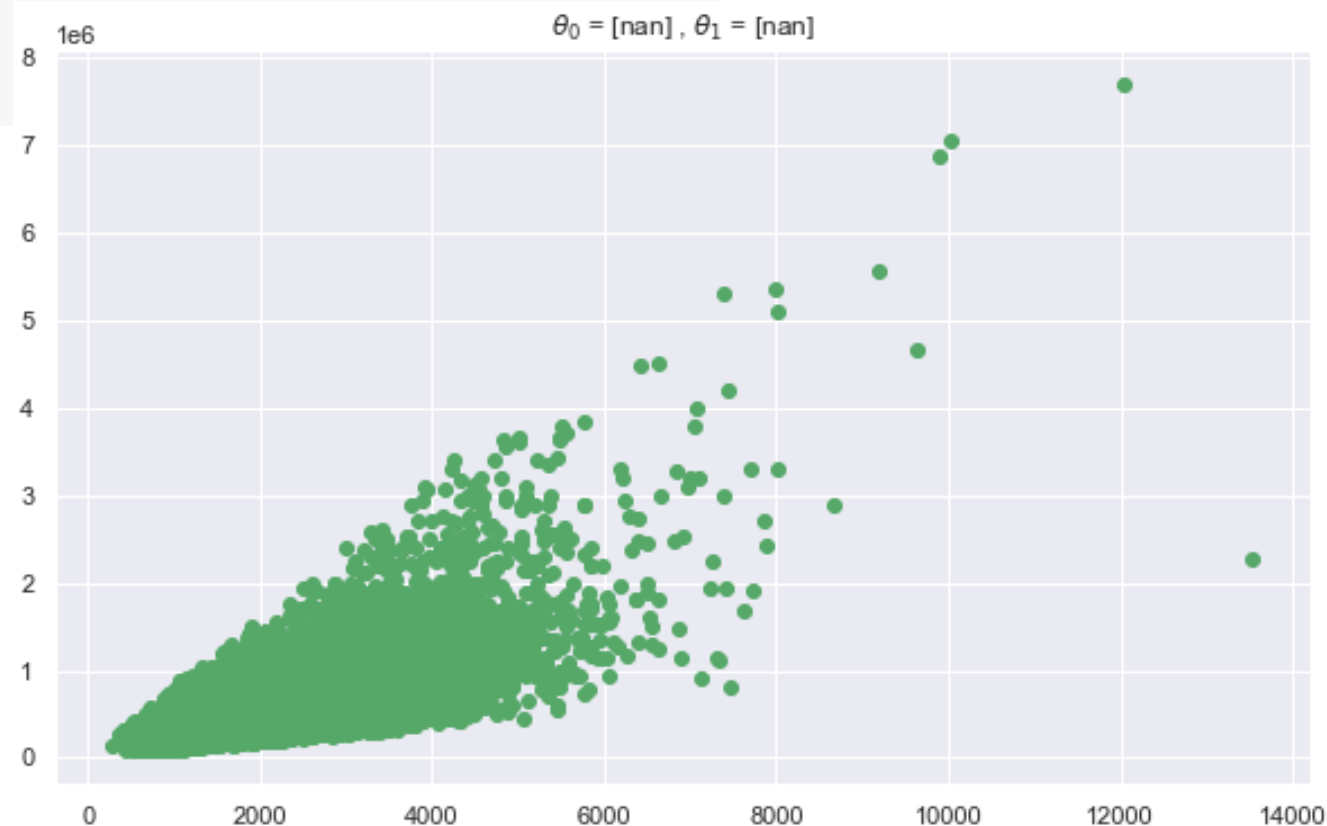
Theta found by Gradient Descent: slope = [nan] and intercept [nan]

Plotting the linear fit

```
[ ] 1 theta.shape
```

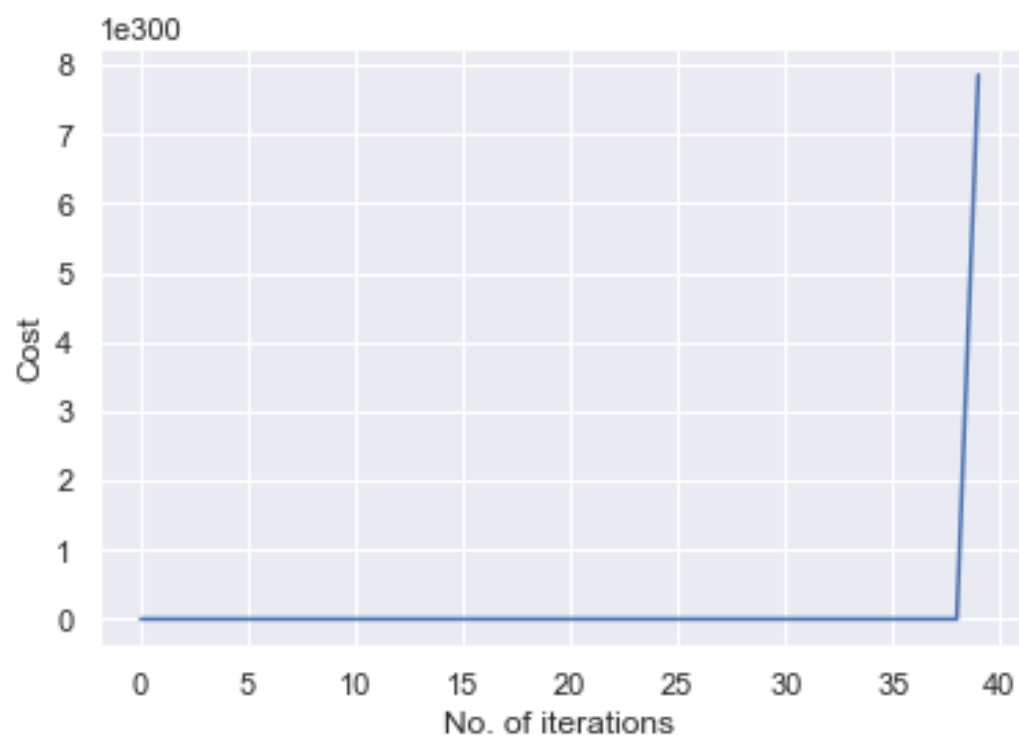
```
(2, 1)
```

```
[ ] 1 plt.figure(figsize=(10,6))  
2 plt.title('$\\theta_0$ = {} , $\\theta_1$ = {}'.format(theta[0], theta[1]))  
3 plt.scatter(x,y, marker='o', color='g')  
4 plt.plot(x,np.dot(x.values, theta.T))  
5 plt.show()
```



```
[ ] 1 plt.plot(cost)
    2 plt.xlabel('No. of iterations')
    3 plt.ylabel('Cost')
```

Text(0, 0.5, 'Cost')



4. Implement with using Scipy

```
[ ] 1 from scipy import stats
    2
    3 xs = x.iloc[:,0]
    4 ys = y.iloc[:,0]
    5 #xs = np.concatenate((np.ones(len(x)).reshape(-1,1), x), axis=1)
    6
    7 slope, intercept, r_value, p_value, std_err = stats.linregress(xs, ys)
```

```
[ ] 1 print('Slope = {} and Intercept = {}'.format(slope, intercept))
    2 print('y = x({}) + {}'.format(slope, intercept))
```

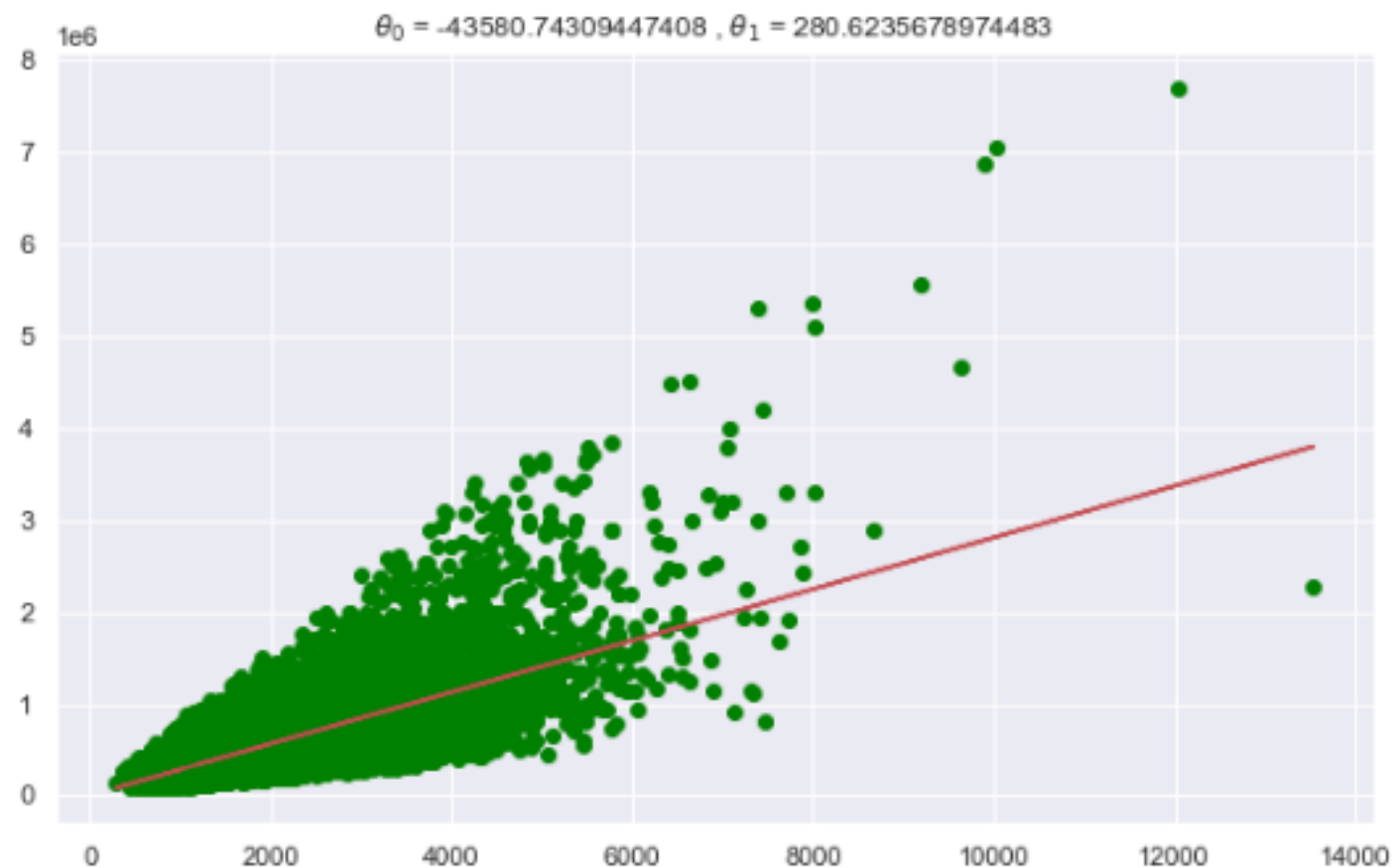
Slope = 280.6235678974483 and Intercept = -43580.74309447408

y = x(280.6235678974483) + -43580.74309447408

Plot the linear fit using the slop and intercept values from scipy

```
[ ] 1 plt.figure(figsize=(10,6))
    2 plt.title('$\\theta_0$ = {} , $\\theta_1$ = {}'.format(intercept, slope))
    3 plt.scatter(xs,y, marker='o', color='green')
    4 plt.plot(xs, np.dot(x, slope), 'r')
```

[<matplotlib.lines.Line2D at 0x7f9eb6b8af10>]



5. Implement using Scikit-Learn

```
[ ] 1 xsl = x.values.reshape(-1,1)
    2 ysl = y.values.reshape(-1,1)
    3 xsl = np.concatenate((np.ones(len(xsl)).reshape(-1,1), xsl), axis=1)
    4
    5 from sklearn.linear_model import LinearRegression
    6
    7 slr = LinearRegression()
    8 slr.fit(xsl[:,1].reshape(-1,1), ysl.reshape(-1,1))
    9 y_hat = slr.predict(xsl[:,1].reshape(-1,1))
   10
   11 print('theta[0] = ', slr.intercept_)
   12 print('theta[1] = ', slr.coef_)
   13
   14 thetas = np.array((slr.intercept_, slr.coef_)).squeeze()
```

```
theta[0] = [-43580.74309447]
theta[1] = [[280.6235679]]
```

```
[ ] 1 plt.figure(figsize=(10,6))
    2 plt.title('$\\theta_0$ = {} , $\\theta_1$ = {}'.format(thetas[0], thetas[1]))
    3 plt.scatter(xsl[:,1],y, marker='x', color='g')
    4 plt.plot(xsl[:,1], np.dot(xsl, thetas), 'r')
```

