# Task-C: Regression outlier effect.

Objective:Visualization best fit linear regression line for different scenarios

```
1   # you should not import any other packages
2   import matplotlib.pyplot as plt
3   import warnings
4   warnings.filterwarnings("ignore")
5   import numpy as np
6   from sklearn.linear_model import SGDRegressor
```

```
1   #Mounting Google drive folder
2   from google.colab import drive
3   drive.mount('/content/drive')
4   %cd /content/drive/My Drive/Appliedai colab/Assignment 8 - Linear models/
```

```
1    import numpy as np
2    import scipy as sp
3    import scipy.optimize
4
5    def angles_in_ellipse(num,a,b):
6        assert(num > 0)
7        assert(a < b)
8        angles = 2 * np.pi * np.arange(num) / num
9        if a != b:
10           e = (1.0 - a ** 2.0 / b ** 2.0) ** 0.5
11           tot_size = sp.special.ellipeinc(2.0 * np.pi, e)
```

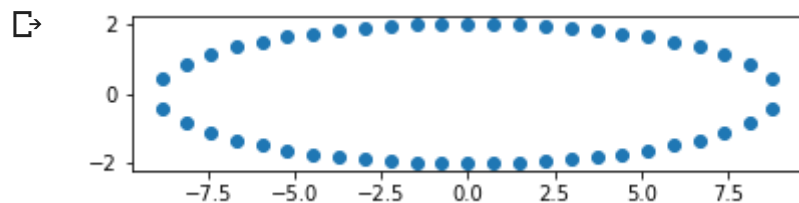```
12              arc_size = tot_size / num
13              arcs = np.arange(num) * arc_size
14              res = sp.optimize.root(
15                  lambda x: (sp.special.ellipeinc(x, e) - arcs), angles)
16              angles = res.x
17          return angles
```

```
1   a = 2
2   b = 9
3   n = 50
4
5   phi = angles_in_ellipse(n, a, b)
6   e = (1.0 - a ** 2.0 / b ** 2.0) ** 0.5
7   arcs = sp.special.ellipeinc(phi, e)
8
9   fig = plt.figure()
10  ax = fig.gca()
11  ax.axes.set_aspect('equal')
12  ax.scatter(b * np.sin(phi), a * np.cos(phi))
13  plt.show()
```



```
1   X = b * np.sin(phi)
2   Y = a * np.cos(phi)
```
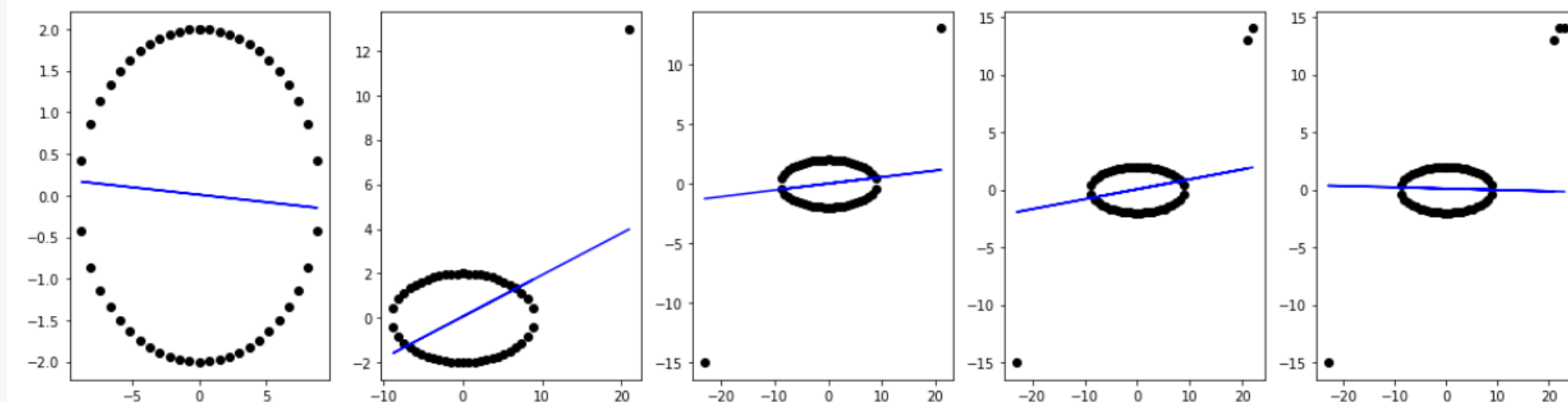
*1. As a part of this assignment you will be working the regression problem and how regularization helps to get rid of outliers*

*2. Use the above created X, Y for this experiment.*

*3. to do this task you can either implement your own SGDRegression(prefered excatly similar to "SGD assignment" with mean sequared error or you can use the SGDRegression of sklearn, for example "SGDRegressor(alpha=0.001, eta0=0.001, learning_rate='constant',random_state=0)" note that you have to use the constant learning rate and learning rate **eta0** initialized.*

*4. as a part of this experiment you will train your linear regression on the data (X, Y) with different regularizations alpha=[0.0001, 1, 100] and observe how prediction hyper plan moves with respect to the outliers*

*5. This the results of one of the experiment we did (title of the plot was not metioned intentionally)*



*in each iteration we were adding single outlier and observed the movement of the hyper plane.*

*6. please consider this list of outliers: [(0,2),(21, 13), (-23, -15), (22,14), (23, 14)] in each of tuple the first elemet is the input feature(X) and the second element is the output(Y)*

*7. for each regularizer, you need to add these outliers one at time to data and then train your model again on the updated data.*

8. you should plot a 3*5 grid of subplots, where each row corresponds to results of model with a single regularizer.

9. Algorithm:

for each regularizer:
  for each outlier:
    #add the outlier to the data
    #fit the linear regression to the updated data
    #get the hyper plane
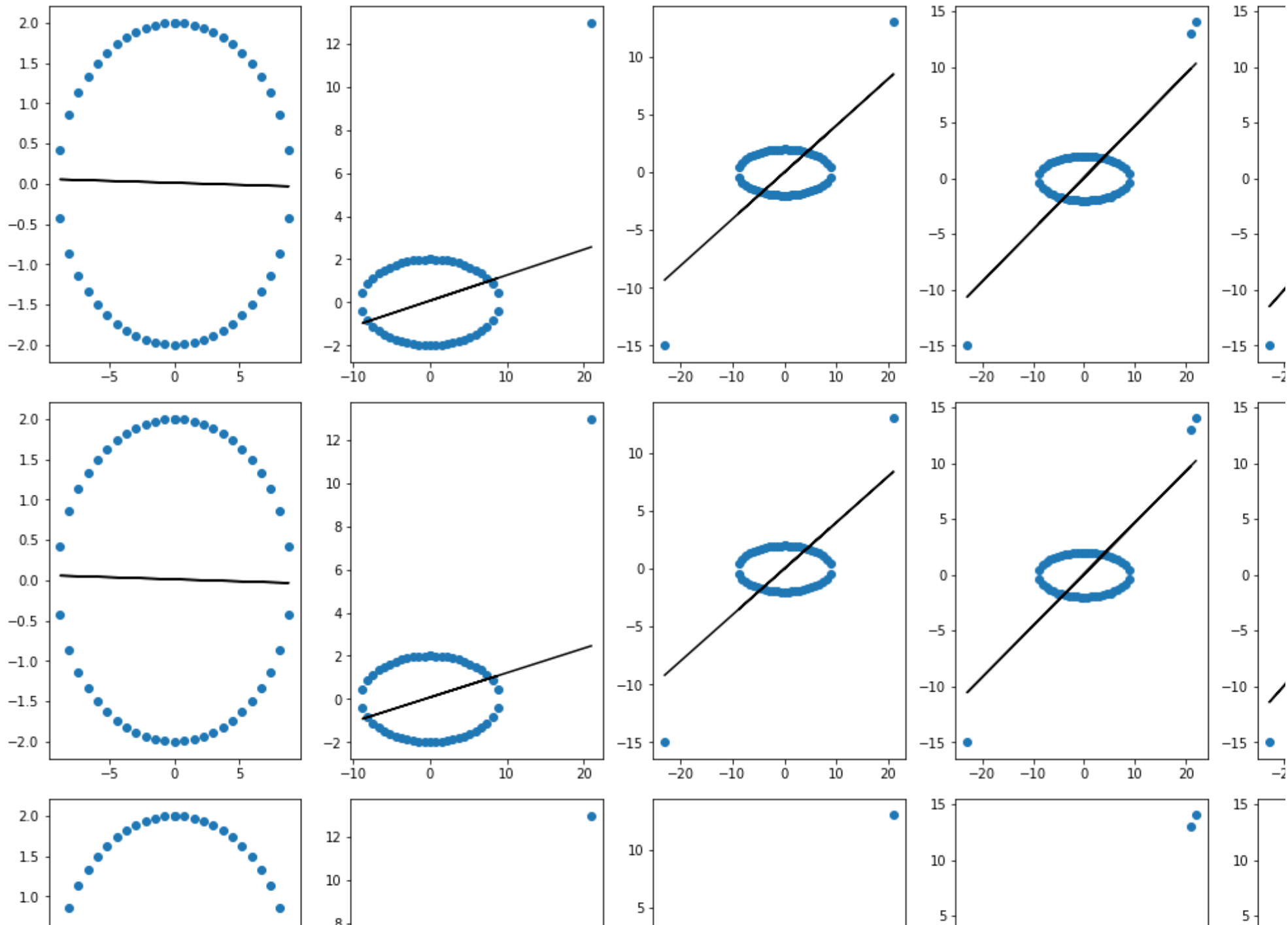    #plot the hyperplane along with the data points

10. MAKE SURE YOU WRITE THE DETAILED OBSERVATIONS, PLEASE CHECK THE LOSS FUNCTION IN THE SKLEARN DOCUMENTATION (please do search for it).
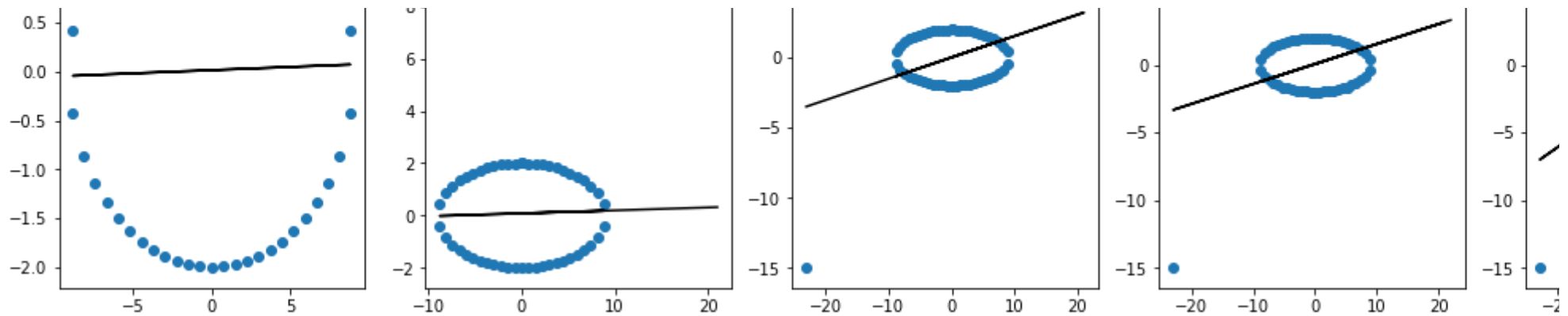
```
1   def plot_decision_boundary(coef, intercept, x):
2     y = coef * x + intercept
3     plt.plot(x, y, color = 'black')
```

```
1   reg_alphas = [0.0001, 1, 100]
2   outliers = [(0,2), (21, 13), (-23, -15), (22,14), (23, 14)]
3
4
5   for alpha in reg_alphas:  #for each regularizer
6     c = 0
7     plt.figure(figsize=(20,5))
8     X = b * np.sin(phi)
9     Y = a * np.cos(phi)
10    for o in outliers:  #for each outlier
11      plt.subplot(1, 5, c+1)
12
```

```
12        c += 1

13
14        #adding outlier to the data
15        X = np.append(X,o[0]).reshape(-1, 1)
16        Y = np.append(Y,o[1]).reshape(-1, 1)

17
18        plt.scatter(X, Y)

19
20        clf = SGDRegressor(alpha = alpha, eta0=0.001, learning_rate='constant', random_state = 0,\
21                           max_iter = 1000, tol = 1e-3)
22        clf.fit(X, Y)
23        coef = clf.coef_
24        intercept = clf.intercept_
25        plot_decision_boundary(coef, intercept, X)
26 plt.show();
```

In this assignment we increased the alpha for SGDRegressor to see how sensitive are the parameters to outliers in data.

Here Alpha is the constant that multiplies the regularization term.

The regularization term(penalty term) defaults to l2 (squared euclidean norm L2) for linear SVM models.

## ▾ Observations

- When outliers makes the absolute value of parameters to be very large, the L2 norm regularization will increase the loss function more and thus it will detect unusual large parameters earlier, which means it is more useful for reducing the sensitivity of parameters to outliers.
- L-2 norm error usually produces much great output value (due to square exponent) and therefore can be logically considered as much more sensitive to outlier data.
- For an even better understanding I read this blog - http://www.chioka.in/differences-between-l1-and-l2-as-loss-function-and-regularization/ and learned how L2 regularization acts as a more stable function as compared to L1.

1