

## Task-C: Regression outlier effect.

Objective: Visualization best fit linear regression line for different scenarios

```
In [13]: # you should not import any other packages
import matplotlib.pyplot as plt
import warnings
warnings.filterwarnings("ignore")
import numpy as np
from sklearn.linear_model import SGDRegressor
```

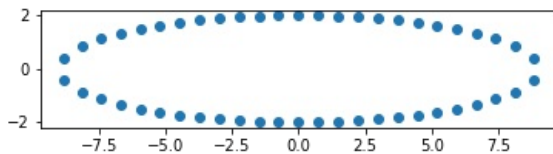
```
In [14]: import numpy as np
import scipy as sp
import scipy.optimize

def angles_in_ellipse(num,a,b):
    assert(num > 0)
    assert(a < b)
    angles = 2 * np.pi * np.arange(num) / num
    if a != b:
        e = (1.0 - a ** 2.0 / b ** 2.0) ** 0.5
        tot_size = sp.special.ellipeinc(2.0 * np.pi, e)
        arc_size = tot_size / num
        arcs = np.arange(num) * arc_size
        res = sp.optimize.root(
            lambda x: (sp.special.ellipeinc(x, e) - arcs), angles)
        angles = res.x
    return angles
```

```
In [15]: a = 2
b = 9
n = 50

phi = angles_in_ellipse(n, a, b)
e = (1.0 - a ** 2.0 / b ** 2.0) ** 0.5
arcs = sp.special.ellipeinc(phi, e)

fig = plt.figure()
ax = fig.gca()
ax.axes.set_aspect('equal')
ax.scatter(b * np.sin(phi), a * np.cos(phi))
plt.show()
```



```
In [16]: X = b * np.sin(phi)
Y = a * np.cos(phi)
```

```
In [18]: print(X)

[ 0.00000000e+00  7.44742410e-01  1.48935470e+00  2.23369584e+00
  2.97760060e+00  3.72086049e+00  4.46319267e+00  5.20418351e+00
  5.94317435e+00  6.67899598e+00  7.40922283e+00  8.12729576e+00
  8.80003723e+00  8.80003723e+00  8.12729576e+00  7.40922283e+00
  6.67899598e+00  5.94317435e+00  5.20418351e+00  4.46319267e+00
  3.72086049e+00  2.97760060e+00  2.23369584e+00  1.48935470e+00
  7.44742410e-01  1.10218212e-15 -7.44742410e-01 -1.48935470e+00
 -2.23369584e+00 -2.97760060e+00 -3.72086049e+00 -4.46319267e+00
 -5.20418351e+00 -5.94317435e+00 -6.67899598e+00 -7.40922283e+00
 -8.12729576e+00 -8.80003723e+00 -8.80003723e+00 -8.12729576e+00
 -7.40922283e+00 -6.67899598e+00 -5.94317435e+00 -5.20418351e+00
 -4.46319267e+00 -3.72086049e+00 -2.97760060e+00 -2.23369584e+00
 -1.48935470e+00 -7.44742410e-01]
```

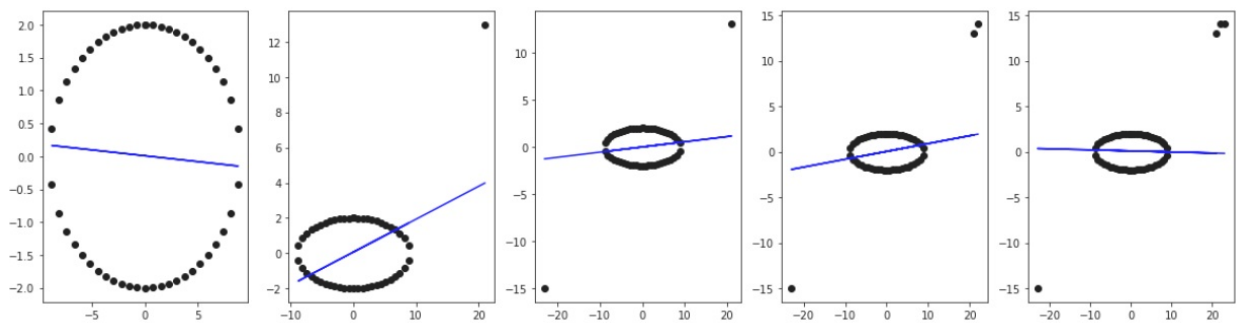
```
In [19]: print(Y)

[ 2.          1.99314081  1.972425   1.93742355  1.88737056  1.82107277
 1.73674751  1.63172973  1.50191119  1.34055475  1.13536674  0.85914295
 0.41924946 -0.41924946 -0.85914295 -1.13536674 -1.34055475 -1.50191119
 -1.63172973 -1.73674751 -1.82107277 -1.88737056 -1.93742355 -1.972425
 -1.99314081 -2.          -1.99314081 -1.972425   -1.93742355 -1.88737056
 -1.82107277 -1.73674751 -1.63172973 -1.50191119 -1.34055475 -1.13536674]
```

```
-0.85914295 -0.41924946 0.41924946 0.85914295 1.13536674 1.34055475
1.50191119 1.63172973 1.73674751 1.82107277 1.88737056 1.93742355
1.972425 1.99314081]
```

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`(preferred) exactly similar to "SGD assignment" with mean squared error or you can use the `SGDRegression` of `sklearn`, for example `"SGDRegressor(alpha=0.001, eta=0.001, learning_rate='constant', random_state=0)"`  
note that you have to use the constant learning rate and learning rate **eta** 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 plane moves with respect to the outliers

5. This the results of one of the experiment we did (title of the plot was not mentioned 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 element 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  $3 \times 5$  grid of subplots, where each row corresponds to results of model with a single regularizer.

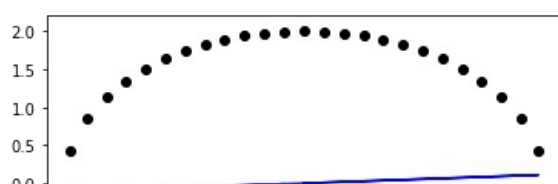
9. Algorithm:

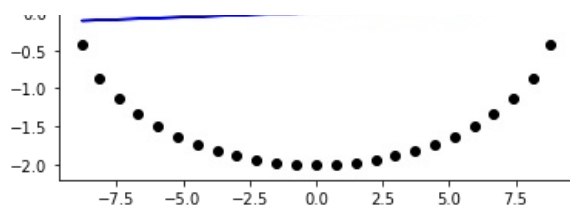
```
for each regularizer:
    for each outlier:
        #add the outlier to the data
        #fit the linear regression to the updated data
        #get the hyper plane
        #plot the hyperplane along with the data points
```

10. MAKE SURE YOU WRITE THE DETAILED OBSERVATIONS, PLEASE CHECK THE LOSS FUNCTION IN THE SKLEARN DOCUMENTATION (please do search for it).

```
In [39]: X = b * np.sin(phi)
Y = a * np.cos(phi)
model = SGDRegressor(alpha = alpha[i], eta=0.001, learning_rate='constant', random_state=0)
model.fit(X.reshape(-1, 1), Y.reshape(-1, 1))
y_pred = model.predict(X.reshape(-1,1))
plt.scatter(X,Y, color = 'black')
plt.plot(X,y_pred, color = 'blue')
```

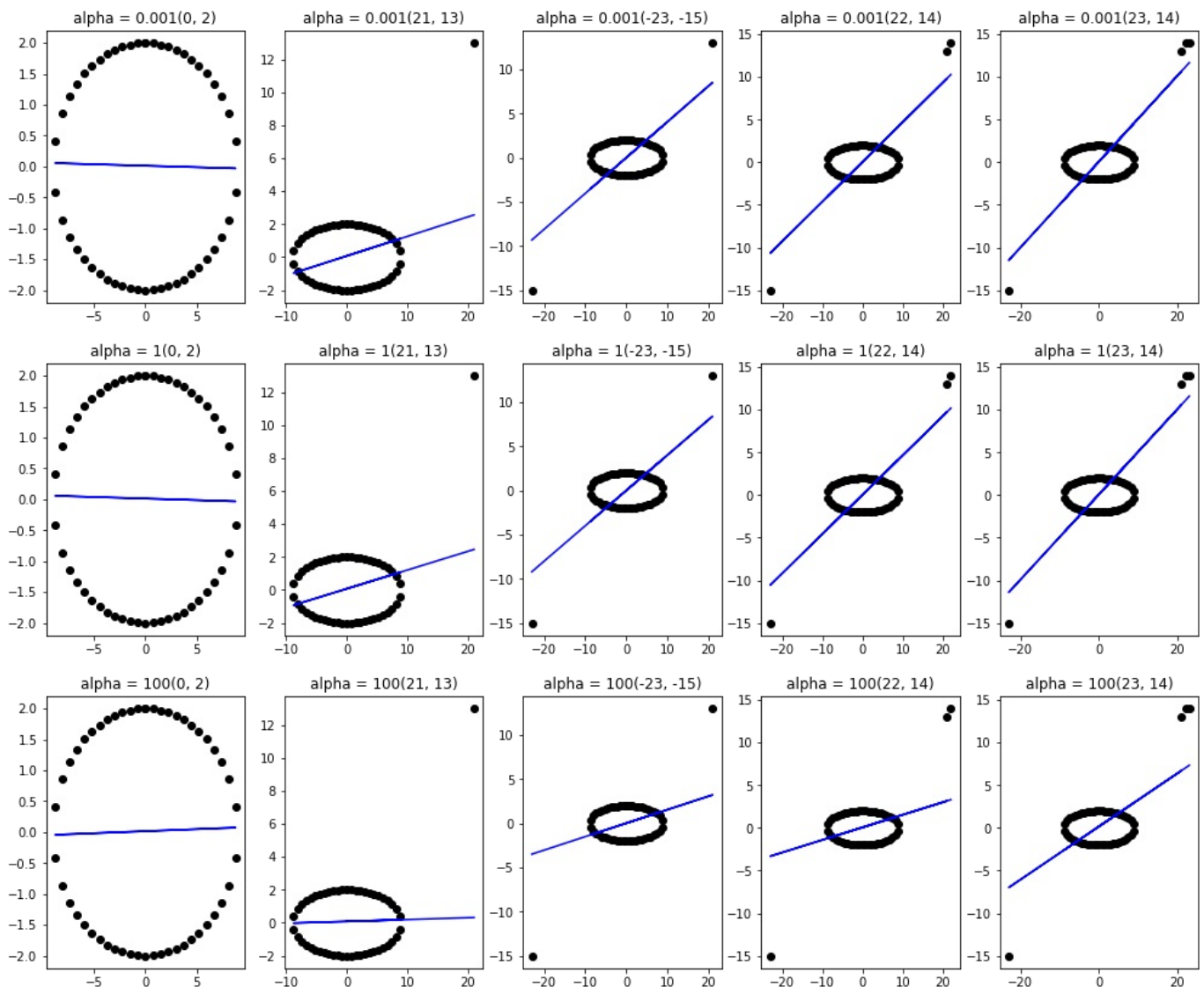
```
Out[39]: [<matplotlib.lines.Line2D at 0x2303e89eca0>]
```





```
In [36]: # https://stackoverflow.com/questions/59367939/how-to-correct-the-position-of-hyper-plane-in-python
alpha = [0.001, 1, 100]
outlier = [(0,2), (21,13), (-23,-15), (22,14), (23,14)]
for i in range(len(alpha)):
    plt.figure(figsize = (17,14))
    n = 1
    X = b * np.sin(phi)
    Y = a * np.cos(phi)
    for j in outlier:
        plt.subplot(3, 5, n)
        n += 1
        X = np.append(X, j[0]).reshape(-1,1)
        Y = np.append(Y, j[1]).reshape(-1,1)
        model = SGDRegressor(alpha = alpha[i], eta0=0.001, learning_rate='constant', random_state=0)
        model.fit(X, Y)
        y_pred = model.predict(X)
        plt.scatter(X,Y, color = 'black')
        plt.plot(X,y_pred, color = 'blue')
        plt.title("alpha = " + str(alpha[i]) + str(j))
plt.show
```

```
Out[36]: <function matplotlib.pyplot.show(close=None, block=None)>
```



Observation Based on outliers plot1 : this outlier not influence the plane plot2 : one outlier point on right side top corner not influencing the plane that much it while alpha = 100 other wise some misclassification is there, by slightly pull up plane plot3 : one outlier point on right side top corner and another outlier point on left side bottom corner which simultaneously pulling the plane to their direction, more misclassifications are there. plot4 : more weight on right side of the

top corner dominates pulling up and down the plane leads more misclassification. But misclassification is reduced while  $\alpha = 100$ . plot5 : same point mentioned in plot4 but comparatively misclassification more here because of the cluster of outlier in one side

In [ ]:

Loading [MathJax]/jax/output/CommonHTML/fonts/TeX/fontdata.js