

ASSIGNMENT 1 : LINEAR REGRESSION

Link for dataset:

Name: _____

Read the whole assignment before starting, Part 2 has the code template for doing Part 1.

Part 1: The Assignment Tasks

1. **Define Loss Functions:** Create two Python functions that calculate the total error for a given slope (m) and intercept (c).
 2. **Optimization:** Use `scipy.optimize.minimize` to find the best-fit parameters for Datasets A and B (using MSE).
 3. **The Outlier Challenge:** For Dataset C, find two sets of parameters:
 - o One that minimizes the **MSE**.
 - o One that minimizes the **MAE**.
 4. **Visualization:** Plot Dataset C with both the MSE line and the MAE line. Highlight the outliers in a different color.
 5. **Statistical Analysis:** For both lines in Dataset C, calculate:
 - o The number of points **above** vs **below** the line.
 - o The **Mean Distance** of points on either side.
-

Part 3: Implementation Template

Required Tools

- **Library:** `from scipy.optimize import minimize`
- **Core Function:** `minimize(fun, x0, args=(x, y))`
 - o `fun`: The loss function you want to minimize.
 - o `x0`: Your initial guess for $[m, c]$, e.g., `[0, 0]`.
 - o `args`: The actual data (x, y) passed to your loss function
 - o This function changes the values of `x0` such that it minimizes the function `fun` for the given `args`.

Code Skeleton

Python

1. `#import matplotlib.pyplot numpy, pandas and the above library`
2. `# --- STEP 1: DEFINE LOSS FUNCTIONS —`
3. `# save the given csv file as df`
4. `x=df['x'].values`
5. `data_c=df['y'].values`
6. `def mse_loss(params, x, y):`
7. `m, c = params`

```
8.     y_hat = m * x + c
9.     # Calculate Mean Squared Error
10.    # YOUR CODE HERE: return np.mean(...)
11.
12. def mae_loss(params, x, y):
13.     m, c = params
14.     y_hat = m * x + c
15.     # Calculate Mean Absolute Error
16.     # YOUR CODE HERE: return np.mean(np.abs(...))
17.
18. # --- STEP 2: OPTIMIZATION ---
19. # Example for Dataset C using MSE
20. initial_guess = [0, 0]
21. result_mse = minimize(mse_loss, initial_guess, args=(x, data_c))
22. m_mse, c_mse = result_mse.x
23.
24. # YOUR CODE HERE: Repeat for MAE on Dataset C
25. # result_mae = minimize(...)
26.
27. # --- STEP 3: RESIDUAL ANALYSIS ---
28. def analyze_line(x, y, m, c):
29.     y_pred = m * x + c
30.     residuals = y - y_pred
31.
32.
33.     # YOUR CODE HERE:
34.     # 1. Count points where residuals > 0 (Above) and residuals < 0 (Below)
35.     # 2. Calculate mean of residuals for Above group
36.     # 3. Calculate mean of ABSOLUTE residuals for Below group
37.
38.     print(f"Above: {count_above} (Avg Dist: {dist_above:.2f})")
39.     print(f"Below: {count_below} (Avg Dist: {dist_below:.2f})"
40.
41. # CODE FOR PLOTTING
42.
43. import matplotlib.pyplot as plt
44. #assuming that your ideal m and c values are stored as m_opt and c_opt
45. # 1. Generate predictions using the optimized parameters
46. # We use the original x_data to create the y-values for the line
47. y_pred = m_opt * x_data + c_opt
48.
49. # 2. Create the plot
50. plt.figure(figsize=(10, 6))
51.
```

```
52. # Plot original data points as a scatter plot
53. plt.scatter(x_data, y_data, color='blue', alpha=0.5, label='Original Data')
54.
55. # Plot the regression line
56. plt.plot(x_data, y_pred, color='red', linewidth=2, label='Fit: y = {m_opt:.2f}x + {c_opt:.2f}')
57.
58. # 3. Add labels and styling
59. plt.title('Linear Regression using scipy.optimize.minimize', fontsize=14)
60. plt.xlabel('x', fontsize=12)
61. plt.ylabel('y', fontsize=12)
62. plt.legend()
63. plt.grid(True, linestyle='--', alpha=0.7)
64.
65. # Show the plot
66. plt.show()
```

Part 4: Discussion Questions for Student

- **Balance:** Which line (MSE or MAE) split the 2,000 points more evenly (1000 above / 1000 below)? Why? Write the count of split of points above and below the line for dataset C.
- **Comment on the sensitivity of the two loss functions in dataset C and write the observations and reasoning for the observations**