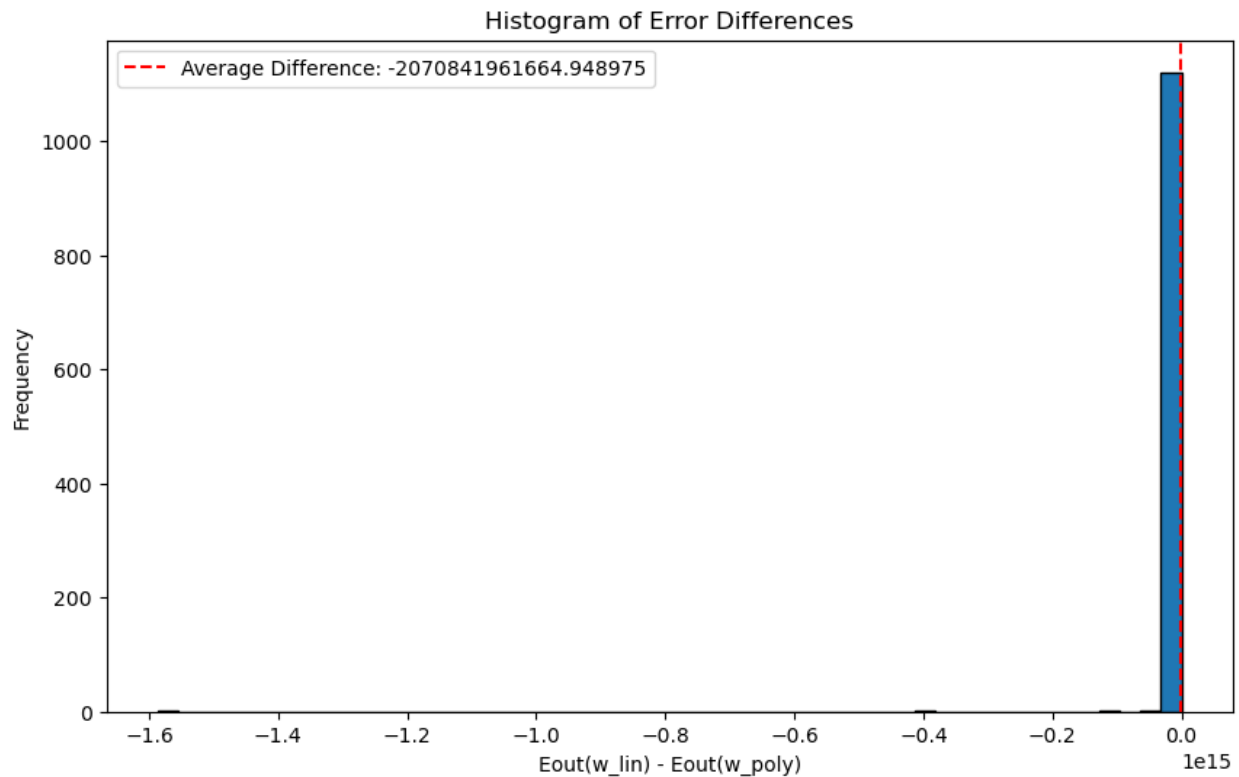


ML homework 4: question 12



Meaning:

From the plot we can see that the differences are all negative, indicating that we're getting way more bigger $E_{out}(w_{poly})$ than $E_{out}(w_{LIN})$, this may due to we're using too many features to fit our data. The original data we have might lie in a smaller space, so overfitting occurs, hence generating enormous out of sample error.

Code:

transform function $\Phi()$

We use this function to map the original data into a higher dimensional space, and save the augmented input vectors in `X_aug_arr`. Which means that each of the element in `X_aug_arr` is a 37-dimensional vector (12 features + 1 constant + 12 squared features + 12 cubed features).

```
1 def Phi(X_arr):
2     X_aug_arr = []
3     for arr in X_arr:
4         Phi_arr = np.concatenate((np.array([1]), arr, arr**2, arr**3))
5         X_aug_arr.append(Phi_arr)
6     return np.array(X_aug_arr)
7
```

[41] ✓ 0.0s

get weight vector

+ Code

+ Markdown

```

1 def weight_vector(X_in_sample_mat_lin, y_in_sample_array, X_in_sample_mat_poly):
2
3     w_lin = np.linalg.inv(X_in_sample_mat_lin.T @ X_in_sample_mat_lin) @ X_in_sample_mat_lin.T @ y_in_sample_array
4     w_poly = np.linalg.inv(X_in_sample_mat_poly.T @ X_in_sample_mat_poly) @ X_in_sample_mat_poly.T @ y_in_sample_array
5     return w_lin, w_poly
6

```

[44] ✓ 0.0s

out of sample error

```

1 def out_of_sample_error(X_out_of_sample_mat_lin, X_out_of_sample_mat_poly, y_out_of_sample_array, w_lin, w_poly):
2
3     out_of_sample_error_lin = np.mean((X_out_of_sample_mat_lin @ w_lin - y_out_of_sample_array) ** 2)
4     out_of_sample_error_poly = np.mean((X_out_of_sample_mat_poly @ w_poly - y_out_of_sample_array) ** 2)
5     return out_of_sample_error_lin, out_of_sample_error_poly
6
7

```

[45] ✓ 0.0s

```

1 for experiment in range(1126):
2
3     seed = experiment
4     random_sample_indices = generate_random_sample(seed)
5
6     X_sample = [X[i] for i in random_sample_indices]
7     y_sample = [y[i] for i in random_sample_indices]
8
9     X_in_sample_mat_lin = np.array(convert_dtype(X_sample))
10    X_in_sample_mat_poly = Phi(convert_dtype(X_sample))
11
12    y_in_sample_array = np.array(y_sample)
13
14    out_ind = generate_out_of_sample_ind(random_sample_indices)
15
16    X_out_of_sample = [X[i] for i in out_ind]
17    y_out_of_sample = [y[i] for i in out_ind]
18
19    X_out_of_sample_mat_lin = np.array(convert_dtype(X_out_of_sample))
20    X_out_of_sample_mat_poly = Phi(convert_dtype(X_out_of_sample))
21
22    y_out_of_sample_array = np.array(y_out_of_sample)
23
24    w_lin, w_poly = weight_vector(X_in_sample_mat_lin, y_in_sample_array, X_in_sample_mat_poly)
25    out_of_sample_error_lin, out_of_sample_error_poly = out_of_sample_error(X_out_of_sample_mat_lin, X_out_of_sample_mat_poly, y_out_of_sample_array, w_lin, w_poly)
26
27    out_sample_error_lin.append(out_of_sample_error_lin)
28    out_sample_error_poly.append(out_of_sample_error_poly)
29    lin_sub_poly_error.append(out_of_sample_error_lin - out_of_sample_error_poly)
30    avg_difference = np.mean(lin_sub_poly_error)
31
32
33 plt.figure(figsize=(10, 6))
34 plt.hist(lin_sub_poly_error, bins=50, edgecolor='black')
35 plt.axvline(x=avg_difference, color='r', linestyle='--',
36            label=f'Average Difference: {avg_difference:.6f}')
37 plt.xlabel('Eout(w_lin) - Eout(w_poly)')
38 plt.ylabel('Frequency')
39 plt.title('Histogram of Error Differences')
40 plt.legend()
41 plt.show()

```