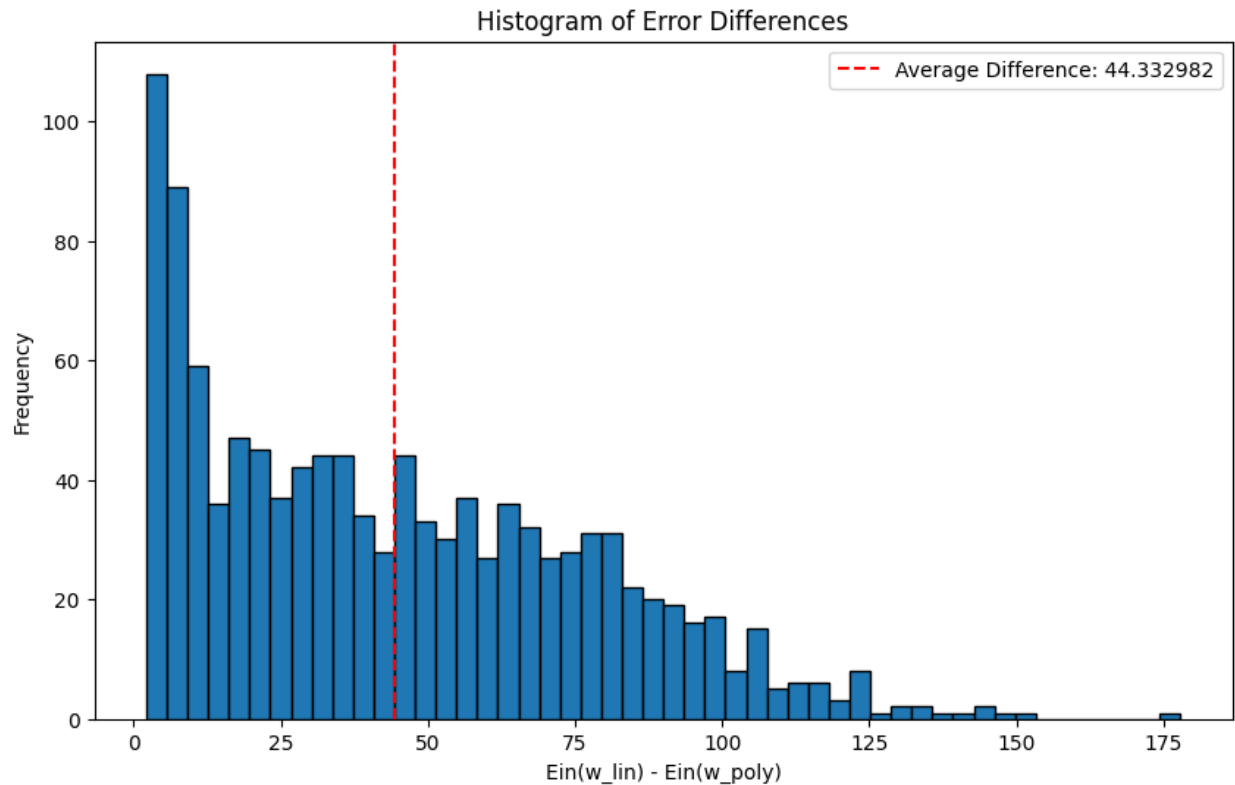


ML homework 4: question 11



Meaning:

From this plot we can see that the differences are mostly small positive values, which means that the error we got using linear regression is slightly greater than the error we got using polynomial transform.

This is because by using polynomial transform, we can use more parameters to fit our data, making our model more flexible to capture the nonlinear relationships among the features, thus, we can draw a function that is nearer to fitting all the examples over the 64 points.

However, the difference between the two approaches are not large, owing to our small dataset and small number of features, since under this situation, it is possible to find a line that nearly fits the points in a low dimensional space.

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

```

[64] ✓ 0.0s

get weight vector

Use the pseudo inverse to get the weight vector.

```

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

```

[26] ✓ 0.0s

in sample error

```

1 def in_sample_error(X_sample_mat_lin, X_sample_mat_poly, y_sample_array, w_lin, w_poly):
2
3     in_sample_error_lin = np.mean((X_sample_mat_lin @ w_lin - y_sample_array) ** 2)
4     in_sample_error_poly = np.mean((X_sample_mat_poly @ w_poly - y_sample_array) ** 2)
5     return in_sample_error_lin, in_sample_error_poly

```

[27] ✓ 0.0s

Main function

```

1 for experiment in range(1126):
2
3     seed = experiment
4     random_sample_indices = generate_random_sample(seed)
5     X_sample = [X[i] for i in random_sample_indices]
6     y_sample = [y[i] for i in random_sample_indices]
7
8     X_sample_mat_lin = np.array(convert_dtype(X_sample))
9     X_sample_mat_poly = Phi(convert_dtype(X_sample))
10
11     y_sample_array = np.array(y_sample)
12
13     w_lin, w_poly = weight_vector(X_sample_mat_lin, y_sample_array, X_sample_mat_poly)
14     in_sample_error_lin, in_sample_error_poly = in_sample_error(X_sample_mat_lin, X_sample_mat_poly, y_sample_array, w_lin, w_poly)
15
16     lin_sub_poly_error.append(in_sample_error_lin - in_sample_error_poly)
17     avg_difference = np.mean(lin_sub_poly_error)
18
19
20

```

[28] ✓ 0.7s

```

1 import matplotlib.pyplot as plt
2 plt.figure(figsize=(10, 6))
3 plt.hist(lin_sub_poly_error, bins=50, edgecolor='black')
4 plt.axvline(x=avg_difference, color='r', linestyle='--',
5            label=f'Average Difference: {avg_difference:.6f}')
6 plt.xlabel('Ein(w_lin) - Ein(w_poly)')
7 plt.ylabel('Frequency')
8 plt.title('Histogram of Error Differences')
9 plt.legend()
10 plt.show()

```

[29] ✓ 0.1s