



Design and Analysis of Algorithms

SOFE3770U Final Project: Part 1: Data Cleaning & Linear Regression

Course Name	Design and Analysis of Algorithms
Course Code	SOFE 3770U
Course CRN	43513
Date	Oct 24, 2025
Group Number	Final Project 20
Group Member Names	Vinujen Dilogen - 100870390 Rukshan Baskaran - 100864410 Khushi Patel - 100940709 Prabhnoor Saini - 100946515

Objective

The goal of this project is to predict the overall State of Health (SOH) of a battery pack using a linear regression model, and then use that prediction to inform a chatbot system that can communicate the battery's health condition and answer user queries. The goal of this part is to develop a linear regression model to predict the State of Health of battery cells based on voltage readings. This model will aim to estimate the SOH and classify the battery as Healthy or Unhealthy, depending on a user-defined threshold (default 0.6)

Dataset Description

The dataset used in this experiment is a PulseBat Dataset, which contains experimental measurements collected from a 21-cell lithium-ion battery pack. It consists of 670 samples and 30 columns, combining metadata, charge-discharge characteristics, voltage readings, and the targeted output variable, State of Health (SOH).

Model and Evaluation Metrics

A linear regression model was trained using an 80-20 train-test split. Three evaluation metrics were used:

R²: Measures how much variance in SOH is explained by the model

MSE (Mean Squared Error): Penalizes larger prediction errors more heavily

MAE (Mean Absolute Error): Average of the absolute errors, interpretable in the same units as SOH.

Classification

A classification rule was applied for the given SOH, if the SOH was 0.8:

- **Healthy:** SOH \geq 0.8, **Unhealthy:** SOH $<$ 0.8

Actual vs Predicted SOH Scatter Plot

This scatter plot compares the actual SOH value from the dataset against the predicted SOH values produced by the Linear Regression model. Each point in the plot represents one test sample from the dataset, x-axis showing the actual SOH value, while the y-axis shows the predicted SOH value. Points close to the red dashed diagonal line, where $y=x$, represent accurate predictions where the model output closely matches the real SOH. The color coding is based on the threshold value. Green points are predicted as healthy with SOH being greater than or equal to the threshold, while red points are predicted as unhealthy where the SOH is less than the threshold. This plot helps visualize how well the model predicts battery health, with points lying close to the diagonal indicating good correlation between predicted and actual SOH.

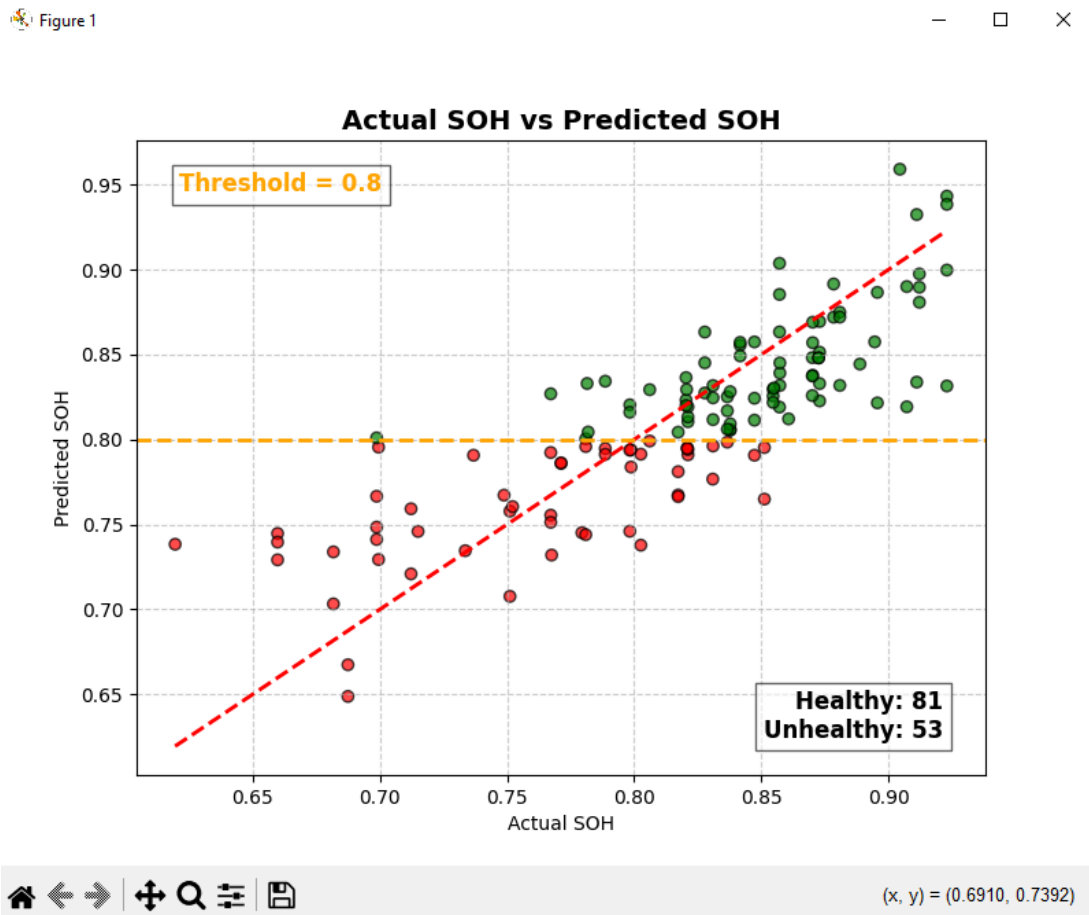


Figure 1: Actual vs Predicted SOH Scatter Plot

Residuals (Error) Distribution Histogram

This histogram of the residuals represents the difference between the actual and predicted State of Health (SOH) values.

$$\text{Residual} = \text{Actual SOH} - \text{Predicted SOH}$$

The x-axis displays the range of residual values, showing how far each prediction deviates from the true SOH. The y-axis represents the frequency, the number of test samples that fall within each residual. The black dashed vertical line at zero marks the point of perfect prediction, where actual and predicted SOH values are equal. A smooth density curve (KDE_) is drawn on top of the bars to show the overall distribution trends. The annotation box ($|\text{Residual}| > 0.05$) counts the number of samples with an absolute error greater than 0.05, identifying them as potential outliers or less accurate predictions.

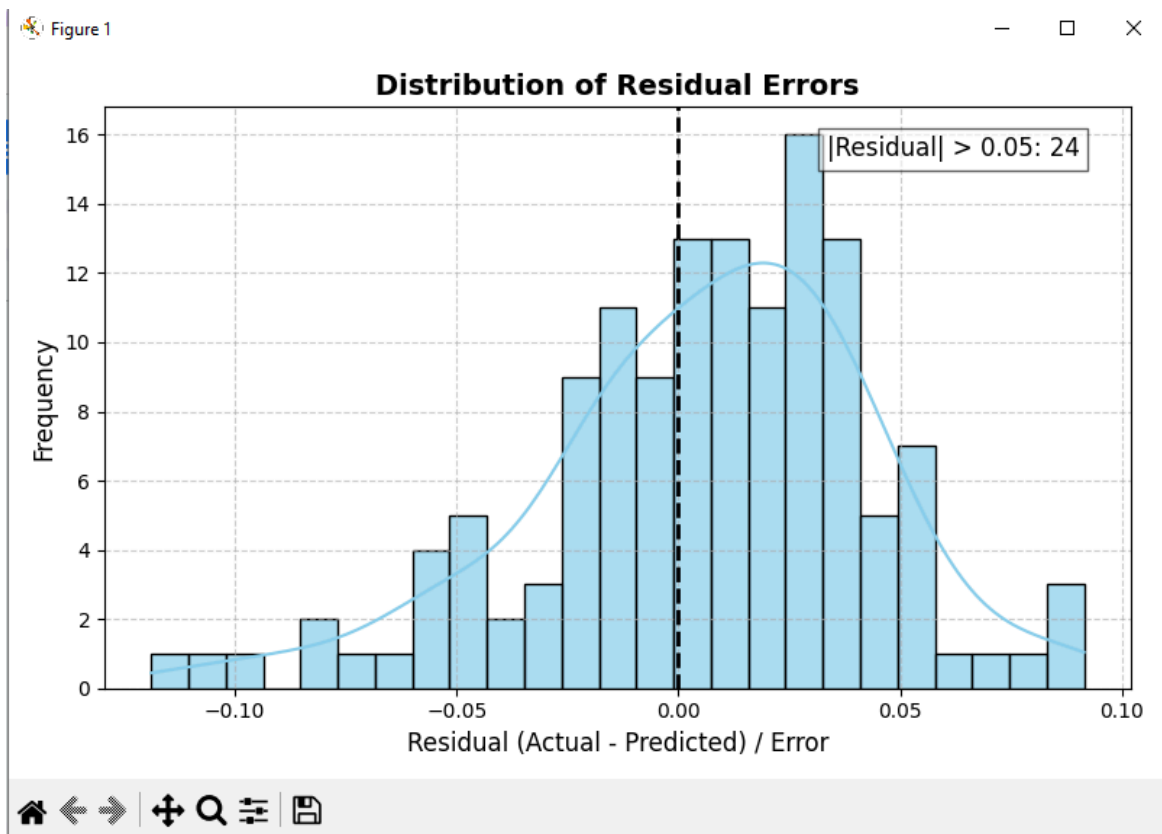


Figure 2: Distribution of Residual Errors

Data Preprocessing and Linear Regression Model Performance on Unsorted, Ascending, and Descending Voltage Data

The output shows that the dataset was successfully loaded, cleaned, and analyzed. Only the voltage readings (U1-U21) and the target variable (SOH) were selected from the original file. While unnecessary columns were removed. No missing values were found, confirming that the dataset is complete. To test whether the order of voltages affects prediction accuracy, three dataset versions were created: unsorted, ascending, and descending. Each was used to train a Linear Regression model, and their performance was evaluated using R^2 , MSE, and MAE metrics. All three versions produced similar accuracy, showing that sorting has minimal effect on results.

```
Columns in dataset: ['Mat', 'No.', 'ID', 'Qn', 'Q', 'Pt', 'SOC', 'SOE', 'U1', 'U2', 'U3', 'U4', 'U5', 'U6', 'U7', 'U8', 'U9', 'U10', 'U11', 'U12', 'U13', 'U14', 'U15', 'U16', 'U17', 'U18', 'U19', 'U20', 'U21', 'SOH']
Selected columns:
  U1    U2    U3    U4    U5    U6    U7    U8    U9    U10   U11   U12   U13   U14   U15   U16   U17   U18   U19   U20   U21   SOH
0  3.4858  3.5072  3.5246  3.5035  3.4877  3.4681  3.4452  3.4654  3.4849  3.5277  3.5599  3.5187  3.4889  3.4530  3.4009  3.4381  3.4840  3.5636  3.6241  3.5469  3.4923  0.912143
1  3.4877  3.5047  3.5317  3.5106  3.4895  3.4706  3.4406  3.4598  3.4864  3.5249  3.5711  3.5320  3.4914  3.4530  3.3903  3.4312  3.4855  3.5692  3.6396  3.5596  3.4945  0.880905
2  3.4858  3.5013  3.5305  3.5131  3.4877  3.4678  3.4387  3.4588  3.4849  3.5258  3.5701  3.5292  3.4892  3.4483  3.3891  3.4309  3.4836  3.5565  3.6399  3.5658  3.4929  0.857333
3  3.4898  3.5084  3.5385  3.5199  3.4917  3.4728  3.4375  3.4576  3.4886  3.5267  3.5810  3.5410  3.4936  3.4539  3.3832  3.4226  3.4874  3.5689  3.6532  3.5773  3.4973  0.831048
4  3.4979  3.5162  3.5494  3.5292  3.5001  3.4799  3.4424  3.4635  3.4967  3.5329  3.5940  3.5556  3.5016  3.4672  3.3860  3.4269  3.4957  3.5658  3.6684  3.5968  3.5060  0.781952

Dropped 0 rows with missing values. Final shape: (670, 22)
Cleaned dataset saved as 'Cleaned_PulseBat_Dataset.xlsx'

Original first row:
[3.4858 3.5072 3.5246 3.5035 3.4877 3.4681 3.4452 3.4654 3.4849 3.5277
 3.5599 3.5187 3.4889 3.453 3.4009 3.4381 3.484 3.5636 3.6241 3.5469
 3.4923]

Ascending sorted row:
[3.4009 3.4381 3.4452 3.453 3.4654 3.4681 3.484 3.4849 3.4858 3.4877
 3.4889 3.4923 3.5035 3.5072 3.5187 3.5246 3.5277 3.5469 3.5599 3.5636
 3.6241]

Descending sorted row:
[3.6241 3.5636 3.5599 3.5469 3.5277 3.5246 3.5187 3.5072 3.5035 3.4923
 3.4889 3.4877 3.4858 3.4849 3.484 3.4681 3.4654 3.453 3.4452 3.4381
 3.4009]

Unsorted Data Model Evaluation Metrics:
R² Score: 0.6561
Mean Squared Error (MSE): 0.0015
Mean Absolute Error (MAE): 0.0303

Ascending Data Model Evaluation Metrics:
R² Score: 0.6588
Mean Squared Error (MSE): 0.0015
Mean Absolute Error (MAE): 0.0304
```

Figure 3: Data Preprocessing and Linear Regression Model Performance

Descending Data Model Evaluation Metrics

This figure displays the performance metrics and sample classifications for the linear regression model trained on descending-sorted battery voltage data. The first value present is the R^2 , which is the coefficient of determination. It shows how much of the variation in the actual SOH can be explained by the regression model. In this model, 65.88% of the variation can be explained, which is a moderate correlation and not perfect, but reasonably strong for battery data. MSE is the average of squared differences between predicted and actual SOH. The smaller the MSE, the better. The value shown, 0.0015, represents that most predictions are close to the actual value. Lastly, the MAE is the average absolute difference between predicted and actual SOH. It shows that on average, there is a 3% margin of error, which is extremely good. Moving on, the total execution time shows how long our Python script took to run, which includes data preprocessing, model training, evaluation, and plot saving. Finally, with the SOH threshold for classification entered as 0.8, with $\text{SOH} \geq 0.8 \rightarrow \text{Healthy}$ and $\text{SOH} < 0.8 \rightarrow \text{Unhealthy}$, the sample predictions with the classification table show how well your predictions matched with reality. You can see that most predictions match with only a few that are slightly off the near 0.8 threshold boundary.

```
Descending Data Model Evaluation Metrics:
R² Score: 0.6588
Mean Squared Error (MSE): 0.0015
Mean Absolute Error (MAE): 0.0304

Total script execution time: 1.7131 seconds

Enter SOH threshold for classification (e.g., 0.6): 0.8

Sample predictions with classification:

Actual SOH Predicted SOH Actual Class Predicted Class
1 0.779476 0.745177 Unhealthy Unhealthy
2 0.872952 0.869526 Healthy Healthy
3 0.767238 0.826942 Unhealthy Healthy
4 0.830952 0.776730 Healthy Unhealthy
5 0.878476 0.872055 Healthy Healthy
6 0.733619 0.734584 Unhealthy Unhealthy
7 0.821190 0.810496 Healthy Healthy
8 0.698762 0.766622 Unhealthy Unhealthy
9 0.817333 0.804345 Healthy Healthy
10 0.904524 0.959155 Healthy Healthy
PS C:\Users\User\Desktop\Battery Pack SOH Prediction>
```

Figure 4: Descending Data Model Evaluation Metrics in Terminal