

Comprehensive Hands-On Guide: Data Cleaning and Preprocessing for IoT Well Log Sensor Data

1 Objective

This guide covers the steps required to clean and preprocess well log data collected via IoT sensors. We will handle common data issues such as missing values, outliers, noise, and normalization. After this process, the data will be ready for further analysis or machine learning models.

2 Setup and Required Libraries

Before starting, install the necessary Python libraries:

```
1 pip install numpy pandas matplotlib seaborn scipy scikit-learn statsmodels
```

3 1. Loading and Inspecting the Well Log Data

Load the well log dataset from a CSV file containing columns like Depth, Gamma Ray, Resistivity, Sonic, and Temperature logs.

```
1 import pandas as pd
2
3 # Load the dataset (replace 'well_log_iot.csv' with your
  actual file path)
4 data = pd.read_csv('well_log_iot.csv')
5
6 # Display the first few rows of the dataset
7 print(data.head())
8
9 # Display basic statistics
10 print(data.describe())
```

4 2. Handling Missing Data

Missing data can arise due to IoT sensor failures or communication interruptions. Let's first identify and handle missing data.

4.1 2.1. Identifying Missing Data

We can check for missing values in the dataset and visualize them using a heatmap.

```
1 # Checking for missing values
2 missing_values = data.isnull().sum()
3 print("Missing values per column:\n", missing_values)
4
5 # Visualizing missing data
6 import seaborn as sns
7 import matplotlib.pyplot as plt
8
9 plt.figure(figsize=(10,6))
10 sns.heatmap(data.isnull(), cbar=False, cmap='magma')
11 plt.title("Heatmap of Missing Values")
12 plt.show()
```

4.2 2.2. Handling Missing Data

We can handle missing data using several methods such as dropping rows, filling with median, or interpolating.

```
1 # Option 1: Dropping rows with too many missing values
2 data_cleaned = data.dropna(thresh=data.shape[1] - 1) # Keep
   rows with at least N-1 non-null values
3
4 # Option 2: Fill missing values with the median of each
   column
5 data_cleaned = data.fillna(data.median())
6
7 # Option 3: Interpolating missing values (useful for time-
   series or depth-indexed data)
8 data_cleaned = data.interpolate()
9
10 # Option 4: KNN Imputation for more advanced imputation
11 from sklearn.impute import KNNImputer
12
13 imputer = KNNImputer(n_neighbors=3)
14 data_cleaned = pd.DataFrame(imputer.fit_transform(data),
   columns=data.columns)
15
16 # Check the cleaned dataset
17 print(data_cleaned.head())
```

5 3. Detecting and Treating Outliers

Outliers in IoT well log data can be due to sensor malfunctions or geological anomalies. We will detect them using visualizations and statistical methods, and then treat them.

5.1 3.1. Visualizing Outliers with Boxplots

```
1 # Visualizing outliers using boxplots
2 plt.figure(figsize=(12,6))
3 sns.boxplot(data=data[['Gamma Ray', 'Resistivity', 'Sonic',
4   'Temperature']])
5 plt.title('Boxplot for Outlier Detection')
6 plt.show()
```

5.2 3.2. Detecting Outliers Using Z-Score and IQR

We can detect outliers statistically using Z-scores and the interquartile range (IQR) method.

```
1 # Z-Score method to detect outliers
2 from scipy import stats
3 import numpy as np
4
5 z_scores = np.abs(stats.zscore(data[['Gamma Ray', '
6   Resistivity', 'Sonic']]))
7 outliers_z = np.where(z_scores > 3)
8 print("Outliers detected using Z-score at rows:", outliers_z
9   )
10
11 # IQR Method to detect outliers
12 Q1 = data[['Gamma Ray', 'Resistivity', 'Sonic']].quantile
13   (0.25)
14 Q3 = data[['Gamma Ray', 'Resistivity', 'Sonic']].quantile
15   (0.75)
16 IQR = Q3 - Q1
17
18 outliers_iqr = (data < (Q1 - 1.5 * IQR)) | (data > (Q3 + 1.5
19   * IQR))
20 print("Outliers detected using IQR method:\n", outliers_iqr.
21   sum())
```

5.3 3.3. Treating Outliers

You can remove outliers or cap them to a certain range.

```
1 # Option 1: Removing outliers (using Z-Score or IQR
2   detection)
```

```

2 data_cleaned = data[(z_scores < 3).all(axis=1)]
3
4 # Option 2: Capping outliers to the 1st and 99th percentile
5 def cap_outliers(df, column):
6     lower = df[column].quantile(0.01)
7     upper = df[column].quantile(0.99)
8     df[column] = np.clip(df[column], lower, upper)
9     return df
10
11 for col in ['Gamma Ray', 'Resistivity', 'Sonic']:
12     data_cleaned = cap_outliers(data_cleaned, col)
13
14 print(data_cleaned.head())

```

6 4. Noise Reduction

Sensor data from IoT systems often contains noise, which can be reduced using smoothing techniques like moving averages or filters.

6.1 4.1. Moving Average Smoothing

```

1 # Apply moving average smoothing with a window of 5
2 data_cleaned['Gamma Ray (Smoothed)'] = data_cleaned['Gamma
   Ray'].rolling(window=5).mean()
3 data_cleaned['Resistivity (Smoothed)'] = data_cleaned['
   Resistivity'].rolling(window=5).mean()
4
5 # Plot the original vs smoothed data
6 plt.figure(figsize=(12,6))
7 plt.plot(data_cleaned['Depth'], data_cleaned['Gamma Ray'],
   label='Original Gamma Ray', alpha=0.5)
8 plt.plot(data_cleaned['Depth'], data_cleaned['Gamma Ray (
   Smoothed)'], label='Smoothed Gamma Ray', color='red')
9 plt.xlabel('Depth')
10 plt.ylabel('Gamma Ray')
11 plt.legend()
12 plt.title('Gamma Ray Before and After Smoothing')
13 plt.show()

```

6.2 4.2. Low-Pass Filter

A low-pass filter removes high-frequency noise while preserving low-frequency signals.

```

1 from scipy.signal import butter, filtfilt
2

```

```

3 # Define a Butterworth low-pass filter
4 def butter_lowpass(cutoff, fs, order=5):
5     nyquist = 0.5 * fs
6     normal_cutoff = cutoff / nyquist
7     b, a = butter(order, normal_cutoff, btype='low', analog=
8         False)
9     return b, a
10
11 def lowpass_filter(data, cutoff, fs, order=5):
12     b, a = butter_lowpass(cutoff, fs, order=order)
13     y = filtfilt(b, a, data)
14     return y
15
16 # Apply the low-pass filter
17 filtered_resistivity = lowpass_filter(data_cleaned['
18     Resistivity'], cutoff=0.1, fs=1.0)
19
20 # Plot original vs filtered resistivity
21 plt.figure(figsize=(12,6))
22 plt.plot(data_cleaned['Depth'], data_cleaned['Resistivity'],
23     label='Original Resistivity', alpha=0.5)
24 plt.plot(data_cleaned['Depth'], filtered_resistivity, label=
25     'Filtered Resistivity', color='red')
26 plt.xlabel('Depth')
27 plt.ylabel('Resistivity')
28 plt.legend()
29 plt.title('Resistivity Before and After Low-Pass Filtering')
30 plt.show()

```

7 5. Data Normalization and Standardization

To ensure the data is on the same scale, we normalize or standardize it.

7.1 5.1. Normalization (Min-Max Scaling)

```

1 from sklearn.preprocessing import MinMaxScaler
2
3 # Normalize to the range [0, 1]
4 scaler = MinMaxScaler()
5 data_normalized = pd.DataFrame(scaler.fit_transform(
6     data_cleaned[['Gamma Ray', 'Resistivity', 'Sonic']]),
7     columns=['Gamma Ray', 'Resistivity', 'Sonic'])
8
9 print(data_normalized.head())

```

7.2 5.2. Standardization

```
1 from sklearn.preprocessing import StandardScaler
2
3 # Standardize to have a mean of 0 and standard deviation of
  1
4 scaler = StandardScaler()
5 data_standardized = pd.DataFrame(scaler.fit_transform(
    data_cleaned[['Gamma Ray', 'Resistivity', 'Sonic']],
    columns=['Gamma Ray', 'Resistivity', 'Sonic'])
6
7 print(data_standardized.head())
```

8 6. Feature Engineering

We can create new features based on existing logs, such as estimating lithology from gamma ray readings or calculating porosity from sonic logs.

8.1 6.1. Deriving Lithology from Gamma Ray

Lithology is often inferred from gamma ray values, where high values indicate shale and low values indicate sand.

```
1 # Creating a simple lithology feature based on Gamma Ray
  threshold
2 data_cleaned['Lithology'] = data_cleaned['Gamma Ray'].apply(
    lambda x: 'Shale' if x > 75 else 'Sand')
3
4 print(data_cleaned[['Depth', 'Gamma Ray', 'Lithology']].head
  ())
```

8.2 6.2. Estimating Porosity from Sonic Log

You can estimate porosity using sonic log data with empirical formulas.

```
1 # Simple porosity calculation from the sonic log
2 data_cleaned['Porosity'] = (data_cleaned['Sonic'] - 50) /
  150 # Simplified formula
3
4 print(data_cleaned[['Depth', 'Sonic', 'Porosity']].head())
```

9 7. Visualization of Cleaned and Preprocessed Data

After cleaning and preprocessing the data, you can visualize it.

```

1 # Plot final cleaned data
2 plt.figure(figsize=(12,8))
3
4 plt.subplot(3, 1, 1)
5 plt.plot(data_cleaned['Depth'], data_cleaned['Gamma Ray'],
6          label='Gamma Ray', color='blue')
7 plt.xlabel('Depth')
8 plt.ylabel('Gamma Ray')
9 plt.title('Gamma Ray vs Depth')
10
11 plt.subplot(3, 1, 2)
12 plt.plot(data_cleaned['Depth'], data_cleaned['Resistivity'],
13          label='Resistivity', color='green')
14 plt.xlabel('Depth')
15 plt.ylabel('Resistivity')
16 plt.title('Resistivity vs Depth')
17
18 plt.subplot(3, 1, 3)
19 plt.plot(data_cleaned['Depth'], data_cleaned['Sonic'], label
20          ='Sonic', color='red')
21 plt.xlabel('Depth')
22 plt.ylabel('Sonic')
23 plt.title('Sonic vs Depth')
24
25 plt.tight_layout()
26 plt.show()

```

10 Conclusion

In this comprehensive hands-on session, we:

- Loaded and inspected IoT well log data.
- Handled missing values using imputation and interpolation.
- Detected and treated outliers using statistical methods.
- Reduced noise with smoothing and low-pass filtering.
- Normalized and standardized the data.
- Engineered new features such as lithology and porosity.
- Visualized the final cleaned data.

These steps prepare IoT well log data for machine learning models or further geophysical analysis.