

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

1 Objective

In this hands-on guide, we will apply data preprocessing and cleaning techniques to well log data using Python. The key tasks include handling missing data, detecting and treating outliers, performing noise reduction, normalizing the data, and engineering features for further analysis.

2 Setup and Required Libraries

To begin, install the necessary Python libraries:

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

3 1. Loading the Well Log Data

We will load a well log dataset using `pandas`. This dataset should contain columns such as Depth, Gamma Ray, Resistivity, and Sonic logs.

```
1 import pandas as pd
2
3 # Load the dataset
4 data = pd.read_csv('well_log.csv')
5
6 # Display the first few rows of the dataset
7 print(data.head())
```

4 2. Handling Missing Data

Missing data is common in well logs due to various reasons such as sensor failures. We will first identify the missing data and then apply techniques to handle it.

4.1 2.1. Identifying Missing Data

We can check for missing values and visualize their presence using a heatmap:

```
1 # Check for missing values
2 missing_data = data.isnull().sum()
3 print("Missing values:\n", missing_data)
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='viridis')
11 plt.title("Missing Values Heatmap")
12 plt.show()
```

4.2 2.2. Imputation Techniques

There are different methods to handle missing values, such as filling with median, using interpolation, or applying K-Nearest Neighbors (KNN) imputation.

```
1 # Option 1: Fill missing values with the median of the
   column
2 data_filled = data.fillna(data.median())
3
4 # Option 2: Interpolate missing values
5 data_interpolated = data.interpolate()
6
7 # Option 3: KNN Imputation (requires scikit-learn)
8 from sklearn.impute import KNNImputer
9
10 imputer = KNNImputer(n_neighbors=5)
11 data_knn = pd.DataFrame(imputer.fit_transform(data), columns
   =data.columns)
12
13 # Visualize the imputed data
14 print(data_knn.head())
```

5 3. Outlier Detection and Treatment

Outliers can significantly affect the analysis. We will detect outliers using box-plots and the Z-score method, and then treat them using different strategies.

5.1 3.1. Detecting Outliers

We can use statistical methods and visualizations to detect outliers.

```

1 # Using boxplots to visualize outliers
2 plt.figure(figsize=(12,6))
3 sns.boxplot(data=data[['Gamma Ray', 'Resistivity', 'Sonic'
4                           ]])
5 plt.title('Boxplot to Detect Outliers')
6 plt.show()
7
8 # Z-Score method
9 from scipy import stats
10 import numpy as np
11
12 z_scores = np.abs(stats.zscore(data[['Gamma Ray', '
13      Resistivity', 'Sonic']]))
14 outliers = np.where(z_scores > 3)
15 print("Outliers detected at rows:", outliers)

```

5.2 3.2. Treating Outliers

You can choose to either remove or cap the outliers using the following code:

```

1 # Option 1: Remove rows with outliers
2 data_cleaned = data[(z_scores < 3).all(axis=1)]
3
4 # Option 2: Cap outliers to the 1st and 99th percentile
5 def cap_outliers(df, column):
6     lower_percentile = df[column].quantile(0.01)
7     upper_percentile = df[column].quantile(0.99)
8     df[column] = np.where(df[column] < lower_percentile,
9                           lower_percentile, df[column])
10    df[column] = np.where(df[column] > upper_percentile,
11                          upper_percentile, df[column])
12    return df
13
14 for col in ['Gamma Ray', 'Resistivity', 'Sonic']:
15     data_cleaned = cap_outliers(data_cleaned, col)
16
17 print(data_cleaned.head())

```

6 4. Noise Reduction

Noise in well log data can be addressed using smoothing techniques like moving averages or low-pass filters.

6.1 4.1. Moving Average Smoothing

```

1 # Apply moving average smoothing
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 Filtering

You can apply a low-pass filter to remove high-frequency noise from the data.

```

1 from scipy.signal import butter, filtfilt
2
3 # Design low-pass filter
4 def butter_lowpass_filter(data, cutoff, fs, order=5):
5     nyquist = 0.5 * fs
6     normal_cutoff = cutoff / nyquist
7     b, a = butter(order, normal_cutoff, btype='low', analog=
   False)
8     y = filtfilt(b, a, data)
9     return y
10
11 # Apply the filter to the resistivity log
12 filtered_resistivity = butter_lowpass_filter(data_cleaned['
   Resistivity'], cutoff=0.1, fs=1.0)
13
14 # Plot original vs filtered data
15 plt.figure(figsize=(12,6))
16 plt.plot(data_cleaned['Depth'], data_cleaned['Resistivity'],
   label='Original Resistivity', alpha=0.5)
17 plt.plot(data_cleaned['Depth'], filtered_resistivity, label=
   'Filtered Resistivity', color='red')
18 plt.xlabel('Depth')
19 plt.ylabel('Resistivity')
20 plt.legend()
21 plt.title('Resistivity Before and After Low-Pass Filtering')
22 plt.show()

```

7 5. Normalization and Scaling

Normalization or standardization is required to scale the data, especially when different logs have varying ranges.

7.1 5.1. Normalization

```
1 # Normalize the data to the range [0, 1]
2 from sklearn.preprocessing import MinMaxScaler
3
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 print(data_normalized.head())
```

7.2 5.2. Standardization

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

8 6. Feature Engineering

You can create new features from the existing logs to enhance your analysis.

```
1 # Example: Deriving a simplified lithology feature based on
2   Gamma Ray values
3 data_cleaned['Lithology'] = data_cleaned['Gamma Ray'].apply(
4     lambda x: 'Shale' if x > 75 else 'Sand')
5
6 print(data_cleaned[['Depth', 'Gamma Ray', 'Lithology']].head(
7     ))
```

9 7. Visualization of Preprocessed Data

Visualize the cleaned and processed data to ensure the changes were applied correctly.

```

1 # Plot final cleaned and processed 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 hands-on session, we:

- Loaded and visualized well log data.
- Handled missing values using imputation techniques.
- Detected and treated outliers using statistical methods.
- Applied noise reduction using moving average and low-pass filtering.
- Normalized and standardized the data for further analysis.
- Engineered new features like lithology based on gamma ray logs.

These preprocessing steps are essential for preparing well log data for machine learning models or further geophysical analysis.