# Data Prep

Great! Let's expand the data cleaning guide to include **every essential step** that will prepare your dataset for machine learning. This will cover all the necessary aspects to ensure your data is fully optimized for training models, enhancing accuracy, performance, and efficiency.

Below, we'll add additional sections such as **dimensionality reduction**, **feature selection**, and **handling class imbalance** in depth. These are critical steps in preparing data for machine learning models.

---

# 🛠 Comprehensive Data Preparation Guide for Machine Learning

This expanded guide will include everything you need to prepare your dataset, ensuring it's ready for robust machine learning model training. We'll cover data cleaning, feature engineering, transformations, and additional preparation steps.

---

## Step 1: Loading and Inspecting the Dataset

```python
import pandas as pd

# Load the dataset
df = pd.read_csv('/kaggle/input/ugransome-dataset/final(2).csv')

# Initial inspection
print(df.head())
print(df.info())
print(df.describe())
print("Number of missing values:\n", df.isnull().sum())
```

---

## Step 2: Renaming Columns for Clarity

```
df.columns = ['Time', 'Protocol', 'Flag', 'Family', 'Clusters',
              'SeedAddress', 'ExpAddress', 'BTC', 'USD',
              'Netflow_Bytes', 'IPaddress', 'Threats', 'Port', 'Prediction']
```

## Step 3: Handling Missing Values

- Fill missing numerical values using mean/median.
- Fill missing categorical values using mode.
- For sequential or time-series data, use interpolation.

```python
df['Netflow_Bytes'].fillna(df['Netflow_Bytes'].median(), inplace=True)
df['Threats'].fillna(df['Threats'].mode()[0], inplace=True)
df['Time'] = pd.to_datetime(df['Time']).interpolate(method='linear')
```

## Step 4: Removing Duplicates and Incorrect Values

```python
# Remove duplicates
df.drop_duplicates(inplace=True)

# Correct misspelled entries
df['Threats'] = df['Threats'].str.replace('Bonet', 'Botnet')
```

## Step 5: Handling Outliers

Remove outliers using the Interquartile Range (IQR) method.

```python
def remove_outliers(df, column):
    Q1 = df[column].quantile(0.25)
    Q3 = df[column].quantile(0.75)
    IQR = Q3 - Q1
    lower_bound = Q1 - 1.5 * IQR
    upper_bound = Q3 + 1.5 * IQR
    return df[(df[column] >= lower_bound) & (df[column] <= upper_bound)]
```

```python
df = remove_outliers(df, 'Netflow_Bytes')
df = remove_outliers(df, 'BTC')
df = remove_outliers(df, 'USD')
```

# Step 6: Feature Engineering

## 6.1 Extracting Date and Time Components

```python
df['Time'] = pd.to_datetime(df['Time'])
df['Hour'] = df['Time'].dt.hour
df['Day'] = df['Time'].dt.day
df['Month'] = df['Time'].dt.month
```

## 6.2 Creating New Features

```python
df['High_BTC'] = df['BTC'].apply(lambda x: 1 if x > df['BTC'].mean() else 0)
```

# Step 7: Data Transformation

## 7.1 Handling Skewed Data

```python
import numpy as np
from scipy import stats

# Log transformation
df['Netflow_Bytes'] = np.log1p(df['Netflow_Bytes'])

# Yeo-Johnson Transformation (handles negative values)
df['BTC'], _ = stats.yeojohnson(df['BTC'])

# Square root transformation
df['USD'] = np.sqrt(df['USD'])
```

## 7.2 Scaling and Normalization

```python
from sklearn.preprocessing import StandardScaler

scaler = StandardScaler()
df[['Netflow_Bytes', 'BTC', 'USD']] =
scaler.fit_transform(df[['Netflow_Bytes', 'BTC', 'USD']])
```

# Step 8: Encoding Categorical Variables

## 8.1 Label Encoding for Ordinal Data

```python
from sklearn.preprocessing import LabelEncoder

encoder = LabelEncoder()
df['Protocol'] = encoder.fit_transform(df['Protocol'])
df['Flag'] = encoder.fit_transform(df['Flag'])
df['Family'] = encoder.fit_transform(df['Family'])
```

## 8.2 One-Hot Encoding for Nominal Data

```python
df = pd.get_dummies(df, columns=['Clusters', 'Port'])
```

## Label Encoding Explained

Label Encoding is a technique used to convert categorical data into numerical values. This method is especially useful for variables that have an inherent ordinal relationship (e.g., low, medium, high). However, it can also be applied to non-ordinal categorical data when the order does not matter.

Label Encoding assigns a unique integer to each category, effectively transforming a categorical column into numerical format.

## When to Use Label Encoding?

- **Ordinal Data**: When there is a natural ordering between the categories.

- Example: "Low", "Medium", "High" → 0, 1, 2.
- **Non-Ordinal Data**: Can also be used for non-ordinal data, but it may introduce unintended biases in some machine learning models (e.g., decision trees are fine, but linear models might assume an ordinal relationship).

---

# Label Encoding Example: Step-by-Step Code

Suppose we have a dataset with a categorical column named `"Color"`.

## Step 1: Load the Dataset

```python
import pandas as pd

# Sample dataset
data = {
    'Color': ['Red', 'Green', 'Blue', 'Green', 'Red', 'Blue', 'Red']
}
df = pd.DataFrame(data)
print("Original DataFrame:")
print(df)
```

**Output:**

```
   Color
0    Red
1  Green
2   Blue
3  Green
4    Red
5   Blue
6    Red
```

---

## Step 2: Apply Label Encoding

```python
from sklearn.preprocessing import LabelEncoder

# Initialize the label encoder
```

```
label_encoder = LabelEncoder()

# Fit and transform the "Color" column
df['Color_Encoded'] = label_encoder.fit_transform(df['Color'])

print("\nLabel Encoded DataFrame:")
print(df)
```

**Output:**

```
   Color  Color_Encoded
0    Red              2
1  Green              1
2   Blue              0
3  Green              1
4    Red              2
5   Blue              0
6    Red              2
```

# Explanation of the Output:

- The Label Encoder assigns numerical values to each unique category in the column:
  - `Red` → 2
  - `Green` → 1
  - `Blue` → 0

# Step 3: Inverse Transform (Optional)

If you want to convert the encoded values back to their original categories:

```
# Inverse transform the encoded labels back to the original values
df['Color_Decoded'] = label_encoder.inverse_transform(df['Color_Encoded'])

print("\nDecoded DataFrame:")
print(df)
```

**Output:**

```
   Color  Color_Encoded Color_Decoded
0   Red              2           Red
1 Green              1         Green
2  Blue              0          Blue
3 Green              1         Green
4   Red              2           Red
5  Blue              0          Blue
6   Red              2           Red
```

## Notes:

1. **Pros**: Simple and efficient for models that can handle categorical variables.
2. **Cons**: Can introduce unintended ordinal relationships between categories, which might mislead some algorithms (e.g., linear regression).
3. **Alternative**: For non-ordinal data, consider using One-Hot Encoding instead.

## Conclusion

Label Encoding is a simple yet powerful tool for converting categorical data into numerical values. However, it's important to choose the encoding method that best suits your dataset and the machine learning model you plan to use.

# Step 9: Feature Selection

## 9.1 Removing Low-Variance Features

```python
from sklearn.feature_selection import VarianceThreshold

selector = VarianceThreshold(threshold=0.01)
df = pd.DataFrame(selector.fit_transform(df),
columns=df.columns[selector.get_support()])
```

## 9.2 Correlation-Based Feature Selection

```python
import seaborn as sns
import matplotlib.pyplot as plt
```

```python
corr_matrix = df.corr()
plt.figure(figsize=(12, 8))
sns.heatmap(corr_matrix, annot=True, cmap='coolwarm')
plt.title('Correlation Matrix')
plt.show()

# Remove highly correlated features
threshold = 0.9
correlated_features = [col for col in corr_matrix if any(corr_matrix[col] >
threshold)]
df.drop(columns=correlated_features, inplace=True)
```

## 9.3 Feature Importance Using Random Forest

```python
from sklearn.ensemble import RandomForestClassifier

X = df.drop('Prediction', axis=1)
y = df['Prediction']
model = RandomForestClassifier(random_state=42)
model.fit(X, y)
importances = model.feature_importances_

# Plot feature importance
plt.figure(figsize=(10, 6))
plt.barh(X.columns, importances)
plt.title('Feature Importance')
plt.show()
```

# Step 10: Dimensionality Reduction

## 10.1 Principal Component Analysis (PCA)

```python
from sklearn.decomposition import PCA

pca = PCA(n_components=10)
X_pca = pca.fit_transform(X)
```

**Explanation**:

- **PCA** reduces dimensionality while retaining variance, which speeds up training and reduces overfitting.

---

# Step 11: Handling Class Imbalance

## 11.1 Using SMOTE for Oversampling

```python
from imblearn.over_sampling import SMOTE

smote = SMOTE(random_state=42)
X_resampled, y_resampled = smote.fit_resample(X, y)
```

**Explanation**:

- SMOTE generates synthetic samples for the minority class to balance the dataset.

---

# Step 12: Splitting Data into Training and Testing Sets

```python
from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(X_resampled, y_resampled,
train_size=0.8, random_state=42)
```

---

# Step 13: Final Checks and Validation

## 13.1 Checking Data Distributions

```python
df.hist(bins=20, figsize=(20, 15))
plt.show()
```

## 13.2 Verifying Class Distribution

```
df['Prediction'].value_counts().plot(kind='bar', color='skyblue')
plt.title('Class Distribution')
plt.xlabel('Class')
plt.ylabel('Count')
plt.show()
```

# Step 14: Saving the Cleaned Dataset

```
df.to_csv('cleaned_ugransome_dataset.csv', index=False)
```

Great! Let's extend the notes to include **encoding techniques**, **data transformations**, and the various strategies for **filling missing values**. This will make your exam prep comprehensive for data preprocessing.

# Section 2: Encoding, Transformations, and Handling Missing Values

## 2.1 Encoding Categorical Variables

Machine learning models require numerical inputs. Therefore, converting categorical variables into numerical formats is essential.

## Types of Encoding

**1. Label Encoding**

- Converts categorical labels into numeric values (0, 1, 2, ...).
- Useful for ordinal variables (where there is a natural order).

**Code Example: Label Encoding**

```
from sklearn.preprocessing import LabelEncoder

# Initialize the encoder
label_encoder = LabelEncoder()
```

```
# Apply to a categorical column
df['category_encoded'] = label_encoder.fit_transform(df['category_column'])

print(df[['category_column', 'category_encoded']].head())
```

## 2. One-Hot Encoding

- Converts categorical variables into a set of binary (0 or 1) columns.
- Suitable for nominal variables (no intrinsic order).

### Code Example: One-Hot Encoding

```
# One-hot encode using pandas
df = pd.get_dummies(df, columns=['category_column'], drop_first=True)

print(df.head())
```

## 3. Ordinal Encoding

- Assigns numerical values to categories based on their order.
- Useful for ordinal data (e.g., low, medium, high).

### Code Example: Ordinal Encoding

```
from sklearn.preprocessing import OrdinalEncoder

# Define the order
categories = [['low', 'medium', 'high']]
ordinal_encoder = OrdinalEncoder(categories=categories)

df['priority_encoded'] = ordinal_encoder.fit_transform(df[['priority']])

print(df[['priority', 'priority_encoded']].head())
```

## 4. Binary Encoding

- Converts categories into binary digits, then encodes them as binary numbers.
- Useful for high-cardinality categorical features.

### Code Example: Binary Encoding

```
!pip install category_encoders
import category_encoders as ce

# Initialize the binary encoder
binary_encoder = ce.BinaryEncoder(cols=['category_column'])
df = binary_encoder.fit_transform(df)

print(df.head())
```

## 5. Frequency Encoding

- Replaces each category with its frequency in the dataset.
- Useful for high-cardinality features where preserving the frequency information is beneficial.

### Code Example: Frequency Encoding

```
df['category_freq'] =
df['category_column'].map(df['category_column'].value_counts())
print(df[['category_column', 'category_freq']].head())
```

---

# 2.2 Data Transformations

Transformations are used to modify data distributions, handle skewness, and make variables more suitable for machine learning algorithms.

## Types of Transformations

### 1. Log Transformation

- Reduces skewness for data that has a long right tail (positive skew).
- Useful for features like income, prices, or any positively skewed data.

### Code Example: Log Transformation

```
df['log_transformed'] = np.log1p(df['skewed_feature'])
```

### 2. Square Root Transformation

- Used to reduce right skew, especially when the data has zeros.

**Code Example: Square Root Transformation**

```
df['sqrt_transformed'] = np.sqrt(df['feature'])
```

### 3. Box-Cox Transformation

- Transforms data to approximate normal distribution.
- Requires all values to be positive.

**Code Example: Box-Cox Transformation**

```
from scipy.stats import boxcox

df['boxcox_transformed'], _ = boxcox(df['positive_feature'] + 1)
```

### 4. Yeo-Johnson Transformation

- Similar to Box-Cox but can handle zero and negative values.

**Code Example: Yeo-Johnson Transformation**

```
from sklearn.preprocessing import PowerTransformer

power_transformer = PowerTransformer(method='yeo-johnson')
df['yeojohnson_transformed'] =
power_transformer.fit_transform(df[['feature']])
```

### 5. Exponential Transformation

- Used to handle left-skewed data.

**Code Example: Exponential Transformation**

```
df['exp_transformed'] = np.exp(df['feature'])
```

---

# 2.3 Handling Missing Values

Missing values are common in real-world data, and there are different strategies for handling them based on the type of data and the extent of missingness.

# Techniques to Fill Missing Values

**1. Dropping Missing Values**

- Use this when the dataset is large and missing values are negligible.

**Code Example: Dropping Missing Values**

```python
df.dropna(inplace=True)
```

**2. Imputation Strategies**

- **Mean Imputation**: For numerical columns with normally distributed data.

```python
df['column'] = df['column'].fillna(df['column'].mean())
```

- **Median Imputation**: For numerical columns with skewed data.

```python
df['column'] = df['column'].fillna(df['column'].median())
```

- **Mode Imputation**: For categorical columns.

```python
df['category'] = df['category'].fillna(df['category'].mode()[0])
```

**3. Forward Fill/Backward Fill**

- Useful for time-series data where trends are preserved.

**Code Example: Forward/Backward Fill**

```python
df.fillna(method='ffill', inplace=True)  # Forward fill
df.fillna(method='bfill', inplace=True)  # Backward fill
```

**4. K-Nearest Neighbors Imputation**

- Imputes missing values based on the similarity of features.

**Code Example: KNN Imputer**

```python
from sklearn.impute import KNNImputer
```

```
imputer = KNNImputer(n_neighbors=5)
df_imputed = pd.DataFrame(imputer.fit_transform(df), columns=df.columns)
```

**5. Interpolation**

- Useful for filling missing values in continuous data (e.g., time series).

**Code Example: Interpolation**

```
df['feature'] = df['feature'].interpolate(method='linear')
```

---

# 2.4 Putting It All Together

Let's consolidate all the steps into a streamlined workflow:

## Comprehensive Data Preprocessing Pipeline

```python
import pandas as pd
import numpy as np
from sklearn.preprocessing import StandardScaler, MinMaxScaler,
PowerTransformer, LabelEncoder
from sklearn.impute import KNNImputer
import category_encoders as ce

# Step 1: Load the Data
df = pd.read_csv('data.csv')

# Step 2: Handle Duplicates
df.drop_duplicates(inplace=True)

# Step 3: Handle Missing Values
df.dropna(subset=['important_column'], inplace=True)  # Drop rows where
crucial data is missing
df['numeric_column'] =
df['numeric_column'].fillna(df['numeric_column'].mean())
df['category_column'] =
df['category_column'].fillna(df['category_column'].mode()[0])

# Step 4: Encoding Categorical Variables
label_encoder = LabelEncoder()
df['label_encoded'] = label_encoder.fit_transform(df['category_column'])
df = pd.get_dummies(df, columns=['one_hot_column'], drop_first=True)
```

```python
# Step 5: Normalize/Standardize Features
scaler = StandardScaler()
df[['col1', 'col2']] = scaler.fit_transform(df[['col1', 'col2']])

# Step 6: Data Transformation
df['log_transformed'] = np.log1p(df['skewed_feature'])
power_transformer = PowerTransformer(method='yeo-johnson')
df['yeo_transformed'] = power_transformer.fit_transform(df[['feature']])

# Step 7: Final Check
print(df.head())
print(df.info())
```

To **fully prepare data for machine learning**, it's essential to go beyond just cleaning, encoding, and normalizing. The data preparation phase is one of the most crucial steps to ensure a model performs optimally. Below, I've outlined a **comprehensive list of techniques and steps** to help prepare your data for machine learning, with code examples and explanations.

# Step 1: Data Loading and Initial Exploration

```python
import pandas as pd

# Load your dataset
df = pd.read_csv("your_dataset.csv")

# Display first few rows
print(df.head())

# Summary statistics
print(df.describe())

# Check for missing values
print(df.isnull().sum())
```

# Step 2: Data Cleaning

## 2.1 Handling Missing Values

- **Strategy 1: Remove rows with missing values**

```python
df.dropna(inplace=True)
```

- **Strategy 2: Fill missing values with mean/median/mode**

```python
df['Age'].fillna(df['Age'].mean(), inplace=True)  # For numerical columns
df['City'].fillna(df['City'].mode()[0], inplace=True)  # For categorical columns
```

- **Strategy 3: Interpolation for time-series data**

```python
df['Temperature'] = df['Temperature'].interpolate(method='linear')
```

## 2.2 Removing Duplicates

```python
df.drop_duplicates(inplace=True)
```

## 2.3 Correcting Data Types

```python
# Convert data types if needed
df['Date'] = pd.to_datetime(df['Date'])
df['Price'] = df['Price'].astype(float)
```

## 2.4 Handling Outliers

```python
import numpy as np

# Remove outliers using the IQR method
Q1 = df['Age'].quantile(0.25)
Q3 = df['Age'].quantile(0.75)
IQR = Q3 - Q1
df = df[~((df['Age'] < (Q1 - 1.5 * IQR)) | (df['Age'] > (Q3 + 1.5 * IQR)))]
```

# Step 3: Feature Engineering

## 3.1 Encoding Categorical Variables

- **One-Hot Encoding**

```python
df = pd.get_dummies(df, columns=['Category'], drop_first=True)
```

- **Label Encoding**

```python
from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()
df['Category'] = le.fit_transform(df['Category'])
```

- **Ordinal Encoding (for ordinal data)**

```python
df['Education_Level'] = df['Education_Level'].map({'High School': 0,
'Bachelor': 1, 'Master': 2, 'PhD': 3})
```

---

## 3.2 Feature Scaling

- **Normalization (Min-Max Scaling)**

```python
from sklearn.preprocessing import MinMaxScaler
scaler = MinMaxScaler()
df[['Age', 'Income']] = scaler.fit_transform(df[['Age', 'Income']])
```

- **Standardization (Z-Score Scaling)**

```python
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
df[['Age', 'Income']] = scaler.fit_transform(df[['Age', 'Income']])
```

## 3.3 Data Transformation

- **Log Transformation (to handle skewed data)**

```
df['Income'] = np.log(df['Income'] + 1)
```

- **Box-Cox Transformation (for normalizing distributions)**

```
from scipy.stats import boxcox
df['Income'], _ = boxcox(df['Income'] + 1)
```

- **Square Root Transformation**

```
df['Income'] = np.sqrt(df['Income'])
```

# 3.4 Handling Missing Values in Categorical Data

- **Using `fillna` with a new category**

```
df['City'] = df['City'].fillna('Unknown')
```

- **Imputation with mode**

```
df['City'].fillna(df['City'].mode()[0], inplace=True)
```

---

# Step 4: Feature Selection

- **Correlation Matrix**

```
import seaborn as sns
import matplotlib.pyplot as plt

corr_matrix = df.corr()
sns.heatmap(corr_matrix, annot=True, cmap='coolwarm')
plt.show()
```

- **Removing Features with Low Variance**

```
from sklearn.feature_selection import VarianceThreshold
selector = VarianceThreshold(threshold=0.1)
```

```
    df = selector.fit_transform(df)
```

- **Feature Importance using Random Forest**

```
from sklearn.ensemble import RandomForestClassifier
model = RandomForestClassifier()
model.fit(X, y)
importances = model.feature_importances_
```

# Step 5: Train-Test Split

```
from sklearn.model_selection import train_test_split

X = df.drop('target', axis=1)
y = df['target']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)
```

# Step 6: Handling Imbalanced Datasets

- **Using SMOTE (Synthetic Minority Over-sampling Technique)**

```
from imblearn.over_sampling import SMOTE
smote = SMOTE()
X_resampled, y_resampled = smote.fit_resample(X_train, y_train)
```

- **Using Class Weights in Algorithms**

```
from sklearn.ensemble import RandomForestClassifier
model = RandomForestClassifier(class_weight='balanced')
```