# Complete Guide to Dataset Preprocessing for Machine Learning

## Contents

# 1 What is Data Preprocessing?

Data preprocessing involves preparing raw data for machine learning by cleaning, transforming, and organizing it. Machine learning models rely on numerical data, so preprocessing ensures the data is consistent, complete, and scaled appropriately. Poor preprocessing can lead to inaccurate models, while good preprocessing improves model performance.

Think of preprocessing as preparing ingredients before cooking. You need to wash, chop, and measure ingredients (data) to ensure the dish (model) turns out well.

# 2 Common Data Issues

Before diving into techniques, lets understand the common problems in datasets:

- **Missing Values**: Some data points are absent (e.g., empty cells in a CSV file).
- **Inconsistent Data**: Typos, mixed formats (e.g., "USA" vs. "United States"), or duplicates.
- **Categorical Data**: Non-numeric data (e.g., "Male"/"Female") that models cant directly process.
- **Outliers**: Extreme values that dont align with the rest of the data.
- **Unscaled Features**: Features with different ranges (e.g., age: 0-100, salary: 0-100000) can bias models.
- **Irrelevant Features**: Columns that dont contribute to the prediction.
- **Imbalanced Data**: When one class dominates (e.g., 90% "No" vs. 10% "Yes" in a dataset).

Preprocessing addresses these issues to make the data model-ready.

# 3 Preprocessing Techniques

## 3.1 Handling Missing Values

Missing values (e.g., `NaN` in Python) can cause errors in machine learning models. Here are common strategies to handle them:

### 3.1.1 Remove Rows/Columns with Missing Values

- **When to Use**: If only a small percentage of data is missing, or if a column has too many missing values to be useful.
- **Pros**: Simple and avoids introducing bias.
- **Cons**: Can lose valuable data if too many rows are removed.

### 3.1.2 Imputation

Replace missing values with a statistic like mean, median, or mode.

- **Mean/Median Imputation**: Use for numerical data (mean for normally distributed data, median for skewed data).
- **Mode Imputation**: Use for categorical data.
- **Pros**: Preserves data size.

- **Cons**: Can introduce bias if missing data is not random.

### 3.1.3 Using Algorithms to Predict Missing Values

Use machine learning (e.g., KNN) to predict missing values based on other features.

- **Pros**: More accurate than simple imputation.

- **Cons**: Computationally expensive.

**Code Example**:

```python
import pandas as pd
import numpy as np
from sklearn.impute import SimpleImputer

# Sample dataset with missing values
data = {
    'Age': [25, np.nan, 30, 35, np.nan],
    'Salary': [50000, 60000, np.nan, 80000, 70000],
    'Gender': ['Male', 'Female', np.nan, 'Male', 'Female']
}
df = pd.DataFrame(data)

# 1. Remove rows with missing values
df_dropped = df.dropna()
print("After dropping rows with missing values:\n", df_dropped)

# 2. Impute numerical columns with mean
num_imputer = SimpleImputer(strategy='mean')
df[['Age', 'Salary']] = num_imputer.fit_transform(df[['Age',
    'Salary']])
print("\nAfter mean imputation for Age and Salary:\n", df)

# 3. Impute categorical column with mode
cat_imputer = SimpleImputer(strategy='most_frequent')
df[['Gender']] = cat_imputer.fit_transform(df[['Gender']])
print("\nAfter mode imputation for Gender:\n", df)
```

**Explanation**:

- `SimpleImputer` from Scikit-learn replaces missing values.

- For `Age` and `Salary`, we use the mean to fill missing values.

- For `Gender`, we use the mode (most frequent value).

- Output shows how the dataset changes after each step.

## 3.2 Data Cleaning

Data cleaning involves fixing inconsistencies, duplicates, and errors in the data.

### 3.2.1 Remove Duplicates

Duplicate rows can bias the model by over-representing certain data points. Use Pandas `drop_duplicates()` to remove them.

### 3.2.2 Fix Inconsistent Formats

Standardize text (e.g., convert "USA" and "United States" to one format). Correct typos or case sensitivity (e.g., "male" vs. "Male").

### 3.2.3 Handle Incorrect Data

Replace invalid entries (e.g., negative age) with valid ones or remove them.

**Code Example**:

```python
# Sample dataset with duplicates and inconsistencies
data = {
    'Name': ['John', 'john', 'Alice', 'Bob', 'Bob'],
    'Age': [25, 25, 30, -5, 35],
    'Country': ['USA', 'United States', 'Canada', 'USA', 'USA']
}
df = pd.DataFrame(data)

# 1. Remove duplicates
df = df.drop_duplicates()
print("After removing duplicates:\n", df)

# 2. Standardize 'Country' column
df['Country'] = df['Country'].replace('United States', 'USA')
print("\nAfter standardizing Country:\n", df)

# 3. Fix invalid ages (e.g., negative values)
df.loc[df['Age'] < 0, 'Age'] = np.nan  # Replace negative ages with NaN
df['Age'] = df['Age'].fillna(df['Age'].mean())  # Impute with mean
print("\nAfter fixing invalid ages:\n", df)
```

**Explanation**:

- `drop_duplicates()` removes identical rows (e.g., two "Bob" entries).
- `replace()` standardizes "United States" to "USA".
- Negative ages are replaced with `NaN`, then imputed with the mean.

### 3.3 Encoding Categorical Variables

Machine learning models require numerical inputs, but datasets often contain categorical data (e.g., "Red", "Blue"). Encoding converts these into numbers.

### 3.3.1 Label Encoding

Assigns a unique integer to each category (e.g., "Male" → 0, "Female" → 1).

- **When to Use**: For ordinal data (where order matters, e.g., "Low", "Medium", "High").
- **Cons**: Can mislead models if used on non-ordinal data.

### 3.3.2 One-Hot Encoding

Creates binary columns for each category (e.g., "Color" with "Red", "Blue" becomes "Color_Red", "Color_Blue").

- **When to Use**: For nominal data (no order, e.g., colors, countries).

- **Cons**: Increases dataset size with many categories.

### 3.3.3  Target Encoding

Replaces categories with the mean of the target variable for that category.

- **When to Use**: For high-cardinality categorical features (many unique values).

**Code Example**:

```python
from sklearn.preprocessing import LabelEncoder, OneHotEncoder

# Sample dataset
data = {
    'Color': ['Red', 'Blue', 'Green', 'Red'],
    'Size': ['Small', 'Medium', 'Large', 'Medium']
}
df = pd.DataFrame(data)

# 1. Label Encoding for 'Size' (ordinal)
label_encoder = LabelEncoder()
df['Size_Encoded'] = label_encoder.fit_transform(df['Size'])
print("After label encoding Size:\n", df)

# 2. One-Hot Encoding for 'Color' (nominal)
one_hot_encoder = OneHotEncoder(sparse_output=False)
color_encoded = one_hot_encoder.fit_transform(df[['Color']])
color_df = pd.DataFrame(color_encoded,
    columns=one_hot_encoder.get_feature_names_out(['Color']))
df = pd.concat([df, color_df], axis=1)
print("\nAfter one-hot encoding Color:\n", df)
```

**Explanation**:

- `LabelEncoder` assigns numbers to "Small", "Medium", "Large" (e.g., 0, 1, 2).

- `OneHotEncoder` creates binary columns for "Red", "Blue", "Green".

- The resulting DataFrame includes both encoded features.

## 3.4  Feature Scaling

Features with different scales (e.g., age: 0-100, salary: 0-100000) can bias models, as larger values may dominate. Scaling standardizes feature ranges.

### 3.4.1  Standardization (Z-Score Normalization)

Scales features to have a mean of 0 and a standard deviation of 1.

- **Formula**: $z = \frac{x - \text{mean}}{\text{std}}$

- **When to Use**: For algorithms like SVM, logistic regression, or neural networks.

### 3.4.2  Min-Max Scaling

Scales features to a fixed range, typically [0, 1].

- **Formula**: $x' = \frac{x - \min}{\max - \min}$

- **When to Use**: For algorithms sensitive to feature ranges, like k-nearest neighbors.

**Code Example**:

```python
from sklearn.preprocessing import StandardScaler, MinMaxScaler

# Sample dataset
data = {
    'Age': [25, 30, 35, 40, 45],
    'Salary': [50000, 60000, 75000, 80000, 90000]
}
df = pd.DataFrame(data)

# 1. Standardization
std_scaler = StandardScaler()
df[['Age_Std', 'Salary_Std']] = std_scaler.fit_transform(df[['Age',
    'Salary']])
print("After standardization:\n", df)

# 2. Min-Max Scaling
minmax_scaler = MinMaxScaler()
df[['Age_MinMax', 'Salary_MinMax']] =
    minmax_scaler.fit_transform(df[['Age', 'Salary']])
print("\nAfter Min-Max scaling:\n", df)
```

**Explanation**:

- `StandardScaler` transforms `Age` and `Salary` to have mean=0 and std=1.

- `MinMaxScaler` scales them to [0, 1].

- The output shows both scaled versions for comparison.

### 3.5   Handling Outliers

Outliers are extreme values that can skew model performance. Common methods to handle them:

#### 3.5.1   Remove Outliers

Use statistical methods like the Interquartile Range (IQR) to detect and remove outliers.

- **IQR Method**: Outliers are values below $Q1 - 1.5 \cdot \text{IQR}$ or above $Q3 + 1.5 \cdot \text{IQR}$.

#### 3.5.2   Cap Outliers

Replace outliers with the nearest acceptable value (e.g., cap at the 95th percentile).

#### 3.5.3   Transform Features

Apply transformations like log or square root to reduce the impact of outliers.

**Code Example**:

```python
# Sample dataset with outliers
data = {
    'Salary': [50000, 60000, 75000, 1000000, 80000]
}
df = pd.DataFrame(data)

# 1. Detect and remove outliers using IQR
```

```
8  Q1 = df['Salary'].quantile(0.25)
9  Q3 = df['Salary'].quantile(0.75)
10 IQR = Q3 - Q1
11 lower_bound = Q1 - 1.5 * IQR
12 upper_bound = Q3 + 1.5 * IQR
13 df_no_outliers = df[(df['Salary'] >= lower_bound) & (df['Salary'] <=
       upper_bound)]
14 print("After removing outliers:\n", df_no_outliers)
15
16 # 2. Cap outliers
17 df['Salary_Capped'] = df['Salary'].clip(lower=lower_bound,
       upper=upper_bound)
18 print("\nAfter capping outliers:\n", df)
```

**Explanation**:

- The IQR method identifies `1000000` as an outlier and removes it.

- `clip()` caps extreme values to the upper/lower bounds.

### 3.6   Feature Selection

Feature selection removes irrelevant or redundant features to improve model performance and reduce training time.

#### 3.6.1   Filter Methods

Select features based on statistical measures (e.g., correlation with the target).

- **Example**: Remove features with low variance or high correlation.

#### 3.6.2   Wrapper Methods

Evaluate subsets of features using a model (e.g., recursive feature elimination).

#### 3.6.3   Embedded Methods

Use algorithms that inherently perform feature selection (e.g., Lasso regression).

**Code Example**:

```
1  from sklearn.feature_selection import VarianceThreshold
2
3  # Sample dataset
4  data = {
5      'Feature1': [1, 1, 1, 1],  # Low variance
6      'Feature2': [2, 4, 6, 8],  # High variance
7      'Feature3': [1, 2, 3, 4]
8  }
9  df = pd.DataFrame(data)
10
11 # Remove low-variance features
12 selector = VarianceThreshold(threshold=0.1)
13 df_selected = selector.fit_transform(df)
14 print("Selected features (after removing low variance):\n",
       pd.DataFrame(df_selected,
       columns=df.columns[selector.get_support()]))
```

**Explanation**:

- `VarianceThreshold` removes `Feature1` because it has low variance (all values are 1).

### 3.7 Feature Engineering

Feature engineering creates new features or transforms existing ones to improve model performance.

#### 3.7.1 Creating New Features

Combine features (e.g., create "Age_Salary_Ratio" from `Age` and `Salary`). Extract features (e.g., extract "Day" from a date column).

#### 3.7.2 Transforming Features

Apply mathematical transformations (e.g., log, square). Bin continuous features into categories (e.g., age groups).

**Code Example**:

```python
# Sample dataset
data = {
    'Age': [25, 30, 35, 40],
    'Salary': [50000, 60000, 75000, 80000],
    'Date': ['2023-01-01', '2023-02-01', '2023-03-01', '2023-04-01']
}
df = pd.DataFrame(data)

# 1. Create new feature: Age_Salary_Ratio
df['Age_Salary_Ratio'] = df['Age'] / df['Salary']

# 2. Extract month from Date
df['Date'] = pd.to_datetime(df['Date'])
df['Month'] = df['Date'].dt.month

# 3. Bin Age into categories
df['Age_Group'] = pd.cut(df['Age'], bins=[0, 30, 40, 100],
    labels=['Young', 'Middle', 'Senior'])
print("After feature engineering:\n", df)
```

**Explanation**:

- `Age_Salary_Ratio` is a new feature combining `Age` and `Salary`.

- `dt.month` extracts the month from the `Date` column.

- `pd.cut` bins `Age` into categories like "Young" and "Middle".

### 3.8 Data Splitting

Before training a model, split the dataset into training, validation, and test sets to evaluate performance.

- **Training Set**: Used to train the model (e.g., 70% of data).

- **Validation Set**: Used to tune hyperparameters (e.g., 15% of data).

- **Test Set**: Used to evaluate final performance (e.g., 15% of data).

**Code Example**:

```python
from sklearn.model_selection import train_test_split

# Sample dataset
data = {
    'Feature1': [1, 2, 3, 4, 5, 6, 7, 8, 9, 10],
    'Target': [0, 1, 0, 1, 0, 1, 0, 1, 0, 1]
}
df = pd.DataFrame(data)

# Split data
X = df[['Feature1']]
y = df['Target']
X_train, X_temp, y_train, y_temp = train_test_split(X, y,
    test_size=0.3, random_state=42)
X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp,
    test_size=0.5, random_state=42)

print("Training set size:", len(X_train))
print("Validation set size:", len(X_val))
print("Test set size:", len(X_test))
```

**Explanation**:

- `train_test_split` splits the data into 70% training and 30% temporary sets.

- The temporary set is further split into 15% validation and 15% test sets.

## 4    Putting It All Together: A Complete Example

Heres a complete preprocessing pipeline combining all techniques.

**Code Example**:

```python
import pandas as pd
import numpy as np
from sklearn.impute import SimpleImputer
from sklearn.preprocessing import StandardScaler, OneHotEncoder
from sklearn.model_selection import train_test_split

# Sample dataset
data = {
    'Age': [25, np.nan, 30, 35, 1000, 40],
    'Salary': [50000, 60000, np.nan, 80000, 90000, 85000],
    'Gender': ['Male', 'Female', np.nan, 'Male', 'Female', 'Male'],
    'Country': ['USA', 'Canada', 'USA', 'United States', 'Canada',
        'USA'],
    'Target': [0, 1, 0, 1, 0, 1]
}
df = pd.DataFrame(data)

# 1. Handle missing values
num_imputer = SimpleImputer(strategy='mean')
df[['Age', 'Salary']] = num_imputer.fit_transform(df[['Age',
    'Salary']])
cat_imputer = SimpleImputer(strategy='most_frequent')
df[['Gender']] = cat_imputer.fit_transform(df[['Gender']])

```

```python
23  # 2. Clean data (standardize Country, handle outliers)
24  df['Country'] = df['Country'].replace('United States', 'USA')
25  Q1 = df['Age'].quantile(0.25)
26  Q3 = df['Age'].quantile(0.75)
27  IQR = Q3 - Q1
28  df['Age'] = df['Age'].clip(lower=Q1 - 1.5 * IQR, upper=Q3 + 1.5 * IQR)
29
30  # 3. Encode categorical variables
31  df = pd.get_dummies(df, columns=['Gender', 'Country'], drop_first=True)
32
33  # 4. Scale features
34  scaler = StandardScaler()
35  df[['Age', 'Salary']] = scaler.fit_transform(df[['Age', 'Salary']])
36
37  # 5. Split data
38  X = df.drop('Target', axis=1)
39  y = df['Target']
40  X_train, X_test, y_train, y_test = train_test_split(X, y,
        test_size=0.2, random_state=42)
41
42  print("Preprocessed training data:\n", X_train)
43  print("\nPreprocessed test data:\n", X_test)
```

**Explanation**:

- The pipeline handles missing values, cleans data, encodes categories, scales features, and splits the data.

- The final output shows the preprocessed training and test sets, ready for model training.

## 5  Best Practices and Tips

1. **Understand Your Data**:

   - Explore the dataset using `df.describe()`, `df.info()`, and visualizations to identify issues.

   - Check for missing values, outliers, and data types.

2. **Avoid Data Leakage**:

   - Apply preprocessing (e.g., scaling, imputation) on the training set first, then apply the same transformation to the test set. Never fit preprocessors on the test set.

3. **Choose Appropriate Techniques**:

   - Use mean imputation for normally distributed data, median for skewed data.

   - Use one-hot encoding for nominal data, label encoding for ordinal data.

4. **Automate Preprocessing**:

   - Use Scikit-learns `Pipeline` to combine preprocessing steps and ensure consistency.

5. **Document Your Steps**:

   - Keep track of preprocessing decisions to reproduce results or debug issues.

6. **Handle Imbalanced Data** (if applicable):

   - Use techniques like oversampling (SMOTE) or undersampling for imbalanced datasets.

**Code Example (Pipeline)**:

```python
from sklearn.pipeline import Pipeline
from sklearn.compose import ColumnTransformer

# Define numeric and categorical columns
numeric_features = ['Age', 'Salary']
categorical_features = ['Gender', 'Country']

# Create preprocessor
preprocessor = ColumnTransformer(
    transformers=[
        ('num', Pipeline([
            ('imputer', SimpleImputer(strategy='mean')),
            ('scaler', StandardScaler())
        ]), numeric_features),
        ('cat', Pipeline([
            ('imputer', SimpleImputer(strategy='most_frequent')),
            ('onehot', OneHotEncoder(drop='first'))
        ]), categorical_features)
    ])

# Apply preprocessor
X_transformed = preprocessor.fit_transform(df.drop('Target', axis=1))
print("Preprocessed data using pipeline:\n",
    pd.DataFrame(X_transformed))
```

**Explanation**:

- `ColumnTransformer` applies different preprocessing steps to numeric and categorical columns.

- `Pipeline` ensures consistent application of imputation and scaling.

This guide covers all essential preprocessing techniques with detailed explanations and code. By following these steps, you can prepare any dataset for machine learning. Save this guide, and youll have a complete reference for preprocessing!