

Exploratory Data Analysis: Data Preprocessing Methods

Akhilesh Joshi [02FE22BCS013]

March 6, 2025

Abstract

This document provides a comprehensive guide to data preprocessing techniques essential for exploratory data analysis. Data preprocessing is a critical step that transforms raw data into a clean, structured format suitable for analysis and modeling. The guide covers methods for handling missing values, outlier detection and treatment, feature scaling and transformation, encoding categorical variables, feature selection, dimensionality reduction, and data balancing. Each section includes practical Python code examples using popular libraries such as pandas, scikit-learn, and numpy.

Contents

1	Introduction to Data Preprocessing	2
2	Data Exploration and Profiling	2
3	Handling Missing Values	4
4	Outlier Detection and Treatment	5
5	Feature Scaling and Transformation	7
6	Encoding Categorical Variables	9
7	Feature Selection	11
8	Dimensionality Reduction	14
8.1	Inconsistent Formatting	16
8.2	Irrelevant Features	16
8.3	High Dimensionality	17
8.4	Data Imbalance	17
8.5	Need for Feature Engineering	17

1 Introduction to Data Preprocessing

Data preprocessing is a crucial step in the data science pipeline that transforms raw data into a clean, structured format suitable for analysis and modeling. Properly preprocessed data can significantly improve the performance of machine learning models and the reliability of analytical insights.

Python Code

```
# Common libraries for data preprocessing
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn import preprocessing
from sklearn.impute import SimpleImputer
from sklearn.decomposition import PCA
from sklearn.feature_selection import SelectKBest, f_classif
from sklearn.model_selection import train_test_split
from sklearn.datasets import load_breast_cancer

# Set visualization style
sns.set(style="whitegrid")
plt.rcParams['figure.figsize'] = (10, 6)

# For reproducibility
np.random.seed(42)

# Load sample dataset for examples
data = load_breast_cancer()
df = pd.DataFrame(data.data, columns=data.feature_names)
df['target'] = data.target

# Quick look at the data
print(f"Dataset shape: {df.shape}")
print("\nFirst 5 rows:")
print(df.head())
```

2 Data Exploration and Profiling

Before applying preprocessing techniques, it's essential to understand the dataset's characteristics through exploratory data analysis.

Python Code

```
# Basic dataset information
def explore_dataset(df):
    """Generate basic exploratory statistics for a dataframe"""
```

```
print(f"Dataset shape: {df.shape} (rows, columns)")
print("\nData types:")
print(df.dtypes.value_counts())

print("\nMissing values per column:")
missing = df.isnull().sum()
print(missing[missing > 0] if missing.any() > 0 else "No
missing values")

print("\nNumeric column statistics:")
print(df.describe().T)

# Check for duplicate rows
dup_count = df.duplicated().sum()
print(f"\nDuplicate rows: {dup_count}")

# For categorical columns
cat_cols = df.select_dtypes(include=['object', 'category']).
columns
if len(cat_cols) > 0:
    print("\nCategorical column statistics:")
    for col in cat_cols:
        print(f"\n{col}:")
        print(df[col].value_counts().head(5))

# Create a dataset with some issues for demonstration
demo_df = pd.DataFrame({
    'age': [25, 30, np.nan, 40, 35, 30],
    'income': [50000, 65000, 75000, np.nan, 55000, 65000],
    'gender': ['M', 'F', 'M', 'F', 'M', 'F'],
    'education': ['Bachelor', 'Master', np.nan, 'PhD', 'Bachelor',
    'Master']
})

# Explore the dataset
explore_dataset(demo_df)

# Visual exploration
fig, axes = plt.subplots(1, 2, figsize=(14, 5))

# Distribution of numeric variables
demo_df[['age', 'income']].boxplot(ax=axes[0])
axes[0].set_title('Boxplot of Numeric Variables')

# Count of categorical variables
demo_df['gender'].value_counts().plot(kind='bar', ax=axes[1])
axes[1].set_title('Gender Distribution')
plt.tight_layout()
```

3 Handling Missing Values

Missing values can significantly impact analysis and model performance. This section covers various techniques to detect and handle missing data.

Python Code

```
# Create a dataset with missing values
missing_df = pd.DataFrame({
    'A': [1, 2, np.nan, 4, 5],
    'B': [np.nan, 2, 3, 4, 5],
    'C': [1, 2, 3, np.nan, np.nan],
    'D': [1, np.nan, np.nan, np.nan, 5]
})

# 1. Visualize missing values
plt.figure(figsize=(10, 6))
sns.heatmap(missing_df.isnull(), cbar=False, cmap='viridis')
plt.title('Missing Value Heatmap')
plt.tight_layout()

# 2. Calculate missing value statistics
missing_count = missing_df.isnull().sum()
missing_percent = missing_df.isnull().mean() * 100
missing_stats = pd.DataFrame({
    'Count': missing_count,
    'Percent': missing_percent
})
print("Missing value statistics:")
print(missing_stats)

# 3. Handle missing values

# a) Remove rows with missing values
removed_rows = missing_df.dropna()
print("\nAfter removing rows with any missing values:")
print(f"Original shape: {missing_df.shape}, New shape: {removed_rows.shape}")

# b) Remove columns with too many missing values (e.g., >50%)
removed_cols = missing_df.dropna(axis=1, thresh=len(missing_df) * 0.5)
print("\nAfter removing columns with >50% missing values:")
print(f"Original columns: {list(missing_df.columns)}")
print(f"Remaining columns: {list(removed_cols.columns)}")

# c) Fill missing values with a constant
filled_constant = missing_df.fillna(0)
print("\nAfter filling with constant (0):")
print(filled_constant)
```

```
# d) Fill missing values with column mean/median/mode
filled_mean = missing_df.fillna(missing_df.mean())
print("\nAfter filling with column means:")
print(filled_mean)

# e) Fill values using forward fill (propagate last valid value)
filled_ffill = missing_df.fillna(method='ffill')
print("\nAfter forward fill:")
print(filled_ffill)

# f) Using scikit-learn's SimpleImputer
from sklearn.impute import SimpleImputer

# Impute with mean
mean_imputer = SimpleImputer(strategy='mean')
imputed_mean = pd.DataFrame(
    mean_imputer.fit_transform(missing_df),
    columns=missing_df.columns
)
print("\nAfter imputing with SimpleImputer (mean):")
print(imputed_mean)

# g) Use KNN imputation for more complex datasets
from sklearn.impute import KNNImputer

knn_imputer = KNNImputer(n_neighbors=2)
imputed_knn = pd.DataFrame(
    knn_imputer.fit_transform(missing_df),
    columns=missing_df.columns
)
print("\nAfter KNN imputation:")
print(imputed_knn)
```

4 Outlier Detection and Treatment

Outliers are extreme values that deviate significantly from other observations. They can distort statistical analyses and model performance.

Python Code

```
# Create a dataset with outliers
np.random.seed(42)
normal_data = np.random.normal(0, 1, 100)
outliers = np.array([10, -10, 8, -8])
data_with_outliers = np.concatenate([normal_data, outliers])

# Convert to DataFrame for easier analysis
```

```
outlier_df = pd.DataFrame({'value': data_with_outliers})

# 1. Visualize the data with outliers
plt.figure(figsize=(12, 5))

plt.subplot(1, 2, 1)
sns.histplot(outlier_df['value'], kde=True)
plt.title('Histogram of Data with Outliers')

plt.subplot(1, 2, 2)
sns.boxplot(x=outlier_df['value'])
plt.title('Boxplot of Data with Outliers')
plt.tight_layout()

# 2. Detect outliers using different methods

# a) Z-score method (for normally distributed data)
from scipy import stats

z_scores = stats.zscore(outlier_df)
z_outliers = outlier_df[abs(z_scores) > 3] # Values more than 3
standard deviations away
print(f"Z-score detected outliers: {len(z_outliers)} values")
print(z_outliers)

# b) IQR (Interquartile Range) method
Q1 = outlier_df['value'].quantile(0.25)
Q3 = outlier_df['value'].quantile(0.75)
IQR = Q3 - Q1
lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR

iqr_outliers = outlier_df[(outlier_df['value'] < lower_bound) |
                           (outlier_df['value'] > upper_bound)]
print(f"\nIQR detected outliers: {len(iqr_outliers)} values")
print(iqr_outliers)

# 3. Handle outliers

# a) Remove outliers
no_outliers_df = outlier_df[(outlier_df['value'] >= lower_bound) &
                             (outlier_df['value'] <= upper_bound)]
print(f"\nOriginal data: {len(outlier_df)} rows")
print(f"After removing outliers: {len(no_outliers_df)} rows")

# b) Cap outliers (Winsorization)
capped_df = outlier_df.copy()
capped_df['value'] = np.where(
    capped_df['value'] > upper_bound,
```

```
        upper_bound,
        capped_df['value']
    )
    capped_df['value'] = np.where(
        capped_df['value'] < lower_bound,
        lower_bound,
        capped_df['value']
    )

# c) Log transformation to handle skewed data with outliers
# Only for positive values
positive_df = outlier_df[outlier_df['value'] > 0].copy()
positive_df['log_value'] = np.log1p(positive_df['value']) # log
(1+x)

# Visualize effect of log transformation
if len(positive_df) > 0:
    plt.figure(figsize=(12, 5))
    plt.subplot(1, 2, 1)
    sns.histplot(positive_df['value'], kde=True)
    plt.title('Original Values')

    plt.subplot(1, 2, 2)
    sns.histplot(positive_df['log_value'], kde=True)
    plt.title('Log Transformed Values')
    plt.tight_layout()

# d) Use robust models less affected by outliers
# Example with RANSAC regressor (demo only - not executed)
from sklearn.linear_model import RANSACRegressor
# ransac = RANSACRegressor()
# ransac.fit(X, y) # X is features, y is target
```

5 Feature Scaling and Transformation

Feature scaling is essential when features have different ranges. It ensures that no feature dominates the others due to its scale.

Python Code

```
# Create a dataset with features of different scales
scaling_df = pd.DataFrame({
    'income': np.random.randint(30000, 150000, 100), # Higher
    range
    'age': np.random.randint(20, 65, 100), # Medium range
    'years_experience': np.random.randint(0, 20, 100) # Lower
    range
})
```

```
# Show original data statistics
print("Original data statistics:")
print(scaling_df.describe().T[['min', 'max', 'mean', 'std']])

# 1. Min-Max Scaling (Normalization) - scales to range [0,1]
from sklearn.preprocessing import MinMaxScaler

min_max_scaler = MinMaxScaler()
scaled_minmax = pd.DataFrame(
    min_max_scaler.fit_transform(scaling_df),
    columns=scaling_df.columns
)

print("\nMin-Max scaled data statistics:")
print(scaled_minmax.describe().T[['min', 'max', 'mean', 'std']])

# 2. Standardization (Z-score) - scales to mean=0, std=1
from sklearn.preprocessing import StandardScaler

std_scaler = StandardScaler()
scaled_std = pd.DataFrame(
    std_scaler.fit_transform(scaling_df),
    columns=scaling_df.columns
)

print("\nStandardized data statistics:")
print(scaled_std.describe().T[['min', 'max', 'mean', 'std']])

# 3. Robust Scaling - uses median and IQR (less affected by outliers)
from sklearn.preprocessing import RobustScaler

robust_scaler = RobustScaler()
scaled_robust = pd.DataFrame(
    robust_scaler.fit_transform(scaling_df),
    columns=scaling_df.columns
)

print("\nRobust scaled data statistics:")
print(scaled_robust.describe().T[['min', 'max', 'mean', 'std']])

# 4. Log Transformation - handle skewed data
# Add a column with skewed distribution
scaling_df['skewed_feature'] = np.random.exponential(scale=2, size=100)

# Apply log transformation
scaling_df['log_skewed'] = np.log1p(scaling_df['skewed_feature'])
```



```
# Visualize effect
plt.figure(figsize=(12, 5))
plt.subplot(1, 2, 1)
sns.histplot(scaling_df['skewed_feature'], kde=True)
plt.title('Original Skewed Distribution')

plt.subplot(1, 2, 2)
sns.histplot(scaling_df['log_skewed'], kde=True)
plt.title('Log Transformed Distribution')
plt.tight_layout()

# 5. Box-Cox Transformation - another way to handle skewed data
from scipy import stats

# Box-Cox requires positive values
positive_vals = scaling_df['skewed_feature'].copy()
boxcox_vals, lambda_val = stats.boxcox(positive_vals)

print(f"\nBox-Cox transformation lambda value: {lambda_val:.4f}")
plt.figure(figsize=(8, 5))
sns.histplot(boxcox_vals, kde=True)
plt.title('Box-Cox Transformed Distribution')
plt.tight_layout()

# 6. Power Transformation - broader family of transformations
from sklearn.preprocessing import PowerTransformer

power_transformer = PowerTransformer(method='yeo-johnson')
power_transformed = power_transformer.fit_transform(scaling_df[['skewed_feature']])

plt.figure(figsize=(8, 5))
sns.histplot(power_transformed, kde=True)
plt.title('Yeo-Johnson Power Transformed Distribution')
plt.tight_layout()
```

6 Encoding Categorical Variables

Categorical variables need to be converted to numeric form for most machine learning algorithms.

Python Code

```
# Create a dataset with categorical variables
cat_df = pd.DataFrame({
    'city': ['New York', 'London', 'Paris', 'Tokyo', 'New York', 'Paris'],
```

```
'color': ['red', 'blue', 'green', 'blue', 'red', 'green'],
'size': ['small', 'medium', 'large', 'medium', 'large', 'small'],
'rating': [4, 5, 3, 5, 2, 3]
})

print("Original categorical data:")
print(cat_df.head())

# 1. Label Encoding - convert categories to numeric values (0, 1, 2, ...)
from sklearn.preprocessing import LabelEncoder

le = LabelEncoder()
cat_df_le = cat_df.copy()
cat_df_le['city_encoded'] = le.fit_transform(cat_df_le['city'])
cat_df_le['color_encoded'] = le.fit_transform(cat_df_le['color'])
cat_df_le['size_encoded'] = le.fit_transform(cat_df_le['size'])

print("\nAfter Label Encoding:")
print(cat_df_le.head())
print("\nLabel mappings for size:")
for i, category in enumerate(le.classes_):
    print(f"{category} -> {i}")

# 2. One-Hot Encoding - create binary columns for each category
# Using pandas get_dummies
onehot_df = pd.get_dummies(cat_df, columns=['city', 'color', 'size'],
    prefix=['city', 'color', 'size'])
print("\nAfter One-Hot Encoding (with pandas):")
print(onehot_df.head())

# Using scikit-learn
from sklearn.preprocessing import OneHotEncoder

ohe = OneHotEncoder(sparse_output=False)
encoded_cols = ohe.fit_transform(cat_df[['city', 'color']])
encoded_df = pd.DataFrame(
    encoded_cols,
    columns=ohe.get_feature_names_out(['city', 'color'])
)

# Combine with original dataframe
cat_df_ohe = pd.concat([cat_df.reset_index(drop=True), encoded_df.
    reset_index(drop=True)], axis=1)
print("\nAfter One-Hot Encoding (with scikit-learn):")
print(cat_df_ohe.head())

# 3. Ordinal Encoding - for categorical variables with clear
```

```
ordering
from sklearn.preprocessing import OrdinalEncoder

# Define the ordering of categories
size_ordering = ['small', 'medium', 'large'] # ordered by size
ordinal_encoder = OrdinalEncoder(categories=[size_ordering])
cat_df['size_ordinal'] = ordinal_encoder.fit_transform(cat_df[['size']])

print("\nAfter Ordinal Encoding for 'size':")
print(cat_df[['size', 'size_ordinal']].head())

# 4. Target Encoding - replace categories with target mean
def target_encode(df, col, target):
    """Simple target encoding implementation"""
    mapping = df.groupby(col)[target].mean().to_dict()
    return df[col].map(mapping)

cat_df['city_target_encoded'] = target_encode(cat_df, 'city', 'rating')
print("\nAfter Target Encoding for 'city' based on 'rating':")
print(cat_df[['city', 'rating', 'city_target_encoded']].head())

# 5. Count/Frequency Encoding
count_map = cat_df['city'].value_counts().to_dict()
cat_df['city_count'] = cat_df['city'].map(count_map)
print("\nAfter Count Encoding for 'city':")
print(cat_df[['city', 'city_count']].head())
```

7 Feature Selection

Feature selection helps identify the most important features, improving model performance and reducing complexity.

Python Code

```
# Load a dataset with many features
from sklearn.datasets import load_breast_cancer
data = load_breast_cancer()
X = pd.DataFrame(data.data, columns=data.feature_names)
y = data.target

print(f"Original dataset shape: {X.shape}")

# 1. Filter Methods

# a) Variance Threshold - remove features with low variance
from sklearn.feature_selection import VarianceThreshold
```

```
var_threshold = VarianceThreshold(threshold=0.1) # features with
variance <= 0.1 will be removed
X_var = var_threshold.fit_transform(X)
print(f"\nAfter Variance Threshold: {X_var.shape[1]} features
retained")

# b) Univariate Feature Selection - ANOVA F-test
from sklearn.feature_selection import SelectKBest, f_classif

k_best = SelectKBest(score_func=f_classif, k=10) # select top 10
features
X_kbest = k_best.fit_transform(X, y)
selected_features = X.columns[k_best.get_support()]

print("\nTop 10 features based on ANOVA F-test:")
feature_scores = pd.DataFrame({
    'Feature': X.columns,
    'Score': k_best.scores_
})
print(feature_scores.sort_values('Score', ascending=False).head
(10))

# c) Correlation-based Feature Selection
# Calculate correlation with the target
X_with_target = X.copy()
X_with_target['target'] = y
correlation_with_target = X_with_target.corr()['target'].
sort_values(ascending=False)

print("\nFeatures with highest correlation with target:")
print(correlation_with_target.head(5))

# Correlation matrix for feature-feature relationships
plt.figure(figsize=(12, 10))
corr_matrix = X.corr().abs()
sns.heatmap(corr_matrix, cmap='viridis', annot=False)
plt.title('Feature Correlation Matrix')
plt.tight_layout()

# Find highly correlated feature pairs
high_corr_pairs = []
for i in range(len(corr_matrix.columns)):
    for j in range(i):
        if abs(corr_matrix.iloc[i, j]) > 0.8: # Correlation
threshold
            high_corr_pairs.append((corr_matrix.columns[i],
corr_matrix.columns[j],
corr_matrix.iloc[i, j]))
```

```
print("\nHighly correlated feature pairs:")
for pair in sorted(high_corr_pairs, key=lambda x: abs(x[2]),
    reverse=True)[:5]:
    print(f"{pair[0]} and {pair[1]}: {pair[2]:.3f}")

# 2. Wrapper Methods

# a) Recursive Feature Elimination (RFE)
from sklearn.feature_selection import RFE
from sklearn.linear_model import LogisticRegression

model = LogisticRegression(max_iter=1000)
rfe = RFE(estimator=model, n_features_to_select=10)
X_rfe = rfe.fit_transform(X, y)

print(f"\nRFE selected {X_rfe.shape[1]} features:")
rfe_features = X.columns[rfe.support_]
print(rfe_features.tolist())

# 3. Embedded Methods

# a) L1 Regularization (Lasso)
from sklearn.linear_model import Lasso
from sklearn.preprocessing import StandardScaler

# Standardize for Lasso
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

# Fit Lasso model
lasso = Lasso(alpha=0.01)
lasso.fit(X_scaled, y)

# Get feature importance
lasso_importance = pd.DataFrame({
    'Feature': X.columns,
    'Coefficient': lasso.coef_
})
lasso_importance = lasso_importance.sort_values('Coefficient', key
    =abs, ascending=False)

print("\nTop features based on Lasso coefficients:")
print(lasso_importance.head(10))

# b) Tree-based feature importance
from sklearn.ensemble import RandomForestClassifier

rf = RandomForestClassifier(n_estimators=100, random_state=42)
```

```
rf.fit(X, y)

# Get feature importance
rf_importance = pd.DataFrame({
    'Feature': X.columns,
    'Importance': rf.feature_importances_
})
rf_importance = rf_importance.sort_values('Importance', ascending=False)

print("\nTop features based on Random Forest importance:")
print(rf_importance.head(10))

plt.figure(figsize=(10, 6))
sns.barplot(x='Importance', y='Feature', data=rf_importance.head(10))
plt.title('Top 10 Features by Random Forest Importance')
plt.tight_layout()
```

8 Dimensionality Reduction

High-dimensional data can lead to overfitting and computational inefficiency. Dimensionality reduction techniques help address these issues.

Python Code

```
# Load a high-dimensional dataset
from sklearn.datasets import fetch_olivetti_faces
faces = fetch_olivetti_faces()
X_faces = faces.data
y_faces = faces.target

print(f"Original data shape: {X_faces.shape}")

# 1. Principal Component Analysis (PCA)
from sklearn.decomposition import PCA

# Apply PCA
pca = PCA(n_components=0.95) # Retain 95% of variance
X_pca = pca.fit_transform(X_faces)

print(f"\nAfter PCA: {X_pca.shape[1]} components retain 95% variance")
print(f"Explained variance ratio: {pca.explained_variance_ratio_.sum():.4f}")

# Visualize cumulative explained variance
plt.figure(figsize=(10, 6))
```

```
plt.plot(np.cumsum(pca.explained_variance_ratio_))
plt.axhline(y=0.95, color='r', linestyle='--')
plt.xlabel('Number of Components')
plt.ylabel('Cumulative Explained Variance')
plt.title('Explained Variance vs. Number of Components')
plt.grid(True)
plt.tight_layout()

# Visualize the first few principal components
if X_faces.shape[1] <= 4096: # Only for face images
    n_row, n_col = 2, 3
    plt.figure(figsize=(2*n_col, 2*n_row))
    for i in range(n_row * n_col):
        if i < len(pca.components_):
            plt.subplot(n_row, n_col, i+1)
            plt.imshow(pca.components_[i].reshape(64, 64), cmap='
gray')
            plt.title(f'PC {i+1}')
            plt.xticks([])
            plt.yticks([])
    plt.tight_layout()

# 2. t-SNE (t-distributed Stochastic Neighbor Embedding)
from sklearn.manifold import TSNE

# Apply t-SNE
tsne = TSNE(n_components=2, random_state=42)
X_tsne = tsne.fit_transform(X_faces)

# Visualize t-SNE results
plt.figure(figsize=(10, 8))
scatter = plt.scatter(X_tsne[:, 0], X_tsne[:, 1], c=y_faces, cmap=
'viridis')
plt.colorbar(scatter, label='Person ID')
plt.title('t-SNE Visualization of Face Images')
plt.xlabel('t-SNE Feature 1')
plt.ylabel('t-SNE Feature 2')
plt.tight_layout()

# 3. UMAP (Uniform Manifold Approximation and Projection)
# Install umap-learn package if not available
try:
    import umap

    # Apply UMAP
    umap_reducer = umap.UMAP(random_state=42)
    X_umap = umap_reducer.fit_transform(X_faces)

    # Visualize UMAP results
```

```
plt.figure(figsize=(10, 8))
scatter = plt.scatter(X_umap[:, 0], X_umap[:, 1], c=y_faces,
                      cmap='viridis')
plt.colorbar(scatter, label='Person ID')
plt.title('UMAP Visualization of Face Images')
plt.xlabel('UMAP Feature 1')
plt.ylabel('UMAP Feature 2')
```

8.1 Inconsistent Formatting

Data collected from various sources often suffers from inconsistent formatting. This includes:

- Different date formats (MM/DD/YYYY vs. DD-MM-YYYY)
- Inconsistent capitalization (Yes/No vs. yes/no)
- Variable spacing and special characters
- Mixed numerical representations (percentages, decimals, scientific notation)

Addressing inconsistent formatting typically requires standardization through string manipulation, regular expressions, and careful validation procedures before analysis can proceed.

8.2 Irrelevant Features

Not all available features contribute meaningfully to the predictive power of a model. Irrelevant features:

- Increase computational complexity
- May introduce noise
- Can lead to overfitting
- Reduce model interpretability

Feature selection techniques such as statistical tests (chi-square, ANOVA), correlation analysis, and wrapper methods help identify and remove irrelevant features, improving model performance and efficiency.

8.3 High Dimensionality

Modern datasets often contain hundreds or thousands of features, creating the "curse of dimensionality":

- Data becomes increasingly sparse in high-dimensional spaces
- Distance metrics become less meaningful
- Computational complexity increases exponentially
- Risk of overfitting grows significantly

Dimension reduction techniques like Principal Component Analysis (PCA), t-SNE, or autoencoders help address this challenge by projecting data into lower-dimensional spaces while preserving important information.

8.4 Data Imbalance

Class imbalance occurs when target classes are not equally represented in the dataset:

- Models tend to favor majority classes
- Performance metrics can be misleading (high accuracy despite poor minority class prediction)
- Rare but important events may be overlooked

Balancing techniques include:

1. Resampling: oversampling minority classes (SMOTE, ADASYN) or undersampling majority classes
2. Class weights: penalizing misclassification of minority samples more heavily
3. Ensemble methods: specialized for imbalanced data (e.g., RUSBoost, EasyEnsemble)

8.5 Need for Feature Engineering

Raw data rarely contains the optimal features for machine learning models:

- Important relationships may be implicit rather than explicit
- Domain knowledge can inform valuable transformations
- Feature interactions often capture critical patterns

Feature engineering approaches include:

1. Creating interaction terms ($x_1 \times x_2$)
2. Polynomial features (e.g., x^2 , x^3)
3. Domain-specific transformations (e.g., log returns in finance)
4. Time-based features from timestamps (hour of day, day of week)
5. Text embeddings for natural language

Automated feature engineering tools like `featuretools` and `tsfresh` can complement manual efforts in complex datasets.