

Data Preprocessing in Machine Learning

A Comprehensive Guide

Dr. ABDELLAOUI

2024/2025

Outline

- 1 Introduction and Motivation
- 2 Data Fundamentals
- 3 Data Cleaning
- 4 Data Transformation
- 5 Dimensionality Reduction
- 6 Sampling and Data Splitting
- 7 Tools and Implementation
- 8 Best Practices and Summary
- 9 Introduction to EDA
- 10 Data Types and Structures
- 11 Descriptive Statistics
- 12 Data Visualization Techniques
- 13 Correlation Analysis
- 14 Data Quality Assessment
- 15 EDA Workflow Summary

Why Data Preprocessing Matters

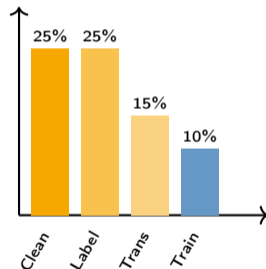
The Reality of ML Projects

- Real-world data is **messy** and **incomplete**
- Raw data rarely meets algorithm requirements
- **“Garbage In, Garbage Out”**
- Quality preprocessing = better performance

Key Insight

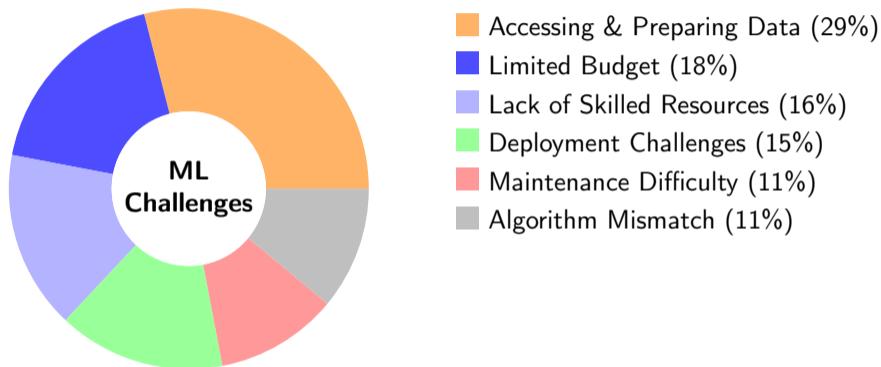
Data preprocessing is the **most time-consuming** but **most impactful** phase!

Time Allocation in ML Projects



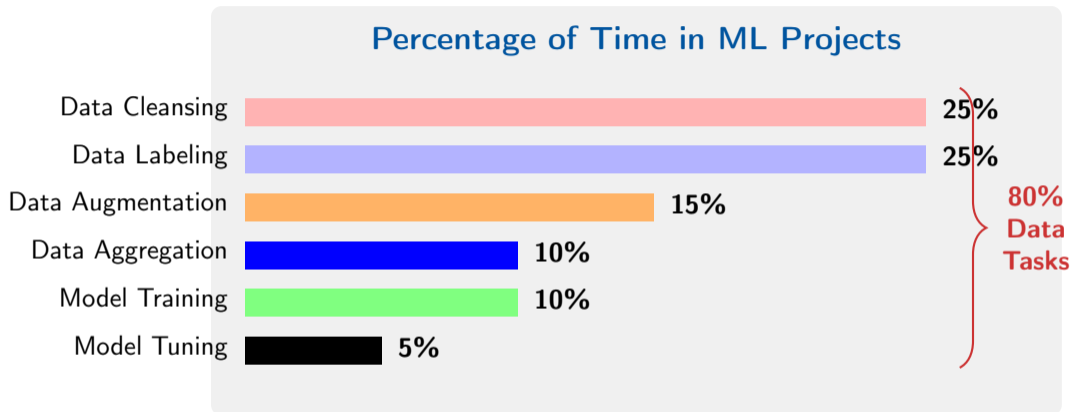
80% on data tasks!

Challenges in Applying Machine Learning



Key Finding: Data-related challenges represent the **#1 obstacle** in ML adoption.

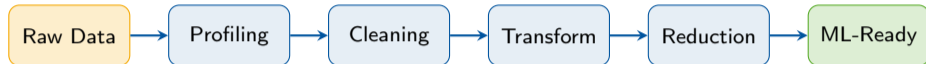
Time Allocation in ML Projects



Industry Statistics

According to multiple surveys, data scientists spend **60-80%** of their time on data preparation tasks, making it the most resource-intensive phase of ML projects.

The Data Preprocessing Pipeline



Preprocessing Objectives

- 1 **Data Quality:** Ensure accuracy and completeness
- 2 **Data Compatibility:** Match algorithm requirements
- 3 **Data Efficiency:** Reduce dimensionality
- 4 **Data Enhancement:** Create meaningful features

What is a Dataset?

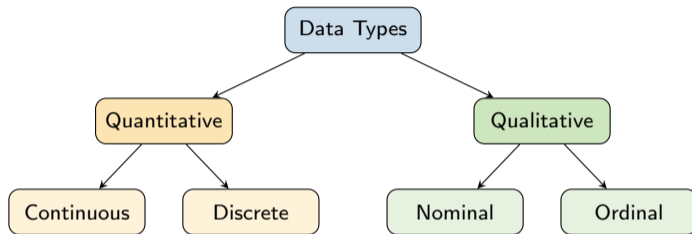
Definition

A **dataset** is an organized collection of data points described by features, formatted for analysis.

		Features (n)				
		x_1	x_2	...	x_n	y
Samples		val	val	...	val	label
						Target

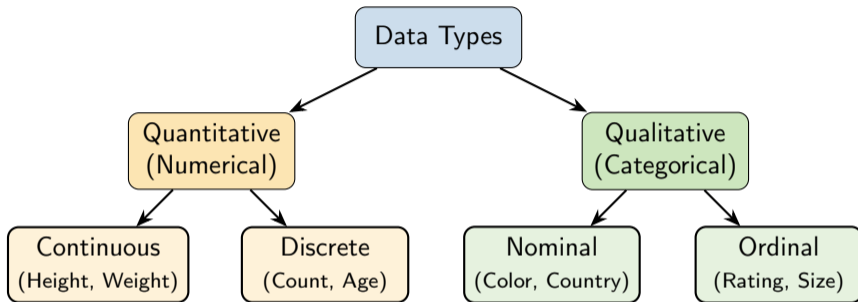
Notation: $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^m$ where $x_i \in \mathbb{R}^n$

Types of Data



Type	Equality	Order	Add	Ratio
Nominal	✓			
Ordinal	✓	✓		
Interval	✓	✓	✓	
Ratio	✓	✓	✓	✓

Types of Data



Data Type	Equality	Order	Addition	Multiplication
Nominal	✓	×	×	×
Ordinal	✓	✓	×	×
Interval	✓	✓	✓	×
Ratio	✓	✓	✓	✓

Common Data Quality Issues

Issue	Description	Example
Missing Values	NULL, NaN, empty	Age = NULL
Duplicates	Repeated records	Same row twice
Outliers	Extreme values	Age = 999
Noise	Random errors	Typos in data
Inconsistency	Format variations	NY vs New York
Errors	Wrong annotations	Mislabeled data

Real-World Example

A customer database: "NY", "New York", "new york", "NYC" — all the same city!

Handling Missing Values

Detection

- Check for NULL, NaN
- `df.isnull().sum()`
- Visualize with heatmaps

Treatment Strategies

- 1 **Deletion** - Remove rows/columns
- 2 **Imputation**
 - Mean/Median/Mode
 - KNN Imputation
 - MICE (Multiple Imputation)
- 3 **Indicator** - Binary flag

Imputation Methods

Method	Best For
Mean	Numerical, normal dist
Median	Numerical, skewed
Mode	Categorical
KNN	Preserves relations
MICE	Uncertainty modeling

Outlier Detection

Statistical Methods

- **Z-score:** $|z| > 3$ is outlier

$$z = \frac{x - \mu}{\sigma}$$

- **IQR Method:**

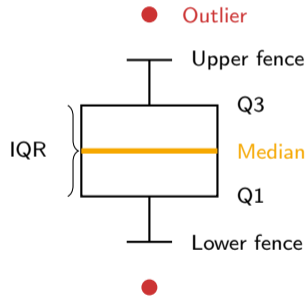
$$x < Q_1 - 1.5 \cdot IQR$$

$$x > Q_3 + 1.5 \cdot IQR$$

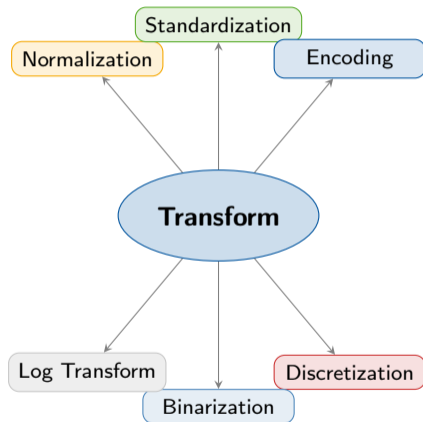
Treatment Options

Remove, Cap (Winsorize), Transform, Keep

Box Plot Anatomy



Transformation Overview



Normalization: Min-Max Scaling

Formula

$$x' = \frac{x - x_{min}}{x_{max} - x_{min}}$$

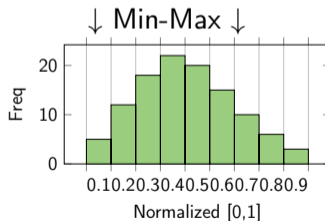
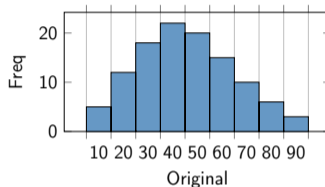
Scales to [0, 1]

When to Use

- Bounded data (age 0-100)
- Uniform distribution
- Neural networks, KNN

Limitation

Sensitive to outliers



Standardization: Z-Score

Formula

$$Z = \frac{X - \mu}{\sigma}$$

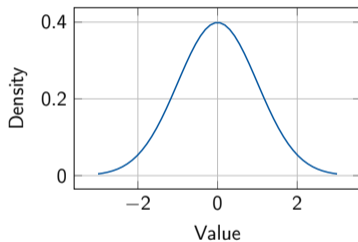
μ =mean, σ =std dev

When to Use

- Data with outliers
- Unknown bounds
- Gaussian assumptions
- SVM, PCA, Gradient descent

Result

Mean=0, Std=1



Standard Normal Distribution

$$\mu = 0, \sigma = 1$$

Normalization vs. Standardization

Aspect	Normalization	Standardization
Range	[0, 1]	Unbounded
Outliers	Sensitive	More robust
Distribution	Preserves	Centers at 0
Use Case	Bounded, NN	Unbounded, SVM
Formula	$\frac{x - x_{min}}{x_{max} - x_{min}}$	$\frac{x - \mu}{\sigma}$

Decision Rule

- **Normalization:** bounded data, no outliers
- **Standardization:** outliers present, Gaussian assumption

Encoding Categorical Variables

One-Hot Encoding

Binary columns per category

Color	R	B	G
Red	1	0	0
Blue	0	1	0
Green	0	0	1

Use for **nominal** data

Label Encoding

Integer per category

Size	Size
Small	0
Medium	1
Large	2

Use for **ordinal** data

Other Methods

Target Encoding, Frequency Encoding, Binary Encoding

Discretization (Binning)

Definition

Converting continuous \rightarrow categorical by grouping into bins

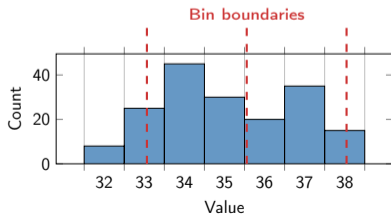
Equal-Width

$$width = \frac{x_{max} - x_{min}}{k}$$

All bins same range

Equal-Frequency (Quantile)

All bins same count



Why Reduce Dimensions?

Curse of Dimensionality

- Data sparse in high dimensions
- Distance metrics lose meaning
- More features = more data needed
- Higher computational cost
- Overfitting risk

Benefits

- Faster training
- Better generalization
- Easier visualization
- Noise reduction

Methods

Feature Selection:

- Filter (correlation, chi-square)
- Wrapper (RFE)
- Embedded (Lasso, trees)

Feature Extraction:

- PCA
- t-SNE
- Autoencoders

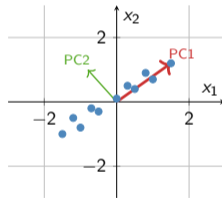
Principal Component Analysis (PCA)

Concept

Find orthogonal directions maximizing variance

Steps

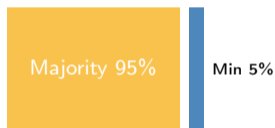
- 1 Standardize data
- 2 Compute covariance matrix
- 3 Get eigenvectors/values
- 4 Sort by eigenvalue
- 5 Select top k components
- 6 Project data



Handling Imbalanced Data

The Problem

Skewed class distribution → model favors majority



Degree	Minority
Mild	20-40%
Moderate	1-20%
Extreme	<1%

Solutions

- **Undersampling:** Reduce majority
- **Oversampling:** Increase minority (SMOTE)
- **Class weights:** Penalize majority errors more

Train-Test Split

Fundamental Rule

Never train and evaluate on the same data!



Key Points

- Test set must be representative
- Use stratified split for imbalanced data
- Temporal split for time series

K-Fold CV

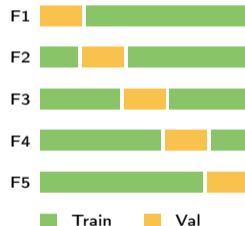
- 1 Split into k folds
- 2 For each fold: use as validation, rest for training
- 3 Average k scores

Benefits

- Reliable estimate
- Uses all data
- Essential for small datasets

Common: $k=5$ or $k=10$

5-Fold CV



Library	Purpose	Key Functions
Pandas	Data manipulation	read_csv, fillna, drop
NumPy	Numerical ops	array, mean, std
Scikit-learn	ML preprocessing	StandardScaler, OneHotEncoder
SciPy	Statistics	zscore, IQR
imbalanced-learn	Resampling	SMOTE, RandomUnderSampler

GUI Tools

Weka, Orange3, RapidMiner, KNIME

Code Example

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.preprocessing import StandardScaler, OneHotEncoder
3 from sklearn.impute import SimpleImputer
4 from sklearn.compose import ColumnTransformer
5 from sklearn.pipeline import Pipeline
6
7 # Numeric preprocessing
8 numeric_transformer = Pipeline([
9     ('imputer', SimpleImputer(strategy='median')),
10    ('scaler', StandardScaler())
11 ])
12
13 # Categorical preprocessing
14 categorical_transformer = Pipeline([
15     ('imputer', SimpleImputer(strategy='most_frequent')),
16     ('encoder', OneHotEncoder(handle_unknown='ignore'))
17 ])
18
19 # Combine
20 preprocessor = ColumnTransformer([
21     ('num', numeric_transformer, numeric_cols),
22     ('cat', categorical_transformer, categorical_cols)
23 ])
24
25 # Apply
26 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
27 X_train_processed = preprocessor.fit_transform(X_train)
28 X_test_processed = preprocessor.transform(X_test) # NO fit!
```

Best Practices

Do's

- Explore data first
- Document transformations
- Use pipelines
- Fit on training only
- Validate assumptions
- Version your data

Don'ts

- Don't leak test data
- Don't over-engineer
- Don't ignore domain
- Don't skip visualization
- Don't forget scaling

Summary

Key Takeaways

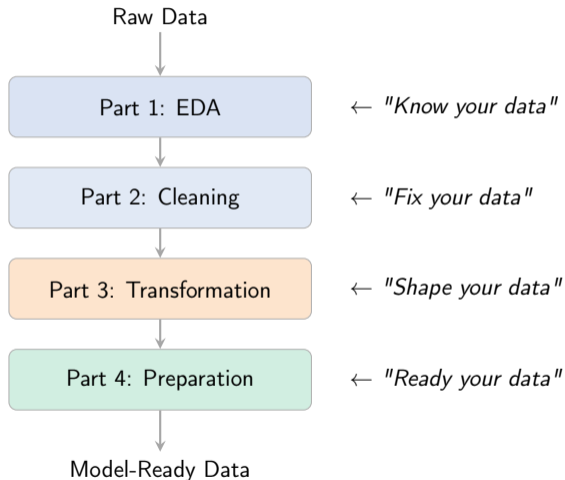
- 1 Data preprocessing is **critical** for ML success
- 2 **80%** of project time on data tasks
- 3 Quality data > Complex algorithms
- 4 Main steps: **Clean, Transform, Reduce, Split**
- 5 Always **fit on training only**
- 6 Use **pipelines** for reproducibility



Questions?

Thank you for your attention!

Pre-Processing Pipeline Overview



Part 1: EDA

"Know your data"

Outline

- 1 Introduction and Motivation
- 2 Data Fundamentals
- 3 Data Cleaning
- 4 Data Transformation
- 5 Dimensionality Reduction
- 6 Sampling and Data Splitting
- 7 Tools and Implementation
- 8 Best Practices and Summary
- 9 Introduction to EDA
- 10 Data Types and Structures
- 11 Descriptive Statistics
- 12 Data Visualization Techniques
- 13 Correlation Analysis
- 14 Data Quality Assessment
- 15 EDA Workflow Summary

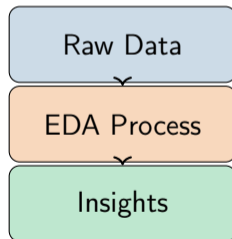
What is Exploratory Data Analysis?

Definition

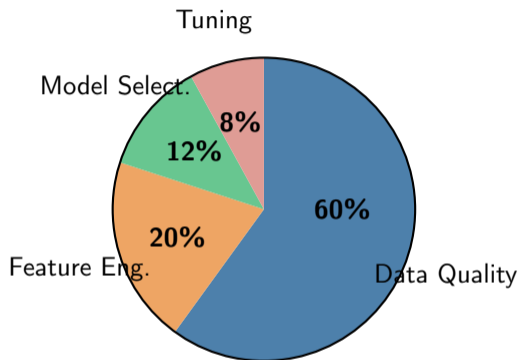
Exploratory Data Analysis (EDA) is an approach to analyzing datasets to summarize their main characteristics, often using statistical graphics and other data visualization methods.

Key Objectives:

- Understand data structure
- Detect patterns and anomalies
- Test hypotheses
- Check assumptions
- Guide feature engineering



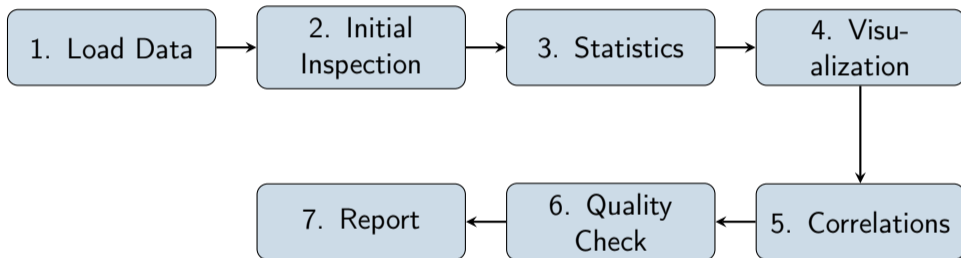
Why is EDA Critical in Machine Learning?



Golden Rule

“Garbage In, Garbage Out” — The quality of your model depends entirely on the quality of your data understanding.

The EDA Workflow



Understanding Data Types

Numerical Data

Continuous:

- Temperature: $T \in \mathbb{R}$
- Price: $P \in \mathbb{R}^+$
- Height: $h \in [0, \infty)$

Discrete:

- Count: $n \in \mathbb{N}$
- Age (years): $a \in \mathbb{Z}^+$

Categorical Data

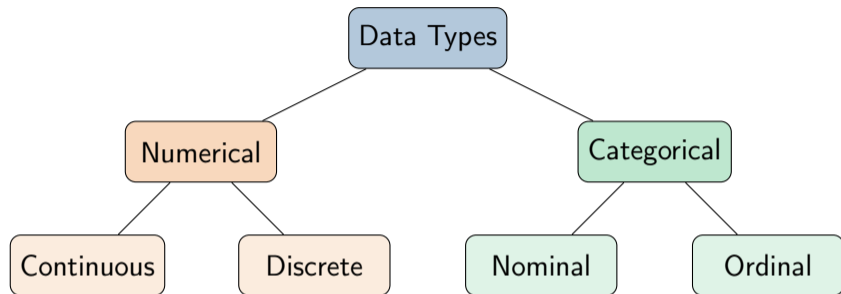
Nominal:

- Color: {Red, Blue, Green}
- Gender: {M, F, Other}

Ordinal:

- Education: {High School < Bachelor < Master < PhD}
- Rating: {1 < 2 < 3 < 4 < 5}

Data Type Hierarchy



Why It Matters

Different data types require different:

- Statistical measures (mean vs. mode)
- Visualization techniques (histogram vs. bar chart)
- Preprocessing methods (scaling vs. encoding)

The Feature Matrix

A dataset can be represented as a matrix $X \in \mathbb{R}^{n \times p}$:

$$X = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{np} \end{pmatrix}$$

Where:

- n = number of samples (rows)
- p = number of features (columns)
- x_{ij} = value of feature j for sample i

Initial Data Inspection Checklist

Questions to Answer:

- 1 What is the shape? $(n, p) = ?$
- 2 What are the column names?
- 3 What are the data types?
- 4 How much memory is used?
- 5 Are there missing values?
- 6 What do the first rows look like?

Python Commands

```
df.shape  
df.columns  
df.dtypes  
df.memory_usage()  
df.isnull().sum()  
df.head()
```

Measures of Central Tendency

Mean (Arithmetic Average)

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{x_1 + x_2 + \cdots + x_n}{n}$$

Median

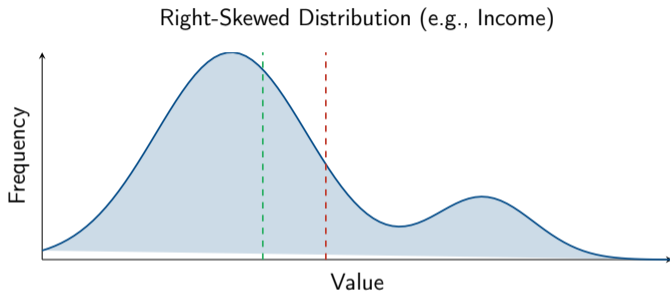
The middle value when data is sorted:

$$\tilde{x} = \begin{cases} x_{\frac{n+1}{2}} & \text{if } n \text{ is odd} \\ \frac{x_{\frac{n}{2}} + x_{\frac{n}{2}+1}}{2} & \text{if } n \text{ is even} \end{cases}$$

Mode

The most frequently occurring value: Mode = $\arg \max_x f(x)$

Mean vs. Median: When to Use Which?



Key Insight

Median is robust to outliers! For skewed data or data with outliers, median provides a better representation of the “typical” value.

Measures of Dispersion

Variance

Population variance:

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$

Sample variance (unbiased estimator):

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

Standard Deviation

$$\sigma = \sqrt{\sigma^2} \quad \text{or} \quad s = \sqrt{s^2}$$

Standard deviation is in the same units as the original data.

Range and Interquartile Range (IQR)

Range

$$R = x_{\max} - x_{\min}$$

Sensitive to outliers

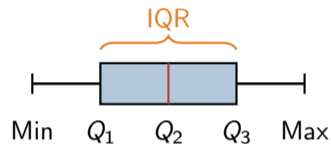
Interquartile Range

$$\text{IQR} = Q_3 - Q_1$$

Where:

- Q_1 = 25th percentile
- Q_3 = 75th percentile

Robust to outliers



Coefficient of Variation (CV)

Definition

The coefficient of variation measures relative variability:

$$CV = \frac{\sigma}{\mu} \times 100\%$$

or for samples:

$$CV = \frac{s}{\bar{x}} \times 100\%$$

Use Case

Comparing variability across different scales:

- Feature A: $\bar{x}_A = 1000$, $s_A = 50 \Rightarrow CV_A = 5\%$
- Feature B: $\bar{x}_B = 10$, $s_B = 2 \Rightarrow CV_B = 20\%$

Feature B has higher *relative* variability despite lower absolute s .

Skewness and Kurtosis

Skewness

Measures asymmetry of distribution:

$$\gamma_1 = \frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s} \right)^3$$

- $\gamma_1 = 0$: Symmetric
- $\gamma_1 > 0$: Right-skewed
- $\gamma_1 < 0$: Left-skewed

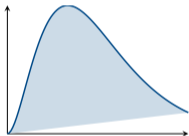
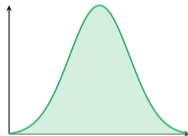
Kurtosis

Measures “tailedness”:

$$\gamma_2 = \frac{1}{n} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s} \right)^4 - 3$$

- $\gamma_2 = 0$: Mesokurtic (normal)
- $\gamma_2 > 0$: Leptokurtic (heavy tails)
- $\gamma_2 < 0$: Platykurtic (light tails)

Visualizing Skewness

Left-Skewed ($\gamma_1 < 0$)Symmetric ($\gamma_1 = 0$)Right-Skewed ($\gamma_1 > 0$)

Summary Statistics Table

Measure	Formula	Robust?
Mean	$\bar{x} = \frac{1}{n} \sum x_i$	No
Median	Middle value	Yes
Mode	Most frequent	Yes
Variance	$s^2 = \frac{1}{n-1} \sum (x_i - \bar{x})^2$	No
Std Dev	$s = \sqrt{s^2}$	No
Range	$R = x_{max} - x_{min}$	No
IQR	$Q_3 - Q_1$	Yes

Tip

Always report both mean/std AND median/IQR when exploring data!

Why Visualize Data?

Anscombe's Quartet

Four datasets with **identical** summary statistics:

- Same mean: $\bar{x} = 9$, $\bar{y} = 7.5$
- Same variance: $s_x^2 = 11$, $s_y^2 = 4.12$
- Same correlation: $r = 0.816$
- Same regression line: $y = 3 + 0.5x$

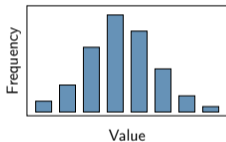
Yet they look completely different when plotted!

Lesson

Never trust statistics alone — always visualize your data!

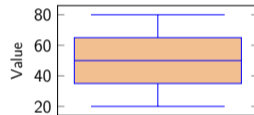
Univariate Visualizations

Histogram



Shows distribution shape, bins continuous data

Box Plot



Shows quartiles, median, outliers

Density Plots and KDE

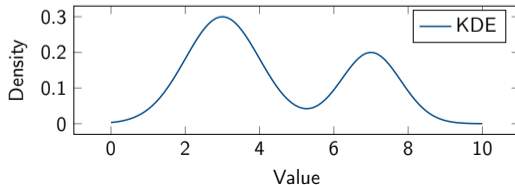
Kernel Density Estimation (KDE)

Smooth estimate of probability density function:

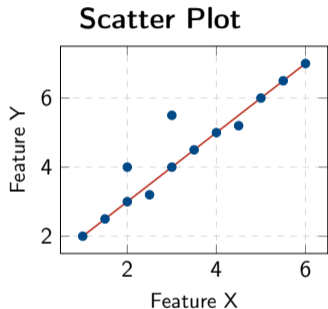
$$\hat{f}(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - x_i}{h}\right)$$

Where:

- $K(\cdot)$ is the kernel function (often Gaussian)
- h is the bandwidth (smoothing parameter)



Bivariate Visualizations



Reveals relationships between two variables

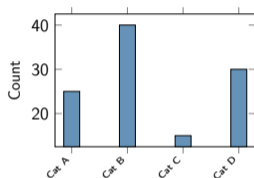
Heatmap (Correlation Matrix)

	A	B	C	D
B	0.7	1.0	0.5	-0.3
C	0.3	0.5	1.0	0.2
D	-0.5	-0.3	0.2	1.0

Shows pairwise correlations

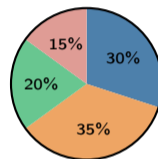
Categorical Data Visualizations

Bar Chart



Compare counts across categories

Pie Chart

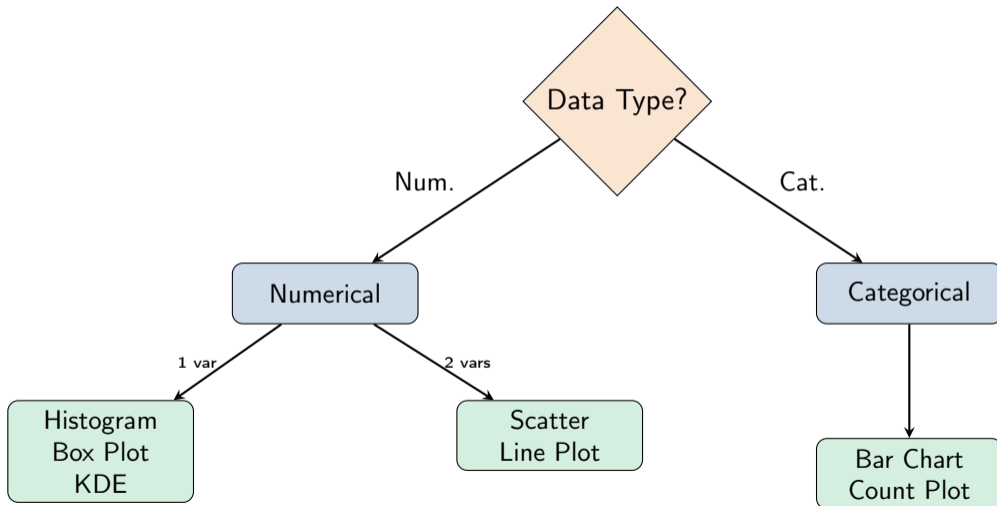


Show proportions (use sparingly!)

Best Practice

Prefer bar charts over pie charts — humans compare lengths better than angles!

Visualization Selection Guide



Pearson Correlation Coefficient

Definition

Measures **linear** relationship between two continuous variables:

$$r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \cdot \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

Or equivalently:

$$r_{xy} = \frac{\text{Cov}(X, Y)}{\sigma_X \cdot \sigma_Y} = \frac{\mathbb{E}[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y}$$

- $r = +1$: Perfect positive linear correlation
- $r = -1$: Perfect negative linear correlation
- $r = 0$: No linear correlation (but may have non-linear relationship!)

Interpreting Correlation Values

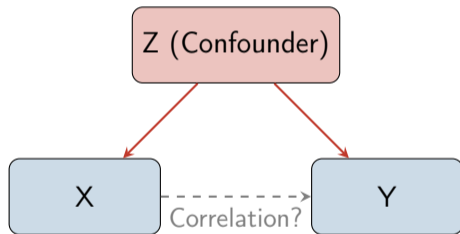


$ r $	Interpretation
0.0 – 0.3	Negligible/Weak
0.3 – 0.5	Moderate
0.5 – 0.7	Strong
0.7 – 1.0	Very Strong

Correlation \neq Causation!

Critical Warning

Correlation does not imply causation!



Example: Ice cream sales and drowning deaths are correlated.

Confounder: Hot weather increases both!

Spearman Rank Correlation

Definition

Measures **monotonic** relationship (not necessarily linear):

$$\rho = 1 - \frac{6 \sum_{i=1}^n d_i^2}{n(n^2 - 1)}$$

Where $d_i = \text{rank}(x_i) - \text{rank}(y_i)$

When to use Spearman:

- Non-linear monotonic relationships
- Ordinal data
- Data with outliers
- Non-normal distributions

When to use Pearson:

- Linear relationships
- Continuous data
- Normally distributed data
- No significant outliers

Correlation for Categorical Variables

Chi-Square Test of Independence

For categorical \times categorical variables:

$$\chi^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

Where:

- O_{ij} = Observed frequency
- $E_{ij} = \frac{(\text{Row Total}_i) \times (\text{Column Total}_j)}{n}$ = Expected frequency

Cramér's V

Strength of association (normalized chi-square):

$$V = \sqrt{\frac{\chi^2}{n}}$$

Point-Biserial Correlation

For Binary \times Continuous Variables

Special case of Pearson correlation:

$$r_{pb} = \frac{\bar{Y}_1 - \bar{Y}_0}{s_Y} \cdot \sqrt{\frac{n_1 n_0}{n^2}}$$

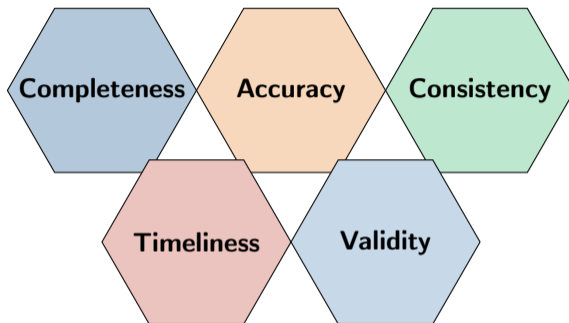
Where:

- \bar{Y}_1, \bar{Y}_0 = means of Y for each binary group
- n_1, n_0 = sample sizes for each group
- s_Y = standard deviation of Y

Example

Correlation between *Gender* (binary) and *Salary* (continuous)

Dimensions of Data Quality



- **Completeness:** Are there missing values?
- **Accuracy:** Are values correct?
- **Consistency:** Are values uniform across the dataset?
- **Timeliness:** Is the data up-to-date?
- **Validity:** Do values conform to expected formats/ranges?

Missing Data Analysis

Types of Missing Data

- **MCAR** (Missing Completely At Random): Missingness unrelated to any data
- **MAR** (Missing At Random): Missingness depends on observed data
- **MNAR** (Missing Not At Random): Missingness depends on unobserved data

Missing Rate Calculation

$$\text{Missing Rate}_j = \frac{\text{Number of missing values in column } j}{n} \times 100\%$$

Rule of Thumb

- < 5% missing: Generally safe to impute
- 5 – 20% missing: Careful analysis required
- > 20% missing: Consider dropping or advanced imputation

Detecting Outliers

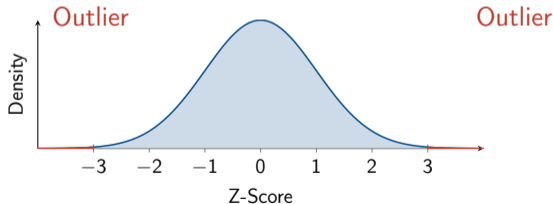
IQR Method

Outlier if: $x < Q_1 - 1.5 \cdot IQR$
or $x > Q_3 + 1.5 \cdot IQR$

Z-Score Method

$$z_i = \frac{x_i - \bar{x}}{s}$$

Outlier if: $|z_i| > 3$



Data Quality Checklist

Structural Checks:

- Correct number of rows/columns
- Appropriate data types
- No duplicate records
- Consistent column names
- Proper encoding (UTF-8)

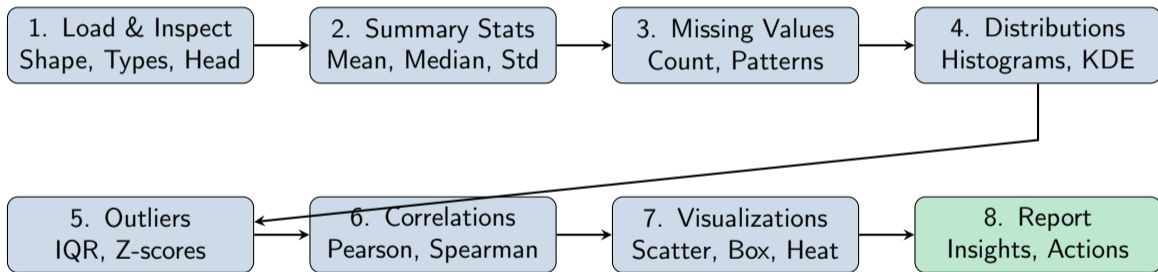
Content Checks:

- Values in valid ranges
- Categorical values expected
- No impossible combinations
- Dates in valid format
- Outliers identified

Quality Score

$$\text{Quality Score} = \frac{\text{Valid Records}}{\text{Total Records}} \times 100\%$$

Complete EDA Workflow



EDA Report Template

What to Include in Your EDA Report

- 1 **Dataset Overview**
 - Source, size, features description
- 2 **Data Quality Summary**
 - Missing values, duplicates, data types
- 3 **Univariate Analysis**
 - Key statistics and distributions for each feature
- 4 **Bivariate Analysis**
 - Important correlations and relationships
- 5 **Key Findings**
 - Patterns, anomalies, insights discovered
- 6 **Recommendations**
 - Suggested preprocessing steps

Python Code: Basic EDA Template

Essential EDA Commands

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

# Load and inspect
df = pd.read_csv('data.csv')
print(df.shape, df.dtypes)
print(df.describe())

# Missing values
print(df.isnull().sum())

# Distributions
df.hist(figsize=(12,8))
```

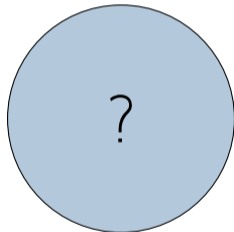
Key Takeaways

- 1 Always start with **EDA** — never skip this step!
- 2 Know your data types — they determine your analysis approach
- 3 Use both statistics **AND visualizations** — they complement each other
- 4 **Correlation \neq Causation** — always remember this
- 5 Document your findings — create reproducible reports
- 6 Quality matters — assess completeness, accuracy, consistency

Next Part





Part 2: Data Cleaning — Handling missing values, duplicates, outliers, and noise.

Questions & Discussion



Contact: `professor@university.edu`

References

-  Tukey, J.W. (1977). *Exploratory Data Analysis*. Addison-Wesley.
-  Wickham, H. & Grolemund, G. (2017). *R for Data Science*. O'Reilly.
-  McKinney, W. (2017). *Python for Data Analysis*. O'Reilly.
-  Géron, A. (2019). *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow*. O'Reilly.

References

Books

- Géron, A. (2022). *Hands-On Machine Learning*
- Zheng & Casari (2018). *Feature Engineering for ML*

Online

- <https://scikit-learn.org/stable/modules/preprocessing.html>
- <https://www.kaggle.com/learn>

Part I

Practice for Part 1 : EDA

What We Will Cover

Practical EDA Skills

- Setting up the Python environment
- Loading and inspecting data
- Understanding data types and structure
- Computing descriptive statistics
- Visualizing distributions and relationships
- Correlation analysis
- Identifying data quality issues

Libraries Used

`pandas, numpy, matplotlib, seaborn, scipy`

Essential Libraries Import

```
1 # Data manipulation
2 import pandas as pd
3 import numpy as np
4
5 # Visualization
6 import matplotlib.pyplot as plt
7 import seaborn as sns
8
9 # Statistical functions
10 from scipy import stats
11
12 # Display settings
13 pd.set_option('display.max_columns', None)
14 pd.set_option('display.precision', 2)
15 plt.style.use('seaborn-v0_8-whitegrid')
16 sns.set_palette("husl")
17
18 # Ignore warnings
19 import warnings
20 warnings.filterwarnings('ignore')
```

Loading Data

```
1 # From CSV file
2 df = pd.read_csv('data.csv')
3
4 # From Excel file
5 df = pd.read_excel('data.xlsx', sheet_name='Sheet1')
6
7 # From URL
8 url = 'https://example.com/data.csv'
9 df = pd.read_csv(url)
10
11 # Common parameters
12 df = pd.read_csv('data.csv',
13                 sep=',',           # Delimiter
14                 header=0,         # Row for column names
15                 index_col=None,   # Column to use as index
16                 na_values=['NA', 'missing', '?'],
17                 parse_dates=['date_column'],
18                 encoding='utf-8')
```

First Look at Data

```
1 # Display first/last rows
2 df.head(10)          # First 10 rows
3 df.tail(5)           # Last 5 rows
4
5 # Random sample
6 df.sample(5)         # 5 random rows
7
8 # Dataset dimensions
9 print(f"Shape: {df.shape}")
10 print(f"Rows: {df.shape[0]}, Columns: {df.shape[1]}")
11
12 # Column names
13 print(df.columns.tolist())
14
15 # Quick overview
16 df.info()
```

Understanding Data Types

```
1 # Check data types
2 print(df.dtypes)
3
4 # Count by data type
5 print(df.dtypes.value_counts())
6
7 # Identify column types
8 numerical_cols = df.select_dtypes(
9     include=['int64', 'float64']).columns.tolist()
10 categorical_cols = df.select_dtypes(
11     include=['object', 'category']).columns.tolist()
12 datetime_cols = df.select_dtypes(
13     include=['datetime64']).columns.tolist()
14
15 print(f"Numerical: {numerical_cols}")
16 print(f"Categorical: {categorical_cols}")
17 print(f"Datetime: {datetime_cols}")
```

Memory Usage

```
1 # Memory usage summary
2 print(df.info(memory_usage='deep'))
3
4 # Memory per column
5 memory_usage = df.memory_usage(deep=True)
6 print(memory_usage)
7
8 # Total memory in MB
9 total_mb = df.memory_usage(deep=True).sum() / 1024**2
10 print(f"Total Memory: {total_mb:.2f} MB")
11
12 # Optimize memory (convert to category for low cardinality)
13 for col in categorical_cols:
14     if df[col].nunique() / len(df) < 0.5:
15         df[col] = df[col].astype('category')
```

Basic Statistics for Numerical Data

```
1 # Complete statistical summary
2 df.describe()
3
4 # Include all columns (categorical too)
5 df.describe(include='all')
6
7 # Specific percentiles
8 df.describe(percentiles=[.01, .05, .25, .5, .75, .95, .99])
9
10 # Individual statistics
11 print(f"Mean:\n{df[numerical_cols].mean()}")
12 print(f"Median:\n{df[numerical_cols].median()}")
13 print(f"Std Dev:\n{df[numerical_cols].std()}")
14 print(f"Min:\n{df[numerical_cols].min()}")
15 print(f"Max:\n{df[numerical_cols].max()}")
```

Advanced Statistics

```
1 # Skewness and Kurtosis
2 from scipy.stats import skew, kurtosis
3
4 for col in numerical_cols:
5     sk = skew(df[col].dropna())
6     kt = kurtosis(df[col].dropna())
7     print(f"{col}: Skewness={sk:.3f}, Kurtosis={kt:.3f}")
8
9 # Quantiles
10 print(df[numerical_cols].quantile([0.25, 0.5, 0.75]))
11
12 # Range and IQR
13 for col in numerical_cols:
14     q1, q3 = df[col].quantile([0.25, 0.75])
15     iqr = q3 - q1
16     range_val = df[col].max() - df[col].min()
17     print(f"{col}: Range={range_val:.2f}, IQR={iqr:.2f}")
```

Statistics for Categorical Data

```
1 # Value counts for each categorical column
2 for col in categorical_cols:
3     print(f"\n{col}:")
4     print(df[col].value_counts())
5     print(f"Unique values: {df[col].nunique()}")
6
7 # Frequency and percentage
8 for col in categorical_cols:
9     freq = df[col].value_counts()
10    pct = df[col].value_counts(normalize=True) * 100
11    summary = pd.DataFrame({'Count': freq, 'Percent': pct})
12    print(f"\n{col}:")
13    print(summary)
14
15 # Mode (most frequent value)
16 print(df[categorical_cols].mode().iloc[0])
```

Creating a Summary Report

```
1 def eda_summary(df):
2     """Generate comprehensive EDA summary"""
3     summary = pd.DataFrame({
4         'dtype': df.dtypes,
5         'non_null': df.count(),
6         'null_count': df.isnull().sum(),
7         'null_pct': (df.isnull().sum() / len(df) * 100).round(2),
8         'unique': df.nunique(),
9         'unique_pct': (df.nunique() / len(df) * 100).round(2)
10    })
11
12    # Add statistics for numeric columns
13    num_stats = df.describe().T[['mean', 'std', 'min', 'max']]
14    summary = summary.join(num_stats)
15
16    return summary
17
18 print(eda_summary(df))
```

Histogram - Distribution of Single Variable

```
1 # Single histogram
2 plt.figure(figsize=(10, 6))
3 plt.hist(df['column'], bins=30, edgecolor='black', alpha=0.7)
4 plt.xlabel('Value')
5 plt.ylabel('Frequency')
6 plt.title('Distribution of Column')
7 plt.show()
8
9 # Multiple histograms with seaborn
10 fig, axes = plt.subplots(2, 3, figsize=(15, 10))
11 for idx, col in enumerate(numerical_cols[:6]):
12     ax = axes[idx // 3, idx % 3]
13     sns.histplot(df[col], kde=True, ax=ax)
14     ax.set_title(f'Distribution of {col}')
15 plt.tight_layout()
16 plt.show()
```

Box Plot - Identifying Outliers

```
1 # Single box plot
2 plt.figure(figsize=(8, 6))
3 sns.boxplot(x=df['column'])
4 plt.title('Box Plot of Column')
5 plt.show()
6
7 # Multiple box plots
8 plt.figure(figsize=(14, 6))
9 df[numerical_cols].boxplot()
10 plt.xticks(rotation=45)
11 plt.title('Box Plots of All Numerical Variables')
12 plt.tight_layout()
13 plt.show()
14
15 # Box plot by category
16 plt.figure(figsize=(10, 6))
17 sns.boxplot(x='category_col', y='numeric_col', data=df)
18 plt.title('Numeric by Category')
19 plt.show()
```

Violin Plot - Distribution Shape

```
1 # Violin plot shows distribution shape
2 plt.figure(figsize=(12, 6))
3 sns.violinplot(x='category', y='value', data=df)
4 plt.title('Distribution by Category')
5 plt.show()
6
7 # Combined violin and box plot
8 plt.figure(figsize=(12, 6))
9 sns.violinplot(x='category', y='value', data=df, inner='box')
10 plt.show()
11
12 # Split violin plot for binary comparison
13 plt.figure(figsize=(12, 6))
14 sns.violinplot(x='category', y='value', hue='binary_col',
15               data=df, split=True)
16 plt.show()
```

Bar Plot - Categorical Data

```
1 # Count plot for categorical variables
2 plt.figure(figsize=(10, 6))
3 sns.countplot(x='category_col', data=df,
4               order=df['category_col'].value_counts().index)
5 plt.xticks(rotation=45)
6 plt.title('Count by Category')
7 plt.show()
8
9 # Bar plot with mean values
10 plt.figure(figsize=(10, 6))
11 sns.barplot(x='category', y='value', data=df,
12            estimator=np.mean, ci=95)
13 plt.title('Mean Value by Category (with 95% CI)')
14 plt.show()
15
16 # Horizontal bar plot
17 plt.figure(figsize=(8, 10))
18 df['category'].value_counts().plot(kind='barh')
19 plt.xlabel('Count')
20 plt.show()
```

Scatter Plot - Relationships

```
1 # Basic scatter plot
2 plt.figure(figsize=(10, 6))
3 plt.scatter(df['x'], df['y'], alpha=0.5)
4 plt.xlabel('X Variable')
5 plt.ylabel('Y Variable')
6 plt.title('Scatter Plot: X vs Y')
7 plt.show()
8
9 # Scatter with regression line
10 plt.figure(figsize=(10, 6))
11 sns.regplot(x='x', y='y', data=df, scatter_kws={'alpha':0.5})
12 plt.title('Scatter Plot with Regression Line')
13 plt.show()
14
15 # Scatter colored by category
16 plt.figure(figsize=(10, 6))
17 sns.scatterplot(x='x', y='y', hue='category', data=df)
18 plt.title('Scatter Plot by Category')
19 plt.show()
```

Pair Plot - Multiple Relationships

```
1 # Basic pair plot (all numerical variables)
2 sns.pairplot(df[numerical_cols[:5]]) # Limit to 5 columns
3 plt.suptitle('Pair Plot of Numerical Variables', y=1.02)
4 plt.show()
5
6 # Pair plot with hue (colored by category)
7 sns.pairplot(df, vars=numerical_cols[:4], hue='target')
8 plt.show()
9
10 # Pair plot with different plot types
11 sns.pairplot(df[numerical_cols[:4]],
12             diag_kind='kde',          # KDE on diagonal
13             plot_kws={'alpha': 0.5},
14             corner=True)             # Only lower triangle
15 plt.show()
```

Heatmap - Visualizing Matrices

```
1 # Correlation heatmap
2 plt.figure(figsize=(12, 10))
3 correlation_matrix = df[numerical_cols].corr()
4 sns.heatmap(correlation_matrix,
5             annot=True,           # Show values
6             fmt='.2f',           # Format
7             cmap='RdBu_r',       # Color map
8             center=0,           # Center at 0
9             square=True,
10            linewidths=0.5)
11 plt.title('Correlation Heatmap')
12 plt.tight_layout()
13 plt.show()
14
15 # Missing values heatmap
16 plt.figure(figsize=(12, 6))
17 sns.heatmap(df.isnull(), cbar=True, yticklabels=False)
18 plt.title('Missing Values Pattern')
19 plt.show()
```

Distribution Comparison (KDE)

```
1 # KDE plot
2 plt.figure(figsize=(10, 6))
3 sns.kdeplot(df['column'], fill=True, alpha=0.5)
4 plt.title('Kernel Density Estimation')
5 plt.show()
6
7 # Compare distributions by group
8 plt.figure(figsize=(10, 6))
9 for category in df['group'].unique():
10     subset = df[df['group'] == category]['value']
11     sns.kdeplot(subset, label=category, fill=True, alpha=0.3)
12 plt.legend()
13 plt.title('Distribution by Group')
14 plt.show()
15
16 # Compare with normal distribution
17 from scipy.stats import norm
18 sns.kdeplot(df['column'], label='Data')
19 x = np.linspace(df['column'].min(), df['column'].max(), 100)
20 plt.plot(x, norm.pdf(x, df['column'].mean(), df['column'].std()),
21         label='Normal', linestyle='--')
22 plt.legend()
```

Pearson Correlation

```
1 # Correlation matrix (Pearson - default)
2 corr_matrix = df[numerical_cols].corr(method='pearson')
3 print(corr_matrix)
4
5 # Correlation with specific target
6 target_corr = df[numerical_cols].corr()['target'].sort_values(
7     ascending=False)
8 print("Correlation with target:")
9 print(target_corr)
10
11 # Statistical significance
12 from scipy.stats import pearsonr
13 for col in numerical_cols:
14     if col != 'target':
15         corr, p_value = pearsonr(df[col].dropna(),
16                                 df['target'].dropna())
17         sig = "***" if p_value < 0.001 else \
18             "**" if p_value < 0.01 else \
19             "*" if p_value < 0.05 else ""
20         print(f"{col}: r={corr:.3f}, p={p_value:.4f} {sig}")
```

Spearman and Kendall Correlation

```
1 # Spearman correlation (rank-based, non-linear monotonic)
2 spearman_corr = df[numerical_cols].corr(method='spearman')
3
4 # Kendall correlation (also rank-based)
5 kendall_corr = df[numerical_cols].corr(method='kendall')
6
7 # Compare all three methods
8 from scipy.stats import spearmanr, kendalltau
9
10 col1, col2 = 'var1', 'var2'
11
12 pearson_r, pearson_p = pearsonr(df[col1], df[col2])
13 spearman_r, spearman_p = spearmanr(df[col1], df[col2])
14 kendall_r, kendall_p = kendalltau(df[col1], df[col2])
15
16 print(f"Pearson: r={pearson_r:.3f}, p={pearson_p:.4f}")
17 print(f"Spearman: r={spearman_r:.3f}, p={spearman_p:.4f}")
18 print(f"Kendall: r={kendall_r:.3f}, p={kendall_p:.4f}")
```

Correlation for Categorical Variables

```
1 # Chi-Square test for categorical vs categorical
2 from scipy.stats import chi2_contingency
3
4 def cramers_v(x, y):
5     """Calculate Cramer's V for categorical association"""
6     contingency = pd.crosstab(x, y)
7     chi2, p, dof, expected = chi2_contingency(contingency)
8     n = contingency.sum().sum()
9     min_dim = min(contingency.shape) - 1
10    return np.sqrt(chi2 / (n * min_dim)), p
11
12 # Example usage
13 v, p = cramers_v(df['cat1'], df['cat2'])
14 print(f"Cramer's V: {v:.3f}, p-value: {p:.4f}")
15
16 # Contingency table
17 contingency = pd.crosstab(df['cat1'], df['cat2'],
18                          margins=True, normalize='all')
19 print(contingency)
```

Point-Biserial Correlation

```
1 # Point-biserial: numerical vs binary categorical
2 from scipy.stats import pointbiserialr
3
4 # Convert binary category to 0/1
5 df['binary_encoded'] = (df['binary_col'] == 'Yes').astype(int)
6
7 # Calculate correlation
8 for col in numerical_cols:
9     corr, p_value = pointbiserialr(df['binary_encoded'], df[col])
10    print(f"{col}: r={corr:.3f}, p={p_value:.4f}")
11
12 # ANOVA for numerical vs multi-class categorical
13 from scipy.stats import f_oneway
14
15 groups = [group['numeric_col'].values
16           for name, group in df.groupby('category_col')]
17 f_stat, p_value = f_oneway(*groups)
18 print(f"ANOVA: F={f_stat:.3f}, p={p_value:.4f}")
```

Correlation Visualization

```
1 # Clustered heatmap (hierarchical clustering)
2 sns.clustermap(df[numerical_cols].corr(),
3               annot=True, fmt='.2f', cmap='RdBu_r',
4               figsize=(12, 12))
5 plt.show()
6
7 # Top correlations (excluding self-correlation)
8 def get_top_correlations(corr_matrix, n=10):
9     pairs = []
10    for i in range(len(corr_matrix.columns)):
11        for j in range(i+1, len(corr_matrix.columns)):
12            pairs.append((corr_matrix.columns[i],
13                        corr_matrix.columns[j],
14                        corr_matrix.iloc[i, j]))
15    return sorted(pairs, key=lambda x: abs(x[2]), reverse=True)[:n]
16
17 top_corr = get_top_correlations(corr_matrix, n=10)
18 for col1, col2, corr in top_corr:
19     print(f"{col1} <-> {col2}: {corr:.3f}")
```

Missing Values Analysis

```
1 # Count missing values
2 print("Missing Values:")
3 print(df.isnull().sum())
4
5 # Percentage of missing values
6 missing_pct = (df.isnull().sum() / len(df) * 100).round(2)
7 missing_df = pd.DataFrame({
8     'Missing Count': df.isnull().sum(),
9     'Missing %': missing_pct
10 }).sort_values('Missing %', ascending=False)
11 print(missing_df[missing_df['Missing Count'] > 0])
12
13 # Visualize missing values
14 import missingno as msno # pip install missingno
15 msno.matrix(df)
16 plt.show()
17
18 msno.bar(df)
19 plt.show()
```

Missing Values Pattern

```
1 # Missing values heatmap (correlation of missingness)
2 msno.heatmap(df)
3 plt.show()
4
5 # Dendrogram of missing values
6 msno.dendrogram(df)
7 plt.show()
8
9 # Check if missing values are related to other columns
10 def missing_analysis(df, col_with_missing, other_cols):
11     """Compare statistics when value is missing vs not missing"""
12     missing_mask = df[col_with_missing].isnull()
13
14     for col in other_cols:
15         if df[col].dtype in ['int64', 'float64']:
16             mean_missing = df.loc[missing_mask, col].mean()
17             mean_not_missing = df.loc[~missing_mask, col].mean()
18             print(f"{col}: Missing={mean_missing:.2f}, "
19                 f"Not Missing={mean_not_missing:.2f}")
20
21 missing_analysis(df, 'col_with_na', numerical_cols)
```

Duplicate Detection

```
1 # Check for duplicate rows
2 print(f"Duplicate rows: {df.duplicated().sum()}")
3
4 # View duplicate rows
5 duplicates = df[df.duplicated(keep=False)]
6 print(duplicates.sort_values(by=df.columns.tolist()))
7
8 # Duplicates based on specific columns
9 key_cols = ['id', 'date', 'category']
10 print(f"Duplicates by key: {df.duplicated(subset=key_cols).sum()}")
11
12 # Identify which rows are duplicates
13 df['is_duplicate'] = df.duplicated(keep=False)
14 df['duplicate_group'] = df.groupby(list(df.columns[:-1])).ngroup()
15
16 # Count duplicates per group
17 dup_counts = df.groupby('duplicate_group').size()
18 print(dup_counts[dup_counts > 1])
```

Outlier Detection - IQR Method

```
1 def detect_outliers_iqr(df, columns):
2     """Detect outliers using IQR method"""
3     outlier_info = {}
4
5     for col in columns:
6         Q1 = df[col].quantile(0.25)
7         Q3 = df[col].quantile(0.75)
8         IQR = Q3 - Q1
9         lower_bound = Q1 - 1.5 * IQR
10        upper_bound = Q3 + 1.5 * IQR
11
12        outliers = df[(df[col] < lower_bound) |
13                      (df[col] > upper_bound)]
14
15        outlier_info[col] = {
16            'count': len(outliers),
17            'percentage': len(outliers) / len(df) * 100,
18            'lower_bound': lower_bound,
19            'upper_bound': upper_bound
20        }
21
22    return pd.DataFrame(outlier_info).T
23
24 print(detect_outliers_iqr(df, numerical_cols))
```

Outlier Detection - Z-Score Method

```
1 from scipy.stats import zscore
2
3 def detect_outliers_zscore(df, columns, threshold=3):
4     """Detect outliers using Z-score method"""
5     outlier_info = {}
6
7     for col in columns:
8         z_scores = np.abs(zscore(df[col].dropna()))
9         outlier_mask = z_scores > threshold
10
11         outlier_info[col] = {
12             'count': outlier_mask.sum(),
13             'percentage': outlier_mask.sum() / len(df) * 100,
14             'max_zscore': z_scores.max()
15         }
16
17     return pd.DataFrame(outlier_info).T
18
19 print(detect_outliers_zscore(df, numerical_cols))
20
21 # Mark outliers in dataframe
22 for col in numerical_cols:
23     df[f'{col}_zscore'] = np.abs(zscore(df[col].fillna(df[col].mean())))
24     df[f'{col}_is_outlier'] = df[f'{col}_zscore'] > 3
```

Outlier Visualization

```
1 # Box plots for all numerical columns
2 fig, axes = plt.subplots(2, 3, figsize=(15, 10))
3 for idx, col in enumerate(numerical_cols[:6]):
4     ax = axes[idx // 3, idx % 3]
5     sns.boxplot(x=df[col], ax=ax)
6     ax.set_title(f'{col}')
7 plt.tight_layout()
8 plt.show()
9
10 # Scatter plot highlighting outliers
11 plt.figure(figsize=(10, 6))
12 outlier_mask = np.abs(zscore(df['column'])) > 3
13 plt.scatter(df.index[~outlier_mask], df.loc[~outlier_mask, 'column'],
14            alpha=0.5, label='Normal')
15 plt.scatter(df.index[outlier_mask], df.loc[outlier_mask, 'column'],
16            color='red', alpha=0.7, label='Outlier')
17 plt.legend()
18 plt.title('Outlier Detection')
19 plt.show()
```

Pandas Profiling / YData Profiling

```
1 # Install: pip install ydata-profiling
2
3 from ydata_profiling import ProfileReport
4
5 # Generate comprehensive report
6 profile = ProfileReport(df,
7                         title="EDA Report",
8                         explorative=True,
9                         minimal=False)
10
11 # Save to HTML file
12 profile.to_file("eda_report.html")
13
14 # Display in Jupyter notebook
15 profile.to_notebook_iframe()
16
17 # Minimal report (faster for large datasets)
18 profile_minimal = ProfileReport(df, minimal=True)
19 profile_minimal.to_file("eda_report_minimal.html")
```

```
1 # Install: pip install sweetviz
2
3 import sweetviz as sv
4
5 # Analyze single dataset
6 report = sv.analyze(df)
7 report.show_html("sweetviz_report.html")
8
9 # Compare two datasets (e.g., train vs test)
10 report = sv.compare([df_train, "Train"], [df_test, "Test"])
11 report.show_html("comparison_report.html")
12
13 # Analyze with target variable
14 report = sv.analyze(df, target_feat='target_column')
15 report.show_html("sweetviz_target_report.html")
16
17 # Compare subsets
18 report = sv.compare_intra(df, df['category'] == 'A',
19                          names=['Group A', 'Others'])
20 report.show_html("intra_comparison.html")
```

D-Tale (Interactive EDA)

```
1 # Install: pip install dtale
2
3 import dtale
4
5 # Launch interactive dashboard
6 d = dtale.show(df)
7
8 # Opens in browser automatically
9 # Or get the URL
10 print(d._url)
11
12 # Features available in D-Tale:
13 # - Column analysis
14 # - Correlation matrix
15 # - Charts builder
16 # - Missing values analysis
17 # - Duplicate detection
18 # - Code export
19
20 # For Jupyter notebooks
21 d.notebook()
```

EDA Function Template (Part 1)

```
1 def complete_eda(df, target=None):
2     """Comprehensive EDA function"""
3     print("="*50)
4     print("DATASET OVERVIEW")
5     print("="*50)
6     print(f"Shape: {df.shape}")
7     print(f"\nData Types:\n{df.dtypes.value_counts()}")
8     print(f"\nMemory Usage: {df.memory_usage(deep=True).sum()/1024**2:.2f} MB")
9
10    # Identify column types
11    num_cols = df.select_dtypes(include=np.number).columns.tolist()
12    cat_cols = df.select_dtypes(include=['object', 'category']).columns.tolist()
13
14    print(f"\nNumerical columns ({len(num_cols)}): {num_cols}")
15    print(f"Categorical columns ({len(cat_cols)}): {cat_cols}")
16
17    # Missing values
18    print("\n" + "="*50)
19    print("MISSING VALUES")
20    print("="*50)
21    missing = df.isnull().sum()
22    missing_pct = (missing / len(df) * 100).round(2)
23    print(pd.DataFrame({'Count': missing, '%': missing_pct}
24                       ).query('Count > 0').sort_values('%', ascending=False))
```

EDA Function Template (Part 2)

```
1 # Continued from previous slide...
2
3 # Descriptive statistics
4 print("\n" + "="*50)
5 print("DESCRIPTIVE STATISTICS")
6 print("="*50)
7 print(df[num_cols].describe().T)
8
9 # Duplicates
10 print(f"\nDuplicate rows: {df.duplicated().sum()}")
11
12 # Target analysis (if provided)
13 if target and target in df.columns:
14     print("\n" + "="*50)
15     print(f"TARGET ANALYSIS: {target}")
16     print("="*50)
17     if df[target].dtype in ['object', 'category']:
18         print(df[target].value_counts())
19         print(f"\nClass balance: {df[target].value_counts(normalize=True)}")
20     else:
21         print(df[target].describe())
22
23 return num_cols, cat_cols
```

EDA Visualization Function

```
1 def plot_eda(df, num_cols, cat_cols, target=None):
2     """Generate EDA visualizations"""
3     # 1. Distribution of numerical variables
4     n_num = len(num_cols)
5     fig, axes = plt.subplots((n_num+2)//3, 3, figsize=(15, 4*((n_num+2)//3)))
6     axes = axes.flatten()
7     for i, col in enumerate(num_cols):
8         sns.histplot(df[col], kde=True, ax=axes[i])
9         axes[i].set_title(f'{col}')
10    plt.tight_layout()
11    plt.savefig('numerical_distributions.png', dpi=150)
12    plt.show()
13
14    # 2. Correlation heatmap
15    plt.figure(figsize=(12, 10))
16    sns.heatmap(df[num_cols].corr(), annot=True, fmt='.2f', cmap='RdBu_r')
17    plt.title('Correlation Matrix')
18    plt.tight_layout()
19    plt.savefig('correlation_matrix.png', dpi=150)
20    plt.show()
```

Running Complete EDA

```
1 # Load data
2 df = pd.read_csv('your_data.csv')
3
4 # Run EDA
5 num_cols, cat_cols = complete_eda(df, target='target_column')
6
7 # Generate plots
8 plot_eda(df, num_cols, cat_cols, target='target_column')
9
10 # Or use automated tools
11 from ydata_profiling import ProfileReport
12 profile = ProfileReport(df, title="Complete EDA Report")
13 profile.to_file("complete_eda_report.html")
14
15 # Save summary to Excel
16 with pd.ExcelWriter('eda_summary.xlsx') as writer:
17     df.describe().to_excel(writer, sheet_name='Statistics')
18     df.corr().to_excel(writer, sheet_name='Correlations')
19     df.isnull().sum().to_excel(writer, sheet_name='Missing')
```

Summary: EDA Checklist

Data Understanding:

- `df.shape, df.info()`
- `df.head(), df.sample()`
- `df.dtypes`
- `df.describe()`

Data Quality:

- `df.isnull().sum()`
- `df.duplicated().sum()`
- Outlier detection (IQR/Z-score)

Visualizations:

- Histograms / KDE
- Box plots
- Scatter plots
- Correlation heatmap
- Pair plot

Correlations:

- Pearson (numerical)
- Spearman (non-linear)
- Chi-square (categorical)

Key Python Functions Reference

Task	Function
Load data	<code>pd.read_csv()</code> , <code>pd.read_excel()</code>
First look	<code>df.head()</code> , <code>df.info()</code> , <code>df.shape</code>
Statistics	<code>df.describe()</code> , <code>df.mean()</code> , <code>df.std()</code>
Missing values	<code>df.isnull().sum()</code> , <code>msno.matrix()</code>
Duplicates	<code>df.duplicated()</code> , <code>df.drop_duplicates()</code>
Value counts	<code>df['col'].value_counts()</code>
Correlation	<code>df.corr()</code> , <code>pearsonr()</code> , <code>spearmanr()</code>
Histogram	<code>sns.histplot()</code> , <code>plt.hist()</code>
Box plot	<code>sns.boxplot()</code> , <code>df.boxplot()</code>
Scatter	<code>sns.scatterplot()</code> , <code>plt.scatter()</code>
Heatmap	<code>sns.heatmap()</code>
Pair plot	<code>sns.pairplot()</code>
Auto EDA	<code>ProfileReport()</code> , <code>sv.analyze()</code>

Questions & Practice

Next: Part 2 - Data Cleaning with Python

Part II

Part 2: Data Cleaning

"Fix your data"

What is Data Cleaning?

Definition

Data Cleaning (or Data Cleansing) is the process of detecting and correcting (or removing) corrupt, inaccurate, or irrelevant records from a dataset.

Common Data Problems:

- Missing values
- Duplicate records
- Outliers and anomalies
- Inconsistent formatting
- Noisy data
- Invalid entries



Why is Data Cleaning Important?

60-80%
of ML project time

Poor data quality
costs \$15M/year avg.

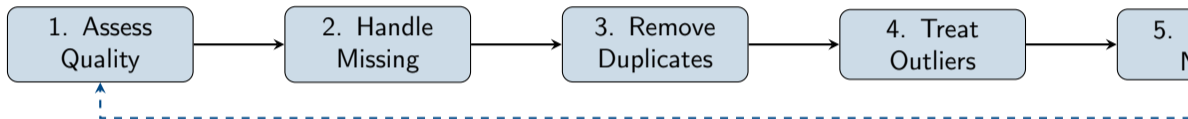
Clean data
improves accuracy 40%+

The Data Quality Principle

“Your model is only as good as your data.”

- Missing values → Biased models
- Duplicates → Overfitting
- Outliers → Skewed predictions
- Noise → Poor generalization

Data Cleaning Workflow



Iterative Process

Data cleaning is often **iterative** — you may need to revisit earlier steps as you discover new issues.

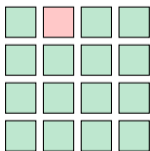
Types of Missing Data

Classification by Rubin (1976)

Missing data mechanisms determine appropriate handling strategies:

MCAR

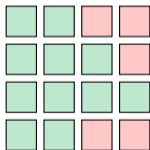
Missing Completely
At Random



Random pattern

MAR

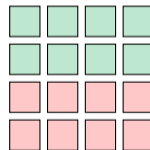
Missing At Random



Depends on observed

MNAR

Missing Not At Random



Depends on missing value

Mathematical Definition of Missingness

Formal Notation

Let X be the complete data matrix and M be the missingness indicator matrix:

$$M_{ij} = \begin{cases} 1 & \text{if } X_{ij} \text{ is missing} \\ 0 & \text{if } X_{ij} \text{ is observed} \end{cases}$$

Missing Data Mechanisms

- **MCAR:** $P(M|X) = P(M)$
- **MAR:** $P(M|X) = P(M|X_{obs})$
- **MNAR:** $P(M|X) = P(M|X_{obs}, X_{miss})$

Key Insight

MCAR and MAR are *ignorable* — standard methods work. MNAR requires specialized techniques.

Detecting Missing Values

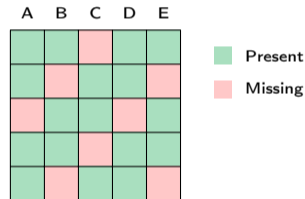
Missing Rate per Feature

$$r_j = \frac{\sum_{i=1}^n M_{ij}}{n} \times 100\%$$

Total Missing Rate

$$r_{total} = \frac{\sum_{i=1}^n \sum_{j=1}^p M_{ij}}{n \times p} \times 100\%$$

Missing Value Heatmap



Strategy 1: Deletion Methods

Listwise Deletion (Complete Case)

Remove entire rows with any missing values:

$$X_{clean} = \{x_i : \sum_{j=1}^p M_{ij} = 0\}$$

Pros:

- Simple to implement
- Works if MCAR

Cons:

- Loses data (reduces n)
- Biased if not MCAR

Pairwise Deletion

Use available data for each analysis:

$\text{Corr}(X_j, X_k)$ uses all rows where both present

Pros:

- Preserves more data
- Flexible

Cons:

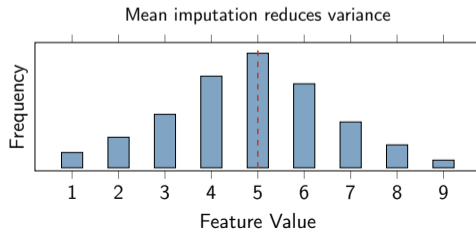
- Different n per calculation
- Can produce inconsistent results

Strategy 2: Simple Imputation Methods

Mean/Median/Mode Imputation

Replace missing values with central tendency measure:

$$\hat{x}_{ij} = \begin{cases} \bar{x}_j = \frac{1}{n_{obs}} \sum_{i: M_{ij}=0} x_{ij} & \text{(Mean)} \\ \tilde{x}_j & \text{(Median)} \\ \text{Mode}(x_j) & \text{(Mode, for categorical)} \end{cases}$$

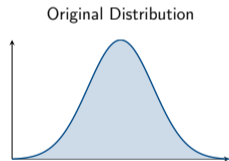


Problems with Simple Imputation

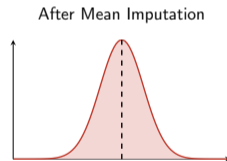
Variance Reduction Problem

Mean imputation artificially reduces variance:

$$\text{Var}(\hat{X}) = \text{Var}(X_{obs}) \cdot \frac{n_{obs}}{n} < \text{Var}(X_{true})$$



Wider spread



Narrower (biased)

Strategy 3: K-Nearest Neighbors (KNN) Imputation

Algorithm

For each missing value x_{ij} :

- 1 Find k nearest neighbors based on other features
- 2 Impute using weighted average of neighbors' values

$$\hat{x}_{ij} = \frac{\sum_{l \in N_k(i)} w_{il} \cdot x_{lj}}{\sum_{l \in N_k(i)} w_{il}}$$

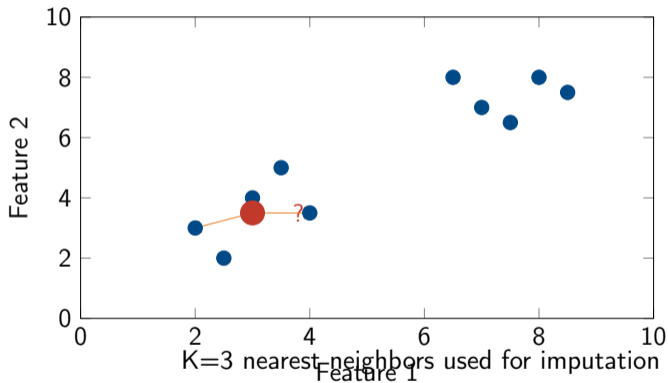
Where weights are typically inverse distance:

$$w_{il} = \frac{1}{d(x_i, x_l) + \epsilon}$$

Distance Metric (Euclidean)

$$d(x_i, x_l) = \sqrt{\sum (x_{ij} - x_{lj})^2}$$

KNN Imputation Visualization



Strategy 4: Regression Imputation

Linear Regression Imputation

Predict missing values using regression on observed features:

$$\hat{x}_{ij} = \beta_0 + \sum_{k \neq j} \beta_k x_{ik}$$

Where β is estimated from complete cases:

$$\hat{\beta} = (X_{obs}^T X_{obs})^{-1} X_{obs}^T y_{obs}$$

Stochastic Regression Imputation

Add random error to preserve variance:

$$\hat{x}_{ij} = \beta_0 + \sum_{k \neq j} \beta_k x_{ik} + \epsilon_i$$

Multivariate Imputation by Chained Equations

Create m complete datasets, analyze each, combine results:

Rubin's Combining Rules

$$\bar{Q} = \frac{1}{m} \sum_{i=1}^m \hat{Q}_i \quad (\text{Combined estimate})$$

Imputation Method Comparison

Method	Preserves Variance	Handles MNAR	Computation	Best For
Deletion	N/A	✗	Low	MCAR, low %
Mean/Median	✗	✗	Low	Quick baseline
KNN	✓	✗	Medium	Local patterns
Regression	Partial	✗	Medium	Linear relations
MICE	✓	Partial	High	Best practice

Rule of Thumb

- < 5% missing: Mean/median usually fine
- 5 – 20% missing: KNN or regression
- > 20% missing: MICE or consider dropping feature

Types of Duplicate Records

Exact Duplicates

All values identical across all columns:

$$x_i = x_j \Leftrightarrow \forall k : x_{ik} = x_{jk}$$

John	25	NYC
Jane	30	LA
John	25	NYC
Bob	22	CHI

Fuzzy/Near Duplicates

Similar but not identical (typos, variations):

John Smith	NYC
Jon Smith	New York
Jane Doe	LA

Different spellings, same entity!

Impact of Duplicates

- Inflate dataset size artificially
- Cause overfitting (same samples in train/test)

Detecting Exact Duplicates

Duplicate Detection Methods

1. Full Row Comparison:

$$\text{isDuplicate}(i, j) = \prod_{k=1}^p 1[x_{ik} = x_{jk}]$$

2. Key-Based Detection:

$$\text{isDuplicate}(i, j) = 1[\text{key}(x_i) = \text{key}(x_j)]$$

Where key might be (ID, timestamp) or composite columns.

Duplicate Rate

$$\text{Duplicate Rate} = \frac{n - n_{\text{unique}}}{n} \times 100\%$$

Fuzzy Matching: String Similarity Metrics

Levenshtein Distance (Edit Distance)

Minimum number of single-character edits (insert, delete, substitute):

$$\text{lev}(a, b) = \begin{cases} |a| & \text{if } |b| = 0 \\ |b| & \text{if } |a| = 0 \\ \text{lev}(\text{tail}(a), \text{tail}(b)) & \text{if } a[0] = b[0] \\ 1 + \min \begin{cases} \text{lev}(\text{tail}(a), b) \\ \text{lev}(a, \text{tail}(b)) \\ \text{lev}(\text{tail}(a), \text{tail}(b)) \end{cases} & \text{otherwise} \end{cases}$$

Example

“kitten” → “sitting”: Distance = 3

① kitten → sitten (substitute k→s)

② sitten → siting (delete t)

③ siting → sitting (insert i)

More Similarity Metrics

Jaro-Winkler Similarity

$$sim_j = \begin{cases} 0 & \text{if } m = 0 \\ \frac{1}{3} \left(\frac{m}{|s_1|} + \frac{m}{|s_2|} + \frac{m-t}{m} \right) & \text{otherwise} \end{cases}$$

Where m = matching characters, t = transpositions

$$sim_{jw} = sim_j + \ell \cdot p \cdot (1 - sim_j)$$

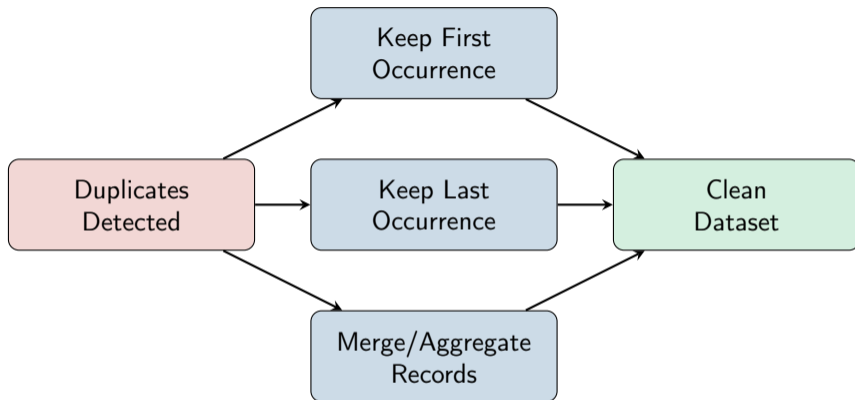
Where ℓ = common prefix length, p = scaling factor (usually 0.1)

Jaccard Similarity (for sets/tokens)

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

Example: {John, Smith} vs {Jon, Smith}

Duplicate Handling Strategies



Aggregation Functions for Merging

For numeric columns: mean, median, sum, min, max

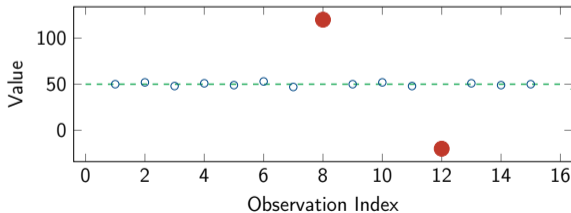
For categorical: mode, concatenate unique values

What are Outliers?

Definition

An **outlier** is a data point that differs significantly from other observations. It may indicate:

- Measurement/data entry errors
- Experimental errors
- Natural variation (rare but valid)
- Fraudulent behavior



Statistical Methods: Z-Score

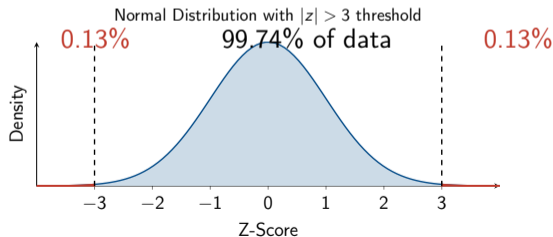
Z-Score Method

Measures how many standard deviations a point is from the mean:

$$z_i = \frac{x_i - \bar{x}}{s}$$

Outlier criterion:

$$|z_i| > \tau \quad (\text{typically } \tau = 3)$$



Modified Z-Score (Robust)

Problem with Standard Z-Score

Mean and standard deviation are sensitive to outliers themselves!

Modified Z-Score (Iglewicz & Hoaglin)

Uses median and MAD (Median Absolute Deviation):

$$M_i = \frac{0.6745 \cdot (x_i - \tilde{x})}{\text{MAD}}$$

Where:

$$\text{MAD} = \text{median}(|x_j - \tilde{x}|)$$

Outlier criterion: $|M_i| > 3.5$

Why 0.6745?

This constant makes MAD consistent with standard deviation for normal distributions:

IQR Method (Tukey's Fences)

Interquartile Range Method

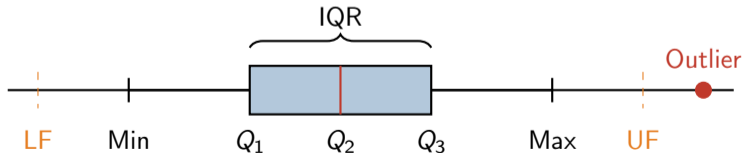
$$\text{IQR} = Q_3 - Q_1$$

Outlier boundaries:

$$\text{Lower Fence} = Q_1 - k \cdot \text{IQR}$$

$$\text{Upper Fence} = Q_3 + k \cdot \text{IQR}$$

Where $k = 1.5$ (outlier) or $k = 3$ (extreme outlier)



Multivariate Outlier Detection: Mahalanobis Distance

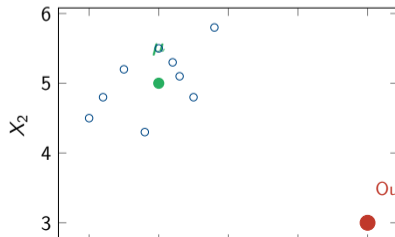
Mahalanobis Distance

Accounts for correlations between variables:

$$D_M(x_i) = \sqrt{(x_i - \mu)^T \Sigma^{-1} (x_i - \mu)}$$

Where:

- μ = mean vector
- Σ = covariance matrix



Machine Learning Methods: Isolation Forest

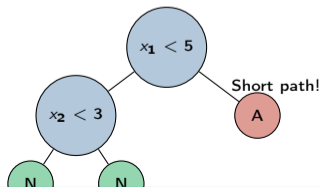
Isolation Forest Algorithm

Anomalies are easier to isolate (fewer splits needed):

$$s(x, n) = 2^{-\frac{E[h(x)]}{c(n)}}$$

Where:

- $h(x)$ = path length to isolate point x
- $c(n)$ = average path length in unsuccessful search in BST
- $s \approx 1$: anomaly, $s \approx 0.5$: normal, $s \approx 0$: very normal



Local Outlier Factor (LOF)

LOF Definition

Compares local density of a point to its neighbors:

$$\text{LOF}_k(x) = \frac{\sum_{o \in N_k(x)} \frac{\text{lrd}_k(o)}{\text{lrd}_k(x)}}{|N_k(x)|}$$

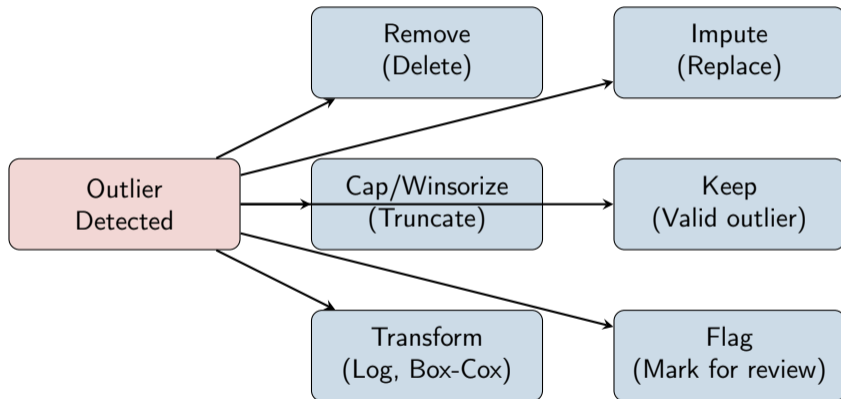
Where local reachability density:

$$\text{lrd}_k(x) = \frac{1}{\frac{\sum_{o \in N_k(x)} \text{reach-dist}_k(x,o)}{|N_k(x)|}}$$

Interpretation

- $\text{LOF} \approx 1$: Similar density to neighbors (normal)
- $\text{LOF} > 1$: Lower density than neighbors (potential outlier)
- $\text{LOF} \gg 1$: Much lower density (likely outlier)

Outlier Treatment Strategies



Winsorization (Capping)

Definition

Replace extreme values with specified percentiles:

$$x_i^{winsor} = \begin{cases} P_\alpha & \text{if } x_i < P_\alpha \\ P_{1-\alpha} & \text{if } x_i > P_{1-\alpha} \\ x_i & \text{otherwise} \end{cases}$$

Typically $\alpha = 0.01$ or $\alpha = 0.05$



Outlier Detection Summary

Method	Univariate	Multivariate	Best For
Z-Score	✓	✗	Normal distributions
Modified Z	✓	✗	Non-normal, robust
IQR/Tukey	✓	✗	Any distribution
Mahalanobis	✗	✓	Correlated features
Isolation Forest	✓	✓	High dimensions
LOF	✓	✓	Local anomalies
DBSCAN	✗	✓	Cluster-based

Important Consideration

Always investigate outliers before removing! They might be:

- Data entry errors (remove/correct)
- Genuine rare events (keep, possibly flag)
- Important signals (keep, they matter!)

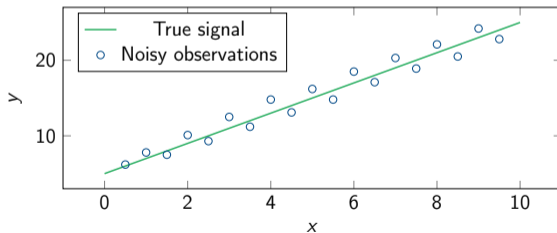
What is Noise in Data?

Definition

Noise is random error or variance in measured variables that obscures the true underlying pattern.

$$X_{observed} = X_{true} + \epsilon$$

Where $\epsilon \sim \mathcal{N}(0, \sigma^2)$ (random noise)



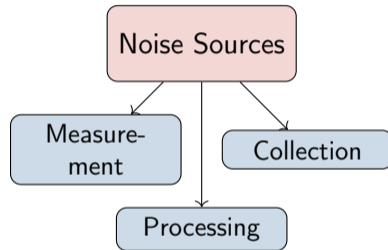
Sources of Noise

Measurement Noise

- Sensor inaccuracy
- Rounding errors
- Timing variations
- Environmental factors

Data Collection Noise

- Human errors
- Inconsistent processes
- Sampling bias



Binning (Discretization)

Smoothing by Binning

Partition data into bins and replace values with bin representative:

1. Bin Means:

$$x_i^{smooth} = \bar{x}_{bin(i)} = \frac{1}{|B_k|} \sum_{j \in B_k} x_j$$

2. Bin Medians:

$$x_i^{smooth} = \tilde{x}_{bin(i)}$$

3. Bin Boundaries:

$$x_i^{smooth} = \begin{cases} \min(B_k) & \text{if } x_i < \text{midpoint}(B_k) \\ \max(B_k) & \text{otherwise} \end{cases}$$

Binning Example

Example: Smoothing by Bin Means

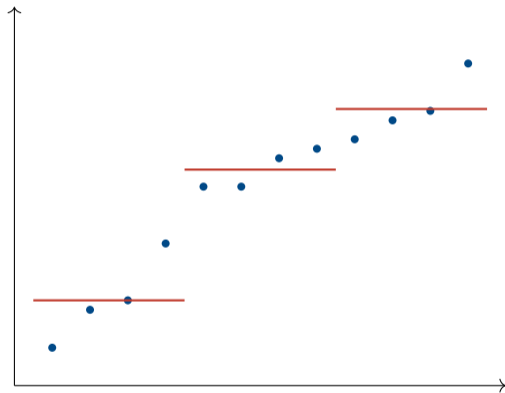
Original sorted data: 4, 8, 9, 15, 21, 21, 24, 25, 26, 28, 29, 34

Partition into bins (size 4):

- Bin 1: 4, 8, 9, 15 \rightarrow Mean = 9
- Bin 2: 21, 21, 24, 25 \rightarrow Mean = 22.75
- Bin 3: 26, 28, 29, 34 \rightarrow Mean = 29.25

Smoothed data: 9, 9, 9, 9, 22.75, 22.75, 22.75, 22.75, 29.25, 29.25, 29.25, 29.25

Binning Example



Moving Average (Rolling Mean)

Simple Moving Average (SMA)

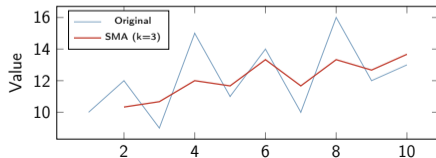
$$SMA_t = \frac{1}{k} \sum_{i=0}^{k-1} x_{t-i}$$

Where k is the window size.

Exponential Moving Average (EMA)

$$EMA_t = \alpha \cdot x_t + (1 - \alpha) \cdot EMA_{t-1}$$

Where $\alpha = \frac{2}{k+1}$ (smoothing factor)



Regression-Based Smoothing

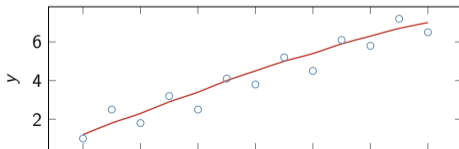
LOWESS/LOESS (Locally Weighted Scatterplot Smoothing)

Fits local weighted regression at each point:

$$\hat{y}_i = \sum_{j=1}^n w_{ij} \cdot y_j$$

Where weights decrease with distance (tricube kernel):

$$w_{ij} = \left(1 - \left| \frac{x_i - x_j}{d_i} \right|^3 \right)^3$$



Signal-to-Noise Ratio (SNR)

Definition

Measure of signal strength relative to background noise:

$$\text{SNR} = \frac{\mu_{\text{signal}}}{\sigma_{\text{noise}}}$$

In decibels:

$$\text{SNR}_{dB} = 10 \cdot \log_{10} \left(\frac{P_{\text{signal}}}{P_{\text{noise}}} \right) = 20 \cdot \log_{10} \left(\frac{A_{\text{signal}}}{A_{\text{noise}}} \right)$$

Quality Interpretation

- SNR > 20 dB: Excellent signal quality
- SNR 10-20 dB: Good signal quality
- SNR < 10 dB: Poor, may need filtering

Data Validation Rules

Types of Validation

- 1 **Type Validation:** Data matches expected type
- 2 **Range Validation:** Values within acceptable bounds
- 3 **Format Validation:** Correct format (dates, emails, etc.)
- 4 **Consistency Validation:** Cross-field logic checks
- 5 **Uniqueness Validation:** Primary keys are unique

Example Rules

- Age: $0 \leq \text{age} \leq 150$ (range)
- Email: matches regex pattern (format)
- Birth date $<$ Hire date (consistency)
- Employee ID is unique (uniqueness)

Key Metrics to Track

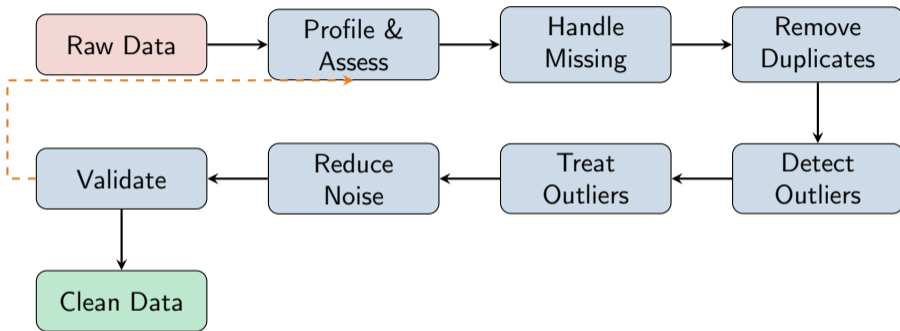
$$\text{Completeness} = \frac{\text{Non-null values}}{\text{Total values}} \times 100\%$$

$$\text{Validity} = \frac{\text{Values passing rules}}{\text{Total values}} \times 100\%$$

$$\text{Uniqueness} = \frac{\text{Distinct values}}{\text{Total values}} \times 100\%$$

$$\text{Consistency} = \frac{\text{Consistent records}}{\text{Total records}} \times 100\%$$

Data Cleaning Pipeline



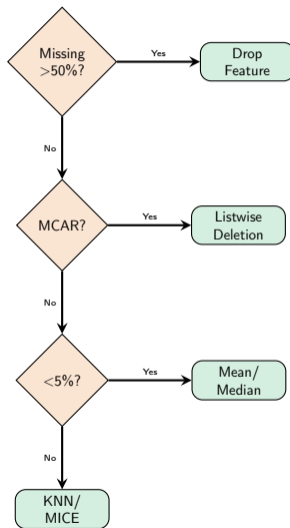
Remember

Data cleaning is **iterative**! Always validate results and be prepared to revisit earlier steps.

Data Cleaning Best Practices

- 1 **Document everything** — Keep logs of all changes
- 2 **Preserve original data** — Never modify source files
- 3 **Automate with scripts** — Reproducible pipelines
- 4 **Validate continuously** — Check quality at each step
- 5 **Understand the domain** — Context matters for decisions
- 6 **Handle edge cases** — Plan for unexpected values
- 7 **Test on samples first** — Verify before full processing
- 8 **Monitor data drift** — Ongoing quality checks

Decision Flowchart: Missing Values



Key Takeaways

Missing Values:

- Understand mechanism (MCAR/MAR/MNAR)
- Choose appropriate imputation
- MICE for best results

Duplicates:

- Check exact & fuzzy
- Use similarity metrics
- Document removal decisions

Outliers:

- Multiple detection methods
- Investigate before removing
- Consider domain context

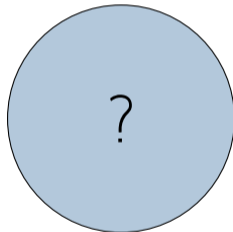
Noise:

- Smoothing techniques
- Balance signal preservation
- Monitor SNR

Next Part

Part 3: Data Transformation — Scaling, Encoding, and Feature Engineering

Questions & Discussion



Contact: professor@university.edu

References

-  Rubin, D.B. (1976). *Inference and Missing Data*. Biometrika.
-  Little, R.J.A. & Rubin, D.B. (2019). *Statistical Analysis with Missing Data*. Wiley.
-  Aggarwal, C.C. (2017). *Outlier Analysis*. Springer.
-  Han, J., Kamber, M., & Pei, J. (2011). *Data Mining: Concepts and Techniques*. Morgan Kaufmann.
-  Wickham, H. (2014). *Tidy Data*. Journal of Statistical Software.

Part 3: Transformation

"Shape your data"

Part III

Practice for Part 2 : Practice for Data Cleaning

What We Will Cover

Practical Data Cleaning Skills

- Detecting and handling missing values
- Identifying and removing duplicates
- Detecting and treating outliers
- Correcting data types
- Cleaning text and string data
- Handling inconsistent data
- Building cleaning pipelines

Libraries Used

pandas, numpy, scikit-learn, scipy, fuzzywuzzy

Essential Libraries Import

```
1 # Data manipulation
2 import pandas as pd
3 import numpy as np
4
5 # Visualization
6 import matplotlib.pyplot as plt
7 import seaborn as sns
8
9 # Scikit-learn imputers
10 from sklearn.impute import SimpleImputer, KNNImputer
11 from sklearn.experimental import enable_iterative_imputer
12 from sklearn.impute import IterativeImputer
13
14 # Statistical functions
15 from scipy import stats
16
17 # String matching (install: pip install fuzzywuzzy python-Levenshtein)
18 from fuzzywuzzy import fuzz, process
19
20 # Display settings
21 pd.set_option('display.max_columns', None)
```

Detecting Missing Values

```
1 # Check for missing values
2 print(df.isnull().sum())           # Count per column
3 print(df.isna().sum())           # Same as isnull()
4
5 # Total missing values
6 print(f"Total missing: {df.isnull().sum().sum()}")
7
8 # Percentage missing per column
9 missing_pct = (df.isnull().sum() / len(df) * 100).round(2)
10 print(missing_pct.sort_values(ascending=False))
11
12 # Rows with any missing value
13 rows_with_missing = df[df.isnull().any(axis=1)]
14 print(f"Rows with missing: {len(rows_with_missing)}")
15
16 # Columns with missing values
17 cols_with_missing = df.columns[df.isnull().any()].tolist()
18 print(f"Columns with missing: {cols_with_missing}")
```

Missing Values Summary Function

```
1 def missing_summary(df):
2     """Generate detailed missing values summary"""
3     missing = df.isnull().sum()
4     missing_pct = (missing / len(df) * 100).round(2)
5
6     summary = pd.DataFrame({
7         'Missing Count': missing,
8         'Missing %': missing_pct,
9         'Dtype': df.dtypes
10    })
11
12    # Filter only columns with missing values
13    summary = summary[summary['Missing Count'] > 0]
14    summary = summary.sort_values('Missing %', ascending=False)
15
16    print(f"Total columns: {len(df.columns)}")
17    print(f"Columns with missing: {len(summary)}")
18    print(f"Total missing values: {missing.sum()}")
19
20    return summary
21
22 print(missing_summary(df))
```

Visualizing Missing Values

```
1 # Using missingno library (pip install missingno)
2 import missingno as msno
3
4 # Matrix visualization
5 msno.matrix(df, figsize=(12, 6))
6 plt.title('Missing Values Matrix')
7 plt.show()
8
9 # Bar chart of missing values
10 msno.bar(df, figsize=(12, 6))
11 plt.show()
12
13 # Heatmap of missing value correlations
14 msno.heatmap(df, figsize=(10, 8))
15 plt.show()
16
17 # Dendrogram showing missing value patterns
18 msno.dendrogram(df, figsize=(12, 6))
19 plt.show()
```

Custom Missing Values Heatmap

```
1 # Create missing values heatmap with seaborn
2 plt.figure(figsize=(14, 8))
3 sns.heatmap(df.isnull(), cbar=True, yticklabels=False,
4             cmap='viridis')
5 plt.title('Missing Values Heatmap')
6 plt.tight_layout()
7 plt.show()
8
9 # Missing values by row (identify problematic rows)
10 row_missing = df.isnull().sum(axis=1)
11 plt.figure(figsize=(12, 4))
12 plt.hist(row_missing, bins=50, edgecolor='black')
13 plt.xlabel('Number of Missing Values per Row')
14 plt.ylabel('Frequency')
15 plt.title('Distribution of Missing Values per Row')
16 plt.show()
17
18 # Rows with too many missing values
19 threshold = len(df.columns) * 0.5 # 50% missing
20 bad_rows = df[row_missing > threshold]
21 print(f"Rows with >50% missing: {len(bad_rows)}")
```

Deletion Methods

```
1 # Drop rows with ANY missing value
2 df_clean = df.dropna()
3
4 # Drop rows with ALL missing values
5 df_clean = df.dropna(how='all')
6
7 # Drop rows with missing in specific columns
8 df_clean = df.dropna(subset=['col1', 'col2'])
9
10 # Drop rows with minimum non-null values (threshold)
11 df_clean = df.dropna(thresh=5) # Keep rows with at least 5 non-null
12
13 # Drop columns with missing values
14 df_clean = df.dropna(axis=1)
15
16 # Drop columns with too many missing values
17 threshold = len(df) * 0.5 # 50% missing threshold
18 cols_to_drop = df.columns[df.isnull().sum() > threshold]
19 df_clean = df.drop(columns=cols_to_drop)
20 print(f"Dropped columns: {cols_to_drop.tolist()}")
```

Simple Imputation - Numerical

```
1 # Fill with constant value
2 df['col'] = df['col'].fillna(0)
3
4 # Fill with mean
5 df['col'] = df['col'].fillna(df['col'].mean())
6
7 # Fill with median (robust to outliers)
8 df['col'] = df['col'].fillna(df['col'].median())
9
10 # Fill with mode (most frequent)
11 df['col'] = df['col'].fillna(df['col'].mode()[0])
12
13 # Fill multiple columns at once
14 numerical_cols = df.select_dtypes(include=np.number).columns
15 df[numerical_cols] = df[numerical_cols].fillna(
16     df[numerical_cols].median()
17 )
18
19 # Fill with column-specific values
20 fill_values = {'col1': 0, 'col2': df['col2'].mean(), 'col3': -1}
21 df = df.fillna(fill_values)
```

Simple Imputation - Categorical

```
1 # Fill with mode (most frequent category)
2 df['category'] = df['category'].fillna(df['category'].mode()[0])
3
4 # Fill with a specific category
5 df['category'] = df['category'].fillna('Unknown')
6 df['category'] = df['category'].fillna('Missing')
7
8 # Fill all categorical columns with mode
9 categorical_cols = df.select_dtypes(include=['object', 'category']).columns
10 for col in categorical_cols:
11     df[col] = df[col].fillna(df[col].mode()[0])
12
13 # Or fill with 'Unknown'
14 df[categorical_cols] = df[categorical_cols].fillna('Unknown')
15
16 # Conditional filling
17 df['category'] = df['category'].fillna(
18     df.groupby('other_col')['category'].transform(
19         lambda x: x.mode()[0] if len(x.mode()) > 0 else 'Unknown'
20     )
21 )
```

Forward and Backward Fill

```
1 # Forward fill (use previous value)
2 df['col'] = df['col'].fillna(method='ffill')
3 # Or using newer syntax:
4 df['col'] = df['col'].ffill()
5
6 # Backward fill (use next value)
7 df['col'] = df['col'].fillna(method='bfill')
8 # Or:
9 df['col'] = df['col'].bfill()
10
11 # Limit the fill (max consecutive NaNs to fill)
12 df['col'] = df['col'].ffill(limit=2)
13
14 # Combine forward and backward fill
15 df['col'] = df['col'].ffill().bfill()
16
17 # Useful for time series data
18 df = df.sort_values('date')
19 df['value'] = df['value'].ffill()
```

Interpolation Methods

```
1 # Linear interpolation (default)
2 df['col'] = df['col'].interpolate(method='linear')
3
4 # Polynomial interpolation
5 df['col'] = df['col'].interpolate(method='polynomial', order=2)
6
7 # Spline interpolation (smooth curves)
8 df['col'] = df['col'].interpolate(method='spline', order=3)
9
10 # Time-based interpolation (for datetime index)
11 df['col'] = df['col'].interpolate(method='time')
12
13 # Index-based interpolation
14 df['col'] = df['col'].interpolate(method='index')
15
16 # Limit interpolation direction
17 df['col'] = df['col'].interpolate(limit=2, limit_direction='forward')
18
19 # Only interpolate inside (not at edges)
20 df['col'] = df['col'].interpolate(limit_area='inside')
```

Scikit-learn SimpleImputer

```
1 from sklearn.impute import SimpleImputer
2
3 # Mean imputation
4 imputer_mean = SimpleImputer(strategy='mean')
5 df[numerical_cols] = imputer_mean.fit_transform(df[numerical_cols])
6
7 # Median imputation
8 imputer_median = SimpleImputer(strategy='median')
9 df[numerical_cols] = imputer_median.fit_transform(df[numerical_cols])
10
11 # Most frequent (mode) - works for categorical too
12 imputer_mode = SimpleImputer(strategy='most_frequent')
13 df[categorical_cols] = imputer_mode.fit_transform(df[categorical_cols])
14
15 # Constant value imputation
16 imputer_const = SimpleImputer(strategy='constant', fill_value=0)
17 df[numerical_cols] = imputer_const.fit_transform(df[numerical_cols])
18
19 # For categorical with constant
20 imputer_cat = SimpleImputer(strategy='constant', fill_value='Missing')
21 df[categorical_cols] = imputer_cat.fit_transform(df[categorical_cols])
```

KNN Imputation

```
1 from sklearn.impute import KNNImputer
2
3 # KNN imputer (uses k nearest neighbors)
4 knn_imputer = KNNImputer(n_neighbors=5, weights='uniform')
5 df[numerical_cols] = knn_imputer.fit_transform(df[numerical_cols])
6
7 # Weighted by distance (closer neighbors have more influence)
8 knn_imputer = KNNImputer(n_neighbors=5, weights='distance')
9 df[numerical_cols] = knn_imputer.fit_transform(df[numerical_cols])
10
11 # Custom number of neighbors
12 knn_imputer = KNNImputer(n_neighbors=10)
13
14 # Note: KNNImputer only works with numerical data
15 # For mixed data, encode categoricals first or use separate imputers
16
17 # Scale data before KNN imputation for better results
18 from sklearn.preprocessing import StandardScaler
19 scaler = StandardScaler()
20 df_scaled = scaler.fit_transform(df[numerical_cols])
21 df_imputed = knn_imputer.fit_transform(df_scaled)
22 df[numerical_cols] = scaler.inverse_transform(df_imputed)
```

Iterative Imputer (MICE)

```
1 from sklearn.experimental import enable_iterative_imputer
2 from sklearn.impute import IterativeImputer
3
4 # Basic iterative imputer (MICE - Multiple Imputation by Chained Equations)
5 iterative_imputer = IterativeImputer(max_iter=10, random_state=42)
6 df[numerical_cols] = iterative_imputer.fit_transform(df[numerical_cols])
7
8 # With specific estimator
9 from sklearn.linear_model import BayesianRidge
10 iterative_imputer = IterativeImputer(
11     estimator=BayesianRidge(),
12     max_iter=10,
13     random_state=42
14 )
15 df[numerical_cols] = iterative_imputer.fit_transform(df[numerical_cols])
16
17 # With Random Forest for non-linear relationships
18 from sklearn.ensemble import RandomForestRegressor
19 iterative_imputer = IterativeImputer(
20     estimator=RandomForestRegressor(n_estimators=10, random_state=42),
21     max_iter=10,
22     random_state=42
23 )
```

Group-Based Imputation

```
1 # Impute based on group statistics
2 # Example: Fill missing salary with median salary by department
3
4 # Using groupby transform
5 df['salary'] = df.groupby('department')['salary'].transform(
6     lambda x: x.fillna(x.median())
7 )
8
9 # Multiple group columns
10 df['salary'] = df.groupby(['department', 'level'])['salary'].transform(
11     lambda x: x.fillna(x.median())
12 )
13
14 # Fallback to global median if group has all NaN
15 def group_impute(group):
16     if group.isnull().all():
17         return group.fillna(df['salary'].median())
18     return group.fillna(group.median())
19
20 df['salary'] = df.groupby('department')['salary'].transform(group_impute)
21
22 # Custom function for complex logic
23 df['value'] = df.groupby('category')['value'].apply(
24     lambda x: x.fillna(x.median() if x.notna().sum() > 0 else 0)
25 )
```

Imputation with Indicator

```
1 # Create indicator column for missing values (useful for modeling)
2 df['col_was_missing'] = df['col'].isnull().astype(int)
3
4 # Then impute
5 df['col'] = df['col'].fillna(df['col'].median())
6
7 # Function to add missing indicators
8 def add_missing_indicators(df, columns):
9     for col in columns:
10         if df[col].isnull().sum() > 0:
11             df[f'{col}_missing'] = df[col].isnull().astype(int)
12     return df
13
14 df = add_missing_indicators(df, numerical_cols)
15
16 # Using SimpleImputer with add_indicator
17 from sklearn.impute import SimpleImputer
18 imputer = SimpleImputer(strategy='median', add_indicator=True)
19 result = imputer.fit_transform(df[numerical_cols])
20 # This returns both imputed values and indicator columns
```

Detecting Duplicates

```
1 # Check for duplicate rows
2 print(f"Duplicate rows: {df.duplicated().sum()}")
3
4 # Check duplicates based on specific columns
5 print(f"Duplicate keys: {df.duplicated(subset=['id', 'date']).sum()}")
6
7 # View duplicate rows
8 duplicates = df[df.duplicated(keep=False)] # keep=False shows all duplicates
9 print(duplicates.sort_values(by=df.columns.tolist()))
10
11 # First occurrence vs duplicates
12 first_occurrence = df[df.duplicated(keep='first') == False]
13 duplicate_rows = df[df.duplicated(keep='first') == True]
14
15 # Count duplicates per group
16 dup_counts = df.groupby(df.columns.tolist()).size().reset_index(name='count')
17 dup_counts = dup_counts[dup_counts['count'] > 1]
18 print(f"Groups with duplicates: {len(dup_counts)}")
```

Removing Duplicates

```
1 # Remove exact duplicate rows (keep first occurrence)
2 df_clean = df.drop_duplicates()
3
4 # Keep last occurrence instead
5 df_clean = df.drop_duplicates(keep='last')
6
7 # Remove all duplicates (don't keep any)
8 df_clean = df.drop_duplicates(keep=False)
9
10 # Remove duplicates based on specific columns
11 df_clean = df.drop_duplicates(subset=['id', 'date'])
12
13 # Remove duplicates and reset index
14 df_clean = df.drop_duplicates().reset_index(drop=True)
15
16 # In-place removal
17 df.drop_duplicates(inplace=True)
18
19 # Keep row with most non-null values among duplicates
20 df_clean = df.loc[df.groupby(['id'])['col'].transform('count').idxmax()]
```

Handling Duplicates with Aggregation

```
1 # Instead of dropping, aggregate duplicate rows
2 # Keep first value for some columns, mean for others
3
4 agg_rules = {
5     'name': 'first',          # Keep first name
6     'value': 'mean',         # Average the values
7     'count': 'sum',          # Sum the counts
8     'date': 'max',           # Keep latest date
9     'category': lambda x: x.mode()[0] # Most frequent
10 }
11
12 df_agg = df.groupby('id').agg(agg_rules).reset_index()
13
14 # Custom aggregation for complex cases
15 def resolve_duplicates(group):
16     return pd.Series({
17         'name': group['name'].iloc[0],
18         'value': group['value'].mean(),
19         'latest_date': group['date'].max(),
20         'all_categories': ', '.join(group['category'].unique())
21     })
22
23 df_resolved = df.groupby('id').apply(resolve_duplicates).reset_index()
```

Fuzzy Duplicate Detection

```
1 # For near-duplicates (e.g., "John Smith" vs "Jon Smith")
2 from fuzzywuzzy import fuzz, process
3
4 # Compare two strings
5 similarity = fuzz.ratio("John Smith", "Jon Smith")
6 print(f"Similarity: {similarity}%") # Output: 91%
7
8 # Find potential duplicates in a column
9 def find_fuzzy_duplicates(df, column, threshold=85):
10     duplicates = []
11     values = df[column].unique()
12
13     for i, val1 in enumerate(values):
14         for val2 in values[i+1:]:
15             score = fuzz.ratio(str(val1), str(val2))
16             if score >= threshold:
17                 duplicates.append((val1, val2, score))
18
19     return pd.DataFrame(duplicates, columns=['Value1', 'Value2', 'Score'])
20
21 fuzzy_dups = find_fuzzy_duplicates(df, 'name', threshold=85)
22 print(fuzzy_dups)
```

Standardizing Fuzzy Duplicates

```
1 # Standardize similar values to one canonical form
2 from fuzzywuzzy import process
3
4 def standardize_column(df, column, threshold=85):
5     """Standardize similar values to most frequent form"""
6     value_counts = df[column].value_counts()
7     mapping = {}
8
9     for value in value_counts.index:
10        if value not in mapping:
11            # Find best match among more frequent values
12            choices = [v for v in value_counts.index if v not in mapping.values()]
13            match, score = process.extractOne(str(value), choices)
14            if score >= threshold:
15                mapping[value] = match
16            else:
17                mapping[value] = value
18
19    return df[column].map(mapping)
20
21 df['name_clean'] = standardize_column(df, 'name', threshold=85)
22
23 # Check the mapping
24 print(df[['name', 'name_clean']].drop_duplicates())
```

IQR Method

```
1 def detect_outliers_iqr(df, column, multiplier=1.5):
2     """Detect outliers using IQR method"""
3     Q1 = df[column].quantile(0.25)
4     Q3 = df[column].quantile(0.75)
5     IQR = Q3 - Q1
6
7     lower_bound = Q1 - multiplier * IQR
8     upper_bound = Q3 + multiplier * IQR
9
10    outliers = (df[column] < lower_bound) | (df[column] > upper_bound)
11
12    return outliers, lower_bound, upper_bound
13
14 # Apply to column
15 outlier_mask, lower, upper = detect_outliers_iqr(df, 'value')
16 print(f"Outliers: {outlier_mask.sum()}")
17 print(f"Bounds: [{lower:.2f}, {upper:.2f}]")
18
19 # View outliers
20 print(df[outlier_mask])
21
22 # Apply to all numerical columns
23 for col in numerical_cols:
24     mask, _, _ = detect_outliers_iqr(df, col)
25     print(f"{col}: {mask.sum()} outliers ({mask.mean()*100:.1f}%")
```

Z-Score Method

```
1 from scipy import stats
2
3 def detect_outliers_zscore(df, column, threshold=3):
4     """Detect outliers using Z-score method"""
5     z_scores = np.abs(stats.zscore(df[column].dropna()))
6     outliers = z_scores > threshold
7     return outliers
8
9 # Apply to column
10 outlier_mask = detect_outliers_zscore(df, 'value', threshold=3)
11 print(f"Outliers (|z| > 3): {outlier_mask.sum()}")
12
13 # Calculate z-scores for all rows
14 df['z_score'] = np.abs(stats.zscore(df['value']))
15 df['is_outlier'] = df['z_score'] > 3
16
17 # Modified Z-score (more robust, uses median)
18 def modified_zscore(x):
19     median = np.median(x)
20     mad = np.median(np.abs(x - median))
21     return 0.6745 * (x - median) / mad
22
23 df['modified_z'] = np.abs(modified_zscore(df['value']))
24 df['is_outlier_robust'] = df['modified_z'] > 3.5
```

Isolation Forest

```
1 from sklearn.ensemble import IsolationForest
2
3 # Isolation Forest for multivariate outlier detection
4 iso_forest = IsolationForest(
5     n_estimators=100,
6     contamination=0.05, # Expected proportion of outliers
7     random_state=42
8 )
9
10 # Fit and predict (-1 for outliers, 1 for inliers)
11 df['outlier_if'] = iso_forest.fit_predict(df[numerical_cols])
12 df['is_outlier_if'] = df['outlier_if'] == -1
13
14 print(f"Outliers detected: {(df['outlier_if'] == -1).sum()}")
15
16 # Get anomaly scores (lower = more anomalous)
17 df['anomaly_score'] = iso_forest.decision_function(df[numerical_cols])
18
19 # Visualize
20 plt.figure(figsize=(10, 6))
21 plt.scatter(df.index, df['anomaly_score'], c=df['outlier_if'], cmap='RdYlGn')
22 plt.xlabel('Index')
23 plt.ylabel('Anomaly Score')
24 plt.title('Isolation Forest Anomaly Scores')
25 plt.show()
```

Local Outlier Factor (LOF)

```
1 from sklearn.neighbors import LocalOutlierFactor
2
3 # LOF for density-based outlier detection
4 lof = LocalOutlierFactor(
5     n_neighbors=20,
6     contamination=0.05
7 )
8
9 # Fit and predict
10 df['outlier_lof'] = lof.fit_predict(df[numerical_cols])
11 df['is_outlier_lof'] = df['outlier_lof'] == -1
12
13 print(f"LOF outliers: {(df['outlier_lof'] == -1).sum()}")
14
15 # Get LOF scores (negative_outlier_factor_)
16 df['lof_score'] = -lof.negative_outlier_factor_
17
18 # Visualize 2D projection with outliers
19 plt.figure(figsize=(10, 8))
20 colors = ['red' if x == -1 else 'blue' for x in df['outlier_lof']]
21 plt.scatter(df['col1'], df['col2'], c=colors, alpha=0.5)
22 plt.xlabel('Column 1')
23 plt.ylabel('Column 2')
24 plt.title('LOF Outlier Detection')
25 plt.show()
```

DBSCAN for Outlier Detection

```
1 from sklearn.cluster import DBSCAN
2 from sklearn.preprocessing import StandardScaler
3
4 # Scale data first
5 scaler = StandardScaler()
6 X_scaled = scaler.fit_transform(df[numerical_cols])
7
8 # DBSCAN clustering
9 dbscan = DBSCAN(eps=0.5, min_samples=5)
10 df['cluster'] = dbscan.fit_predict(X_scaled)
11
12 # Points labeled -1 are outliers (noise)
13 df['is_outlier_dbscan'] = df['cluster'] == -1
14
15 print(f"DBSCAN outliers: {(df['cluster'] == -1).sum()}")
16 print(f"Number of clusters: {df['cluster'].max() + 1}")
17
18 # Visualize
19 plt.figure(figsize=(10, 8))
20 scatter = plt.scatter(df['col1'], df['col2'], c=df['cluster'], cmap='viridis')
21 plt.colorbar(scatter)
22 plt.title('DBSCAN Clustering (-1 = Outliers)')
23 plt.show()
```

Removing Outliers

```
1 # Remove outliers detected by IQR
2 outlier_mask, lower, upper = detect_outliers_iqr(df, 'value')
3 df_clean = df[~outlier_mask]
4 print(f"Removed {outlier_mask.sum()} outliers")
5
6 # Remove outliers from multiple columns
7 def remove_outliers_iqr(df, columns, multiplier=1.5):
8     mask = pd.Series([False] * len(df))
9     for col in columns:
10         outliers, _, _ = detect_outliers_iqr(df, col, multiplier)
11         mask = mask | outliers
12     return df[~mask]
13
14 df_clean = remove_outliers_iqr(df, numerical_cols)
15
16 # Remove based on Z-score
17 z_scores = np.abs(stats.zscore(df[numerical_cols]))
18 df_clean = df[(z_scores < 3).all(axis=1)]
19
20 # Remove based on percentiles
21 lower = df['value'].quantile(0.01)
22 upper = df['value'].quantile(0.99)
23 df_clean = df[(df['value'] >= lower) & (df['value'] <= upper)]
```

Capping / Winsorization

```
1 # Cap outliers to boundary values (Winsorization)
2 def cap_outliers_iqr(df, column, multiplier=1.5):
3     Q1 = df[column].quantile(0.25)
4     Q3 = df[column].quantile(0.75)
5     IQR = Q3 - Q1
6     lower = Q1 - multiplier * IQR
7     upper = Q3 + multiplier * IQR
8
9     df[column] = df[column].clip(lower=lower, upper=upper)
10    return df
11
12 df = cap_outliers_iqr(df, 'value')
13
14 # Cap using percentiles
15 def cap_percentile(df, column, lower_pct=0.01, upper_pct=0.99):
16     lower = df[column].quantile(lower_pct)
17     upper = df[column].quantile(upper_pct)
18     df[column] = df[column].clip(lower=lower, upper=upper)
19     return df
20
21 # Using scipy's winsorize
22 from scipy.stats.mstats import winsorize
23 df['value_winsorized'] = winsorize(df['value'], limits=[0.05, 0.05])
```

Transformation Methods

```
1 # Log transformation (for right-skewed data)
2 df['value_log'] = np.log1p(df['value']) # log(1+x) handles zeros
3
4 # Square root transformation
5 df['value_sqrt'] = np.sqrt(df['value'])
6
7 # Box-Cox transformation (requires positive values)
8 from scipy.stats import boxcox
9 df['value_boxcox'], lambda_param = boxcox(df['value'] + 1)
10 print(f"Optimal lambda: {lambda_param}")
11
12 # Yeo-Johnson transformation (handles negative values)
13 from sklearn.preprocessing import PowerTransformer
14 pt = PowerTransformer(method='yeo-johnson')
15 df['value_yeojohnson'] = pt.fit_transform(df[['value']])
16
17 # Compare distributions
18 fig, axes = plt.subplots(1, 3, figsize=(15, 4))
19 df['value'].hist(ax=axes[0], bins=30)
20 axes[0].set_title('Original')
21 df['value_log'].hist(ax=axes[1], bins=30)
22 axes[1].set_title('Log Transform')
23 df['value_boxcox'].hist(ax=axes[2], bins=30)
24 axes[2].set_title('Box-Cox Transform')
25 plt.tight_layout()
```

Imputing Outliers

```
1 # Replace outliers with median
2 def impute_outliers_median(df, column, multiplier=1.5):
3     Q1 = df[column].quantile(0.25)
4     Q3 = df[column].quantile(0.75)
5     IQR = Q3 - Q1
6     lower = Q1 - multiplier * IQR
7     upper = Q3 + multiplier * IQR
8     median = df[column].median()
9
10    df[column] = df[column].apply(
11        lambda x: median if (x < lower or x > upper) else x
12    )
13    return df
14
15 # Replace outliers with NaN, then impute
16 def outliers_to_nan(df, column, multiplier=1.5):
17     mask, lower, upper = detect_outliers_iqr(df, column, multiplier)
18     df.loc[mask, column] = np.nan
19     return df
20
21 df = outliers_to_nan(df, 'value')
22 df['value'] = df['value'].fillna(df['value'].median())
```

Converting Data Types

```
1 # Check current data types
2 print(df.dtypes)
3
4 # Convert to numeric (coerce errors to NaN)
5 df['col'] = pd.to_numeric(df['col'], errors='coerce')
6
7 # Convert to integer (must handle NaN first)
8 df['col'] = df['col'].fillna(0).astype(int)
9
10 # Nullable integer type (handles NaN)
11 df['col'] = df['col'].astype('Int64') # Note: capital I
12
13 # Convert to float
14 df['col'] = df['col'].astype(float)
15
16 # Convert to string
17 df['col'] = df['col'].astype(str)
18
19 # Convert to category (memory efficient)
20 df['category_col'] = df['category_col'].astype('category')
21
22 # Convert to boolean
23 df['bool_col'] = df['bool_col'].astype(bool)
24 df['bool_col'] = df['col'].map({'Yes': True, 'No': False})
```

Converting Datetime

```
1 # Convert to datetime
2 df['date'] = pd.to_datetime(df['date'])
3
4 # Specify format for faster parsing
5 df['date'] = pd.to_datetime(df['date'], format='%Y-%m-%d')
6
7 # Handle different formats
8 df['date'] = pd.to_datetime(df['date'], format='mixed', dayfirst=True)
9
10 # Handle errors
11 df['date'] = pd.to_datetime(df['date'], errors='coerce') # Invalid -> NaT
12
13 # Extract components
14 df['year'] = df['date'].dt.year
15 df['month'] = df['date'].dt.month
16 df['day'] = df['date'].dt.day
17 df['dayofweek'] = df['date'].dt.dayofweek
18 df['hour'] = df['date'].dt.hour
19
20 # Convert Unix timestamp to datetime
21 df['date'] = pd.to_datetime(df['timestamp'], unit='s')
22
23 # Convert datetime to string
24 df['date_str'] = df['date'].dt.strftime('%Y-%m-%d')
```

Automatic Type Inference

```
1 # Infer better dtypes automatically
2 df = df.infer_objects()
3
4 # Convert object columns to best possible types
5 df = df.convert_dtypes()
6
7 # Custom type conversion function
8 def optimize_dtypes(df):
9     for col in df.columns:
10         col_type = df[col].dtype
11
12         if col_type == 'object':
13             # Try to convert to numeric
14             try:
15                 df[col] = pd.to_numeric(df[col])
16             except:
17                 # Try to convert to datetime
18                 try:
19                     df[col] = pd.to_datetime(df[col])
20                 except:
21                     # Convert to category if low cardinality
22                     if df[col].nunique() / len(df) < 0.5:
23                         df[col] = df[col].astype('category')
24
25     return df
26 df = optimize_dtypes(df)
```

Basic String Cleaning

```
1 # Remove leading/trailing whitespace
2 df['col'] = df['col'].str.strip()
3
4 # Remove extra whitespace
5 df['col'] = df['col'].str.replace(r'\s+', ' ', regex=True)
6
7 # Convert to lowercase/uppercase
8 df['col'] = df['col'].str.lower()
9 df['col'] = df['col'].str.upper()
10 df['col'] = df['col'].str.title() # Title Case
11
12 # Remove specific characters
13 df['col'] = df['col'].str.replace('[^\w\s]', '', regex=True) # Keep only alphanumeric
14
15 # Remove digits
16 df['col'] = df['col'].str.replace(r'\d+', '', regex=True)
17
18 # Remove specific strings
19 df['col'] = df['col'].str.replace('unwanted', '', regex=False)
20
21 # Replace multiple values
22 df['col'] = df['col'].replace({'old1': 'new1', 'old2': 'new2'})
```

Advanced String Operations

```
1 # Extract patterns with regex
2 df['phone'] = df['text'].str.extract(r'(\d{3}-\d{3}-\d{4})')
3 df['email'] = df['text'].str.extract(r'([\w.-]+@[w.-]+)')
4
5 # Split strings
6 df[['first_name', 'last_name']] = df['full_name'].str.split(' ', n=1, expand=True)
7
8 # Check for patterns
9 df['has_email'] = df['text'].str.contains(r'@', regex=True)
10 df['starts_with_a'] = df['col'].str.startswith('A')
11
12 # Pad strings
13 df['code'] = df['code'].str.zfill(5) # Pad with zeros: '42' -> '00042'
14
15 # Slice strings
16 df['first_3'] = df['col'].str[:3]
17 df['last_4'] = df['col'].str[-4:]
18
19 # Get string length
20 df['length'] = df['col'].str.len()
21
22 # Count occurrences
23 df['word_count'] = df['text'].str.count(r'\w+')
```

Standardizing Categorical Values

```
1 # Map inconsistent values to standard forms
2 mapping = {
3     'yes': 'Yes', 'YES': 'Yes', 'y': 'Yes', 'Y': 'Yes', '1': 'Yes',
4     'no': 'No', 'NO': 'No', 'n': 'No', 'N': 'No', '0': 'No'
5 }
6 df['response'] = df['response'].map(mapping).fillna(df['response'])
7
8 # Using replace with regex
9 df['col'] = df['col'].replace({
10     r'^[Yy]es.*': 'Yes',
11     r'^[Nn]o.*': 'No'
12 }, regex=True)
13
14 # Standardize case and strip
15 def standardize_text(text):
16     if pd.isna(text):
17         return text
18     return str(text).strip().lower().title()
19
20 df['name'] = df['name'].apply(standardize_text)
21
22 # Fix common typos
23 typo_fixes = {'teh': 'the', 'recieve': 'receive', 'seperate': 'separate'}
24 for typo, fix in typo_fixes.items():
25     df['text'] = df['text'].str.replace(typo, fix, case=False)
```

Complete Cleaning Function

```
1 def clean_dataframe(df):
2     """Comprehensive data cleaning pipeline"""
3     df = df.copy()
4
5     # 1. Remove duplicate rows
6     initial_rows = len(df)
7     df = df.drop_duplicates()
8     print(f"Removed {initial_rows - len(df)} duplicate rows")
9
10    # 2. Handle column names
11    df.columns = df.columns.str.strip().str.lower().str.replace(' ', '_')
12
13    # 3. Identify column types
14    num_cols = df.select_dtypes(include=np.number).columns.tolist()
15    cat_cols = df.select_dtypes(include='object').columns.tolist()
16
17    # 4. Clean string columns
18    for col in cat_cols:
19        df[col] = df[col].str.strip()
20
21    # 5. Handle missing values
22    # Numerical: median imputation
23    for col in num_cols:
24        df[col] = df[col].fillna(df[col].median())
25
26    return df
```

Sklearn Pipeline for Cleaning

```
1 from sklearn.pipeline import Pipeline
2 from sklearn.compose import ColumnTransformer
3 from sklearn.impute import SimpleImputer
4 from sklearn.preprocessing import StandardScaler, OneHotEncoder
5
6 # Define preprocessing for different column types
7 numerical_pipeline = Pipeline([
8     ('imputer', SimpleImputer(strategy='median')),
9     ('scaler', StandardScaler())
10 ])
11
12 categorical_pipeline = Pipeline([
13     ('imputer', SimpleImputer(strategy='constant', fill_value='Missing')),
14     ('encoder', OneHotEncoder(handle_unknown='ignore', sparse_output=False))
15 ])
16
17 # Combine pipelines
18 preprocessor = ColumnTransformer([
19     ('num', numerical_pipeline, numerical_cols),
20     ('cat', categorical_pipeline, categorical_cols)
21 ])
22
23 # Fit and transform
24 X_clean = preprocessor.fit_transform(df)
```

Custom Transformer

```
1 from sklearn.base import BaseEstimator, TransformerMixin
2
3 class OutlierRemover(BaseEstimator, TransformerMixin):
4     def __init__(self, multiplier=1.5):
5         self.multiplier = multiplier
6         self.bounds_ = {}
7
8     def fit(self, X, y=None):
9         for col in X.columns:
10            Q1, Q3 = X[col].quantile([0.25, 0.75])
11            IQR = Q3 - Q1
12            self.bounds_[col] = (Q1 - self.multiplier * IQR,
13                               Q3 + self.multiplier * IQR)
14
15        return self
16
17    def transform(self, X):
18        X = X.copy()
19        for col in X.columns:
20            lower, upper = self.bounds_[col]
21            X[col] = X[col].clip(lower=lower, upper=upper)
22
23        return X
24
25 # Use in pipeline
26 pipeline = Pipeline([
27     ('outlier', OutlierRemover(multiplier=1.5)),
28     ('imputer', SimpleImputer(strategy='median'))
29 ])
```

Saving and Loading Clean Data

```
1 # Save cleaned data
2 df_clean.to_csv('data_cleaned.csv', index=False)
3 df_clean.to_parquet('data_cleaned.parquet', index=False) # Better for large files
4 df_clean.to_pickle('data_cleaned.pkl') # Preserves dtypes
5
6 # Save with compression
7 df_clean.to_csv('data_cleaned.csv.gz', compression='gzip', index=False)
8
9 # Load cleaned data
10 df = pd.read_csv('data_cleaned.csv')
11 df = pd.read_parquet('data_cleaned.parquet')
12 df = pd.read_pickle('data_cleaned.pkl')
13
14 # Save preprocessing pipeline (for production)
15 import joblib
16 joblib.dump(preprocessor, 'preprocessor.pkl')
17
18 # Load and apply to new data
19 preprocessor = joblib.load('preprocessor.pkl')
20 X_new_clean = preprocessor.transform(df_new)
```

Data Cleaning Checklist

Missing Values:

- Identify missing patterns
- Decide: delete or impute
- Apply appropriate imputation
- Add missing indicators if needed

Duplicates:

- Check exact duplicates
- Check key-based duplicates
- Check fuzzy duplicates
- Remove or aggregate

Outliers:

- Detect with IQR/Z-score
- Verify if real or errors
- Remove, cap, or transform

Data Types & Text:

- Fix incorrect types
- Parse dates properly
- Clean string columns
- Standardize categories

Key Python Functions Reference

Task	Function
Detect missing	<code>df.isnull().sum(), msno.matrix()</code>
Drop missing	<code>df.dropna(), df.dropna(subset=[...])</code>
Fill missing	<code>df.fillna(), df.interpolate()</code>
Impute (sklearn)	<code>SimpleImputer, KNNImputer</code>
Detect duplicates	<code>df.duplicated(), fuzz.ratio()</code>
Remove duplicates	<code>df.drop_duplicates()</code>
Detect outliers	<code>zscore(), IsolationForest</code>
Cap outliers	<code>df.clip(), winsorize()</code>
Convert types	<code>astype(), pd.to_numeric()</code>
Clean strings	<code>str.strip(), str.replace()</code>

Questions & Practice

Next: Part 3 - Data Transformation with Python

Part IV

Part 3 : Data Transformation

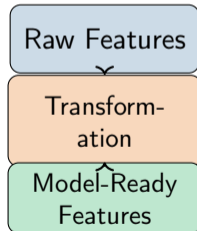
What is Data Transformation?

Definition

Data Transformation is the process of converting data from one format, structure, or value range to another to make it suitable for machine learning algorithms.

Main Objectives:

- Normalize feature scales
- Encode categorical variables
- Create new informative features
- Reduce dimensionality
- Handle non-linear relationships



Why Transform Data?

Scale Difference

Age: 0-100
Income: 0-1M

Categorical Data

Color: Red, Blue
ML needs numbers

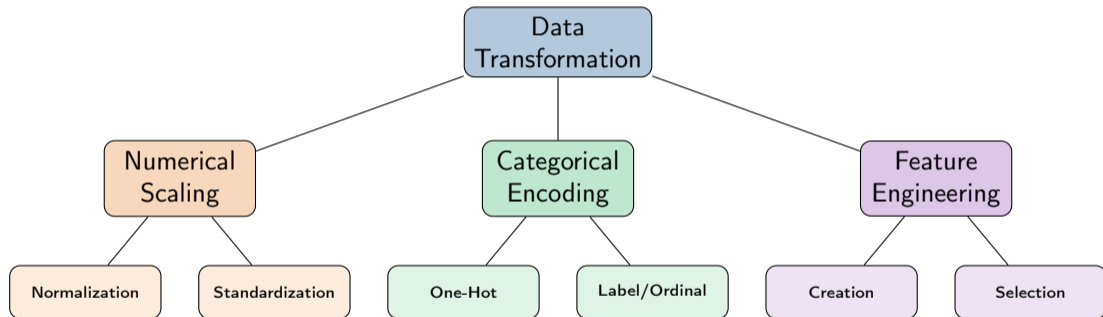
Non-linearity

Exponential growth
Linear models fail

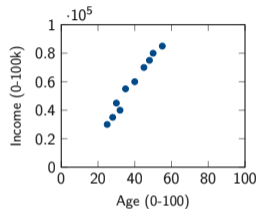
Impact on ML Algorithms

- **Gradient-based methods** (Neural Networks, Logistic Regression): Sensitive to feature scales
- **Distance-based methods** (KNN, K-Means, SVM): Dominated by large-scale features
- **Tree-based methods** (Random Forest, XGBoost): Generally scale-invariant

Transformation Categories Overview

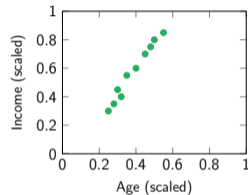


Why Scale Numerical Features?



Before Scaling

Income dominates distance!



After Scaling

Equal contribution!

Euclidean Distance Example

$$d = \sqrt{(100 - 0)^2 + (100000 - 0)^2} \approx 100000 \quad (\text{Income dominates!})$$

Min-Max Normalization (Rescaling)

Formula

Scales features to a fixed range $[a, b]$ (typically $[0, 1]$):

$$x'_i = \frac{x_i - x_{min}}{x_{max} - x_{min}} \cdot (b - a) + a$$

For $[0, 1]$ range:

$$x'_i = \frac{x_i - x_{min}}{x_{max} - x_{min}}$$

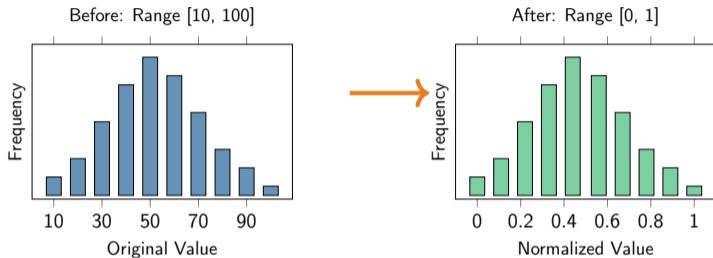
Properties:

- Output range: $[0, 1]$ or $[a, b]$
- Preserves zero values (if $a = 0$)
- Preserves relationships

Limitations:

- Sensitive to outliers
- New data may fall outside $[0, 1]$
- Doesn't center data

Min-Max Normalization: Visual Example



Calculation Example

$$x = 55, x_{min} = 10, x_{max} = 100$$

$$x' = \frac{55 - 10}{100 - 10} = \frac{45}{90} = 0.5$$

Z-Score Standardization

Formula

Transforms data to have mean $\mu = 0$ and standard deviation $\sigma = 1$:

$$z_i = \frac{x_i - \mu}{\sigma} = \frac{x_i - \bar{x}}{s}$$

Where:

- $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ (sample mean)
- $s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$ (sample std dev)

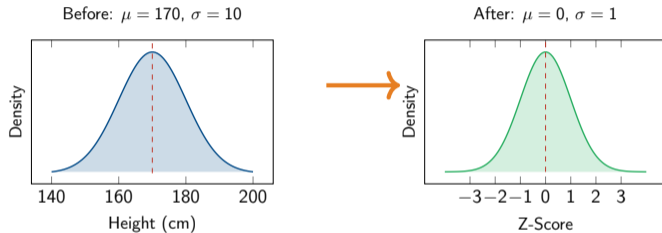
Properties:

- Centered at zero
- Unit variance
- Unbounded range

Best for:

- Gaussian-like distributions
- Algorithms assuming $\mathcal{N}(0, 1)$
- When outliers present

Z-Score Standardization: Visual



Interpretation

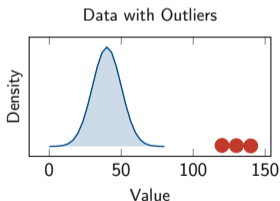
A z-score of $z = 2$ means the value is 2 standard deviations above the mean.

Robust Scaling

Formula

Uses median and IQR (robust to outliers):

$$x'_i = \frac{x_i - \text{median}(x)}{\text{IQR}} = \frac{x_i - Q_2}{Q_3 - Q_1}$$



Z-Score: Affected by outliers
 σ inflated \rightarrow poor scaling

Robust: Uses median/IQR
Outliers have minimal impact

MaxAbs Scaling

Formula

Scales by maximum absolute value (preserves sparsity):

$$x'_i = \frac{x_i}{\max(|x|)} = \frac{x_i}{|x_{max}|}$$

Result range: $[-1, 1]$

Use Case: Sparse Data

- Preserves zero entries (important for sparse matrices)
- No centering (doesn't destroy sparsity)
- Useful for data already centered at zero

Example

Data: $[-10, -5, 0, 3, 8]$, $\max(|x|) = 10$

Power Transformations: Box-Cox

Box-Cox Transformation

Makes data more Gaussian-like:

$$x'_i = \begin{cases} \frac{x_i^\lambda - 1}{\lambda} & \text{if } \lambda \neq 0 \\ \ln(x_i) & \text{if } \lambda = 0 \end{cases}$$

Constraint: $x_i > 0$ (all values must be positive)

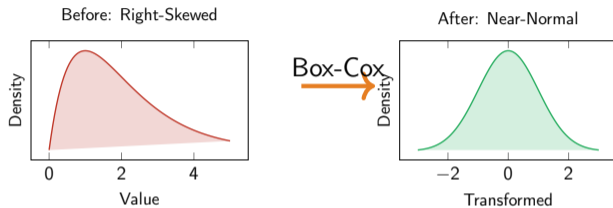
Special Cases:

- $\lambda = 1$: No transformation
- $\lambda = 0$: Log transform
- $\lambda = 0.5$: Square root
- $\lambda = -1$: Reciprocal

λ Selection:

- Maximum Likelihood Estimation
- Optimizes normality
- Typically $\lambda \in [-5, 5]$

Box-Cox Transformation: Visual



When to Use

- Linear regression (assumes normal residuals)
- Algorithms sensitive to distribution shape
- Reducing heteroscedasticity

Yeo-Johnson Transformation

Formula (handles negative values)

$$x'_i = \begin{cases} \frac{(x_i+1)^\lambda - 1}{\lambda} & \text{if } \lambda \neq 0, x_i \geq 0 \\ \ln(x_i + 1) & \text{if } \lambda = 0, x_i \geq 0 \\ \frac{-((-x_i+1)^{2-\lambda} - 1)}{2-\lambda} & \text{if } \lambda \neq 2, x_i < 0 \\ -\ln(-x_i + 1) & \text{if } \lambda = 2, x_i < 0 \end{cases}$$

Advantage over Box-Cox

- Works with zero and negative values
- No need to shift data first
- More flexible for real-world data

Log and Square Root Transformations

Log Transformation

$$x' = \log(x + c)$$

Where c ensures $x + c > 0$

Use for:

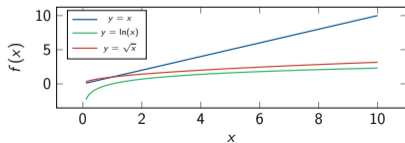
- Right-skewed data
- Multiplicative relationships
- Data spanning orders of magnitude

Square Root Transformation

$$x' = \sqrt{x}$$

Use for:

- Count data
- Moderate right skew
- Variance stabilization



Caution

Non-parametric, so doesn't preserve relationships between features. Best for single features.

Definition

Maps data to a specified distribution (uniform or normal) using quantiles:

$$x'_i = \Phi^{-1}(F(x_i))$$

Where:

- $F(x_i)$ = empirical CDF (rank / n)
- Φ^{-1} = inverse CDF of target distribution

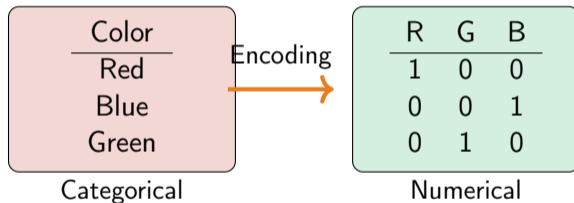
Scaling Methods Comparison

Method	Range	Centered	Outlier Robust	Sparse Safe	Use Case
Min-Max	$[0, 1]$	No	✗	✗	Bounded data
Z-Score	$(-\infty, \infty)$	Yes	✗	✗	Gaussian data
Robust	$(-\infty, \infty)$	Yes	✓	✗	Outliers present
MaxAbs	$[-1, 1]$	No	✗	✓	Sparse matrices
Box-Cox	Varies	Varies	Partial	✗	Non-normal
Quantile	Varies	Varies	✓	✗	Any distribution

Why Encode Categorical Features?

The Problem

Machine learning algorithms work with **numerical data**. Categorical variables must be converted to numbers.



Important

The encoding method matters! Wrong encoding can introduce false ordinal relationships.

Label Encoding (Integer Encoding)

Definition

Assigns a unique integer to each category:

$$f : \mathcal{C} \rightarrow \{0, 1, 2, \dots, k - 1\}$$

Where k = number of unique categories

Example

Category	Encoded
Red	0
Blue	1
Green	2

Problem!

Implies false ordering:

$$\text{Red} < \text{Blue} < \text{Green}$$

And false distances:

$$d(\text{Red}, \text{Green}) = 2$$

Ordinal Encoding

Definition

Like label encoding, but respects natural ordering:

$$f : \mathcal{C}_{ordered} \rightarrow \{0, 1, 2, \dots, k - 1\}$$

Where ordering is meaningful

Appropriate Examples

Education Level:

Level	Code
High School	0
Bachelor's	1
Master's	2
PhD	3

Size:

Size	Code
Small	0
Medium	1
Large	2
XL	3

One-Hot Encoding (Dummy Variables)

Definition

Creates binary columns for each category:

$$x \in \mathcal{C} \rightarrow e_x \in \{0, 1\}^k$$

Where e_x is a one-hot vector with 1 at position x

Color	is_Blue	
Red	0	
Blue	1	
Green	0	
Red	0	

One-Hot Encoding: Pros and Cons

Advantages

- No ordinal assumption
- Equal distance between categories:

$$d(e_i, e_j) = \sqrt{2} \quad \forall i \neq j$$

- Works with all algorithms
- Interpretable coefficients

Disadvantages

- **Curse of dimensionality:**
 k categories $\rightarrow k$ (or $k - 1$) columns
- High cardinality problems
- Sparse matrices
- Multicollinearity (dummy trap)

Dummy Variable Trap

Drop one column to avoid perfect multicollinearity:

$$\text{is_Red} + \text{is_Blue} + \text{is_Green} = 1 \quad (\text{always})$$

Binary Encoding

Definition

Encode categories as binary numbers, then split into columns:

$$\text{Category} \xrightarrow{\text{integer}} n \xrightarrow{\text{binary}} (b_1, b_2, \dots, b_m)$$

Where $m = \lceil \log_2(k) \rceil$ columns needed for k categories

Example with 5 categories

Category	Integer	Binary	b_1	b_2	b_3
A	0	000	0	0	0
B	1	001	0	0	1
C	2	010	0	1	0
D	3	011	0	1	1
E	4	100	1	0	0

5 categories → 3 columns (vs 5 for one-hot)

Target Encoding (Mean Encoding)

Definition

Replace category with mean of target variable for that category:

$$x'_c = \frac{1}{n_c} \sum_{i:x_i=c} y_i = \bar{y}_c$$

Where n_c = count of samples with category c

Example: Predicting Salary

Department	Avg Salary	Encoded
Engineering	\$95,000	95000
Marketing	\$75,000	75000
Sales	\$65,000	65000

Danger: Data Leakage!

Target Encoding: Smoothing

Smoothed Target Encoding

Blend category mean with global mean to handle rare categories:

$$x'_c = \lambda(n_c) \cdot \bar{y}_c + (1 - \lambda(n_c)) \cdot \bar{y}_{global}$$

Where smoothing factor:

$$\lambda(n) = \frac{n}{n + m}$$

m = smoothing parameter (e.g., 10-100)

Interpretation

- Large n_c : $\lambda \approx 1$, trust category mean
- Small n_c : $\lambda \approx 0$, use global mean

Frequency Encoding

Definition

Replace category with its frequency (count or proportion):

$$x'_c = \frac{n_c}{n} \quad \text{or} \quad x'_c = n_c$$

Example

City	Count	Freq
NYC	500	0.50
LA	300	0.30
Chicago	200	0.20

Pros:

- No target leakage
- Single column output
- Captures rarity

Cons:

- Same encoding for different categories with same frequency

Encoding Methods Comparison

Method	Output Dims	High Cardinality	Target Leakage	Best For
Label	1	✓	✓	Trees, Ordinal
One-Hot	k or $k - 1$	✗	✓	Low cardinality
Binary	$\log_2(k)$	✓	✓	Medium cardinality
Target	1	✓	✗	Supervised (with CV)
Frequency	1	✓	✓	Any

Rule of Thumb

- $k < 10$: One-Hot encoding
- $10 \leq k < 100$: Binary or Target encoding
- $k \geq 100$: Target or Frequency encoding

Andrew Ng's Quote

"Coming up with features is difficult, time-consuming, requires expert knowledge. Applied machine learning is basically feature engineering."

Definition

Feature Engineering is the process of using domain knowledge to create new features from raw data that better represent the underlying problem to the model.

Mathematical Feature Creation

Arithmetic Operations

Combine features using math operations:

$$x_{new} = f(x_1, x_2, \dots, x_n)$$

Common Operations

- **Sum:** $x_1 + x_2$
- **Difference:** $x_1 - x_2$
- **Product:** $x_1 \times x_2$
- **Ratio:** $\frac{x_1}{x_2}$
- **Power:** x^2, x^3
- **Root:** \sqrt{x}
- **Log:** $\log(x)$
- **Exponential:** e^x

Real Example: BMI Calculation

$$\text{BMI} = \frac{\text{Weight (kg)}}{\text{Height (m)}^2}$$

Polynomial Features

Definition

Generate polynomial and interaction terms up to degree d :

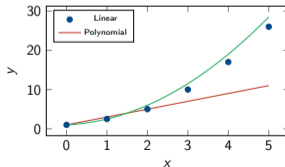
$$\phi(x_1, x_2) = (1, x_1, x_2, x_1^2, x_1x_2, x_2^2, x_1^3, \dots)$$

Example: Degree 2 with 2 features

Input: (x_1, x_2)

Output: $(1, x_1, x_2, x_1^2, x_1x_2, x_2^2)$

Number of features: $\binom{n+d}{d} = \binom{2+2}{2} = 6$



Datetime Feature Engineering

Extracting Components from Datetime

`datetime` \rightarrow $\{year, month, day, hour, minute, dayofweek, \dots\}$

Common Datetime Features

Direct Extraction:

- Year, Month, Day
- Hour, Minute, Second
- Day of week (0-6)
- Day of year (1-365)
- Week of year
- Quarter

Derived Features:

- Is weekend (binary)
- Is holiday (binary)
- Time since event
- Part of day (morning/afternoon)
- Season
- Is month start/end

Cyclical Encoding for Time Features

Problem with Linear Encoding

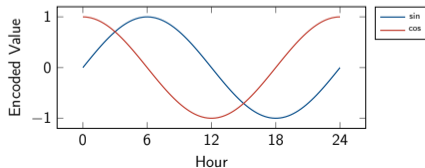
Hour 23 and Hour 0 should be close, but:

$$|23 - 0| = 23 \quad (\text{Wrong! They're 1 hour apart})$$

Solution: Sine-Cosine Encoding

$$x_{\sin} = \sin\left(\frac{2\pi \cdot t}{T}\right), \quad x_{\cos} = \cos\left(\frac{2\pi \cdot t}{T}\right)$$

Where T = period (24 for hours, 7 for days, 12 for months)



Binning (Discretization)

Definition

Convert continuous features into discrete bins/categories:

$$x \in \mathbb{R} \rightarrow b \in \{1, 2, \dots, k\}$$

Equal-Width Binning

$$\text{width} = \frac{x_{max} - x_{min}}{k}$$

Bin boundaries evenly spaced

Equal-Frequency (Quantile)

Each bin has $\approx \frac{n}{k}$ samples
Uses percentiles as boundaries

Example: Age Binning

Age	Equal-Width	Domain-Based
25	20-40	Young Adult
45	40-60	Middle Age
70	60-80	Old

Text Feature Engineering

Basic Text Features

- **Length:** Character count, word count
- **Statistics:** Average word length, unique word ratio
- **Patterns:** Number of digits, special chars, uppercase ratio

Bag of Words (BoW)

Count occurrences of each word:

$$\text{BoW}(d) = (c_1, c_2, \dots, c_V)$$

Where c_i = count of word i in document d

TF-IDF (Term Frequency - Inverse Document Frequency)

$$\text{TF-IDF}(t, d) = \text{TF}(t, d) \times \text{IDF}(t)$$

Aggregation Features

Group-Based Statistics

Compute statistics within groups:

$$x_{agg} = \text{Agg}_{group}(x)$$

Common Aggregations

- Mean: \bar{x}_{group}
- Median: \tilde{x}_{group}
- Sum: $\sum x_{group}$
- Count: n_{group}
- Min/Max
- Standard deviation
- Skewness, Kurtosis
- Range

Example: Customer Features

Raw	Engineered
Transaction amounts	Avg transaction per customer
Purchase dates	Days since last purchase

Aggregation Features

Example: Customer Features

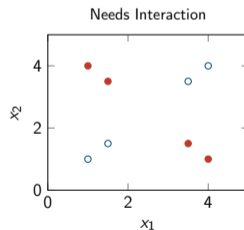
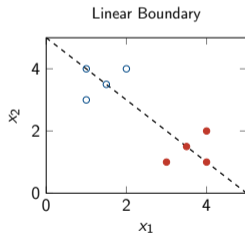
Raw	Engineered
Transaction amounts	Avg transaction per customer
Purchase dates	Days since last purchase
Product categories	Most frequent category

Interaction Features

Definition

Capture relationships between features:

$$x_{interaction} = x_i \times x_j \quad \text{or} \quad x_{interaction} = f(x_i, x_j)$$



XOR pattern requires $x_1 \times x_2$ interaction term to separate linearly

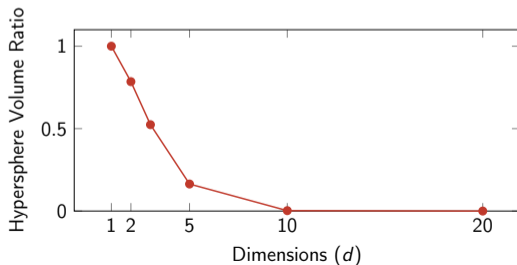
The Curse of Dimensionality

Problem

As dimensions increase, data becomes increasingly sparse:

$$\text{Volume of unit hypercube} = 1^d = 1$$

$$\text{Volume of inscribed hypersphere} = \frac{\pi^{d/2}}{\Gamma(\frac{d}{2} + 1)} \cdot \left(\frac{1}{2}\right)^d \xrightarrow{d \rightarrow \infty} 0$$



Feature Selection vs Feature Extraction

Feature Selection

Select subset of original features:

$$X' = X[:, S]$$

Where $S \subset \{1, \dots, p\}$

Methods:

- Filter (correlation, variance)
- Wrapper (RFE, forward selection)
- Embedded (Lasso, tree importance)

Feature Extraction

Create new features from combinations:

$$X' = X \cdot W$$

Where $W \in \mathbb{R}^{p \times k}$, $k < p$

Methods:

- PCA
- LDA
- Autoencoders
- t-SNE, UMAP

Principal Component Analysis (PCA) Overview

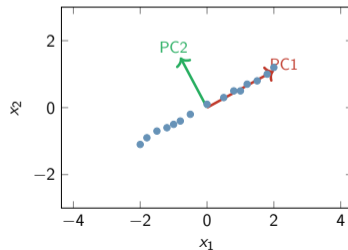
Goal

Find orthogonal directions of maximum variance:

$$z = W^T(x - \mu)$$

Where columns of W are eigenvectors of covariance matrix:

$$\Sigma w_i = \lambda_i w_i$$



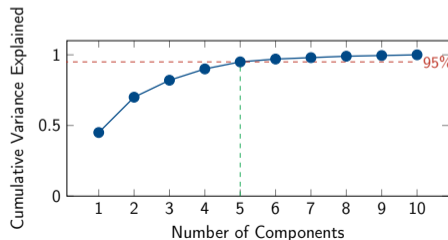
Variance Explained by PCA

Explained Variance Ratio

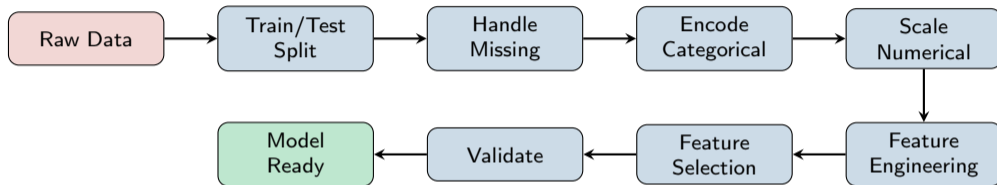
$$\text{EVR}_i = \frac{\lambda_i}{\sum_{j=1}^p \lambda_j}$$

Cumulative:

$$\text{CEVR}_k = \sum_{i=1}^k \text{EVR}_i$$



Complete Transformation Pipeline



Critical Rule: Fit on Train Only!

`scaler.fit(X_{train})` then `scaler.transform(X_{test})`

Never fit transformers on test data — this causes **data leakage!**

Avoiding Data Leakage

What is Data Leakage?

Using information from test set during training, leading to overly optimistic performance estimates.

Wrong (Leakage)

- 1 Scale ALL data
- 2 Split into train/test
- 3 Train model

Test statistics leak into scaling!

Correct

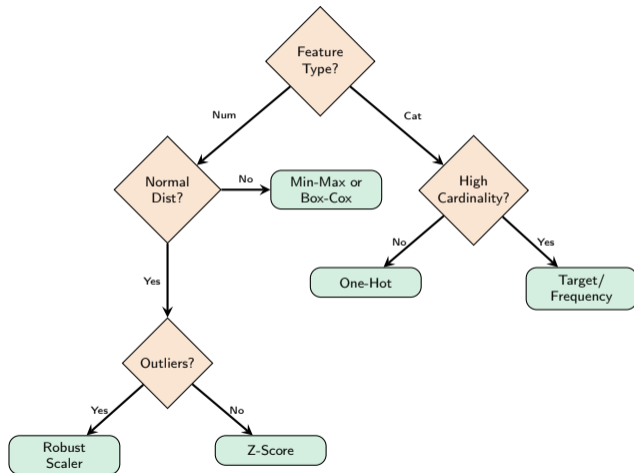
- 1 Split into train/test
- 2 Fit scaler on train
- 3 Transform both sets
- 4 Train model

Test data remains unseen!

Use Pipelines!

Sklearn Pipeline ensures correct order and prevents leakage automatically.

Transformation Decision Guide








Key Takeaways

- 1 **Scaling matters** for distance-based and gradient-based algorithms
- 2 **Choose encoding wisely** — wrong encoding creates false relationships
- 3 **Feature engineering** often has the biggest impact on model performance
- 4 **Avoid data leakage** — always fit on training data only
- 5 **Use pipelines** for reproducible and correct transformations
- 6 **Domain knowledge** is crucial for effective feature engineering

Next Part

Part 4: Data Preparation for Modeling — Train/Test Split, Cross-Validation, Handling Imbalanced Data

References

-  Pedregosa et al. (2011). *Scikit-learn: Machine Learning in Python*. JMLR.
-  Zheng, A. & Casari, A. (2018). *Feature Engineering for Machine Learning*. O'Reilly.
-  Box, G.E.P. & Cox, D.R. (1964). *An Analysis of Transformations*. JRSS.
-  Kuhn, M. & Johnson, K. (2019). *Feature Engineering and Selection*. CRC Press.
-  Géron, A. (2019). *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow*. O'Reilly.

Part V

Part 3 : Practice Data Transformation

What We Will Cover

Practical Data Cleaning Skills

- Detecting and handling missing values
- Identifying and removing duplicates
- Detecting and treating outliers
- Correcting data types
- Cleaning text and string data
- Handling inconsistent data
- Building cleaning pipelines

Libraries Used

pandas, numpy, scikit-learn, scipy, fuzzywuzzy

Essential Libraries Import

```
1 # Data manipulation
2 import pandas as pd
3 import numpy as np
4
5 # Visualization
6 import matplotlib.pyplot as plt
7 import seaborn as sns
8
9 # Scikit-learn imputers
10 from sklearn.impute import SimpleImputer, KNNImputer
11 from sklearn.experimental import enable_iterative_imputer
12 from sklearn.impute import IterativeImputer
13
14 # Statistical functions
15 from scipy import stats
16
17 # String matching (install: pip install fuzzywuzzy python-Levenshtein)
18 from fuzzywuzzy import fuzz, process
19
20 # Display settings
21 pd.set_option('display.max_columns', None)
```

Detecting Missing Values

```
1 # Check for missing values
2 print(df.isnull().sum())           # Count per column
3 print(df.isna().sum())           # Same as isnull()
4
5 # Total missing values
6 print(f"Total missing: {df.isnull().sum().sum()}")
7
8 # Percentage missing per column
9 missing_pct = (df.isnull().sum() / len(df) * 100).round(2)
10 print(missing_pct.sort_values(ascending=False))
11
12 # Rows with any missing value
13 rows_with_missing = df[df.isnull().any(axis=1)]
14 print(f"Rows with missing: {len(rows_with_missing)}")
15
16 # Columns with missing values
17 cols_with_missing = df.columns[df.isnull().any()].tolist()
18 print(f"Columns with missing: {cols_with_missing}")
```

Missing Values Summary Function

```
1 def missing_summary(df):
2     """Generate detailed missing values summary"""
3     missing = df.isnull().sum()
4     missing_pct = (missing / len(df) * 100).round(2)
5
6     summary = pd.DataFrame({
7         'Missing Count': missing,
8         'Missing %': missing_pct,
9         'Dtype': df.dtypes
10    })
11
12    # Filter only columns with missing values
13    summary = summary[summary['Missing Count'] > 0]
14    summary = summary.sort_values('Missing %', ascending=False)
15
16    print(f"Total columns: {len(df.columns)}")
17    print(f"Columns with missing: {len(summary)}")
18    print(f"Total missing values: {missing.sum()}")
19
20    return summary
21
22 print(missing_summary(df))
```

Visualizing Missing Values

```
1 # Using missingno library (pip install missingno)
2 import missingno as msno
3
4 # Matrix visualization
5 msno.matrix(df, figsize=(12, 6))
6 plt.title('Missing Values Matrix')
7 plt.show()
8
9 # Bar chart of missing values
10 msno.bar(df, figsize=(12, 6))
11 plt.show()
12
13 # Heatmap of missing value correlations
14 msno.heatmap(df, figsize=(10, 8))
15 plt.show()
16
17 # Dendrogram showing missing value patterns
18 msno.dendrogram(df, figsize=(12, 6))
19 plt.show()
```

Custom Missing Values Heatmap

```
1 # Create missing values heatmap with seaborn
2 plt.figure(figsize=(14, 8))
3 sns.heatmap(df.isnull(), cbar=True, yticklabels=False,
4             cmap='viridis')
5 plt.title('Missing Values Heatmap')
6 plt.tight_layout()
7 plt.show()
8
9 # Missing values by row (identify problematic rows)
10 row_missing = df.isnull().sum(axis=1)
11 plt.figure(figsize=(12, 4))
12 plt.hist(row_missing, bins=50, edgecolor='black')
13 plt.xlabel('Number of Missing Values per Row')
14 plt.ylabel('Frequency')
15 plt.title('Distribution of Missing Values per Row')
16 plt.show()
17
18 # Rows with too many missing values
19 threshold = len(df.columns) * 0.5 # 50% missing
20 bad_rows = df[row_missing > threshold]
21 print(f"Rows with >50% missing: {len(bad_rows)}")
```

Deletion Methods

```
1 # Drop rows with ANY missing value
2 df_clean = df.dropna()
3
4 # Drop rows with ALL missing values
5 df_clean = df.dropna(how='all')
6
7 # Drop rows with missing in specific columns
8 df_clean = df.dropna(subset=['col1', 'col2'])
9
10 # Drop rows with minimum non-null values (threshold)
11 df_clean = df.dropna(thresh=5) # Keep rows with at least 5 non-null
12
13 # Drop columns with missing values
14 df_clean = df.dropna(axis=1)
15
16 # Drop columns with too many missing values
17 threshold = len(df) * 0.5 # 50% missing threshold
18 cols_to_drop = df.columns[df.isnull().sum() > threshold]
19 df_clean = df.drop(columns=cols_to_drop)
20 print(f"Dropped columns: {cols_to_drop.tolist()}")
```

Simple Imputation - Numerical

```
1 # Fill with constant value
2 df['col'] = df['col'].fillna(0)
3
4 # Fill with mean
5 df['col'] = df['col'].fillna(df['col'].mean())
6
7 # Fill with median (robust to outliers)
8 df['col'] = df['col'].fillna(df['col'].median())
9
10 # Fill with mode (most frequent)
11 df['col'] = df['col'].fillna(df['col'].mode()[0])
12
13 # Fill multiple columns at once
14 numerical_cols = df.select_dtypes(include=np.number).columns
15 df[numerical_cols] = df[numerical_cols].fillna(
16     df[numerical_cols].median()
17 )
18
19 # Fill with column-specific values
20 fill_values = {'col1': 0, 'col2': df['col2'].mean(), 'col3': -1}
21 df = df.fillna(fill_values)
```

Simple Imputation - Categorical

```
1 # Fill with mode (most frequent category)
2 df['category'] = df['category'].fillna(df['category'].mode()[0])
3
4 # Fill with a specific category
5 df['category'] = df['category'].fillna('Unknown')
6 df['category'] = df['category'].fillna('Missing')
7
8 # Fill all categorical columns with mode
9 categorical_cols = df.select_dtypes(include=['object', 'category']).columns
10 for col in categorical_cols:
11     df[col] = df[col].fillna(df[col].mode()[0])
12
13 # Or fill with 'Unknown'
14 df[categorical_cols] = df[categorical_cols].fillna('Unknown')
15
16 # Conditional filling
17 df['category'] = df['category'].fillna(
18     df.groupby('other_col')['category'].transform(
19         lambda x: x.mode()[0] if len(x.mode()) > 0 else 'Unknown'
20     )
21 )
```

Forward and Backward Fill

```
1 # Forward fill (use previous value)
2 df['col'] = df['col'].fillna(method='ffill')
3 # Or using newer syntax:
4 df['col'] = df['col'].ffill()
5
6 # Backward fill (use next value)
7 df['col'] = df['col'].fillna(method='bfill')
8 # Or:
9 df['col'] = df['col'].bfill()
10
11 # Limit the fill (max consecutive NaNs to fill)
12 df['col'] = df['col'].ffill(limit=2)
13
14 # Combine forward and backward fill
15 df['col'] = df['col'].ffill().bfill()
16
17 # Useful for time series data
18 df = df.sort_values('date')
19 df['value'] = df['value'].ffill()
```

Interpolation Methods

```
1 # Linear interpolation (default)
2 df['col'] = df['col'].interpolate(method='linear')
3
4 # Polynomial interpolation
5 df['col'] = df['col'].interpolate(method='polynomial', order=2)
6
7 # Spline interpolation (smooth curves)
8 df['col'] = df['col'].interpolate(method='spline', order=3)
9
10 # Time-based interpolation (for datetime index)
11 df['col'] = df['col'].interpolate(method='time')
12
13 # Index-based interpolation
14 df['col'] = df['col'].interpolate(method='index')
15
16 # Limit interpolation direction
17 df['col'] = df['col'].interpolate(limit=2, limit_direction='forward')
18
19 # Only interpolate inside (not at edges)
20 df['col'] = df['col'].interpolate(limit_area='inside')
```

Scikit-learn SimpleImputer

```
1 from sklearn.impute import SimpleImputer
2
3 # Mean imputation
4 imputer_mean = SimpleImputer(strategy='mean')
5 df[numerical_cols] = imputer_mean.fit_transform(df[numerical_cols])
6
7 # Median imputation
8 imputer_median = SimpleImputer(strategy='median')
9 df[numerical_cols] = imputer_median.fit_transform(df[numerical_cols])
10
11 # Most frequent (mode) - works for categorical too
12 imputer_mode = SimpleImputer(strategy='most_frequent')
13 df[categorical_cols] = imputer_mode.fit_transform(df[categorical_cols])
14
15 # Constant value imputation
16 imputer_const = SimpleImputer(strategy='constant', fill_value=0)
17 df[numerical_cols] = imputer_const.fit_transform(df[numerical_cols])
18
19 # For categorical with constant
20 imputer_cat = SimpleImputer(strategy='constant', fill_value='Missing')
21 df[categorical_cols] = imputer_cat.fit_transform(df[categorical_cols])
```

KNN Imputation

```
1 from sklearn.impute import KNNImputer
2
3 # KNN imputer (uses k nearest neighbors)
4 knn_imputer = KNNImputer(n_neighbors=5, weights='uniform')
5 df[numerical_cols] = knn_imputer.fit_transform(df[numerical_cols])
6
7 # Weighted by distance (closer neighbors have more influence)
8 knn_imputer = KNNImputer(n_neighbors=5, weights='distance')
9 df[numerical_cols] = knn_imputer.fit_transform(df[numerical_cols])
10
11 # Custom number of neighbors
12 knn_imputer = KNNImputer(n_neighbors=10)
13
14 # Note: KNNImputer only works with numerical data
15 # For mixed data, encode categoricals first or use separate imputers
16
17 # Scale data before KNN imputation for better results
18 from sklearn.preprocessing import StandardScaler
19 scaler = StandardScaler()
20 df_scaled = scaler.fit_transform(df[numerical_cols])
21 df_imputed = knn_imputer.fit_transform(df_scaled)
22 df[numerical_cols] = scaler.inverse_transform(df_imputed)
```

Iterative Imputer (MICE)

```
1 from sklearn.experimental import enable_iterative_imputer
2 from sklearn.impute import IterativeImputer
3
4 # Basic iterative imputer (MICE - Multiple Imputation by Chained Equations)
5 iterative_imputer = IterativeImputer(max_iter=10, random_state=42)
6 df[numerical_cols] = iterative_imputer.fit_transform(df[numerical_cols])
7
8 # With specific estimator
9 from sklearn.linear_model import BayesianRidge
10 iterative_imputer = IterativeImputer(
11     estimator=BayesianRidge(),
12     max_iter=10,
13     random_state=42
14 )
15 df[numerical_cols] = iterative_imputer.fit_transform(df[numerical_cols])
16
17 # With Random Forest for non-linear relationships
18 from sklearn.ensemble import RandomForestRegressor
19 iterative_imputer = IterativeImputer(
20     estimator=RandomForestRegressor(n_estimators=10, random_state=42),
21     max_iter=10,
22     random_state=42
23 )
```

Group-Based Imputation

```
1 # Impute based on group statistics
2 # Example: Fill missing salary with median salary by department
3
4 # Using groupby transform
5 df['salary'] = df.groupby('department')['salary'].transform(
6     lambda x: x.fillna(x.median())
7 )
8
9 # Multiple group columns
10 df['salary'] = df.groupby(['department', 'level'])['salary'].transform(
11     lambda x: x.fillna(x.median())
12 )
13
14 # Fallback to global median if group has all NaN
15 def group_impute(group):
16     if group.isnull().all():
17         return group.fillna(df['salary'].median())
18     return group.fillna(group.median())
19
20 df['salary'] = df.groupby('department')['salary'].transform(group_impute)
21
22 # Custom function for complex logic
23 df['value'] = df.groupby('category')['value'].apply(
24     lambda x: x.fillna(x.median() if x.notna().sum() > 0 else 0)
25 )
```

Imputation with Indicator

```
1 # Create indicator column for missing values (useful for modeling)
2 df['col_was_missing'] = df['col'].isnull().astype(int)
3
4 # Then impute
5 df['col'] = df['col'].fillna(df['col'].median())
6
7 # Function to add missing indicators
8 def add_missing_indicators(df, columns):
9     for col in columns:
10         if df[col].isnull().sum() > 0:
11             df[f'{col}_missing'] = df[col].isnull().astype(int)
12     return df
13
14 df = add_missing_indicators(df, numerical_cols)
15
16 # Using SimpleImputer with add_indicator
17 from sklearn.impute import SimpleImputer
18 imputer = SimpleImputer(strategy='median', add_indicator=True)
19 result = imputer.fit_transform(df[numerical_cols])
20 # This returns both imputed values and indicator columns
```

Detecting Duplicates

```
1 # Check for duplicate rows
2 print(f"Duplicate rows: {df.duplicated().sum()}")
3
4 # Check duplicates based on specific columns
5 print(f"Duplicate keys: {df.duplicated(subset=['id', 'date']).sum()}")
6
7 # View duplicate rows
8 duplicates = df[df.duplicated(keep=False)] # keep=False shows all duplicates
9 print(duplicates.sort_values(by=df.columns.tolist()))
10
11 # First occurrence vs duplicates
12 first_occurrence = df[df.duplicated(keep='first') == False]
13 duplicate_rows = df[df.duplicated(keep='first') == True]
14
15 # Count duplicates per group
16 dup_counts = df.groupby(df.columns.tolist()).size().reset_index(name='count')
17 dup_counts = dup_counts[dup_counts['count'] > 1]
18 print(f"Groups with duplicates: {len(dup_counts)}")
```

Removing Duplicates

```
1 # Remove exact duplicate rows (keep first occurrence)
2 df_clean = df.drop_duplicates()
3
4 # Keep last occurrence instead
5 df_clean = df.drop_duplicates(keep='last')
6
7 # Remove all duplicates (don't keep any)
8 df_clean = df.drop_duplicates(keep=False)
9
10 # Remove duplicates based on specific columns
11 df_clean = df.drop_duplicates(subset=['id', 'date'])
12
13 # Remove duplicates and reset index
14 df_clean = df.drop_duplicates().reset_index(drop=True)
15
16 # In-place removal
17 df.drop_duplicates(inplace=True)
18
19 # Keep row with most non-null values among duplicates
20 df_clean = df.loc[df.groupby(['id'])['col'].transform('count').idxmax()]
```

Handling Duplicates with Aggregation

```
1 # Instead of dropping, aggregate duplicate rows
2 # Keep first value for some columns, mean for others
3
4 agg_rules = {
5     'name': 'first',          # Keep first name
6     'value': 'mean',         # Average the values
7     'count': 'sum',          # Sum the counts
8     'date': 'max',           # Keep latest date
9     'category': lambda x: x.mode()[0] # Most frequent
10 }
11
12 df_agg = df.groupby('id').agg(agg_rules).reset_index()
13
14 # Custom aggregation for complex cases
15 def resolve_duplicates(group):
16     return pd.Series({
17         'name': group['name'].iloc[0],
18         'value': group['value'].mean(),
19         'latest_date': group['date'].max(),
20         'all_categories': ', '.join(group['category'].unique())
21     })
22
23 df_resolved = df.groupby('id').apply(resolve_duplicates).reset_index()
```

Fuzzy Duplicate Detection

```
1 # For near-duplicates (e.g., "John Smith" vs "Jon Smith")
2 from fuzzywuzzy import fuzz, process
3
4 # Compare two strings
5 similarity = fuzz.ratio("John Smith", "Jon Smith")
6 print(f"Similarity: {similarity}%") # Output: 91%
7
8 # Find potential duplicates in a column
9 def find_fuzzy_duplicates(df, column, threshold=85):
10     duplicates = []
11     values = df[column].unique()
12
13     for i, val1 in enumerate(values):
14         for val2 in values[i+1:]:
15             score = fuzz.ratio(str(val1), str(val2))
16             if score >= threshold:
17                 duplicates.append((val1, val2, score))
18
19     return pd.DataFrame(duplicates, columns=['Value1', 'Value2', 'Score'])
20
21 fuzzy_dups = find_fuzzy_duplicates(df, 'name', threshold=85)
22 print(fuzzy_dups)
```

Standardizing Fuzzy Duplicates

```
1 # Standardize similar values to one canonical form
2 from fuzzywuzzy import process
3
4 def standardize_column(df, column, threshold=85):
5     """Standardize similar values to most frequent form"""
6     value_counts = df[column].value_counts()
7     mapping = {}
8
9     for value in value_counts.index:
10        if value not in mapping:
11            # Find best match among more frequent values
12            choices = [v for v in value_counts.index if v not in mapping.values()]
13            match, score = process.extractOne(str(value), choices)
14            if score >= threshold:
15                mapping[value] = match
16            else:
17                mapping[value] = value
18
19    return df[column].map(mapping)
20
21 df['name_clean'] = standardize_column(df, 'name', threshold=85)
22
23 # Check the mapping
24 print(df[['name', 'name_clean']].drop_duplicates())
```

IQR Method

```
1 def detect_outliers_iqr(df, column, multiplier=1.5):
2     """Detect outliers using IQR method"""
3     Q1 = df[column].quantile(0.25)
4     Q3 = df[column].quantile(0.75)
5     IQR = Q3 - Q1
6
7     lower_bound = Q1 - multiplier * IQR
8     upper_bound = Q3 + multiplier * IQR
9
10    outliers = (df[column] < lower_bound) | (df[column] > upper_bound)
11
12    return outliers, lower_bound, upper_bound
13
14 # Apply to column
15 outlier_mask, lower, upper = detect_outliers_iqr(df, 'value')
16 print(f"Outliers: {outlier_mask.sum()}")
17 print(f"Bounds: [{lower:.2f}, {upper:.2f}]")
18
19 # View outliers
20 print(df[outlier_mask])
21
22 # Apply to all numerical columns
23 for col in numerical_cols:
24     mask, _, _ = detect_outliers_iqr(df, col)
25     print(f"{col}: {mask.sum()} outliers ({mask.mean()*100:.1f}%")
```

Z-Score Method

```
1 from scipy import stats
2
3 def detect_outliers_zscore(df, column, threshold=3):
4     """Detect outliers using Z-score method"""
5     z_scores = np.abs(stats.zscore(df[column].dropna()))
6     outliers = z_scores > threshold
7     return outliers
8
9 # Apply to column
10 outlier_mask = detect_outliers_zscore(df, 'value', threshold=3)
11 print(f"Outliers (|z| > 3): {outlier_mask.sum()}")
12
13 # Calculate z-scores for all rows
14 df['z_score'] = np.abs(stats.zscore(df['value']))
15 df['is_outlier'] = df['z_score'] > 3
16
17 # Modified Z-score (more robust, uses median)
18 def modified_zscore(x):
19     median = np.median(x)
20     mad = np.median(np.abs(x - median))
21     return 0.6745 * (x - median) / mad
22
23 df['modified_z'] = np.abs(modified_zscore(df['value']))
24 df['is_outlier_robust'] = df['modified_z'] > 3.5
```

Isolation Forest

```
1 from sklearn.ensemble import IsolationForest
2
3 # Isolation Forest for multivariate outlier detection
4 iso_forest = IsolationForest(
5     n_estimators=100,
6     contamination=0.05, # Expected proportion of outliers
7     random_state=42
8 )
9
10 # Fit and predict (-1 for outliers, 1 for inliers)
11 df['outlier_if'] = iso_forest.fit_predict(df[numerical_cols])
12 df['is_outlier_if'] = df['outlier_if'] == -1
13
14 print(f"Outliers detected: {(df['outlier_if'] == -1).sum()}")
15
16 # Get anomaly scores (lower = more anomalous)
17 df['anomaly_score'] = iso_forest.decision_function(df[numerical_cols])
18
19 # Visualize
20 plt.figure(figsize=(10, 6))
21 plt.scatter(df.index, df['anomaly_score'], c=df['outlier_if'], cmap='RdYlGn')
22 plt.xlabel('Index')
23 plt.ylabel('Anomaly Score')
24 plt.title('Isolation Forest Anomaly Scores')
25 plt.show()
```

Local Outlier Factor (LOF)

```
1 from sklearn.neighbors import LocalOutlierFactor
2
3 # LOF for density-based outlier detection
4 lof = LocalOutlierFactor(
5     n_neighbors=20,
6     contamination=0.05
7 )
8
9 # Fit and predict
10 df['outlier_lof'] = lof.fit_predict(df[numerical_cols])
11 df['is_outlier_lof'] = df['outlier_lof'] == -1
12
13 print(f"LOF outliers: {(df['outlier_lof'] == -1).sum()}")
14
15 # Get LOF scores (negative_outlier_factor_)
16 df['lof_score'] = -lof.negative_outlier_factor_
17
18 # Visualize 2D projection with outliers
19 plt.figure(figsize=(10, 8))
20 colors = ['red' if x == -1 else 'blue' for x in df['outlier_lof']]
21 plt.scatter(df['col1'], df['col2'], c=colors, alpha=0.5)
22 plt.xlabel('Column 1')
23 plt.ylabel('Column 2')
24 plt.title('LOF Outlier Detection')
25 plt.show()
```

DBSCAN for Outlier Detection

```
1 from sklearn.cluster import DBSCAN
2 from sklearn.preprocessing import StandardScaler
3
4 # Scale data first
5 scaler = StandardScaler()
6 X_scaled = scaler.fit_transform(df[numerical_cols])
7
8 # DBSCAN clustering
9 dbscan = DBSCAN(eps=0.5, min_samples=5)
10 df['cluster'] = dbscan.fit_predict(X_scaled)
11
12 # Points labeled -1 are outliers (noise)
13 df['is_outlier_dbscan'] = df['cluster'] == -1
14
15 print(f"DBSCAN outliers: {(df['cluster'] == -1).sum()}")
16 print(f"Number of clusters: {df['cluster'].max() + 1}")
17
18 # Visualize
19 plt.figure(figsize=(10, 8))
20 scatter = plt.scatter(df['col1'], df['col2'], c=df['cluster'], cmap='viridis')
21 plt.colorbar(scatter)
22 plt.title('DBSCAN Clustering (-1 = Outliers)')
23 plt.show()
```

Removing Outliers

```
1 # Remove outliers detected by IQR
2 outlier_mask, lower, upper = detect_outliers_iqr(df, 'value')
3 df_clean = df[~outlier_mask]
4 print(f"Removed {outlier_mask.sum()} outliers")
5
6 # Remove outliers from multiple columns
7 def remove_outliers_iqr(df, columns, multiplier=1.5):
8     mask = pd.Series([False] * len(df))
9     for col in columns:
10         outliers, _, _ = detect_outliers_iqr(df, col, multiplier)
11         mask = mask | outliers
12     return df[~mask]
13
14 df_clean = remove_outliers_iqr(df, numerical_cols)
15
16 # Remove based on Z-score
17 z_scores = np.abs(stats.zscore(df[numerical_cols]))
18 df_clean = df[(z_scores < 3).all(axis=1)]
19
20 # Remove based on percentiles
21 lower = df['value'].quantile(0.01)
22 upper = df['value'].quantile(0.99)
23 df_clean = df[(df['value'] >= lower) & (df['value'] <= upper)]
```

Capping / Winsorization

```
1 # Cap outliers to boundary values (Winsorization)
2 def cap_outliers_iqr(df, column, multiplier=1.5):
3     Q1 = df[column].quantile(0.25)
4     Q3 = df[column].quantile(0.75)
5     IQR = Q3 - Q1
6     lower = Q1 - multiplier * IQR
7     upper = Q3 + multiplier * IQR
8
9     df[column] = df[column].clip(lower=lower, upper=upper)
10    return df
11
12 df = cap_outliers_iqr(df, 'value')
13
14 # Cap using percentiles
15 def cap_percentile(df, column, lower_pct=0.01, upper_pct=0.99):
16     lower = df[column].quantile(lower_pct)
17     upper = df[column].quantile(upper_pct)
18     df[column] = df[column].clip(lower=lower, upper=upper)
19     return df
20
21 # Using scipy's winsorize
22 from scipy.stats.mstats import winsorize
23 df['value_winsorized'] = winsorize(df['value'], limits=[0.05, 0.05])
```

Transformation Methods

```
1 # Log transformation (for right-skewed data)
2 df['value_log'] = np.log1p(df['value']) # log(1+x) handles zeros
3
4 # Square root transformation
5 df['value_sqrt'] = np.sqrt(df['value'])
6
7 # Box-Cox transformation (requires positive values)
8 from scipy.stats import boxcox
9 df['value_boxcox'], lambda_param = boxcox(df['value'] + 1)
10 print(f"Optimal lambda: {lambda_param}")
11
12 # Yeo-Johnson transformation (handles negative values)
13 from sklearn.preprocessing import PowerTransformer
14 pt = PowerTransformer(method='yeo-johnson')
15 df['value_yeojohnson'] = pt.fit_transform(df[['value']])
16
17 # Compare distributions
18 fig, axes = plt.subplots(1, 3, figsize=(15, 4))
19 df['value'].hist(ax=axes[0], bins=30)
20 axes[0].set_title('Original')
21 df['value_log'].hist(ax=axes[1], bins=30)
22 axes[1].set_title('Log Transform')
23 df['value_boxcox'].hist(ax=axes[2], bins=30)
24 axes[2].set_title('Box-Cox Transform')
25 plt.tight_layout()
```

Imputing Outliers

```
1 # Replace outliers with median
2 def impute_outliers_median(df, column, multiplier=1.5):
3     Q1 = df[column].quantile(0.25)
4     Q3 = df[column].quantile(0.75)
5     IQR = Q3 - Q1
6     lower = Q1 - multiplier * IQR
7     upper = Q3 + multiplier * IQR
8     median = df[column].median()
9
10    df[column] = df[column].apply(
11        lambda x: median if (x < lower or x > upper) else x
12    )
13    return df
14
15 # Replace outliers with NaN, then impute
16 def outliers_to_nan(df, column, multiplier=1.5):
17     mask, lower, upper = detect_outliers_iqr(df, column, multiplier)
18     df.loc[mask, column] = np.nan
19     return df
20
21 df = outliers_to_nan(df, 'value')
22 df['value'] = df['value'].fillna(df['value'].median())
```

Converting Data Types

```
1 # Check current data types
2 print(df.dtypes)
3
4 # Convert to numeric (coerce errors to NaN)
5 df['col'] = pd.to_numeric(df['col'], errors='coerce')
6
7 # Convert to integer (must handle NaN first)
8 df['col'] = df['col'].fillna(0).astype(int)
9
10 # Nullable integer type (handles NaN)
11 df['col'] = df['col'].astype('Int64') # Note: capital I
12
13 # Convert to float
14 df['col'] = df['col'].astype(float)
15
16 # Convert to string
17 df['col'] = df['col'].astype(str)
18
19 # Convert to category (memory efficient)
20 df['category_col'] = df['category_col'].astype('category')
21
22 # Convert to boolean
23 df['bool_col'] = df['bool_col'].astype(bool)
24 df['bool_col'] = df['col'].map({'Yes': True, 'No': False})
```

Converting Datetime

```
1 # Convert to datetime
2 df['date'] = pd.to_datetime(df['date'])
3
4 # Specify format for faster parsing
5 df['date'] = pd.to_datetime(df['date'], format='%Y-%m-%d')
6
7 # Handle different formats
8 df['date'] = pd.to_datetime(df['date'], format='mixed', dayfirst=True)
9
10 # Handle errors
11 df['date'] = pd.to_datetime(df['date'], errors='coerce') # Invalid -> NaT
12
13 # Extract components
14 df['year'] = df['date'].dt.year
15 df['month'] = df['date'].dt.month
16 df['day'] = df['date'].dt.day
17 df['dayofweek'] = df['date'].dt.dayofweek
18 df['hour'] = df['date'].dt.hour
19
20 # Convert Unix timestamp to datetime
21 df['date'] = pd.to_datetime(df['timestamp'], unit='s')
22
23 # Convert datetime to string
24 df['date_str'] = df['date'].dt.strftime('%Y-%m-%d')
```

Automatic Type Inference

```
1 # Infer better dtypes automatically
2 df = df.infer_objects()
3
4 # Convert object columns to best possible types
5 df = df.convert_dtypes()
6
7 # Custom type conversion function
8 def optimize_dtypes(df):
9     for col in df.columns:
10         col_type = df[col].dtype
11
12         if col_type == 'object':
13             # Try to convert to numeric
14             try:
15                 df[col] = pd.to_numeric(df[col])
16             except:
17                 # Try to convert to datetime
18                 try:
19                     df[col] = pd.to_datetime(df[col])
20                 except:
21                     # Convert to category if low cardinality
22                     if df[col].nunique() / len(df) < 0.5:
23                         df[col] = df[col].astype('category')
24
25     return df
26 df = optimize_dtypes(df)
```

Basic String Cleaning

```
1 # Remove leading/trailing whitespace
2 df['col'] = df['col'].str.strip()
3
4 # Remove extra whitespace
5 df['col'] = df['col'].str.replace(r'\s+', ' ', regex=True)
6
7 # Convert to lowercase/uppercase
8 df['col'] = df['col'].str.lower()
9 df['col'] = df['col'].str.upper()
10 df['col'] = df['col'].str.title() # Title Case
11
12 # Remove specific characters
13 df['col'] = df['col'].str.replace('[^\w\s]', '', regex=True) # Keep only alphanumeric
14
15 # Remove digits
16 df['col'] = df['col'].str.replace(r'\d+', '', regex=True)
17
18 # Remove specific strings
19 df['col'] = df['col'].str.replace('unwanted', '', regex=False)
20
21 # Replace multiple values
22 df['col'] = df['col'].replace({'old1': 'new1', 'old2': 'new2'})
```

Advanced String Operations

```
1 # Extract patterns with regex
2 df['phone'] = df['text'].str.extract(r'(\d{3}-\d{3}-\d{4})')
3 df['email'] = df['text'].str.extract(r'([\w.-]+@[w.-]+)')
4
5 # Split strings
6 df[['first_name', 'last_name']] = df['full_name'].str.split(' ', n=1, expand=True)
7
8 # Check for patterns
9 df['has_email'] = df['text'].str.contains(r'@', regex=True)
10 df['starts_with_a'] = df['col'].str.startswith('A')
11
12 # Pad strings
13 df['code'] = df['code'].str.zfill(5) # Pad with zeros: '42' -> '00042'
14
15 # Slice strings
16 df['first_3'] = df['col'].str[:3]
17 df['last_4'] = df['col'].str[-4:]
18
19 # Get string length
20 df['length'] = df['col'].str.len()
21
22 # Count occurrences
23 df['word_count'] = df['text'].str.count(r'\w+')
```

Standardizing Categorical Values

```
1 # Map inconsistent values to standard forms
2 mapping = {
3     'yes': 'Yes', 'YES': 'Yes', 'y': 'Yes', 'Y': 'Yes', '1': 'Yes',
4     'no': 'No', 'NO': 'No', 'n': 'No', 'N': 'No', '0': 'No'
5 }
6 df['response'] = df['response'].map(mapping).fillna(df['response'])
7
8 # Using replace with regex
9 df['col'] = df['col'].replace({
10     r'^[Yy]es.*': 'Yes',
11     r'^[Nn]o.*': 'No'
12 }, regex=True)
13
14 # Standardize case and strip
15 def standardize_text(text):
16     if pd.isna(text):
17         return text
18     return str(text).strip().lower().title()
19
20 df['name'] = df['name'].apply(standardize_text)
21
22 # Fix common typos
23 typo_fixes = {'teh': 'the', 'recieve': 'receive', 'seperate': 'separate'}
24 for typo, fix in typo_fixes.items():
25     df['text'] = df['text'].str.replace(typo, fix, case=False)
```

Complete Cleaning Function

```
1 def clean_dataframe(df):
2     """Comprehensive data cleaning pipeline"""
3     df = df.copy()
4
5     # 1. Remove duplicate rows
6     initial_rows = len(df)
7     df = df.drop_duplicates()
8     print(f"Removed {initial_rows - len(df)} duplicate rows")
9
10    # 2. Handle column names
11    df.columns = df.columns.str.strip().str.lower().str.replace(' ', '_')
12
13    # 3. Identify column types
14    num_cols = df.select_dtypes(include=np.number).columns.tolist()
15    cat_cols = df.select_dtypes(include='object').columns.tolist()
16
17    # 4. Clean string columns
18    for col in cat_cols:
19        df[col] = df[col].str.strip()
20
21    # 5. Handle missing values
22    # Numerical: median imputation
23    for col in num_cols:
24        df[col] = df[col].fillna(df[col].median())
25
26    return df
```

Sklearn Pipeline for Cleaning

```
1 from sklearn.pipeline import Pipeline
2 from sklearn.compose import ColumnTransformer
3 from sklearn.impute import SimpleImputer
4 from sklearn.preprocessing import StandardScaler, OneHotEncoder
5
6 # Define preprocessing for different column types
7 numerical_pipeline = Pipeline([
8     ('imputer', SimpleImputer(strategy='median')),
9     ('scaler', StandardScaler())
10 ])
11
12 categorical_pipeline = Pipeline([
13     ('imputer', SimpleImputer(strategy='constant', fill_value='Missing')),
14     ('encoder', OneHotEncoder(handle_unknown='ignore', sparse_output=False))
15 ])
16
17 # Combine pipelines
18 preprocessor = ColumnTransformer([
19     ('num', numerical_pipeline, numerical_cols),
20     ('cat', categorical_pipeline, categorical_cols)
21 ])
22
23 # Fit and transform
24 X_clean = preprocessor.fit_transform(df)
```

Custom Transformer

```
1 from sklearn.base import BaseEstimator, TransformerMixin
2
3 class OutlierRemover(BaseEstimator, TransformerMixin):
4     def __init__(self, multiplier=1.5):
5         self.multiplier = multiplier
6         self.bounds_ = {}
7
8     def fit(self, X, y=None):
9         for col in X.columns:
10            Q1, Q3 = X[col].quantile([0.25, 0.75])
11            IQR = Q3 - Q1
12            self.bounds_[col] = (Q1 - self.multiplier * IQR,
13                               Q3 + self.multiplier * IQR)
14
15        return self
16
17    def transform(self, X):
18        X = X.copy()
19        for col in X.columns:
20            lower, upper = self.bounds_[col]
21            X[col] = X[col].clip(lower=lower, upper=upper)
22
23        return X
24
25 # Use in pipeline
26 pipeline = Pipeline([
27     ('outlier', OutlierRemover(multiplier=1.5)),
28     ('imputer', SimpleImputer(strategy='median'))
29 ])
```

Saving and Loading Clean Data

```
1 # Save cleaned data
2 df_clean.to_csv('data_cleaned.csv', index=False)
3 df_clean.to_parquet('data_cleaned.parquet', index=False) # Better for large files
4 df_clean.to_pickle('data_cleaned.pkl') # Preserves dtypes
5
6 # Save with compression
7 df_clean.to_csv('data_cleaned.csv.gz', compression='gzip', index=False)
8
9 # Load cleaned data
10 df = pd.read_csv('data_cleaned.csv')
11 df = pd.read_parquet('data_cleaned.parquet')
12 df = pd.read_pickle('data_cleaned.pkl')
13
14 # Save preprocessing pipeline (for production)
15 import joblib
16 joblib.dump(preprocessor, 'preprocessor.pkl')
17
18 # Load and apply to new data
19 preprocessor = joblib.load('preprocessor.pkl')
20 X_new_clean = preprocessor.transform(df_new)
```

Data Cleaning Checklist

Missing Values:

- Identify missing patterns
- Decide: delete or impute
- Apply appropriate imputation
- Add missing indicators if needed

Duplicates:

- Check exact duplicates
- Check key-based duplicates
- Check fuzzy duplicates
- Remove or aggregate

Outliers:

- Detect with IQR/Z-score
- Verify if real or errors
- Remove, cap, or transform

Data Types & Text:

- Fix incorrect types
- Parse dates properly
- Clean string columns
- Standardize categories

Key Python Functions Reference

Task	Function
Detect missing	<code>df.isnull().sum()</code> , <code>msno.matrix()</code>
Drop missing	<code>df.dropna()</code> , <code>df.dropna(subset=[...])</code>
Fill missing	<code>df.fillna()</code> , <code>df.interpolate()</code>
Impute (sklearn)	<code>SimpleImputer</code> , <code>KNNImputer</code>
Detect duplicates	<code>df.duplicated()</code> , <code>fuzz.ratio()</code>
Remove duplicates	<code>df.drop_duplicates()</code>
Detect outliers	<code>zscore()</code> , <code>IsolationForest</code>
Cap outliers	<code>df.clip()</code> , <code>winsorize()</code>
Convert types	<code>astype()</code> , <code>pd.to_numeric()</code>
Clean strings	<code>str.strip()</code> , <code>str.replace()</code>

Questions & Practice

Next: Part 3 - Data Transformation with Python

Part VI

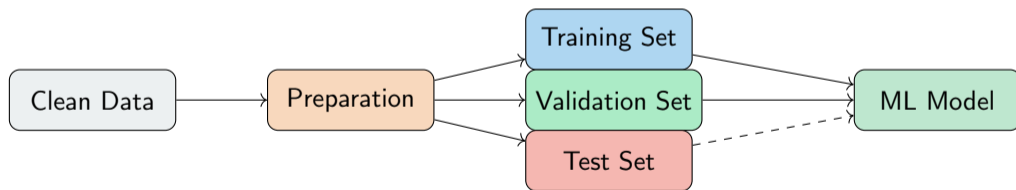
Part 4: Preparation

"Ready your data"

The Final Step: Data Preparation

Definition

Data Preparation for Modeling is the final preprocessing phase where clean, transformed data is organized into appropriate subsets for training, validation, and testing machine learning models.



Why Proper Data Preparation Matters

Without Proper Preparation

- **Data Leakage:** Test info leaks into training
- **Overfitting:** Model memorizes noise
- **Biased Evaluation:** Unrealistic performance
- **Poor Generalization:** Fails on new data

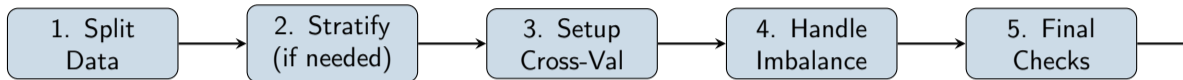
With Proper Preparation

- **Honest Evaluation:** True performance estimate
- **Robust Models:** Generalizes well
- **Reliable Comparisons:** Fair model selection
- **Production Ready:** Works on real data

Key Principle

The test set must simulate truly unseen data — never use it for any decision during model development!

Data Preparation Workflow



Topics Covered

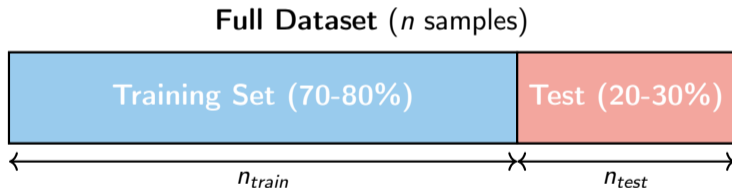
- 1 Train/Test/Validation Splitting
- 2 Cross-Validation Techniques
- 3 Stratified Sampling
- 4 Handling Imbalanced Datasets
- 5 Time Series Considerations
- 6 Final Validation Checklist

The Fundamental Split

Purpose

Divide data into separate sets to:

- **Train:** Learn patterns from data
- **Test:** Evaluate generalization to unseen data



Mathematical Notation

$$\mathcal{D} = \mathcal{D}_{train} \cup \mathcal{D}_{test}, \quad \mathcal{D}_{train} \cap \mathcal{D}_{test} = \emptyset$$

Common Split Ratios

Dataset Size	Train	Validation	Test
Small (< 1,000)	60%	20%	20%
Medium (1,000 – 100,000)	70%	15%	15%
Large (100,000 – 1M)	80%	10%	10%
Very Large (> 1M)	98%	1%	1%

Two-Way Split

$$\mathcal{D} \rightarrow \mathcal{D}_{train}, \mathcal{D}_{test}$$

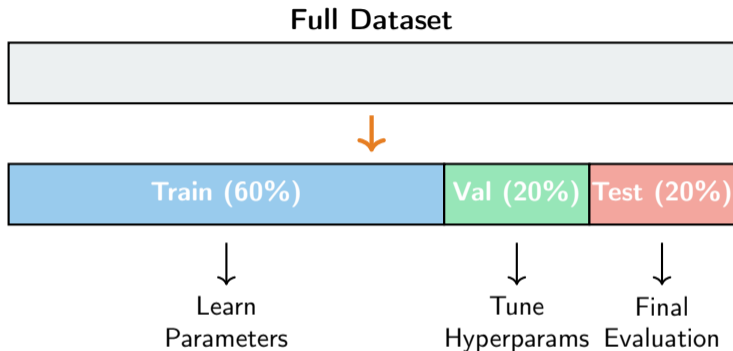
Use cross-validation on train set

Three-Way Split

$$\mathcal{D} \rightarrow \mathcal{D}_{train}, \mathcal{D}_{val}, \mathcal{D}_{test}$$

Explicit validation set

Train-Validation-Test Split



Role of Each Set

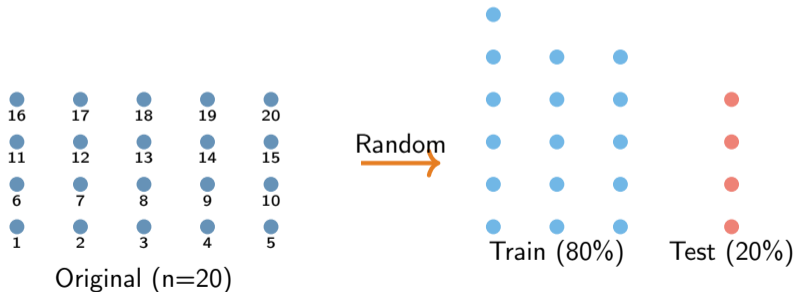
- **Training:** Fit model parameters (θ)
- **Validation:** Tune hyperparameters (λ , learning rate, etc.)
- **Test:** Final evaluation of the model performance (ONCE)

Random Sampling: The Basics

Simple Random Split

Each sample has equal probability of being in test set:

$$P(x_i \in \mathcal{D}_{test}) = \frac{n_{test}}{n} = 1 - \frac{n_{train}}{n}$$



Importance of Random State (Reproducibility)

Setting Random Seed

`random_state = s` \implies Same split every time

Without Fixed Seed

- Different splits each run
- Non-reproducible results
- Can't compare experiments
- Debugging nightmare

With Fixed Seed

- Reproducible experiments
- Fair model comparisons
- Easier debugging
- Scientific rigor

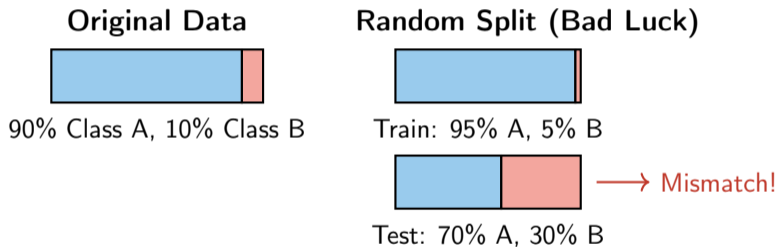
Best Practice

Always set `random_state=42` (or any fixed integer) and document it!

The Problem with Random Sampling

Class Imbalance Issue

Random sampling may not preserve class distribution, especially with imbalanced data!



Consequence

Model trained on different distribution than it's tested on → unreliable evaluation!

Stratified Sampling: The Solution

Definition

Stratified sampling ensures each split has the same proportion of classes as the original dataset:

$$\frac{n_{c,train}}{n_{train}} = \frac{n_{c,test}}{n_{test}} = \frac{n_c}{n} \quad \forall \text{ class } c$$

Original (90%/10%) Train (90%/10%) Test (90%/10%)



Result

Both train and test maintain the original 90%/10% class ratio!

When to Use Stratified Sampling

Always Stratify When:

- Classification problems (especially imbalanced)
- Small datasets (random variation is high)
- Multi-class problems
- Target variable has skewed distribution

Stratification for Regression

For continuous targets, bin the target into quantiles and stratify:

$$y \rightarrow \text{bins}(y) = \{Q_1, Q_2, Q_3, Q_4\}$$

Then stratify on bins to ensure similar target distributions.

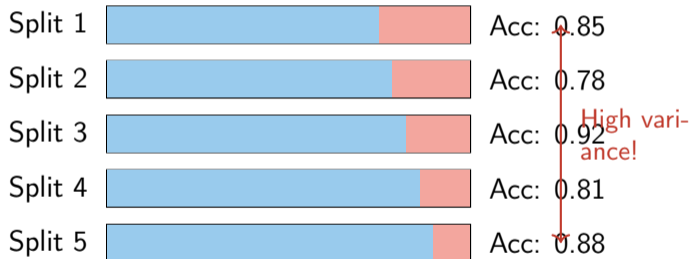
Multi-Label Stratification

For multi-label problems, use iterative stratification algorithms (more complex).

Why Cross-Validation?

Problem with Single Split

A single train/test split may not be representative — performance depends on which samples end up in each set.



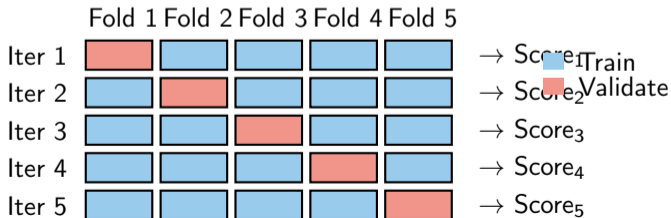
Solution

Cross-validation: Use multiple splits and average the results for a robust estimate.

K-Fold Cross-Validation

Algorithm

- 1 Divide data into K equal-sized folds
- 2 For each fold $k = 1, \dots, K$:
 - Use fold k as validation set
 - Use remaining $K - 1$ folds as training set
 - Train model and record performance
- 3 Average performance across all K folds



K-Fold: Mathematical Formulation

Performance Estimation

Let $\mathcal{L}(\hat{f}^{(-k)}, \mathcal{D}_k)$ be the loss on fold k using model trained without fold k :

$$CV(K) = \frac{1}{K} \sum_{k=1}^K \mathcal{L}(\hat{f}^{(-k)}, \mathcal{D}_k)$$

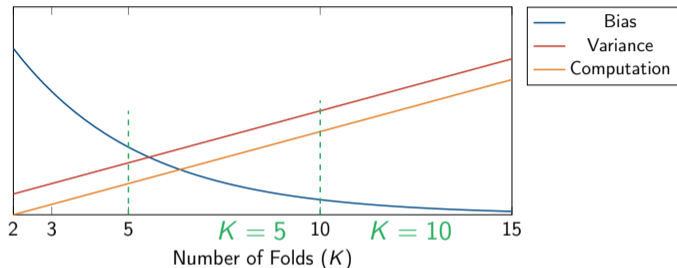
Variance of CV Estimate

$$\text{Var}[CV(K)] = \frac{1}{K} \sigma^2 + \frac{K-1}{K} \rho \sigma^2$$

Where:

- σ^2 = variance of individual fold estimates
- ρ = correlation between fold estimates (due to overlapping training sets)

Typical Values

Choosing K : Trade-offs

Common Choices

- $K = 5$: Fast, good for large datasets
- $K = 10$: Standard choice, good balance
- $K = n$ (LOOCV): Minimum bias, maximum variance, expensive

Leave-One-Out Cross-Validation (LOOCV)

Definition

Special case where $K = n$ (each sample is its own fold):

$$\text{LOOCV} = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(\hat{f}^{(-i)}, (x_i, y_i))$$



...

n iterations total

Pros:

Cons:

Stratified K-Fold Cross-Validation

Definition

K-Fold CV that preserves class distribution in each fold:

$$\frac{n_{c,k}}{n_k} \approx \frac{n_c}{n} \quad \forall \text{ fold } k, \text{ class } c$$

Each fold maintains 80%/20% ratio



Always Use Stratified K-Fold for Classification!

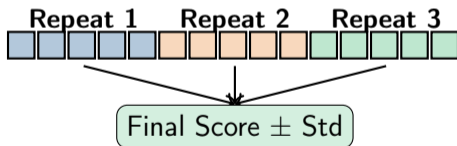
Repeated K-Fold Cross-Validation

Definition

Repeat K-Fold CV multiple times with different random splits:

$$\text{Repeated-CV} = \frac{1}{R \cdot K} \sum_{r=1}^R \sum_{k=1}^K \mathcal{L}(\hat{f}^{(r,-k)}, \mathcal{D}_k^{(r)})$$

Where R = number of repetitions



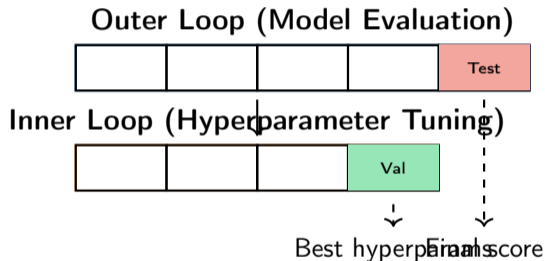
Common Configuration

5-Fold repeated 10 times = 50 total model fits → robust estimate!

Nested Cross-Validation

Purpose

Unbiased evaluation when both model selection AND hyperparameter tuning are needed.



Process

- **Inner loop:** Tune hyperparameters using CV on training folds
- **Outer loop:** Evaluate best model on held-out test fold

Cross-Validation Methods Comparison

Method	Bias	Variance	Cost	Best For
Hold-out	High	High	Low	Quick testing
5-Fold CV	Medium	Medium	Medium	Large datasets
10-Fold CV	Low	Medium	Medium	Standard choice
LOOCV	Very Low	High	High	Small datasets
Stratified K-Fold	Low	Medium	Medium	Classification
Repeated K-Fold	Very Low	Low	High	Final evaluation
Nested CV	Very Low	Low	Very High	Model selection

Recommendation

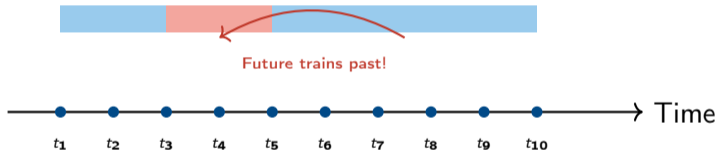
- Default: Stratified 10-Fold CV
- Final reporting: Repeated Stratified K-Fold
- With hyperparameter tuning: Nested CV

Why Standard CV Fails for Time Series

Temporal Leakage Problem

Standard K-Fold uses future data to predict the past — this is **data leakage**!

Wrong: Standard K-Fold



Key Principle

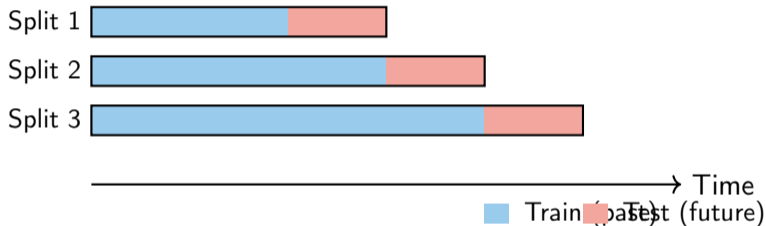
In time series: **Always train on past, test on future!**

$$t_{train} < t_{test} \quad (\text{strictly})$$

Time Series Split (Walk-Forward Validation)

Definition

Expanding window where training set grows over time:



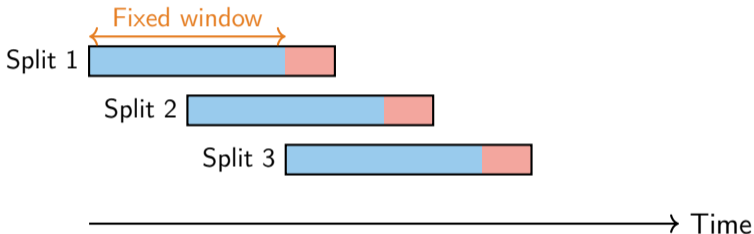
Properties

- Training set always precedes test set
- Training set grows with each split
- Mimics real-world forecasting scenario

Sliding Window Cross-Validation

Definition

Fixed-size training window that slides through time:



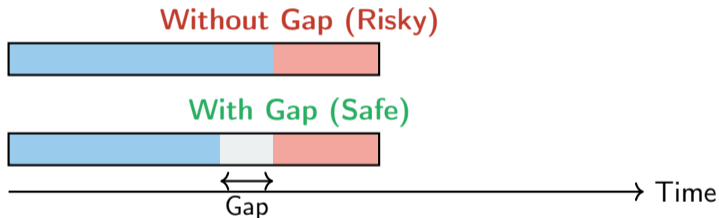
When to Use

- When recent data is more relevant than old data
- Non-stationary time series
- Concept drift scenarios

Gap Between Train and Test

Forecast Horizon Gap

Add gap to prevent information leakage from lag features:



Gap Size

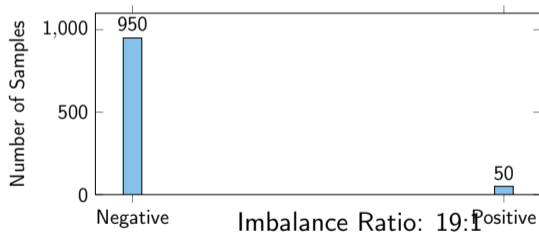
$$\text{Gap} \geq \max(\text{lag features}, \text{forecast horizon})$$

The Class Imbalance Problem

Definition

A dataset where classes have significantly different frequencies:

$$\text{Imbalance Ratio} = \frac{n_{\text{majority}}}{n_{\text{minority}}} \gg 1$$



Real-World Examples

Fraud detection (0.1%), disease diagnosis (1-5%), spam (10-20%)

Why Imbalance is Problematic

Model Behavior on Imbalanced Data

A model can achieve 95% accuracy by simply predicting the majority class!

$$\text{Accuracy} = \frac{950 + 0}{950 + 50} = 95\% \quad (\text{useless model!})$$

	Pred Pos
Actual Neg	0 (FP)
Actual Pos	0 (TP)

Precision = $0/0 = \text{unde-}$
fined

Recall = $0/50 = 0\%$

F1 = 0

Acc = 95% but useless!

Better Metrics for Imbalanced Data

Precision, Recall, F1-Score

$$\text{Precision} = \frac{TP}{TP + FP}, \quad \text{Recall} = \frac{TP}{TP + FN}$$
$$F_1 = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} = \frac{2TP}{2TP + FP + FN}$$

F-beta Score (Weighted)

$$F_\beta = (1 + \beta^2) \cdot \frac{\text{Precision} \cdot \text{Recall}}{\beta^2 \cdot \text{Precision} + \text{Recall}}$$

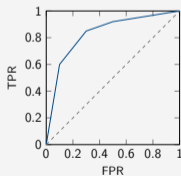
- $\beta = 1$: Equal weight (standard F1)
- $\beta = 2$: Recall is 2x more important than precision
- $\beta = 0.5$: Precision is 2x more important than recall

ROC-AUC and PR-AUC

ROC Curve

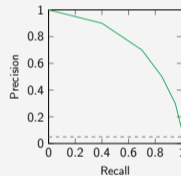
$$TPR = \frac{TP}{TP + FN}$$

$$FPR = \frac{FP}{FP + TN}$$



Precision-Recall Curve

Better for imbalanced data!



Baseline = $\frac{n_{pos}}{n}$ (class ratio)

Resampling: Oversampling Techniques

Random Oversampling

Duplicate minority class samples randomly:

$$n'_{minority} = n_{majority} \quad \text{or} \quad n'_{minority} = k \cdot n_{minority}$$

SMOTE (Synthetic Minority Oversampling Technique)

Generate synthetic samples along line segments between minority samples:

$$x_{new} = x_i + \lambda \cdot (x_{nn} - x_i), \quad \lambda \in [0, 1]$$

Where x_{nn} is a nearest neighbor of x_i



SMOTE Variants for Different Scenarios

- **Borderline-SMOTE:** Only oversample near decision boundary
- **ADASYN:** Generate more samples where density is low
- **SMOTE-NC:** Handles mixed numerical and categorical features
- **SVMSMOTE:** Uses SVM to identify support vectors



Borderline-SMOTE focuses here

Resampling: Undersampling Techniques

Random Undersampling

Remove majority class samples randomly:

$$n'_{majority} = n_{minority} \quad \text{or} \quad n'_{majority} = k \cdot n_{minority}$$

Risk: May lose important information!

Informed Undersampling

- **Tomek Links:** Remove majority samples that form Tomek links
- **ENN (Edited Nearest Neighbors):** Remove samples misclassified by KNN
- **NearMiss:** Keep majority samples closest to minority
- **Cluster Centroids:** Replace majority samples with cluster centers

Tomek Link

Pair (x_i, x_j) where they are each other's nearest neighbor but different classes.

SMOTEENN

Combine SMOTE oversampling with ENN cleaning:

$$\mathcal{D}' = \text{ENN}(\text{SMOTE}(\mathcal{D}))$$

SMOTETomek

Combine SMOTE with Tomek links removal:

$$\mathcal{D}' = \text{Tomek}(\text{SMOTE}(\mathcal{D}))$$

Class Weights

Cost-Sensitive Learning

Assign higher weight to minority class in loss function:

$$\mathcal{L}_{weighted} = \sum_{i=1}^n w_{y_i} \cdot \ell(y_i, \hat{y}_i)$$

Where weights are typically:

$$w_c = \frac{n}{k \cdot n_c}$$

n = total samples, k = number of classes, n_c = samples in class c

Example

- 950 negative, 50 positive samples
- $w_{neg} = \frac{1000}{2 \times 950} = 0.53$
- $w_{pos} = \frac{1000}{2 \times 50} = 10.0$

Threshold Adjustment

Default Threshold

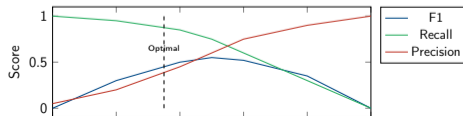
Most classifiers use threshold $\tau = 0.5$:

$$\hat{y} = \begin{cases} 1 & \text{if } P(y = 1|x) \geq 0.5 \\ 0 & \text{otherwise} \end{cases}$$

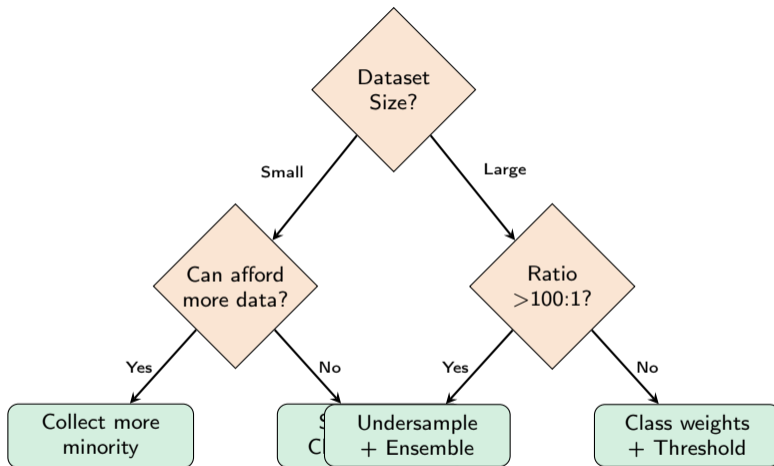
Adjusting Threshold for Imbalance

Lower threshold to increase recall for minority class:

$$\hat{y} = \begin{cases} 1 & \text{if } P(y = 1|x) \geq \tau \\ 0 & \text{otherwise} \end{cases} \quad \text{where } \tau < 0.5$$



Imbalanced Data: Strategy Summary



Important Reminder

Apply resampling ONLY to training data! Test set must remain untouched to reflect

Pre-Modeling Validation Checklist

Data Quality:

- No missing values (or handled)
- No duplicate records
- Outliers addressed
- Data types correct
- No invalid values

Transformations:

- Numerical features scaled
- Categorical features encoded
- Features engineered
- Transformations fitted on train only

Splitting:

- Train/val/test split done
- Stratification applied (if needed)
- No data leakage
- Random state documented

Imbalance:

- Class distribution checked
- Appropriate metrics chosen
- Resampling applied (train only)
- Class weights considered

Data Leakage: Final Check

Common Sources of Data Leakage

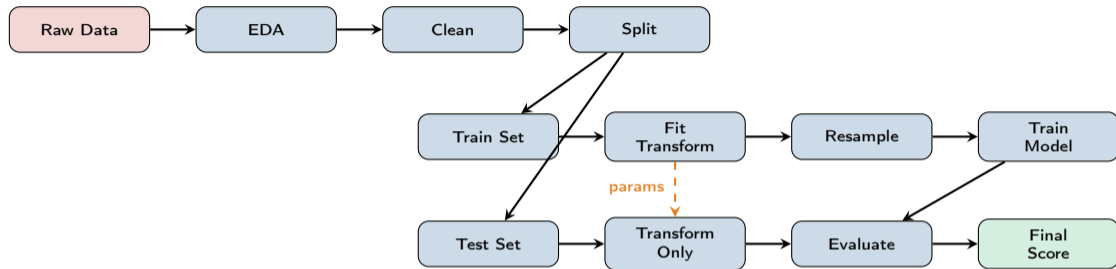
- 1 **Preprocessing on full data:** Scaling, imputation, encoding
- 2 **Feature selection on full data:** Using test info to select features
- 3 **Time series mistakes:** Using future data to predict past
- 4 **Duplicate records:** Same record in train and test
- 5 **Target leakage:** Features that contain target information

Prevention Strategy

Pipeline = Preprocessing → Model

Use sklearn Pipeline to ensure all transformations are fit only on training data!

Complete ML Preparation Pipeline



Key Points

- EDA and cleaning on full data (before split)
- All transformers fitted on training data only
- Resampling applied to training data only
- Test set touched only for final evaluation

Key Takeaways

- 1 **Split first, transform second** — avoid data leakage
- 2 **Stratify for classification** — preserve class distribution
- 3 **Use cross-validation** — robust performance estimates
- 4 **Choose appropriate metrics** — accuracy is not always best
- 5 **Handle imbalance properly** — resample training data only
- 6 **Time series is special** — always respect temporal order
- 7 **Document everything** — reproducibility is key

Course Complete!

You've now covered all four parts of ML Pre-Processing:

- 1 EDA ✓
- 2 Data Cleaning ✓
- 3 Data Transformation ✓
- 4 Data Preparation ✓

Part VII

Part 4: Practice Preparation

"Ready your data"

What We Will Cover

Practical Data Preparation Skills

- Train-test splitting strategies
- Stratified sampling techniques
- Cross-validation methods (K-Fold, Stratified, LOOCV)
- Time series splitting
- Handling imbalanced datasets (SMOTE, undersampling)
- Evaluation metrics for imbalanced data
- Complete preprocessing pipelines
- Avoiding data leakage

Libraries Used

`scikit-learn`, `imbalanced-learn`, `pandas`, `numpy`

Essential Libraries Import

```
1 # Data manipulation
2 import pandas as pd
3 import numpy as np
4
5 # Train-test splitting
6 from sklearn.model_selection import (
7     train_test_split, KFold, StratifiedKFold, LeaveOneOut,
8     RepeatedKFold, RepeatedStratifiedKFold, TimeSeriesSplit,
9     cross_val_score, cross_validate, GridSearchCV
10 )
11
12 # Imbalanced data handling (pip install imbalanced-learn)
13 from imblearn.over_sampling import SMOTE, ADASYN, RandomOverSampler
14 from imblearn.under_sampling import RandomUnderSampler, TomekLinks
15 from imblearn.combine import SMOTETomek, SMOTEENN
16 from imblearn.pipeline import Pipeline as ImbPipeline
17
18 # Metrics
19 from sklearn.metrics import (
20     accuracy_score, precision_score, recall_score, f1_score,
21     classification_report, confusion_matrix, roc_auc_score
22 )
```

Basic Train-Test Split

```
1 from sklearn.model_selection import train_test_split
2
3 # Basic split (80% train, 20% test)
4 X_train, X_test, y_train, y_test = train_test_split(
5     X, y,
6     test_size=0.2,
7     random_state=42
8 )
9
10 print(f"Training set: {X_train.shape[0]} samples")
11 print(f"Test set: {X_test.shape[0]} samples")
12
13 # Alternative: specify train_size
14 X_train, X_test, y_train, y_test = train_test_split(
15     X, y,
16     train_size=0.8,
17     random_state=42
18 )
19
20 # Shuffle data (default is True)
21 X_train, X_test, y_train, y_test = train_test_split(
22     X, y,
23     test_size=0.2,
24     shuffle=True,
25     random_state=42
26 )
```

Train-Validation-Test Split

```
1 # Three-way split: 60% train, 20% validation, 20% test
2
3 # Method 1: Two sequential splits
4 X_train_full, X_test, y_train_full, y_test = train_test_split(
5     X, y, test_size=0.2, random_state=42
6 )
7 X_train, X_val, y_train, y_val = train_test_split(
8     X_train_full, y_train_full, test_size=0.25, random_state=42 # 0.25 * 0.8 = 0.2
9 )
10
11 print(f"Train: {len(X_train)} ({len(X_train)/len(X)*100:.0f}%)")
12 print(f"Val: {len(X_val)} ({len(X_val)/len(X)*100:.0f}%)")
13 print(f"Test: {len(X_test)} ({len(X_test)/len(X)*100:.0f}%)")
14
15 # Method 2: Using numpy split
16 indices = np.arange(len(X))
17 np.random.seed(42)
18 np.random.shuffle(indices)
19 train_idx = indices[:int(0.6*len(X))]
20 val_idx = indices[int(0.6*len(X)):int(0.8*len(X))]
21 test_idx = indices[int(0.8*len(X)):]
22
23 X_train, X_val, X_test = X[train_idx], X[val_idx], X[test_idx]
24 y_train, y_val, y_test = y[train_idx], y[val_idx], y[test_idx]
```

Stratified Split (Classification)

```
1 # Stratified split: preserves class distribution
2 X_train, X_test, y_train, y_test = train_test_split(
3     X, y,
4     test_size=0.2,
5     stratify=y,          # Key parameter!
6     random_state=42
7 )
8
9 # Verify class distribution is preserved
10 print("Original distribution:")
11 print(pd.Series(y).value_counts(normalize=True))
12
13 print("\nTraining distribution:")
14 print(pd.Series(y_train).value_counts(normalize=True))
15
16 print("\nTest distribution:")
17 print(pd.Series(y_test).value_counts(normalize=True))
18
19 # Stratified three-way split
20 X_temp, X_test, y_temp, y_test = train_test_split(
21     X, y, test_size=0.2, stratify=y, random_state=42
22 )
23 X_train, X_val, y_train, y_val = train_test_split(
24     X_temp, y_temp, test_size=0.25, stratify=y_temp, random_state=42
25 )
```

Stratified Split for Regression

```
1 # For regression: stratify by binned target
2
3 # Create bins from continuous target
4 n_bins = 5
5 y_binned = pd.qcut(y, q=n_bins, labels=False, duplicates='drop')
6
7 # Stratified split using binned target
8 X_train, X_test, y_train, y_test = train_test_split(
9     X, y,
10    test_size=0.2,
11    stratify=y_binned,
12    random_state=42
13 )
14
15 # Verify distribution
16 print("Original target stats:")
17 print(f"Mean: {y.mean():.2f}, Std: {y.std():.2f}")
18
19 print("\nTrain target stats:")
20 print(f"Mean: {y_train.mean():.2f}, Std: {y_train.std():.2f}")
21
22 print("\nTest target stats:")
23 print(f"Mean: {y_test.mean():.2f}, Std: {y_test.std():.2f}")
```

Group-Based Split

```
1 from sklearn.model_selection import GroupShuffleSplit
2
3 # When data has groups that should stay together
4 # Example: multiple samples from same patient, same customer, etc.
5
6 groups = df['patient_id'].values # Group identifier
7
8 # Group shuffle split
9 gss = GroupShuffleSplit(n_splits=1, test_size=0.2, random_state=42)
10 train_idx, test_idx = next(gss.split(X, y, groups))
11
12 X_train, X_test = X[train_idx], X[test_idx]
13 y_train, y_test = y[train_idx], y[test_idx]
14
15 # Verify no group appears in both sets
16 train_groups = set(groups[train_idx])
17 test_groups = set(groups[test_idx])
18 print(f"Groups in train: {len(train_groups)}")
19 print(f"Groups in test: {len(test_groups)}")
20 print(f"Overlap: {len(train_groups & test_groups)}") # Should be 0
21
22 # Alternative: GroupKFold for cross-validation
23 from sklearn.model_selection import GroupKFold
24 gkf = GroupKFold(n_splits=5)
25 for train_idx, test_idx in gkf.split(X, y, groups):
26     # Each fold has different groups
27     pass
```

K-Fold Cross-Validation

```
1 from sklearn.model_selection import KFold, cross_val_score
2 from sklearn.ensemble import RandomForestClassifier
3
4 # Create K-Fold splitter
5 kfold = KFold(n_splits=5, shuffle=True, random_state=42)
6
7 # Method 1: Manual iteration
8 model = RandomForestClassifier(random_state=42)
9 scores = []
10 for fold, (train_idx, val_idx) in enumerate(kfold.split(X)):
11     X_train, X_val = X[train_idx], X[val_idx]
12     y_train, y_val = y[train_idx], y[val_idx]
13
14     model.fit(X_train, y_train)
15     score = model.score(X_val, y_val)
16     scores.append(score)
17     print(f"Fold {fold+1}: {score:.4f}")
18
19 print(f"\nMean: {np.mean(scores):.4f} (+/- {np.std(scores)*2:.4f})")
20
21 # Method 2: Using cross_val_score (simpler)
22 scores = cross_val_score(model, X, y, cv=5, scoring='accuracy')
23 print(f"CV Accuracy: {scores.mean():.4f} (+/- {scores.std()*2:.4f})")
```

Stratified K-Fold Cross-Validation

```
1 from sklearn.model_selection import StratifiedKFold, cross_val_score
2
3 # Stratified K-Fold: preserves class distribution in each fold
4 skfold = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
5
6 # Manual iteration
7 for fold, (train_idx, val_idx) in enumerate(skfold.split(X, y)):
8     y_train_fold = y[train_idx]
9     y_val_fold = y[val_idx]
10
11     # Check distribution
12     train_dist = pd.Series(y_train_fold).value_counts(normalize=True)
13     val_dist = pd.Series(y_val_fold).value_counts(normalize=True)
14     print(f"Fold {fold+1} - Train: {dict(train_dist)}, Val: {dict(val_dist)}")
15
16 # Using cross_val_score with stratified CV
17 model = RandomForestClassifier(random_state=42)
18 scores = cross_val_score(model, X, y, cv=skfold, scoring='accuracy')
19 print(f"\nStratified CV Accuracy: {scores.mean():.4f} (+/- {scores.std()*2:.4f})")
20
21 # Shortcut: cv=5 with classifier automatically uses StratifiedKFold
22 scores = cross_val_score(model, X, y, cv=5) # Auto-stratified for classifiers
```

Leave-One-Out Cross-Validation (LOOCV)

```
1 from sklearn.model_selection import LeaveOneOut, cross_val_score
2
3 # LOOCV: each sample is used once as validation
4 loo = LeaveOneOut()
5
6 print(f"Number of splits: {loo.get_n_splits(X)}") # = n_samples
7
8 # Using cross_val_score (can be slow for large datasets!)
9 model = RandomForestClassifier(n_estimators=10, random_state=42)
10 scores = cross_val_score(model, X, y, cv=loo)
11
12 print(f"LOOCV Accuracy: {scores.mean():.4f} (+/- {scores.std()*2:.4f})")
13
14 # Manual iteration (for more control)
15 predictions = []
16 for train_idx, val_idx in loo.split(X):
17     X_train, X_val = X[train_idx], X[val_idx]
18     y_train, y_val = y[train_idx], y[val_idx]
19
20     model.fit(X_train, y_train)
21     pred = model.predict(X_val)
22     predictions.append(pred[0])
23
24 accuracy = np.mean(np.array(predictions) == y)
25 print(f"LOOCV Accuracy (manual): {accuracy:.4f}")
```

Repeated K-Fold Cross-Validation

```
1 from sklearn.model_selection import RepeatedKFold, RepeatedStratifiedKFold
2
3 # Repeated K-Fold: run K-Fold multiple times with different shuffles
4 rkfold = RepeatedKFold(n_splits=5, n_repeats=10, random_state=42)
5
6 print(f"Total iterations: {rkfold.get_n_splits(X)}") # 5 * 10 = 50
7
8 model = RandomForestClassifier(random_state=42)
9 scores = cross_val_score(model, X, y, cv=rkfold, scoring='accuracy')
10
11 print(f"Repeated CV: {scores.mean():.4f} (+/- {scores.std()*2:.4f})")
12
13 # Repeated Stratified K-Fold (recommended for classification)
14 rskfold = RepeatedStratifiedKFold(n_splits=5, n_repeats=10, random_state=42)
15
16 scores = cross_val_score(model, X, y, cv=rskfold, scoring='accuracy')
17 print(f"Repeated Stratified CV: {scores.mean():.4f} (+/- {scores.std()*2:.4f})")
18
19 # Detailed results per fold
20 for i, score in enumerate(scores):
21     repeat = i // 5 + 1
22     fold = i % 5 + 1
23     print(f"Repeat {repeat}, Fold {fold}: {score:.4f}")
```

Cross-Validation with Multiple Metrics

```
1 from sklearn.model_selection import cross_validate
2
3 # Get multiple metrics at once
4 scoring = ['accuracy', 'precision_macro', 'recall_macro', 'f1_macro', 'roc_auc_ovr']
5
6 model = RandomForestClassifier(random_state=42)
7 cv_results = cross_validate(
8     model, X, y,
9     cv=5,
10    scoring=scoring,
11    return_train_score=True
12 )
13
14 # Print results
15 for metric in scoring:
16     train_key = f'train_{metric}'
17     test_key = f'test_{metric}'
18     print(f"\n{metric}:")
19     print(f"  Train: {cv_results[train_key].mean():.4f} (+/- {cv_results[train_key].std()*2:.4f})")
20     print(f"  Test:  {cv_results[test_key].mean():.4f} (+/- {cv_results[test_key].std()*2:.4f})")
21
22 # Also returns fit_time and score_time
23 print(f"\nAverage fit time: {cv_results['fit_time'].mean():.2f}s")
24 print(f"Average score time: {cv_results['score_time'].mean():.2f}s")
```

Cross-Validation with Predictions

```
1 from sklearn.model_selection import cross_val_predict
2
3 # Get predictions for each sample (when it was in validation set)
4 model = RandomForestClassifier(random_state=42)
5 y_pred = cross_val_predict(model, X, y, cv=5)
6
7 # Now can compute any metric
8 from sklearn.metrics import classification_report, confusion_matrix
9
10 print("Classification Report (CV predictions):")
11 print(classification_report(y, y_pred))
12
13 print("\nConfusion Matrix:")
14 print(confusion_matrix(y, y_pred))
15
16 # Get probabilities instead of predictions
17 y_proba = cross_val_predict(model, X, y, cv=5, method='predict_proba')
18 print(f"\nProbability shape: {y_proba.shape}")
19
20 # Compute ROC-AUC from probabilities
21 from sklearn.metrics import roc_auc_score
22 if len(np.unique(y)) == 2: # Binary classification
23     auc = roc_auc_score(y, y_proba[:, 1])
24     print(f"ROC-AUC: {auc:.4f}")
```

Time Series Split

```
1 from sklearn.model_selection import TimeSeriesSplit
2
3 # Time series split: train on past, test on future
4 tscv = TimeSeriesSplit(n_splits=5)
5
6 # Visualize splits
7 for fold, (train_idx, val_idx) in enumerate(tscv.split(X)):
8     print(f"Fold {fold+1}:")
9     print(f"  Train: indices {train_idx[0]} to {train_idx[-1]} ({len(train_idx)} samples)")
10    print(f"  Val:   indices {val_idx[0]} to {val_idx[-1]} ({len(val_idx)} samples)")
11
12 # Use with cross_val_score
13 from sklearn.ensemble import GradientBoostingRegressor
14
15 model = GradientBoostingRegressor(random_state=42)
16 scores = cross_val_score(model, X, y, cv=tscv, scoring='neg_mean_squared_error')
17
18 rmse_scores = np.sqrt(-scores)
19 print(f"\nTime Series CV RMSE: {rmse_scores.mean():.4f} (+/- {rmse_scores.std()*2:.4f})")
```

Time Series Split with Gap

```
1 from sklearn.model_selection import TimeSeriesSplit
2
3 # Add gap between train and test to prevent leakage from lag features
4 tscv = TimeSeriesSplit(n_splits=5, gap=10) # 10 samples gap
5
6 for fold, (train_idx, val_idx) in enumerate(tscv.split(X)):
7     print(f"Fold {fold+1}:")
8     print(f"  Train ends at: {train_idx[-1]}")
9     print(f"  Gap: {val_idx[0] - train_idx[-1] - 1} samples")
10    print(f"  Val starts at: {val_idx[0]}")
11
12 # With maximum train size (sliding window)
13 tscv = TimeSeriesSplit(n_splits=5, max_train_size=1000)
14
15 for fold, (train_idx, val_idx) in enumerate(tscv.split(X)):
16     print(f"Fold {fold+1}: Train size = {len(train_idx)}")
17
18 # Custom time series split with test size
19 tscv = TimeSeriesSplit(n_splits=5, test_size=100)
20
21 for fold, (train_idx, val_idx) in enumerate(tscv.split(X)):
22     print(f"Fold {fold+1}: Val size = {len(val_idx)}")
```

Custom Sliding Window Split

```
1 def sliding_window_split(X, window_size, test_size, step=1):
2     """Custom sliding window time series split"""
3     n_samples = len(X)
4     indices = np.arange(n_samples)
5
6     splits = []
7     start = 0
8     while start + window_size + test_size <= n_samples:
9         train_idx = indices[start:start + window_size]
10        test_idx = indices[start + window_size:start + window_size + test_size]
11        splits.append((train_idx, test_idx))
12        start += step
13
14    return splits
15
16 # Usage
17 splits = sliding_window_split(X, window_size=100, test_size=20, step=20)
18 print(f"Number of splits: {len(splits)}")
19
20 scores = []
21 for train_idx, test_idx in splits:
22     X_train, X_test = X[train_idx], X[test_idx]
23     y_train, y_test = y[train_idx], y[test_idx]
24
25     model.fit(X_train, y_train)
26     scores.append(model.score(X_test, y_test))
27
28 print(f"Sliding Window CV: {np.mean(scores):.4f} (+/- {np.std(scores)*2:.4f})")
```

Blocking Time Series Split

```
1 from sklearn.model_selection import TimeSeriesSplit
2
3 class BlockingTimeSeriesSplit:
4     """Purged time series split to prevent leakage"""
5     def __init__(self, n_splits=5, purge_gap=0, embargo_gap=0):
6         self.n_splits = n_splits
7         self.purge_gap = purge_gap # Gap before test set
8         self.embargo_gap = embargo_gap # Gap after test set
9
10    def split(self, X, y=None, groups=None):
11        n_samples = len(X)
12        fold_size = n_samples // (self.n_splits + 1)
13
14        for i in range(self.n_splits):
15            test_start = (i + 1) * fold_size
16            test_end = test_start + fold_size
17
18            train_end = test_start - self.purge_gap
19            train_idx = np.arange(0, train_end)
20            test_idx = np.arange(test_start, min(test_end, n_samples))
21
22            yield train_idx, test_idx
23
24    # Usage
25    btscv = BlockingTimeSeriesSplit(n_splits=5, purge_gap=5)
26    scores = cross_val_score(model, X, y, cv=btscv)
```

Checking Class Imbalance

```
1 # Check class distribution
2 print("Class distribution:")
3 print(pd.Series(y).value_counts())
4 print(f"\nImbalance ratio: {pd.Series(y).value_counts().max() / pd.Series(y).value_counts().min():.2f}")
5
6 # Visualize
7 import matplotlib.pyplot as plt
8 plt.figure(figsize=(8, 5))
9 pd.Series(y).value_counts().plot(kind='bar')
10 plt.title('Class Distribution')
11 plt.xlabel('Class')
12 plt.ylabel('Count')
13 plt.tight_layout()
14 plt.show()
15
16 # Calculate class weights
17 from sklearn.utils.class_weight import compute_class_weight
18
19 class_weights = compute_class_weight('balanced', classes=np.unique(y), y=y)
20 weight_dict = dict(zip(np.unique(y), class_weights))
21 print(f"\nComputed class weights: {weight_dict}")
```

Random Over-Sampling

```
1 from imblearn.over_sampling import RandomOverSampler
2
3 # Random over-sampling: duplicate minority class samples
4 ros = RandomOverSampler(random_state=42)
5 X_resampled, y_resampled = ros.fit_resample(X_train, y_train)
6
7 print(f"Original: {pd.Series(y_train).value_counts().to_dict()}")
8 print(f"Resampled: {pd.Series(y_resampled).value_counts().to_dict()}")
9
10 # With specific sampling strategy
11 ros = RandomOverSampler(
12     sampling_strategy='minority', # Only oversample minority
13     random_state=42
14 )
15 X_resampled, y_resampled = ros.fit_resample(X_train, y_train)
16
17 # Custom ratio
18 ros = RandomOverSampler(
19     sampling_strategy=0.5, # Minority will be 50% of majority
20     random_state=42
21 )
22 X_resampled, y_resampled = ros.fit_resample(X_train, y_train)
23
24 # Specific counts per class
25 ros = RandomOverSampler(
26     sampling_strategy={0: 1000, 1: 1000}, # Both classes to 1000
27     random_state=42
28 )
```

SMOTE (Synthetic Minority Over-sampling)

```
1 from imblearn.over_sampling import SMOTE
2
3 # SMOTE: create synthetic samples
4 smote = SMOTE(random_state=42)
5 X_resampled, y_resampled = smote.fit_resample(X_train, y_train)
6
7 print(f"Original: {pd.Series(y_train).value_counts().to_dict()}")
8 print(f"After SMOTE: {pd.Series(y_resampled).value_counts().to_dict()}")
9
10 # SMOTE with parameters
11 smote = SMOTE(
12     sampling_strategy='auto', # Resample all classes but majority
13     k_neighbors=5,           # Number of neighbors for synthesis
14     random_state=42
15 )
16
17 # Borderline-SMOTE: focus on borderline samples
18 from imblearn.over_sampling import BorderlineSMOTE
19 bsmote = BorderlineSMOTE(random_state=42, kind='borderline-1')
20 X_resampled, y_resampled = bsmote.fit_resample(X_train, y_train)
21
22 # SVMSMOTE: use SVM to identify samples
23 from imblearn.over_sampling import SVMSMOTE
24 svmsmote = SVMSMOTE(random_state=42)
25 X_resampled, y_resampled = svmsmote.fit_resample(X_train, y_train)
```

ADASYN (Adaptive Synthetic Sampling)

```
1 from imblearn.over_sampling import ADASYN
2
3 # ADASYN: more synthetic samples where density is low
4 adasyn = ADASYN(random_state=42)
5 X_resampled, y_resampled = adasyn.fit_resample(X_train, y_train)
6
7 print(f"Original: {pd.Series(y_train).value_counts().to_dict()}")
8 print(f"After ADASYN: {pd.Series(y_resampled).value_counts().to_dict()}")
9
10 # Note: ADASYN may not produce exact target ratio
11 # because it adapts to local data density
12
13 # ADASYN parameters
14 adasyn = ADASYN(
15     sampling_strategy='minority',
16     n_neighbors=5,
17     random_state=42
18 )
19
20 # Comparison: SMOTE vs ADASYN
21 # SMOTE: uniform synthetic samples
22 # ADASYN: more samples in harder-to-learn regions
```

Random Under-Sampling

```
1 from imblearn.under_sampling import RandomUnderSampler
2
3 # Random under-sampling: remove majority class samples
4 rus = RandomUnderSampler(random_state=42)
5 X_resampled, y_resampled = rus.fit_resample(X_train, y_train)
6
7 print(f"Original: {pd.Series(y_train).value_counts().to_dict()}")
8 print(f"After undersampling: {pd.Series(y_resampled).value_counts().to_dict()}")
9
10 # Warning: may lose important information!
11
12 # With replacement (bootstrap)
13 rus = RandomUnderSampler(replacement=True, random_state=42)
14
15 # Specific ratio
16 rus = RandomUnderSampler(
17     sampling_strategy=0.5, # Majority will be 2x minority
18     random_state=42
19 )
20
21 # NearMiss: keep majority samples close to minority
22 from imblearn.under_sampling import NearMiss
23 nm = NearMiss(version=1) # version 1, 2, or 3
24 X_resampled, y_resampled = nm.fit_resample(X_train, y_train)
```

Tomek Links and ENN

```
1 from imblearn.under_sampling import TomekLinks, EditedNearestNeighbours
2
3 # Tomek Links: remove majority samples that form Tomek links
4 tokek = TomekLinks()
5 X_resampled, y_resampled = tokek.fit_resample(X_train, y_train)
6 print(f"Removed {len(X_train) - len(X_resampled)} samples (Tomek)")
7
8 # ENN: remove samples misclassified by KNN
9 enn = EditedNearestNeighbours(n_neighbors=3)
10 X_resampled, y_resampled = enn.fit_resample(X_train, y_train)
11 print(f"Removed {len(X_train) - len(X_resampled)} samples (ENN)")
12
13 # Repeated ENN
14 from imblearn.under_sampling import RepeatedEditedNearestNeighbours
15 renn = RepeatedEditedNearestNeighbours()
16 X_resampled, y_resampled = renn.fit_resample(X_train, y_train)
17
18 # These methods clean the data by removing noisy samples
19 # Often used in combination with SMOTE
```

Combined Sampling: SMOTE + Tomek/ENN

```
1 from imblearn.combine import SMOTETomek, SMOTEENN
2
3 # SMOTE + Tomek Links: oversample then clean
4 smote_tomek = SMOTETomek(random_state=42)
5 X_resampled, y_resampled = smote_tomek.fit_resample(X_train, y_train)
6
7 print(f"Original: {pd.Series(y_train).value_counts().to_dict()}")
8 print(f"After SMOTE-Tomek: {pd.Series(y_resampled).value_counts().to_dict()}")
9
10 # SMOTE + ENN: oversample then clean with ENN
11 smote_enn = SMOTEENN(random_state=42)
12 X_resampled, y_resampled = smote_enn.fit_resample(X_train, y_train)
13
14 print(f"After SMOTE-ENN: {pd.Series(y_resampled).value_counts().to_dict()}")
15
16 # Custom combination
17 from imblearn.over_sampling import SMOTE
18 from imblearn.under_sampling import TomekLinks
19 from imblearn.pipeline import Pipeline
20
21 resampling_pipeline = Pipeline([
22     ('smote', SMOTE(random_state=42)),
23     ('tomek', TomekLinks())
24 ])
25 X_resampled, y_resampled = resampling_pipeline.fit_resample(X_train, y_train)
```

Using Class Weights

```
1 from sklearn.ensemble import RandomForestClassifier
2 from sklearn.linear_model import LogisticRegression
3
4 # Method 1: class_weight='balanced' (auto-compute weights)
5 model = RandomForestClassifier(class_weight='balanced', random_state=42)
6 model.fit(X_train, y_train)
7
8 # Method 2: Custom class weights
9 class_weights = {0: 1, 1: 10} # Minority class 10x more important
10 model = RandomForestClassifier(class_weight=class_weights, random_state=42)
11
12 # Method 3: Compute balanced weights
13 from sklearn.utils.class_weight import compute_class_weight
14 weights = compute_class_weight('balanced', classes=np.unique(y_train), y=y_train)
15 class_weights = dict(zip(np.unique(y_train), weights))
16 print(f"Computed weights: {class_weights}")
17
18 # For Logistic Regression
19 model = LogisticRegression(class_weight='balanced', random_state=42)
20
21 # Sample weights (per-sample, not per-class)
22 sample_weights = np.where(y_train == 1, 10, 1) # Minority samples weight 10
23 model.fit(X_train, y_train, sample_weight=sample_weights)
```

Imbalanced-Learn Pipeline

```
1 from imblearn.pipeline import Pipeline as ImbPipeline
2 from sklearn.preprocessing import StandardScaler
3 from sklearn.ensemble import RandomForestClassifier
4 from imblearn.over_sampling import SMOTE
5
6 # Pipeline that includes resampling
7 # Note: use imblearn's Pipeline, not sklearn's!
8 pipeline = ImbPipeline([
9     ('scaler', StandardScaler()),
10    ('smote', SMOTE(random_state=42)),
11    ('classifier', RandomForestClassifier(random_state=42))
12 ])
13
14 # Fit the pipeline
15 pipeline.fit(X_train, y_train)
16
17 # Predictions
18 y_pred = pipeline.predict(X_test)
19
20 # Cross-validation with resampling
21 from sklearn.model_selection import cross_val_score
22 scores = cross_val_score(pipeline, X, y, cv=5, scoring='f1')
23 print(f"CV F1: {scores.mean():.4f} (+/- {scores.std()*2:.4f})")
24
25 # Important: Resampling is applied only during fit, not predict!
```

Classification Report

```
1 from sklearn.metrics import classification_report, confusion_matrix
2
3 # Full classification report
4 print(classification_report(y_test, y_pred))
5
6 # With target names
7 print(classification_report(y_test, y_pred,
8                             target_names=['Negative', 'Positive']))
9
10 # As dictionary (for programmatic access)
11 report = classification_report(y_test, y_pred, output_dict=True)
12 print(f"Precision (class 1): {report['1']['precision']:.4f}")
13 print(f"Recall (class 1): {report['1']['recall']:.4f}")
14 print(f"F1 (class 1): {report['1']['f1-score']:.4f}")
15
16 # Confusion matrix
17 cm = confusion_matrix(y_test, y_pred)
18 print("\nConfusion Matrix:")
19 print(cm)
20
21 # As DataFrame for better visualization
22 cm_df = pd.DataFrame(cm, index=['Actual 0', 'Actual 1'],
23                       columns=['Pred 0', 'Pred 1'])
24 print(cm_df)
```

Precision, Recall, F1-Score

```
1 from sklearn.metrics import precision_score, recall_score, f1_score
2
3 # Binary classification
4 precision = precision_score(y_test, y_pred)
5 recall = recall_score(y_test, y_pred)
6 f1 = f1_score(y_test, y_pred)
7
8 print(f"Precision: {precision:.4f}")
9 print(f"Recall: {recall:.4f}")
10 print(f"F1-Score: {f1:.4f}")
11
12 # Multi-class: specify averaging method
13 precision_macro = precision_score(y_test, y_pred, average='macro')
14 precision_weighted = precision_score(y_test, y_pred, average='weighted')
15 precision_micro = precision_score(y_test, y_pred, average='micro')
16
17 print(f"\nMacro Precision: {precision_macro:.4f}") # Unweighted mean
18 print(f"Weighted Precision: {precision_weighted:.4f}") # Weighted by support
19 print(f"Micro Precision: {precision_micro:.4f}") # Global calculation
20
21 # F-beta score (custom beta)
22 from sklearn.metrics import fbeta_score
23 f2 = fbeta_score(y_test, y_pred, beta=2) # Recall 2x more important
24 f05 = fbeta_score(y_test, y_pred, beta=0.5) # Precision 2x more important
```

ROC-AUC Score

```
1 from sklearn.metrics import roc_auc_score, roc_curve
2
3 # Get probabilities
4 y_proba = model.predict_proba(X_test)[: , 1]
5
6 # ROC-AUC score
7 auc = roc_auc_score(y_test, y_proba)
8 print(f"ROC-AUC: {auc:.4f}")
9
10 # ROC curve
11 fpr, tpr, thresholds = roc_curve(y_test, y_proba)
12
13 # Plot ROC curve
14 plt.figure(figsize=(8, 6))
15 plt.plot(fpr, tpr, label=f'ROC (AUC = {auc:.4f})')
16 plt.plot([0, 1], [0, 1], 'k--', label='Random')
17 plt.xlabel('False Positive Rate')
18 plt.ylabel('True Positive Rate')
19 plt.title('ROC Curve')
20 plt.legend()
21 plt.show()
22
23 # Multi-class ROC-AUC
24 y_proba_multi = model.predict_proba(X_test)
25 auc_ovr = roc_auc_score(y_test, y_proba_multi, multi_class='ovr')
26 auc_ovo = roc_auc_score(y_test, y_proba_multi, multi_class='ovo')
```

Precision-Recall Curve (Better for Imbalanced)

```
1 from sklearn.metrics import precision_recall_curve, average_precision_score
2
3 # Get probabilities
4 y_proba = model.predict_proba(X_test)[: , 1]
5
6 # Average Precision (Area under PR curve)
7 ap = average_precision_score(y_test, y_proba)
8 print(f"Average Precision: {ap:.4f}")
9
10 # Precision-Recall curve
11 precision, recall, thresholds = precision_recall_curve(y_test, y_proba)
12
13 # Plot PR curve
14 plt.figure(figsize=(8, 6))
15 plt.plot(recall, precision, label=f'PR (AP = {ap:.4f})')
16 plt.axhline(y=y_test.mean(), color='r', linestyle='--', label='Baseline')
17 plt.xlabel('Recall')
18 plt.ylabel('Precision')
19 plt.title('Precision-Recall Curve')
20 plt.legend()
21 plt.show()
22
23 # Find best threshold (maximize F1)
24 f1_scores = 2 * (precision * recall) / (precision + recall + 1e-8)
25 best_idx = np.argmax(f1_scores[:-1]) # Last element has no threshold
26 best_threshold = thresholds[best_idx]
27 print(f"Best threshold for F1: {best_threshold:.4f}")
```

Threshold Tuning

```
1 from sklearn.metrics import precision_score, recall_score, f1_score
2
3 # Get probabilities
4 y_proba = model.predict_proba(X_test)[: , 1]
5
6 # Test different thresholds
7 thresholds = np.arange(0.1, 0.9, 0.05)
8 results = []
9
10 for thresh in thresholds:
11     y_pred_thresh = (y_proba >= thresh).astype(int)
12     results.append({
13         'threshold': thresh,
14         'precision': precision_score(y_test, y_pred_thresh),
15         'recall': recall_score(y_test, y_pred_thresh),
16         'f1': f1_score(y_test, y_pred_thresh)
17     })
18
19 results_df = pd.DataFrame(results)
20 print(results_df)
21
22 # Find optimal threshold
23 best_f1_idx = results_df['f1'].idxmax()
24 best_threshold = results_df.loc[best_f1_idx, 'threshold']
25 print(f"\nBest threshold (F1): {best_threshold}")
26
27 # Apply best threshold
28 y_pred_optimal = (y_proba >= best_threshold).astype(int)
29 print(f"F1 with optimal threshold: {f1_score(y_test, y_pred_optimal):.4f}")
```

Confusion Matrix Visualization

```
1 from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
2 import matplotlib.pyplot as plt
3
4 # Method 1: ConfusionMatrixDisplay
5 cm = confusion_matrix(y_test, y_pred)
6 disp = ConfusionMatrixDisplay(cm, display_labels=['Negative', 'Positive'])
7 disp.plot(cmap='Blues')
8 plt.title('Confusion Matrix')
9 plt.show()
10
11 # Method 2: Seaborn heatmap (more customizable)
12 plt.figure(figsize=(8, 6))
13 sns.heatmap(cm, annot=True, fmt='d', cmap='Blues',
14             xticklabels=['Pred Neg', 'Pred Pos'],
15             yticklabels=['Actual Neg', 'Actual Pos'])
16 plt.title('Confusion Matrix')
17 plt.ylabel('Actual')
18 plt.xlabel('Predicted')
19 plt.show()
20
21 # Normalized confusion matrix (percentages)
22 cm_normalized = confusion_matrix(y_test, y_pred, normalize='true')
23 sns.heatmap(cm_normalized, annot=True, fmt='.2%', cmap='Blues')
24 plt.title('Normalized Confusion Matrix')
25 plt.show()
```

Full Pipeline with Preprocessing

```
1 from sklearn.pipeline import Pipeline
2 from sklearn.compose import ColumnTransformer
3 from sklearn.preprocessing import StandardScaler, OneHotEncoder
4 from sklearn.impute import SimpleImputer
5 from sklearn.ensemble import RandomForestClassifier
6
7 # Define column types
8 numerical_cols = ['age', 'income', 'balance']
9 categorical_cols = ['gender', 'education']
10
11 # Preprocessing pipelines
12 num_pipeline = Pipeline([
13     ('imputer', SimpleImputer(strategy='median')),
14     ('scaler', StandardScaler())
15 ])
16
17 cat_pipeline = Pipeline([
18     ('imputer', SimpleImputer(strategy='constant', fill_value='Missing')),
19     ('encoder', OneHotEncoder(handle_unknown='ignore', sparse_output=False))
20 ])
21
22 # Column transformer
23 preprocessor = ColumnTransformer([
24     ('num', num_pipeline, numerical_cols),
25     ('cat', cat_pipeline, categorical_cols)
26 ])
27
28 # Full pipeline
29 full_pipeline = Pipeline([
30     ('preprocessor', preprocessor)
```

Pipeline with Resampling

```
1 from imblearn.pipeline import Pipeline as ImbPipeline
2 from imblearn.over_sampling import SMOTE
3
4 # Full pipeline with SMOTE
5 full_pipeline = ImbPipeline([
6     ('preprocessor', preprocessor),
7     ('smote', SMOTE(random_state=42)),
8     ('classifier', RandomForestClassifier(random_state=42))
9 ])
10
11 # Train-test split
12 X_train, X_test, y_train, y_test = train_test_split(
13     df[numerical_cols + categorical_cols], df['target'],
14     test_size=0.2, stratify=df['target'], random_state=42
15 )
16
17 # Fit pipeline
18 full_pipeline.fit(X_train, y_train)
19
20 # Evaluate
21 y_pred = full_pipeline.predict(X_test)
22 print(classification_report(y_test, y_pred))
23
24 # Cross-validation
25 scores = cross_val_score(full_pipeline, X_train, y_train, cv=5, scoring='f1')
26 print(f"CV F1: {scores.mean():.4f} (+/- {scores.std()*2:.4f})")
```

Hyperparameter Tuning with GridSearchCV

```
1 from sklearn.model_selection import GridSearchCV
2
3 # Parameter grid
4 param_grid = {
5     'preprocessor__num__imputer__strategy': ['mean', 'median'],
6     'smote__k_neighbors': [3, 5, 7],
7     'classifier__n_estimators': [50, 100, 200],
8     'classifier__max_depth': [5, 10, None],
9     'classifier__class_weight': [None, 'balanced']
10 }
11
12 # Grid search with stratified CV
13 grid_search = GridSearchCV(
14     full_pipeline,
15     param_grid,
16     cv=5,
17     scoring='f1',
18     n_jobs=-1,
19     verbose=1
20 )
21
22 grid_search.fit(X_train, y_train)
23
24 print(f"Best parameters: {grid_search.best_params_}")
25 print(f"Best CV F1: {grid_search.best_score_:.4f}")
26
27 # Evaluate best model on test set
28 y_pred = grid_search.predict(X_test)
29 print(f"Test F1: {f1_score(y_test, y_pred):.4f}")
```

Avoiding Data Leakage

```
1 # WRONG: Preprocessing before split (DATA LEAKAGE!)
2 # scaler.fit(X) # Sees test data!
3 # X_scaled = scaler.transform(X)
4 # X_train, X_test = train_test_split(X_scaled, ...)
5
6 # CORRECT: Split first, then preprocess
7 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
8
9 # Fit on training data only
10 scaler = StandardScaler()
11 X_train_scaled = scaler.fit_transform(X_train) # fit_transform on train
12 X_test_scaled = scaler.transform(X_test) # transform only on test
13
14 # BEST: Use Pipeline (handles this automatically)
15 pipeline = Pipeline([
16     ('scaler', StandardScaler()),
17     ('model', RandomForestClassifier())
18 ])
19 pipeline.fit(X_train, y_train) # Scaler sees only training data
20 predictions = pipeline.predict(X_test)
21
22 # Same for resampling: ONLY on training data
23 smote = SMOTE(random_state=42)
24 X_train_resampled, y_train_resampled = smote.fit_resample(X_train, y_train)
25 # Never apply SMOTE to test data!
```

Saving and Loading Complete Pipeline

```
1 import joblib
2
3 # Save the complete pipeline
4 joblib.dump(full_pipeline, 'model_pipeline.joblib')
5
6 # Save with compression
7 joblib.dump(full_pipeline, 'model_pipeline.joblib.gz', compress=3)
8
9 # Load pipeline
10 loaded_pipeline = joblib.load('model_pipeline.joblib')
11
12 # Make predictions with loaded pipeline
13 new_data = pd.DataFrame({
14     'age': [35, 45],
15     'income': [50000, 75000],
16     'balance': [10000, 25000],
17     'gender': ['M', 'F'],
18     'education': ['Bachelor', 'Master']
19 })
20
21 predictions = loaded_pipeline.predict(new_data)
22 probabilities = loaded_pipeline.predict_proba(new_data)
23
24 print(f"Predictions: {predictions}")
25 print(f"Probabilities: {probabilities}")
26
27 # Save grid search results
28 joblib.dump(grid_search.best_estimator_, 'best_model.joblib')
```

Data Preparation Checklist

Splitting:

- Train/test split done
- Stratification applied
- Random state documented
- Groups kept together (if applicable)

Cross-Validation:

- CV method chosen
- Stratified for classification
- Time-aware for time series
- Multiple metrics computed

Imbalanced Data:

- Class distribution checked
- Resampling applied (train only!)
- Appropriate metrics chosen
- Threshold tuned

Pipeline:

- No data leakage
- Pipeline tested
- Model saved
- Reproducible

Key Python Functions Reference

Task	Function/Class
Train-test split	<code>train_test_split(stratify=y)</code>
K-Fold CV	<code>KFold</code> , <code>StratifiedKFold</code>
Repeated CV	<code>RepeatedStratifiedKFold</code>
Time series CV	<code>TimeSeriesSplit</code>
CV scores	<code>cross_val_score</code> , <code>cross_validate</code>
CV predictions	<code>cross_val_predict</code>
Oversampling	<code>SMOTE</code> , <code>ADASYN</code> , <code>RandomOverSampler</code>
Undersampling	<code>RandomUnderSampler</code> , <code>TomekLinks</code>
Combined	<code>SMOTETomek</code> , <code>SMOTEENN</code>
Class weights	<code>class_weight='balanced'</code>
Metrics	<code>f1_score</code> , <code>roc_auc_score</code>
Pipeline	<code>Pipeline</code> , <code>imblearn.pipeline.Pipeline</code>

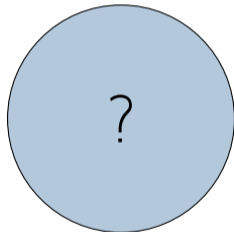
Congratulations!

You have completed all 4 parts of the
Machine Learning Pre-Processing Course









Questions?

Questions & Discussion



Contact: abdellaoui.e@gmail.com

References

-  Kohavi, R. (1995). *A Study of Cross-Validation and Bootstrap*. IJCAI.
-  Chawla, N.V. et al. (2002). *SMOTE: Synthetic Minority Over-sampling Technique*. JAIR.
-  He, H. & Garcia, E.A. (2009). *Learning from Imbalanced Data*. IEEE TKDE.
-  Lemaître, G. et al. (2017). *Imbalanced-learn: A Python Toolbox*. JMLR.
-  Bergstra, J. & Bengio, Y. (2012). *Random Search for Hyper-Parameter Optimization*. JMLR.
-  Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The Elements of Statistical Learning*. Springer.

Thank You!

Good luck with your Machine Learning projects!