# 📊 PYTHON PANDAS LIBRARY

Powerful, Open-Source Tool for Data Manipulation, Analysis, and Transformation in Data Science, Machine Learning, AI, and Business Analytics.

## ◆ CORE CONCEPTS
- SERIES & DATAFRAME
- INDEXING & REINDEXING
- DATA TYPES & CASTING

## ◆ DATA INPUT / OUTPUT
- READ & WRITE CSV, EXCEL, JSON, SQL, HTML
- HANDLING LARGE DATASETS EFFICIENTLY

## ◆ DATA CLEANING & PREPROCESSING
- HANDLING MISSING VALUES (NAN)
- REMOVING DUPLICATES
- DATA NORMALIZATION & FORMATTING
- STRING OPERATIONS

## ◆ DATA SELECTION & FILTERING
- loc[ ], iloc[ ]
- BOOLEAN INDEXING
- CONDITIONAL FILTERING

## ◆ DATA MANIPULATION
- SORTING & RANKING
- GROUPBY & AGGREGATION
- MERGE, JOIN & CONCATENATE
- PIVOT TABLES & CROSSTAB

## ◆ STATISTICAL & ANALYTICAL OPERATIONS
- DESCRIPTIVE STATISTICS
- CORRELATION & COVARIANCE
- APPLY, MAP & LAMBDA FUNCTIONS

## ◆ PERFORMANCE & OPTIMIZATION
- VECTORIZED OPERATIONS
- MEMORY OPTIMIZATION

## ◆ TIME SERIES & DATE HANDLING
- DATETIME INDEXING
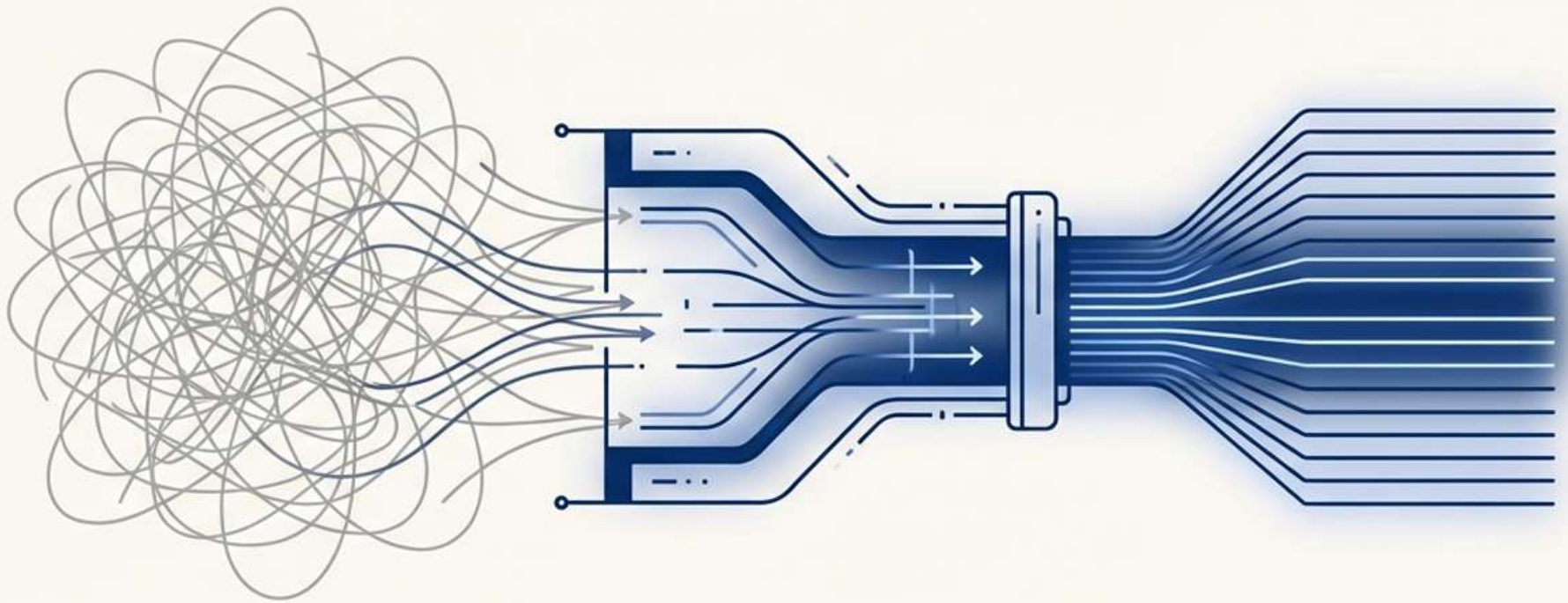- RESAMPLING & ROLLING OPERATIONS

## ◆ INTEGRATION
- WORKS SEAMLESSLY WITH NUMPY, MATPLOTLIB, SEABORN, SCIKIT-LEARN

💡 **USED FOR:** Data analysis, ETL pipelines, dashboards, financial analysis, real-time analytics, ML preprocessing.

# Pandas: The Engine of Data Transformation

A structured guide to manipulating, cleaning, and analyzing
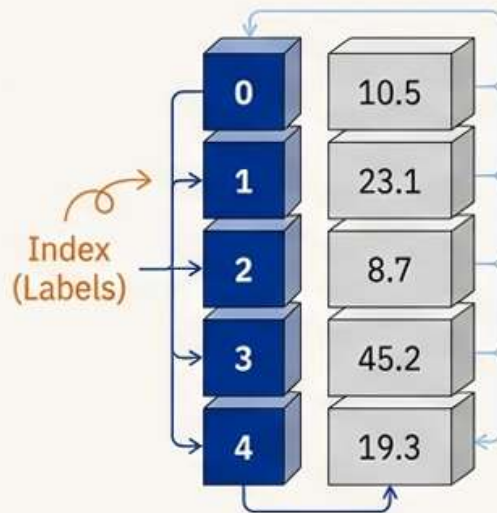data with Python's premier library

# The Foundational Structures: Series and DataFrame

At its core, Pandas organizes data into two primary structures, providing an intuitive and powerful way to manage information.
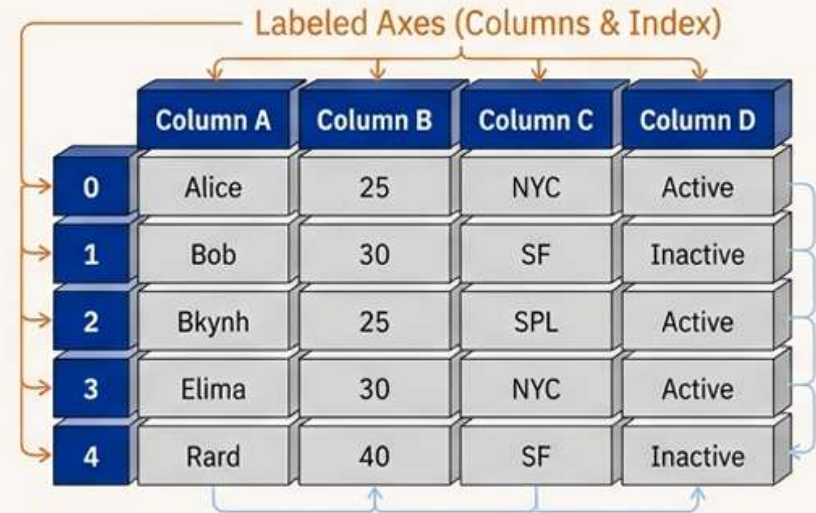
## Series: The 1D Building Block

A one-dimensional labeled array, similar to a column in a spreadsheet. Each element has an associated label, called an index.

| Index (Labels) | |
|:---:|:---:|
| 0 | 10.5 |
| 1 | 23.1 |
| 2 | 8.7 |
| 3 | 45.2 |
| 4 | 19.3 |

## DataFrame: The 2D Workhorse

A two-dimensional labeled data structure with columns of potentially different types. It is the primary Pandas object, analogous to a SQL table or an Excel spreadsheet.

Labeled Axes (Columns & Index)

| | Column A | Column B | Column C | Column D |
|:---:|:---:|:---:|:---:|:---:|
| 0 | Alice | 25 | NYC | Active |
| 1 | Bob | 30 | SF | Inactive |
| 2 | Bkynh | 25 | SPL | Active |
| 3 | Elima | 30 | NYC | Active |
| 4 | Rard | 40 | SF | Inactive |

Also covers essential concepts like **Indexing & Reindexing** (the address system for your data) and IBM Plex Sans.

Values sike **Data Types & Casting** (defining the nature of your data).

int float toan str

# The Data Workflow: Stage 1

## Ingesting the Raw Material: Data Input/Output

A project begins with data. Pandas provides a comprehensive and efficient toolkit for reading from and writing to a wide array of formats.
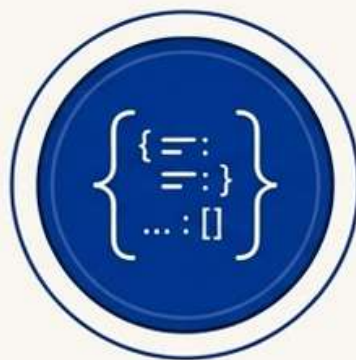
**CSV**

Reading and writing comma-separated values.

**Excel**

Interacting directly with Excel spreadsheets.

**JSON**

Handling JavaScript Object Notation files.

**SQL**

Querying databases and loading results into a DataFrame.

**HTML**

Parsing tables directly from web pages.

**Built for Scale:** Includes optimized readers and techniques for handling datasets that are too large to fit in memory.

# Precise Data Selection and Filtering

Isolate the exact data you need with a powerful and flexible indexing system. Pandas offers multiple methods for targeted data selection.

### `loc[]`: Label-based selection

Access a group of rows and columns by labels or a boolean array.

### `iloc[]`: Integer-position-based selection

Access data by its integer position in the structure.

### Boolean Indexing

The most powerful method. Filter your data based on the values within it, allowing for complex conditional logic.

|   | Column A | Column B | Column C | Column D |
|---|----------|----------|----------|----------|
| 0 | 15 | 42.5 | East | Active |
| 1 | 28 | 18.9 | West | Inactive |
| 2 | 56 | 31.2 | North | Active |
| 3 | 34 | 27.8 | South | Active |
| 4 | 19 | 50.1 | East | Inactive |

# The Essential Craft of Data Cleaning and Preprocessing

Raw data is often messy. Pandas provides a robust suite of tools to handle imperfections, ensuring the quality and integrity of your analysis.

**Before Table**

| | ID | Name | Age | Score | Status |
|---|---|---|---|---|---|
| **1** | 0 | 'john' | NaN | 85 | 'Active' |
| **2** | 1 | 'Jane' | 28 | 92 | 'Inactive' |
| **3** | 2 | 'john' | NaN | 85 | 'Active' |
| **4** | 3 | 'mike' | 32 | NaN | 'Active' |
| **5** | 4 | 'Sarah' | 45 | 78 | 'Inactive' |

**After Table**

| | ID | Name | Age | Score | Status |
|---|---|---|---|---|---|
| **1** | 0 | 'John' | 35.5 | 85 | 'Active' |
| **2** | 1 | 'Jane' | 28 | 92 | 'Inactive' |
| **3** | 3 | 'Mike' | 32 | 35.5 | 'Active' |
| **4** | 4 | 'Sarah' | 45 | 78 | 'Inactive' |

## Core Cleaning Tasks

1. **Handling Missing Values (NaN):** Find, remove, or fill missing data points using various strategies.

2. **Removing Duplicates:** Easily identify and drop duplicate records.

3. **Data Normalization & Formatting:** Standardize data formats for consistency (e.g., text cases, date formats).

4. **Vectorized String Operations:** Perform complex text cleaning and manipulation efficiently across entire columns using the `str` accessor.

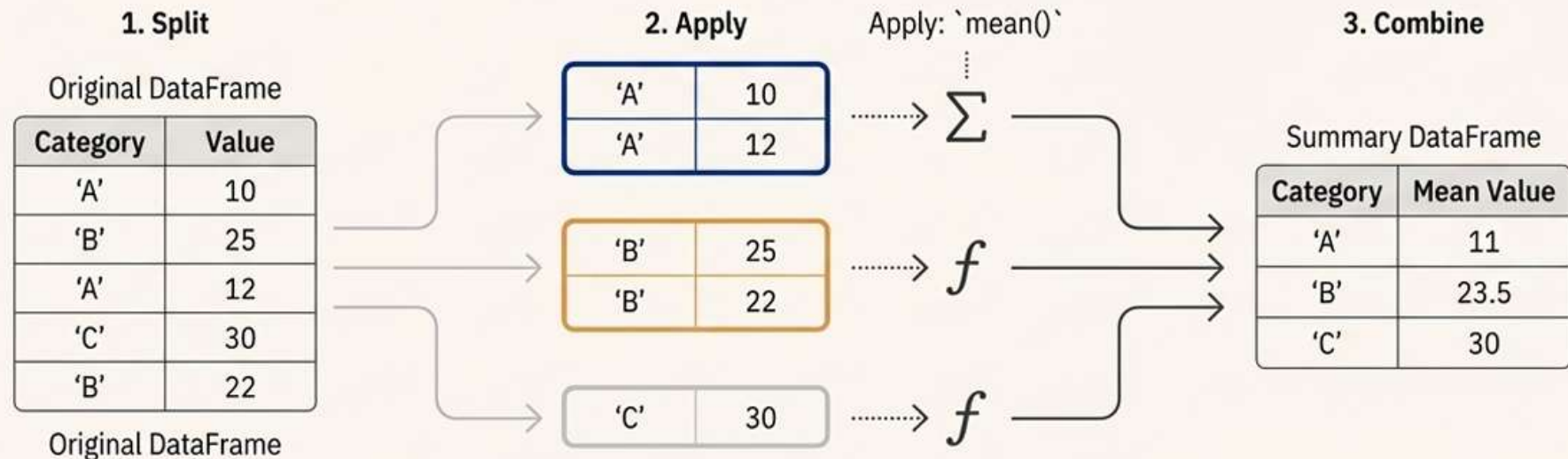# Shaping Data: Sorting, Grouping, and Aggregating

Once clean, data can be reshaped to answer specific questions. Start by imposing order and creating meaningful summaries.

## Sorting & Ranking

Arrange data based on criteria to quickly identify extremes and patterns. Functions for sorting by index or values, and for assigning ranks.

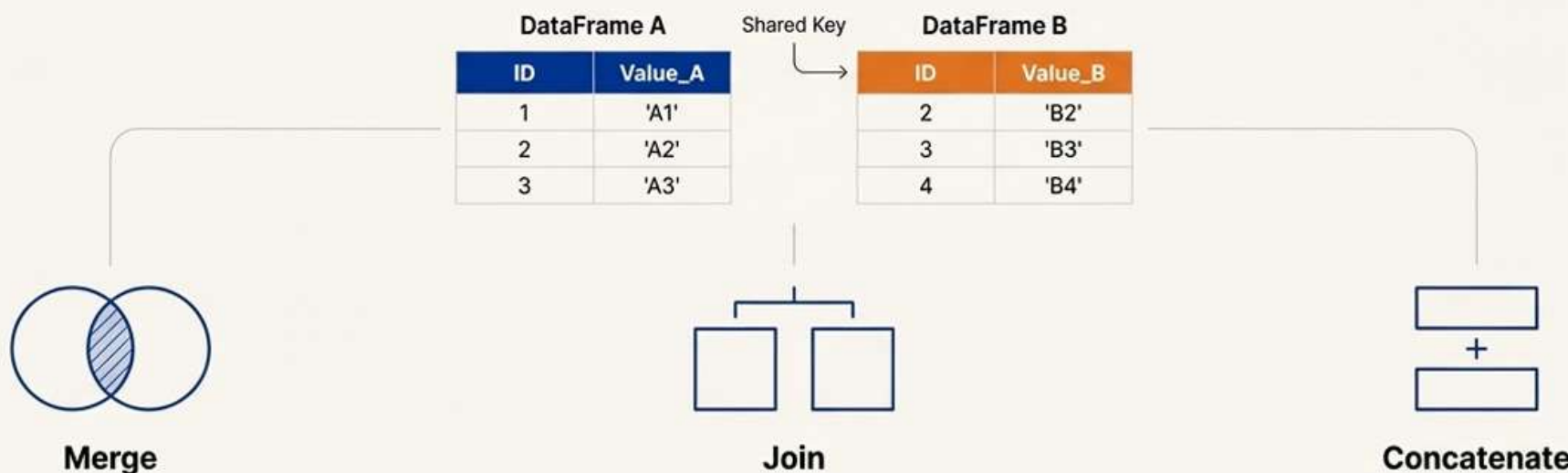## GroupBy & Aggregation: The Split-Apply-Combine Paradigm

The cornerstone of data analysis. This powerful paradigm allows you to:

**1. Split**

Original DataFrame

| Category | Value |
|----------|-------|
| 'A' | 10 |
| 'B' | 25 |
| 'A' | 12 |
| 'C' | 30 |
| 'B' | 22 |

Original DataFrame

**2. Apply**

| 'A' | 10 |
|-----|----|
| 'A' | 12 |

$\Sigma$

| 'B' | 25 |
|-----|----|
| 'B' | 22 |

$f$

| 'C' | 30 |
|-----|----|

$f$

Apply: `mean()`

**3. Combine**

Summary DataFrame

| Category | Mean Value |
|----------|------------|
| 'A' | 11 |
| 'B' | 23.5 |
| 'C' | 30 |

# Synthesizing Information: Merge, Join, and Concatenate

Your insights often live at the intersection of multiple datasets. Pandas provides a high-performance, in-memory toolkit for combining DataFrames.

**DataFrame A**

| ID | Value_A |
|----|---------|
| 1  | 'A1'    |
| 2  | 'A2'    |
| 3  | 'A3'    |

Shared Key →

**DataFrame B**

| ID | Value_B |
|----|---------|
| 2  | 'B2'    |
| 3  | 'B3'    |
| 4  | 'B4'    |

## Merge

For SQL-style, database-like joining operations on columns or indexes.

**Result (Merge)**

| ID | Value_A | Value_B |
|----|---------|---------|
| 2  | 'A2'    | 'B2'    |
| 3  | 'A3'    | 'B3'    |

## Join

A convenient method for combining the columns of two potentially differently-indexed DataFrames.

**Result (Join)**

| Index | Value_A | Value_B |
|-------|---------|---------|
| 0     | 'A1'    | 'B2'    |
| 1     | 'A2'    | 'B3'    |
| 2     | 'A3'    | 'B4'    |

## Concatenate

For stacking DataFrames on top of one another (along rows) or side-by-side (along columns).

**Result (Concatenate - Rows)**

| ID | Value_A | Value_B |
|----|---------|---------|
| 1  | 'A1'    | NaN     |
| 2  | 'A2'    | NaN     |
| 3  | 'A3'    | NaN     |
| 2  | NaN     | 'B2'    |
| 3  | NaN     | 'B3'    |
| 4  | NaN     | 'B4'    |

# Reshaping Layouts with Pivot Tables and Crosstabs

Sometimes, the key to insight is not in the data itself, but in its presentation. Reshape your dataset to summarize and explore relationships between categorical variables.

## Pivot Tables

A familiar concept from spreadsheets, this tool allows you to transform data from a 'long' format to a 'wide' format, summarizing data with an aggregation function.

## Crosstab

A specialized function that computes a frequency map of two (or more) factors.

### Long Format

| 'Date' | City | Sales |
|---|---|---|
| 2023-10-01 | London | 100 |
| 2023-10-01 | Tokyo | 200 |
| 2023-10-01 | NYC | 150 |
| 2023-10-02 | London | 120 |
| 2023-10-02 | Tokyo | 180 |
| 2023-10-02 | NYC | 160 |
| 2023-10-03 | London | 110 |
| 2023-10-03 | Tokyo | 210 |
| 2023-10-03 | NYC | 140 |

### Pivoted 'Wide' Format

| 'Date' | 'London' | 'Tokyo' | 'NYC' |
|---|---|---|---|
| 2023-10-01 | 100 | 200 | 150 |
| 2023-10-02 | 120 | 180 | 160 |
| 2023-10-03 | 110 | 210 | 140 |

# Mastering Time: Specialized Time Series and Date Handling

Time is a critical dimension in many datasets. Pandas has a rich set of tools specifically designed for time series analysis.
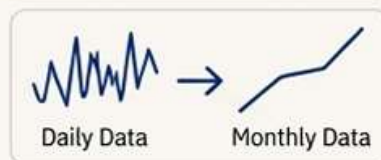
### Datetime Indexing

Use dates and times as the primary index, enabling powerful time-based slicing and selection.
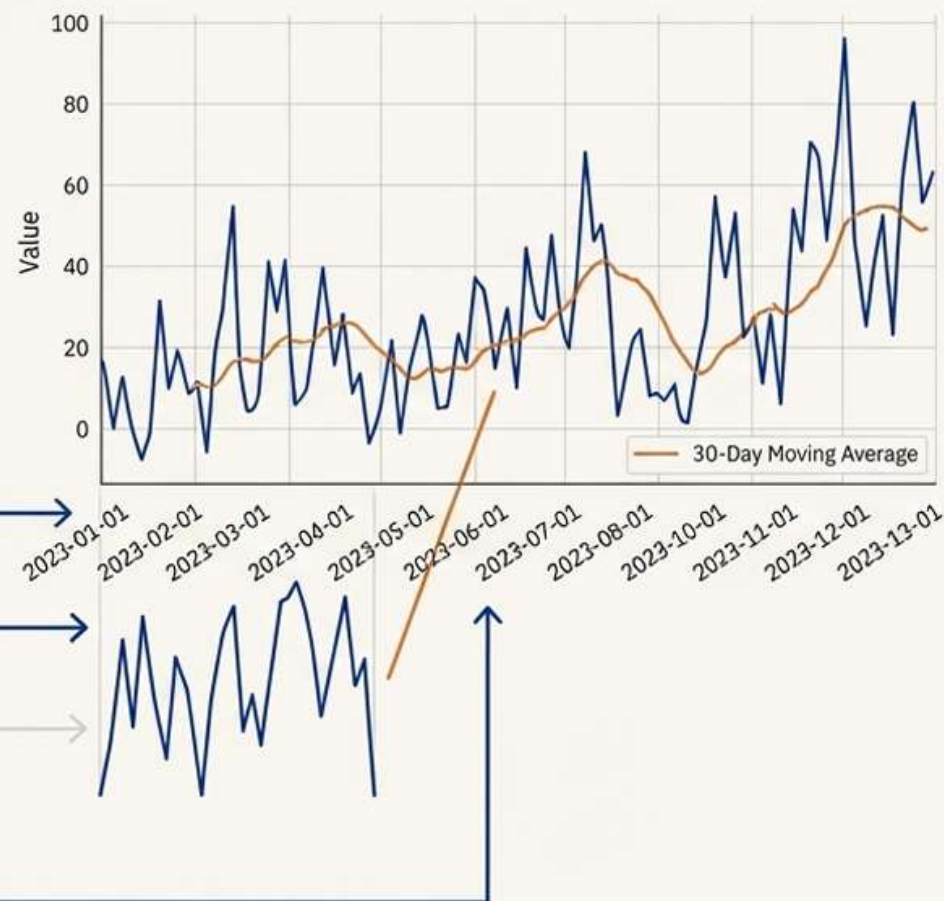
### Resampling

Easily convert time series from one frequency to another (e.g., from daily to monthly data) while applying an aggregation.

Daily Data → Monthly Data

### Rolling Operations

Calculate rolling window statistics, such as as moving averages, to smooth out data and identify trends.

# The Data Workflow: Stage 5

## Uncovering Insights: Statistical and Analytical Operations

With prepared data, you can now perform a wide range of calculations to extract meaning and test hypotheses.

### Descriptive Statistics

Generate summary statistics (count, mean, std, min, max, etc.) for each column with a single command (`.describe()`).

### Correlation & Covariance

Understand the relationships between numerical columns.

### Apply, Map & Lambda Functions

$f(x)$

For ultimate flexibility, apply any custom function to your data, from simple arithmetic to complex algorithms.

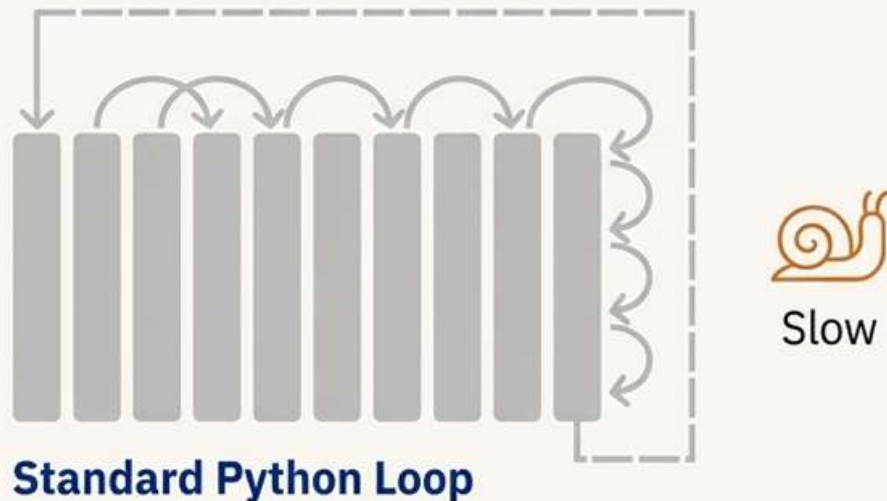| A | B |
|---|---|
| 10 | 15 |
| 20 | 25 |
| 30 | 35 |

```
# Calculate the range for each numeric column
column_range = df.apply(lambda x: x.max() - x.min())
```

# The Power Plant: Ensuring Performance and Optimization

Speed is a feature. Pandas is built on highly optimized C and Cython code, enabling rapid processing of large datasets.

## Vectorized Operations

The secret to Pandas' speed. Instead of looping through rows in Python, operations are applied to entire arrays at once at the C level, which is orders of magnitude faster.
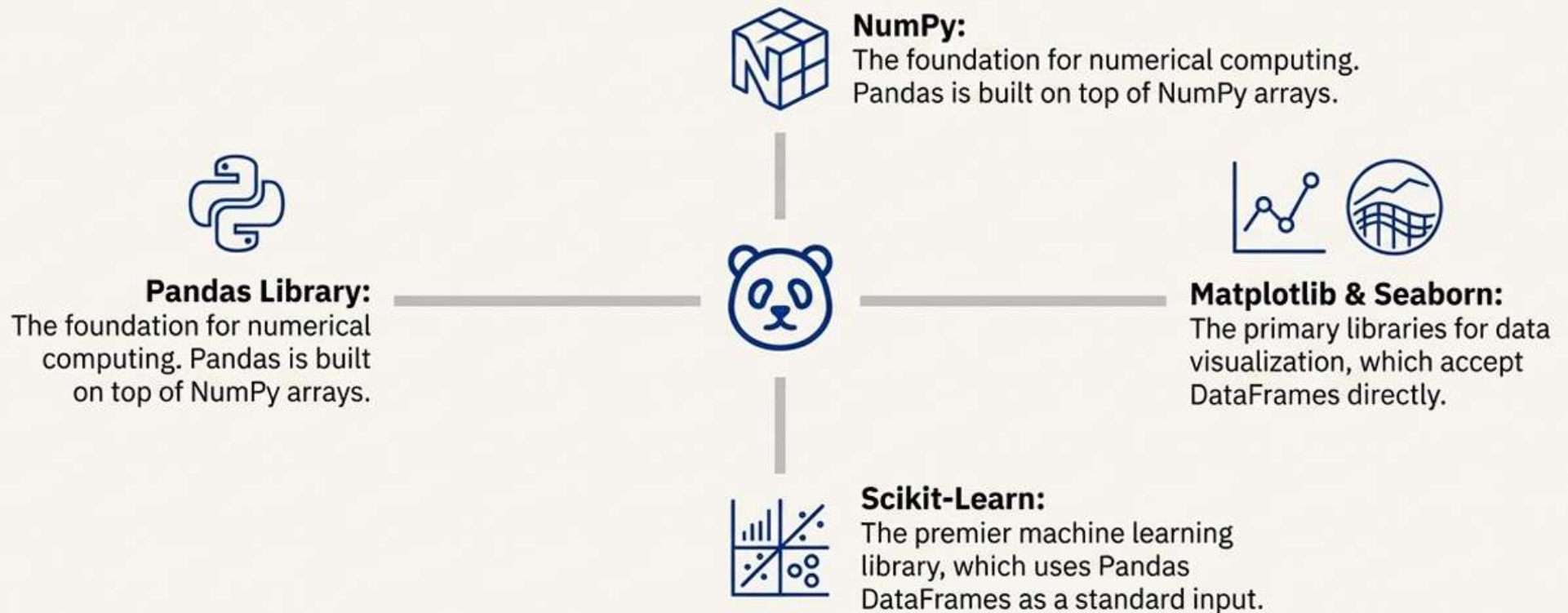
## Memory Optimization

Use appropriate data types and techniques to reduce the memory footprint of your DataFrames, allowing you to work with larger datasets on a given machine.

**Pandas Vectorized Operation**

Slow

Fast

**Standard Python Loop**

**Pandas Vectorized Operation**

# The Central Hub of the Python Data Science Ecosystem

Pandas does not stand alone. It integrates seamlessly with and is the foundational data structure for the most important libraries in data science.

**NumPy:**
The foundation for numerical computing. Pandas is built on top of NumPy arrays.

**Pandas Library:**
The foundation for numerical computing. Pandas is built on top of NumPy arrays.

**Matplotlib & Seaborn:**
The primary libraries for data visualization, which accept DataFrames directly.

**Scikit-Learn:**
The premier machine learning library, which uses Pandas DataFrames as a standard input.

# Real-World Impact: Where Pandas Drives Value

From finance to machine learning, Pandas is the tool of choice for a wide range of data-driven tasks.

**Exploratory Data Analysis**
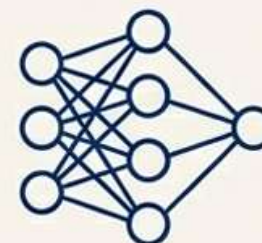
**ETL Pipelines (Extract, Transform, Load)**

**Data Feeds for Dashboards**

**Quantitative Financial Analysis**

**Real-time Analytics**

**Machine Learning Preprocessing**

# The Pandas Philosophy: From Chaos to Clarity

Pandas provides the essential structure that makes insight possible. It is the definitive tool for taking raw, messy, and complex data and transforming it into a clean, orderly, and analyzable format.