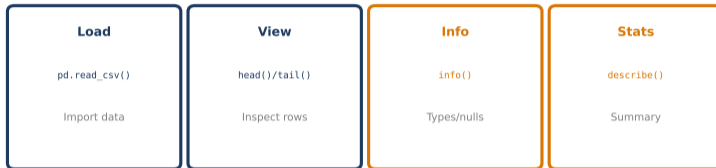


Lesson 05 Summary: DataFrames Introduction

Data Science with Python – Key Concepts

Data Science Program

pandas DataFrame Essentials



DataFrame = Rows (Index) + Columns + Values

Access:

`df.shape` | `df.columns` | `df.index` | `df.values`

Series = 1 column | *DataFrame = multiple Series*

pandas is THE library for data manipulation in Python

DataFrame Structure

Columns (features)

Index	Date	AAPL	MSFT	GOOGL
0	2024-01-02	185.2	376.1	140.9
1	2024-01-03	184.8	374.2	139.5
2	2024-01-04	186.1	378.5	141.2

Rows (observations)

2D labeled data structure with rows and columns

DataFrames are 2D labeled data structures

Two core pandas data structures:

Series (1D):

- Single column with index
- Like a labeled array
- Example: `df['Price']` returns a Series

DataFrame (2D):

- Multiple columns sharing an index
- Collection of Series
- Example: `df[['Price', 'Volume']]` returns a DataFrame

Series = 1 column; DataFrame = multiple Series

Loading CSV Data



Common Parameters:

```
filepath: "data/prices.csv"  
index_col: "Date"  
parse_dates: True  
usecols: ["AAPL", "MSFT"]
```

`pd.read_csv()` handles most CSV formats automatically

Essential inspection methods:

- **df.head(n)**: First n rows (default 5)
- **df.tail(n)**: Last n rows (default 5)
- **df.sample(n)**: Random n rows
- **df.shape**: (rows, columns) tuple

Always inspect after loading:

```
df = pd.read_csv('data.csv')
print(df.head()) # First 5 rows
print(df.shape) # (252, 5)
```

Always inspect data after loading!

DataFrame Info: df.info()

```
<class pandas.DataFrame>
RangeIndex: 252 entries, 0 to 251
Data columns (5 columns):
  Date      252 non-null datetime64
  AAPL     252 non-null float64
  MSFT     252 non-null float64
  GOOGL    250 non-null float64  (2 missing)
memory usage: 10.0 KB
```

info() reveals data types and missing values

df.describe() provides statistical summary:

- **count:** Non-null values
- **mean:** Average value
- **std:** Standard deviation
- **min/max:** Range
- **25%/50%/75%:** Quartiles

Finance insight: Compare std (volatility) across assets

describe() provides statistical summary of numeric columns

Index and Columns



`df.shape: (252, 4) | df.dtypes: column data types`

Index labels rows; columns label data fields

Essential DataFrame Operations:

Operation	Syntax
Import pandas	<code>import pandas as pd</code>
Load CSV	<code>df = pd.read_csv('file.csv')</code>
View first rows	<code>df.head(n)</code>
View last rows	<code>df.tail(n)</code>
Data types	<code>df.info()</code>
Statistics	<code>df.describe()</code>
Shape	<code>df.shape</code>
Column names	<code>df.columns</code>
Index	<code>df.index</code>

Master these for efficient data exploration