

CS109A: Introduction to Data Science

Lecture 03: Introduction to Pandas

Harvard University

Fall 2024

- **Course:** CS109A: Introduction to Data Science
- **Lecture:** Lecture 03: Introduction to Pandas
- **Instructor:** Chris Gumb
- **Objective:** Learn the fundamentals of Pandas for data manipulation and analysis in Python

Key Summary

This lecture introduces **Pandas**, the essential Python library for data manipulation and analysis. You'll learn about the two core data structures: `Series` (1-dimensional) and `DataFrame` (2-dimensional tables). We cover how to create these structures, load data from CSV files, inspect and clean data, select subsets using boolean indexing and the `loc/iloc` accessors, and perform aggregations with `groupby`. By the end of this lecture, you'll have the foundational skills needed to work with tabular data in Python.

Contents

1	Introduction: Why Pandas?	3
1.1	What is Pandas?	3
1.2	Why Use Pandas Instead of Plain Python?	3
1.3	Importing Pandas	3
2	The Series Data Structure	4
2.1	What is a Series?	4
2.2	Creating a Series	4
2.3	Series Attributes	4
2.4	Custom Indexes	6
3	The DataFrame Data Structure	7
3.1	What is a DataFrame?	7
3.2	Creating a DataFrame	7
3.2.1	Method 1: From a List of Dictionaries (Row by Row)	7

3.2.2	Method 2: From a Dictionary of Lists (Column by Column)	7
3.3	DataFrame Attributes	9
4	Loading Data from Files	10
4.1	Reading CSV Files	10
4.2	Inspecting Your Data	10
5	Selecting Data	12
5.1	Selecting Columns	12
5.2	The loc and iloc Accessors	12
6	Boolean Indexing (Filtering)	14
6.1	Creating Boolean Masks	14
6.2	Applying the Mask	14
6.3	Combining Conditions	14
7	Common DataFrame Operations	16
7.1	Sorting	16
7.2	Value Counts	16
7.3	Group By and Aggregation	16
7.4	Handling Missing Data	18
7.5	Renaming Columns	18
8	Method Chaining	19
8.1	What is Method Chaining?	19
9	Practical Example: Survey Data	20
10	Advanced: Exploding Lists	21
11	Key Takeaways	22

1 Introduction: Why Pandas?

1.1 What is Pandas?

Definition:

Pandas **Pandas** is a Python library that provides high-performance, easy-to-use data structures and data analysis tools. The name comes from “Panel Data” (a term from econometrics) and “Python Data Analysis.”

Pandas gives you access to data types that don’t exist in standard Python:

- **Series:** A 1-dimensional labeled array
- **DataFrame:** A 2-dimensional labeled table (like an Excel spreadsheet)

1.2 Why Use Pandas Instead of Plain Python?

You could store data in Python lists and dictionaries, but Pandas offers significant advantages:

- **Speed:** Pandas is built on NumPy, which stores data in contiguous memory blocks for fast operations
- **Convenience:** Methods like `.head()`, `.describe()`, `.groupby()` make common operations trivial
- **Data Alignment:** Pandas automatically aligns data by labels during operations
- **Missing Data Handling:** Built-in support for missing values (`NaN`)
- **File I/O:** Easy reading/writing of CSV, Excel, SQL databases, JSON, and more

1.3 Importing Pandas

The standard convention is to import pandas with the alias `pd`:

```
1 import pandas as pd
2 import numpy as np # Often used alongside pandas
```

Warning

Don’t break conventions! Always use `pd` for pandas and `np` for numpy. Using other aliases will confuse anyone reading your code.

2 The Series Data Structure

2.1 What is a Series?

Definition:

Series A **Series** is a one-dimensional array with labels (called an **index**). Think of it as:

- A single column from a spreadsheet
- A dictionary where keys are the index and values are the data
- A NumPy array with attached labels

2.2 Creating a Series

The simplest way is from a Python list:

```

1 # Create a simple Series
2 my_list = [10, 20, 30, 40, 50]
3 s = pd.Series(my_list)
4 print(s)
```

Listing 1: Creating a Series from a list

Output:

```

0    10
1    20
2    30
3    40
4    50
dtype: int64
```

Notice the output shows:

- **Left column:** The **index** (0, 1, 2, 3, 4)—created automatically
- **Right column:** The **values** from your list
- **Bottom:** The **dtype** (data type)—`int64` means 64-bit integers

2.3 Series Attributes

Every Series has four key attributes:

```

1 s = pd.Series([10, 20, 30], index=['a', 'b', 'c'], name='my_series')
2
3 print(s.values)    # The actual data (as a NumPy array)
4 # Output: array([10, 20, 30])
5
6 print(s.index)    # The labels
7 # Output: Index(['a', 'b', 'c'], dtype='object')
```

```
8  
9 print(s.name)      # The series name (becomes column name in DataFrame)  
10 # Output: 'my_series'  
11  
12 print(s.dtype)     # The data type  
13 # Output: dtype('int64')
```

Listing 2: Series attributes

2.4 Custom Indexes

The index doesn't have to be integers—you can use any labels:

```
1 s = pd.Series([100, 200, 300])
2 s.index = ['apple', 'banana', 'cherry']
3 print(s)
4
5 # Output:
6 # apple      100
7 # banana     200
8 # cherry     300
9 # dtype: int64
10
11 # Access by label
12 print(s['banana'])  # Output: 200
```

Listing 3: Custom string index

Warning

Index Labels Can Be Non-Unique!

Unlike dictionary keys, Series indices can have duplicates:

```
1 s.index = ['a', 'a', 'b']
2 print(s['a'])
3 # Returns BOTH values with index 'a':
4 # a      100
5 # a      200
```

This can cause unexpected behavior if you're not careful!

3 The DataFrame Data Structure

3.1 What is a DataFrame?

Definition:

DataFrame A **DataFrame** is a 2-dimensional labeled data structure—like an Excel spreadsheet or SQL table.

- Each **column** is a Series
- All columns share the same **index** (row labels)
- Different columns can have different data types

3.2 Creating a DataFrame

3.2.1 Method 1: From a List of Dictionaries (Row by Row)

Each dictionary represents one row; keys become column names:

```

1 data = [
2     {'fruit': 'apple', 'color': 'red', 'price': 1.50},
3     {'fruit': 'banana', 'color': 'yellow', 'price': 0.75},
4     {'fruit': 'cherry', 'color': 'red', 'price': 3.00}
5 ]
6
7 df = pd.DataFrame(data)
8 print(df)

```

Listing 4: DataFrame from list of dictionaries

Output:

	fruit	color	price
0	apple	red	1.50
1	banana	yellow	0.75
2	cherry	red	3.00

3.2.2 Method 2: From a Dictionary of Lists (Column by Column)

Each key is a column name; the list becomes that column's values:

```

1 data = {
2     'fruit': ['apple', 'banana', 'cherry'],
3     'color': ['red', 'yellow', 'red'],
4     'price': [1.50, 0.75, 3.00]
5 }
6
7 df = pd.DataFrame(data)
8 # Same result as above

```

Listing 5: DataFrame from dictionary of lists

Key Information

Which Method to Use?

- **List of dictionaries:** Useful when scraping data (you build up rows one at a time)
- **Dictionary of lists:** Useful when you already have column-wise data

3.3 DataFrame Attributes

```
1 print(df.shape)      # (rows, columns)
2 # Output: (3, 3)
3
4 print(df.columns)    # Column names
5 # Output: Index(['fruit', 'color', 'price'], dtype='object')
6
7 print(df.index)      # Row labels
8 # Output: RangeIndex(start=0, stop=3, step=1)
9
10 print(df.dtypes)    # Data type of each column
11 # Output:
12 # fruit      object
13 # color      object
14 # price     float64
15 # dtype: object
```

Listing 6: DataFrame attributes

4 Loading Data from Files

4.1 Reading CSV Files

The most common way to get data into pandas:

```

1 # Read a CSV file into a DataFrame
2 df = pd.read_csv('data/survey_results.csv')
3
4 # Peek at the first few rows
5 df.head()

```

Listing 7: Reading a CSV file

Key Information

Other File Formats

Pandas can read many formats:

- `pd.read_excel('file.xlsx')` — Excel files
- `pd.read_json('file.json')` — JSON files
- `pd.read_html('url')` — Tables from web pages
- `pd.read_sql(query, connection)` — SQL databases
- `pd.read_clipboard()` — From your clipboard!

4.2 Inspecting Your Data

After loading data, always inspect it first:

```

1 # First 5 rows (default) or specify n
2 df.head()
3 df.head(2)
4
5 # Last 5 rows
6 df.tail()
7
8 # Shape: (rows, columns)
9 df.shape
10
11 # Column names
12 df.columns
13
14 # Data types and memory usage
15 df.info()
16
17 # Summary statistics for numeric columns
18 df.describe()
19
20 # Data types only
21 df.dtypes

```

Listing 8: Inspecting a DataFrame

5 Selecting Data

5.1 Selecting Columns

```

1 # Single column (returns a Series)
2 df['fruit']

3
4 # Single column using dot notation (shortcut)
5 df.fruit # Only works if column name has no spaces

6
7 # Multiple columns (returns a DataFrame)
8 df[['fruit', 'price']]

```

Listing 9: Selecting columns

Warning

Be Careful with Dot Notation

`df.fruit` is convenient but dangerous if:

- Column name has spaces (`df.my column` doesn't work)
- Column name matches a DataFrame method (`df.head` returns the method, not a column!)

The bracket notation `df['fruit']` is always safe.

5.2 The loc and iloc Accessors

Pandas provides two ways to select rows:

Definition:

loc vs iloc

- **loc**: Selection by **label** (what you see in the index)
- **iloc**: Selection by **integer position** (0, 1, 2, ...)

```

1 # Create a DataFrame with a non-default index
2 df = pd.DataFrame({
3     'name': ['Alice', 'Bob', 'Charlie'],
4     'age': [25, 30, 35]
5 }, index=['a', 'b', 'c'])

6
7 # loc: by label
8 df.loc['b']           # Row with label 'b'
9 df.loc['a':'b']       # Rows 'a' through 'b' (inclusive!)
10 df.loc['a', 'age']   # Specific cell: row 'a', column 'age'

11
12 # iloc: by position
13 df.iloc[1]            # Second row (position 1)
14 df.iloc[0:2]          # First two rows (exclusive end!)
15 df.iloc[0, 1]         # Specific cell: row 0, column 1

```

Listing 10: Using loc and iloc

Important:

The Index Trap After filtering or sorting, row indices get scrambled:

```
1 df_sorted = df.sort_values('age', ascending=False)
2 print(df_sorted)
3 #      name  age
4 # c    Charlie  35
5 # b      Bob   30
6 # a    Alice   25
7
8 # Trying to get "first row" by position
9 df_sorted.iloc[0]  # Gets Charlie (correct)
10
11 # Trying by label 0
12 df_sorted.loc[0]  # ERROR! No label '0' exists
```

Solution: Use `.reset_index(drop=True)` after sorting/filtering.

6 Boolean Indexing (Filtering)

6.1 Creating Boolean Masks

A **boolean mask** is a Series of True/False values that you use to filter rows:

```

1 df = pd.DataFrame({
2     'name': ['Alice', 'Bob', 'Charlie', 'Diana'],
3     'age': [25, 30, 35, 28],
4     'city': ['NYC', 'LA', 'NYC', 'LA']
5 })
6
7 # Create a mask: which rows have age > 27?
8 mask = df['age'] > 27
9 print(mask)
10 # 0      False
11 # 1      True
12 # 2      True
13 # 3      True
14 # dtype: bool

```

Listing 11: Creating boolean masks

6.2 Applying the Mask

Use the mask inside brackets to filter:

```

1 # Get rows where age > 27
2 df[mask]
3 # or directly:
4 df[df['age'] > 27]
5
6 #      name  age  city
7 # 1      Bob   30   LA
8 # 2  Charlie   35  NYC
9 # 3   Diana   28   LA

```

Listing 12: Filtering with boolean masks

6.3 Combining Conditions

```

1 # Use & for AND, | for OR, ~ for NOT
2 # IMPORTANT: Each condition must be in parentheses!
3
4 # Age > 27 AND city is NYC
5 df[(df['age'] > 27) & (df['city'] == 'NYC')]
6
7 # Age > 30 OR city is LA
8 df[(df['age'] > 30) | (df['city'] == 'LA')]
9

```

```
10 # NOT in NYC  
11 df[~(df['city'] == 'NYC')]
```

Listing 13: Combining conditions

Warning

Parentheses are Required!

Due to Python operator precedence, you must wrap each condition in parentheses:

```
1 # WRONG:  
2 df[df['age'] > 27 & df['city'] == 'NYC'] # Error!  
3  
4 # CORRECT:  
5 df[(df['age'] > 27) & (df['city'] == 'NYC')]
```

7 Common DataFrame Operations

7.1 Sorting

```

1 # Sort by one column
2 df.sort_values('age')

3

4 # Sort descending
5 df.sort_values('age', ascending=False)

6

7 # Sort by multiple columns
8 df.sort_values(['city', 'age'])

9

10 # Reset index after sorting (usually what you want)
11 df.sort_values('age').reset_index(drop=True)

```

Listing 14: Sorting DataFrames

7.2 Value Counts

Count occurrences of unique values:

```

1 # How many people in each city?
2 df['city'].value_counts()
3 # NYC      2
4 # LA       2

5

6 # Get unique values
7 df['city'].unique()
8 # array(['NYC', 'LA'], dtype=object)

9

10 # Number of unique values
11 df['city'].nunique()
12 # 2

```

Listing 15: Counting values

7.3 Group By and Aggregation

Definition:

GroupBy **GroupBy** splits data into groups based on a column's values, then applies an aggregation function to each group.

```

1 # Average age by city
2 df.groupby('city')['age'].mean()
3 # city
4 # LA      29.0
5 # NYC    30.0

```

```
6  
7 # Multiple aggregations  
8 df.groupby('city').agg({  
9     'age': ['mean', 'min', 'max'],  
10    'name': 'count'  
11 })  
12  
13 # Count rows per group  
14 df.groupby('city').size()
```

Listing 16: GroupBy operations

7.4 Handling Missing Data

```

1 # Check for missing values
2 df.isna()                      # Boolean mask of NaN locations
3 df.isna().sum()                  # Count NaNs per column
4
5 # Drop rows with any NaN
6 df.dropna()
7
8 # Fill NaN with a value
9 df.fillna(0)
10 df['column'].fillna('Unknown')
11
12 # Fill with column mean
13 df['age'].fillna(df['age'].mean())

```

Listing 17: Working with missing values

7.5 Renaming Columns

```

1 # Rename specific columns
2 df.rename(columns={'old_name': 'new_name'})
3
4 # Rename all columns at once
5 df.columns = ['col1', 'col2', 'col3', 'col4']
6
7 # Make all column names lowercase
8 df.columns = df.columns.str.lower()

```

Listing 18: Renaming columns

8 Method Chaining

8.1 What is Method Chaining?

Instead of saving intermediate results to variables, you can chain multiple operations together:

```
1 # Without chaining (verbose)
2 df_filtered = df[df['age'] > 25]
3 df_sorted = df_filtered.sort_values('age')
4 df_top = df_sorted.head(3)
5
6 # With chaining (concise)
7 result = (df[df['age'] > 25]
8         .sort_values('age')
9         .head(3))
```

Listing 19: Method chaining example

Key Information

Benefits of Method Chaining

- More concise code
- No intermediate variables cluttering your namespace
- Reads like a pipeline: “Take df, filter it, sort it, take top 3”

Tip: Wrapping in parentheses lets you break across multiple lines for readability.

9 Practical Example: Survey Data

Let's walk through a realistic example using class survey data:

```

1 # 1. Load the data
2 df = pd.read_csv('data/survey_raw.csv')
3
4 # 2. Inspect
5 print(df.shape)
6 print(df.columns)
7 df.head()
8
9 # 3. Clean column names (remove spaces, lowercase)
10 df.columns = ['timestamp', 'program', 'jupyter_exp',
11               'python_exp', 'pandas_skill', 'os',
12               'dark_mode', 'languages', 'dob']
13
14 # 4. Check data types
15 df.dtypes
16 df.info()
17
18 # 5. Convert dark_mode to boolean
19 df['dark_mode'] = df['dark_mode'] == 'Yes'
20
21 # 6. Filter to students with pandas skill > 3
22 skilled = df[df['pandas_skill'] > 3]
23
24 # 7. Count programs
25 df['program'].value_counts()
26
27 # 8. Average pandas skill by program
28 (df.groupby('program')['pandas_skill']
29   .mean()
30   .sort_values(ascending=False))
31
32 # 9. Save cleaned data
33 df.to_csv('data/survey_clean.csv', index=False)

```

Listing 20: Complete example workflow

10 Advanced: Exploding Lists

Sometimes a column contains comma-separated values that you want to split:

```
1 # Sample data
2 df = pd.DataFrame({
3     'name': ['Alice', 'Bob'],
4     'languages': ['English, Spanish', 'English, Mandarin, French']
5 })
6
7 # Step 1: Split the string into a list
8 df['lang_list'] = df['languages'].str.split(', ')
9
10 # Step 2: Explode - each list item becomes its own row
11 df_exploded = df.explode('lang_list')
12
13 #      name          languages    lang_list
14 # 0  Alice        English, Spanish    English
15 # 0  Alice        English, Spanish   Spanish
16 # 1    Bob  English, Mandarin, French    English
17 # 1    Bob  English, Mandarin, French   Mandarin
18 # 1    Bob  English, Mandarin, French    French
19
20 # Now you can count languages, filter, etc.
21 df_exploded['lang_list'].value_counts()
```

Listing 21: Exploding comma-separated values

11 Key Takeaways

Key Summary

Summary of Lecture 03: Pandas Fundamentals

Core Data Structures

- **Series:** 1D labeled array (values + index + name + dtype)
- **DataFrame:** 2D table of Series (columns share an index)

Creating DataFrames

- `pd.DataFrame(list_of_dicts)` — Row by row
- `pd.DataFrame(dict_of_lists)` — Column by column
- `pd.read_csv('file.csv')` — From file

Inspecting Data

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

Selecting Data

- Columns: `df['col']` or `df[['col1', 'col2']]`
- Rows by label: `df.loc['label']`
- Rows by position: `df.iloc[0]`
- Boolean filtering: `df[df['col'] > value]`

Common Operations

- Sort: `df.sort_values('col')`
- Count: `df['col'].value_counts()`
- Group: `df.groupby('col').agg('other': 'mean')`
- Reset: `df.reset_index(drop=True)`

Key Warnings

- `loc` uses labels; `iloc` uses positions
- After sorting/filtering, indices get scrambled—use `reset_index()`
- Use parentheses around each condition when combining with `&/|`
- Mixed types become `object` dtype (slow!)

