

# Week 7: Comprehensive Lecture Notes on Pandas (Part 2)

## 1 Introduction

Pandas is a powerful library for data manipulation and analysis, widely used in data science. In this session, we will explore advanced techniques that can enhance efficiency and performance when working with large datasets.

## 2 Performance Optimization in Pandas

### 2.1 Using Efficient Data Types

- Use `astype()` to optimize memory usage. - Convert categorical data using `pd.Categorical()`. - Convert date-time columns to `datetime64` for efficient operations.

```
1 import pandas as pd
2 import numpy as np
3
4 df = pd.DataFrame({
5     'category': ['A', 'B', 'A', 'C'],
6     'value': [10, 20, 15, 30]
7 })
8
9 df['category'] = df['category'].astype('category') # Optimized memory usage
10 print(df.dtypes)
```

### 2.2 Vectorized Operations vs. Apply

- Use vectorized operations instead of `apply()` or loops for better performance.

```
1 df['new_value'] = df['value'] * 2 # Vectorized operation
```

- Avoid `apply()` for row-wise operations unless necessary.

```
1 df['new_value'] = df.apply(lambda row: row['value'] * 2, axis=1) # Slower
```

### 2.3 Using `eval()` and `query()`

- `eval()` and `query()` offer faster alternatives to standard pandas operations.

```
1 df.eval('new_value = value * 2', inplace=True) # Faster than apply
2 df_filtered = df.query('value > 15') # Efficient filtering
```

## 3 Advanced Data Manipulation

### 3.1 MultiIndex for Hierarchical Data

- MultiIndex allows better organization of hierarchical data.

```

1 arrays = [
2     ['A', 'A', 'B', 'B'],
3     [1, 2, 1, 2]
4 ]
5
6 index = pd.MultiIndex.from_arrays(arrays, names=('Group', 'Num'))
7 df = pd.DataFrame({'Values': [10, 20, 30, 40]}, index=index)
8 print(df)

```

- Accessing specific data using `.loc[]`:

```

1 print(df.loc['A', 1])

```

## 3.2 Pivot Tables

A pivot table is a powerful tool in pandas that allows you to summarize and reorganize data in a tabular format. It is used to aggregate data in a way that makes it easier to analyze, similar to Excel pivot tables.

- Pivot tables summarize data efficiently.

```

1 df = pd.DataFrame({
2     'Category': ['A', 'A', 'B', 'B'],
3     'Type': ['X', 'Y', 'X', 'Y'],
4     'Value': [10, 20, 15, 25]
5 })
6
7 pivot_df = df.pivot_table(values='Value', index='Category', columns='Type', aggfunc='sum')
8 print(pivot_df)

```

## 3.3 Window Functions (Rolling & Expanding)

- Rolling windows compute statistics over moving time windows.

```

1 df['rolling_avg'] = df['Value'].rolling(window=2).mean()
2 print(df)

```

- Expanding applies functions over an expanding window.

```

1 df['expanding_sum'] = df['Value'].expanding().sum()
2 print(df)

```

# 4 Handling Large Datasets

## 4.1 Chunk Processing

- Use `chunksize` when reading large files.

```

1 chunks = pd.read_csv('large_file.csv', chunksize=10000)
2 for chunk in chunks:
3     print(chunk.head())

```

## 4.2 Dask for Parallel Processing

- Dask provides parallel computing for pandas operations.

```

1 import dask.dataframe as dd
2 ddf = dd.read_csv('large_file.csv')
3 print(ddf.head())

```

## 5 Load Datasets in Pandas

Pandas is a powerful Python library for **data analysis and manipulation**. In this tutorial, we will explore different ways to **load datasets** into Pandas from various file formats such as **CSV, Excel, JSON, SQL, Parquet, and more**.

By the end of this tutorial, you will be able to:

- Read data from different file formats into Pandas.
- Work with large datasets efficiently.
- Load data from online sources and databases.

### 5.1 Load a CSV File

CSV (Comma-Separated Values) is the most common file format for datasets.

```
1 import pandas as pd
2
3 # Load a CSV file
4 df = pd.read_csv("dataset.csv")
5
6 # Display the first 5 rows
7 print(df.head())
```

#### 5.1.1 Additional Options for `pd.read_csv()`

- `delimiter=";"` → Use if the file is separated by semicolons.
- `header=None` → Use if the file does not contain column names.
- `index_col=0` → Use if the first column should be used as an index.

### 5.2 Load an Excel File

Pandas also supports reading data from Excel files (`.xlsx`, `.xls`).

```
1 df = pd.read_excel("dataset.xlsx", sheet_name="Sheet1") # Specify sheet name if needed
2 print(df.head())
```

#### 5.2.1 Install Required Library

Excel reading requires `openpyxl` or `xlrd`. If not installed, run:

```
1 pip install openpyxl
```

### 5.3 Load a JSON File

JSON (JavaScript Object Notation) is commonly used for **web APIs and structured data**.

```
1 df = pd.read_json("dataset.json")
2 print(df.head())
```

#### 5.3.1 Read JSON from an API

```
1 import requests
2
3 # Fetch JSON from an API
4 url = "https://api.example.com/data"
5 response = requests.get(url)
6
7 # Convert JSON to DataFrame
8 df = pd.DataFrame(response.json())
9 print(df.head())
```



## 5.4 Load Data from a SQL Database

If your dataset is stored in a **SQL database**, you can directly import it using Pandas.

```
1 import sqlite3
2
3 # Connect to a database
4 conn = sqlite3.connect("database.db")
5
6 # Run a SQL query
7 df = pd.read_sql_query("SELECT * FROM my_table", conn)
8 print(df.head())
```

### 5.4.1 Read from MySQL or PostgreSQL

Install SQLAlchemy if using MySQL/PostgreSQL:

```
1 pip install sqlalchemy
```

## 5.5 Load a Parquet File (Efficient for Large Data)

Parquet is a high-performance file format often used in **big data applications**.

```
1 df = pd.read_parquet("dataset.parquet")
2 print(df.head())
```

### 5.5.1 Install Parquet Support

If you get an error, install pyarrow:

```
1 pip install pyarrow
```

## 5.6 Load a Dataset from a URL

If your dataset is hosted online, you can load it directly using a URL.

```
1 url = "https://example.com/dataset.csv"
2 df = pd.read_csv(url)
3 print(df.head())
```

**Works for:** CSV, JSON, and other file formats.

## 5.7 Load Data from a ZIP File

If your dataset is compressed inside a ZIP file, Pandas can read it directly.

```
1 df = pd.read_csv("dataset.zip", compression="zip")
2 print(df.head())
```

## 5.8 Load a Dataset from Google Sheets

Google Sheets can be **converted into a CSV link** and loaded into Pandas.

```
1 sheet_url = "https://docs.google.com/spreadsheets/d/example/export?format=csv"
2 df = pd.read_csv(sheet_url)
3 print(df.head())
```

# Exercises: Pandas Problems and Solutions

## Problem 1: Optimize Memory Usage with Categorical Data

You have the following dataset:

```
1 import pandas as pd
2
3 df = pd.DataFrame({
4     'Category': ['Apple', 'Banana', 'Apple', 'Orange', 'Banana', 'Apple'],
5     'Value': [10, 20, 30, 40, 50, 60]
6 })
7
8 print(df.dtypes)
```

**Task:** Convert the "Category" column to an optimized data type.

**Solution:**

```
1 df['Category'] = df['Category'].astype('category')
2 print(df.dtypes)
```

## Problem 2: Use Vectorized Operations Instead of apply()

Given the dataset:

```
1 df = pd.DataFrame({'Value': [10, 20, 30, 40, 50]})
```

**Task:** Double the values in the "Value" column using a vectorized operation.

**Solution:**

```
1 df['Double_Value'] = df['Value'] * 2
2 print(df)
```

## Problem 3: Use eval() for Efficient Computation

Compute a new column named "Triple\_Value" using eval().

**Solution:**

```
1 df.eval('Triple_Value = Value * 3', inplace=True)
2 print(df)
```

## Problem 4: Filter Data Using query()

Given:

```
1 df = pd.DataFrame({
2     'Name': ['Alice', 'Bob', 'Charlie', 'David'],
3     'Age': [25, 30, 35, 40]
4 })
```

**Task:** Filter out rows where Age is greater than 30.

**Solution:**

```
1 filtered_df = df.query('Age > 30')
2 print(filtered_df)
```



## Problem 5: Use MultiIndex for Hierarchical Data

Given:

```
1 arrays = [  
2     ['USA', 'USA', 'Canada', 'Canada'],  
3     ['East', 'West', 'East', 'West']  
4 ]  
5  
6 index = pd.MultiIndex.from_arrays(arrays, names=('Country', 'Region'))  
7  
8 df = pd.DataFrame({'Sales': [100, 150, 200, 250]}, index=index)
```

**Task:** Retrieve sales for Canada in the East region.

**Solution:**

```
1 print(df.loc['Canada', 'East'])
```

## Problem 6: Create a Pivot Table

Given:

```
1 df = pd.DataFrame({  
2     'Department': ['HR', 'HR', 'IT', 'IT'],  
3     'Gender': ['Male', 'Female', 'Male', 'Female'],  
4     'Salary': [50000, 52000, 70000, 68000]  
5 })
```

**Task:** Create a pivot table showing the average salary per department and gender.

**Solution:**

```
1 pivot_df = df.pivot_table(values='Salary', index='Department', columns='Gender', aggfunc='mean')  
2 print(pivot_df)
```

## Problem 7: Apply Rolling Window Mean

Given:

```
1 df = pd.DataFrame({'Values': [10, 20, 30, 40, 50]})
```

**Task:** Calculate the rolling average with a window size of 2.

**Solution:**

```
1 df['Rolling_Mean'] = df['Values'].rolling(window=2).mean()  
2 print(df)
```

## Problem 8: Process Large CSV Files in Chunks

Suppose you have a large CSV file called `large_file.csv`.

**Task:** Read and process the file in chunks of 5000 rows.

**Solution:**

```
1 chunks = pd.read_csv('large_file.csv', chunksize=5000)  
2  
3 for chunk in chunks:  
4     print(chunk.head()) # Process each chunk separately
```

## Problem 9: Merge Two DataFrames

Given:

```
1 df1 = pd.DataFrame({'ID': [1, 2, 3], 'Value': [100, 200, 300]})
2 df2 = pd.DataFrame({'ID': [2, 3, 4], 'Score': [90, 80, 70]})
```

**Task:** Merge the DataFrames on "ID" using an inner join.

**Solution:**

```
1 merged_df = pd.merge(df1, df2, on='ID', how='inner')
2 print(merged_df)
```

## Problem 10: Resample Time Series Data

Given:

```
1 df = pd.DataFrame({
2     'Date': pd.date_range(start='2024-01-01', periods=10, freq='D'),
3     'Sales': [100, 120, 130, 90, 150, 160, 200, 180, 210, 190]
4 })
5 df.set_index('Date', inplace=True)
```

**Task:** Resample the data to compute the mean sales per week.

**Solution:**

```
1 df_weekly = df.resample('W').mean()
2 print(df_weekly)
```

# Problems to be Solved

## Problem 1: Memory Optimization

You are given the following DataFrame:

```
1 import pandas as pd
2
3 df = pd.DataFrame({
4     'Product': ['Laptop', 'Tablet', 'Laptop', 'Smartphone'],
5     'Price': [1200, 600, 1100, 900],
6     'Stock': [30, 50, 20, 80]
7 })
```

Optimize this DataFrame for memory usage without losing any information.

**Solution:**

```
1 df['Product'] = df['Product'].astype('category')
2 df['Price'] = df['Price'].astype('int16')
3 df['Stock'] = df['Stock'].astype('int16')
```

## Problem 2: MultiIndexing

Create a MultiIndex DataFrame where products belong to different categories and cities. Then, retrieve all the products available in "New York."

**Solution:**

```
1 index = pd.MultiIndex.from_tuples([
2     ('Electronics', 'New York', 'Laptop'),
3     ('Electronics', 'New York', 'Smartphone'),
4     ('Electronics', 'Chicago', 'Tablet')
5 ], names=['Category', 'City', 'Product'])
6
7 df = pd.DataFrame({'Stock': [20, 30, 40]}, index=index)
8 df.xs('New York', level='City')
```

## Problem 3: Merging and Joining

You have two DataFrames:

```
1 df1 = pd.DataFrame({'ID': [1, 2, 3], 'Name': ['Alice', 'Bob', 'Charlie']})
2 df2 = pd.DataFrame({'ID': [2, 3, 4], 'Salary': [50000, 60000, 70000]})
```

Perform an outer join on the "ID" column.

**Solution:**

```
1 df_merged = pd.merge(df1, df2, on='ID', how='outer')
```

## Problem 4: Pivot Tables

Create a pivot table from the following DataFrame to show total sales per product per region.

```
1 df = pd.DataFrame({
2     'Region': ['North', 'South', 'North', 'South'],
3     'Product': ['A', 'B', 'A', 'B'],
4     'Sales': [100, 200, 150, 250]
5 })
```

**Solution:**

```
1 df.pivot_table(values='Sales', index='Region', columns='Product', aggfunc='sum')
```





## Problem 5: Custom Aggregations

Calculate the range (max - min) of sales for each region.

```
1 df = pd.DataFrame({  
2     'Region': ['East', 'East', 'West', 'West'],  
3     'Sales': [100, 150, 200, 250]  
4 })
```

**Solution:**

```
1 df.groupby('Region')['Sales'].agg(lambda x: x.max() - x.min())
```

## Problem 6: Loading and Working with Datasets

- Try loading a dataset using different formats.
- Explore **data cleaning and transformation** with Pandas.
- Practice with **real-world datasets** (Kaggle, UCI, etc.).



