# Data Wrangling

## ◆ Reminder

Yesterday we cleaned the Titanic dataset (handled missing values, duplicates, datatypes).
👉 Today we'll **transform** it so it becomes usable for insights and analysis.

## ◆ Notes on Data Wrangling

- **Definition** → Data wrangling is the process of **transforming raw data into a structured, usable format**.

- **Why important?**

    - Raw datasets are messy.

    - Analysis/ML requires tidy data (rows = observations, columns = features).

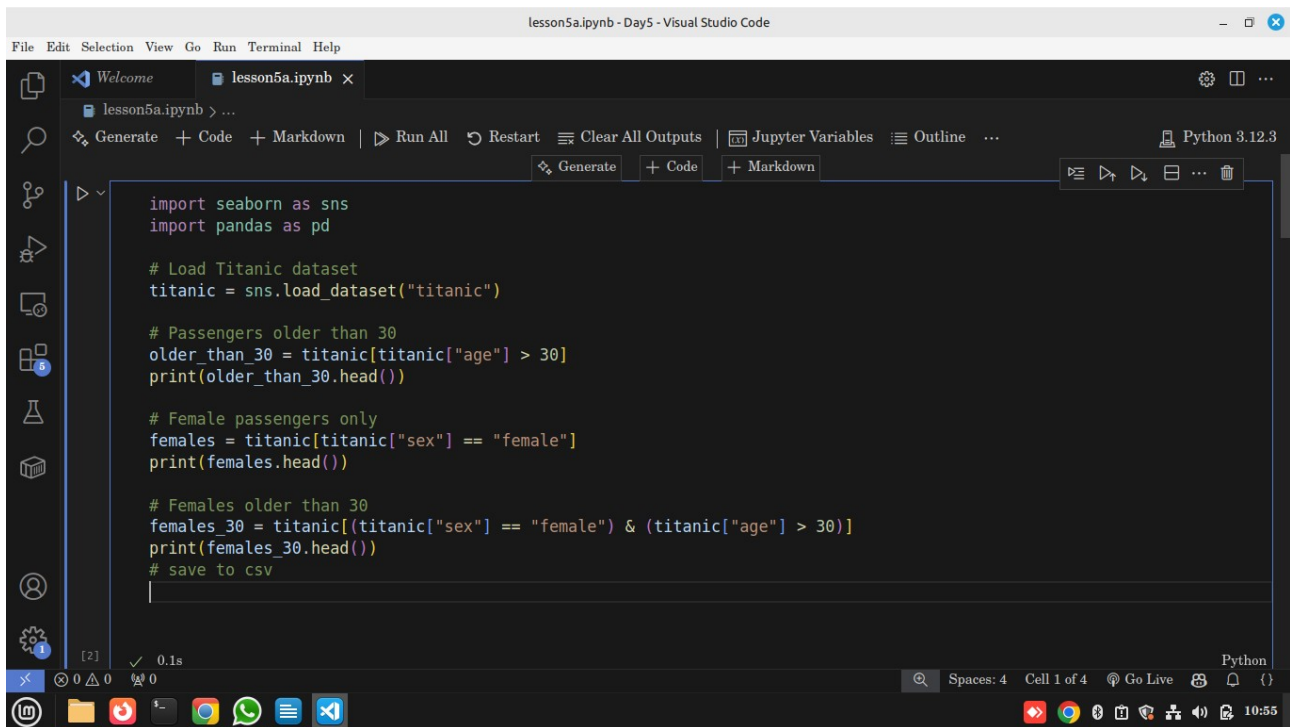    - Enables pattern discovery and insights.

Common Wrangling Operations:

1. Filtering → Selecting specific rows.

2. Grouping → Summarizing by categories.

3. Merging → Combining datasets.

4. Pivoting → Reshaping data.

📌 **Use in Data Science**: Wrangling prepares datasets for visualization, dashboards, and machine learning models.

## ◆ Filtering

# Practical



```python
import seaborn as sns
import pandas as pd

# Load Titanic dataset
titanic = sns.load_dataset("titanic")

# Passengers older than 30
older_than_30 = titanic[titanic["age"] > 30]
print(older_than_30.head())

# Female passengers only
females = titanic[titanic["sex"] == "female"]
print(females.head())

# Females older than 30
females_30 = titanic[(titanic["sex"] == "female") & (titanic["age"] > 30)]
print(females_30.head())
# save to csv
```
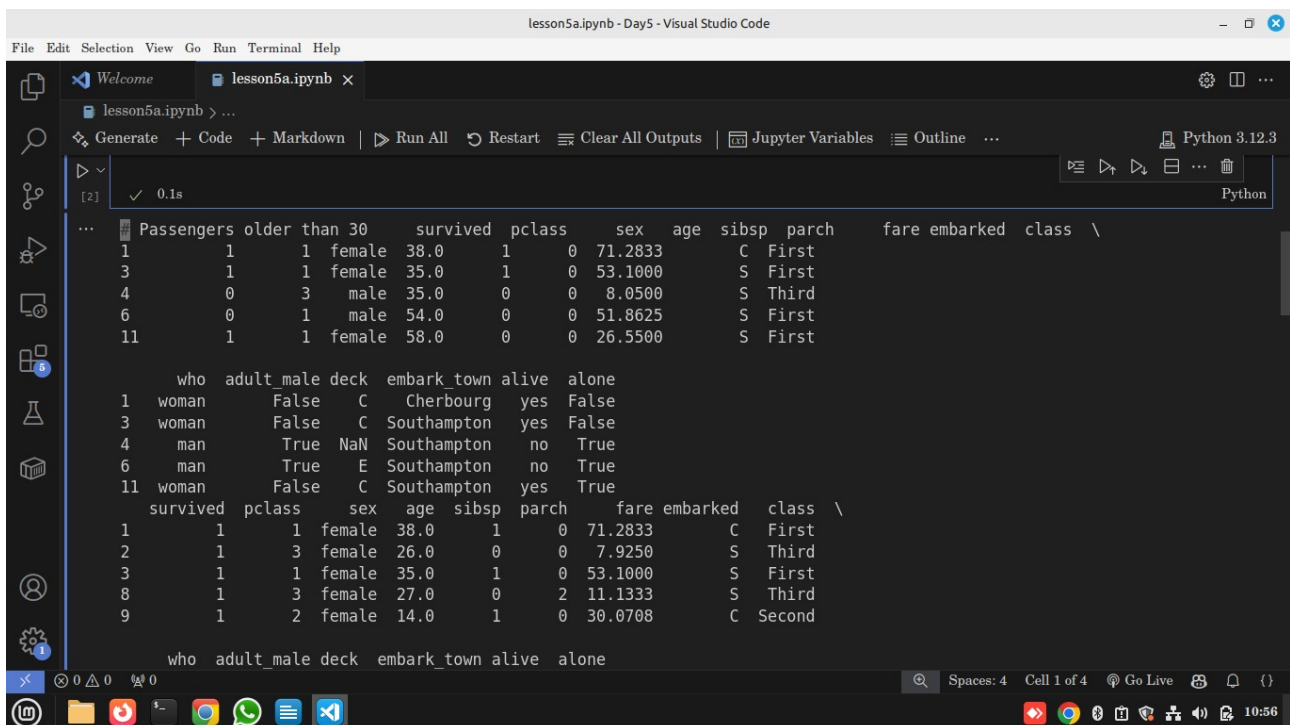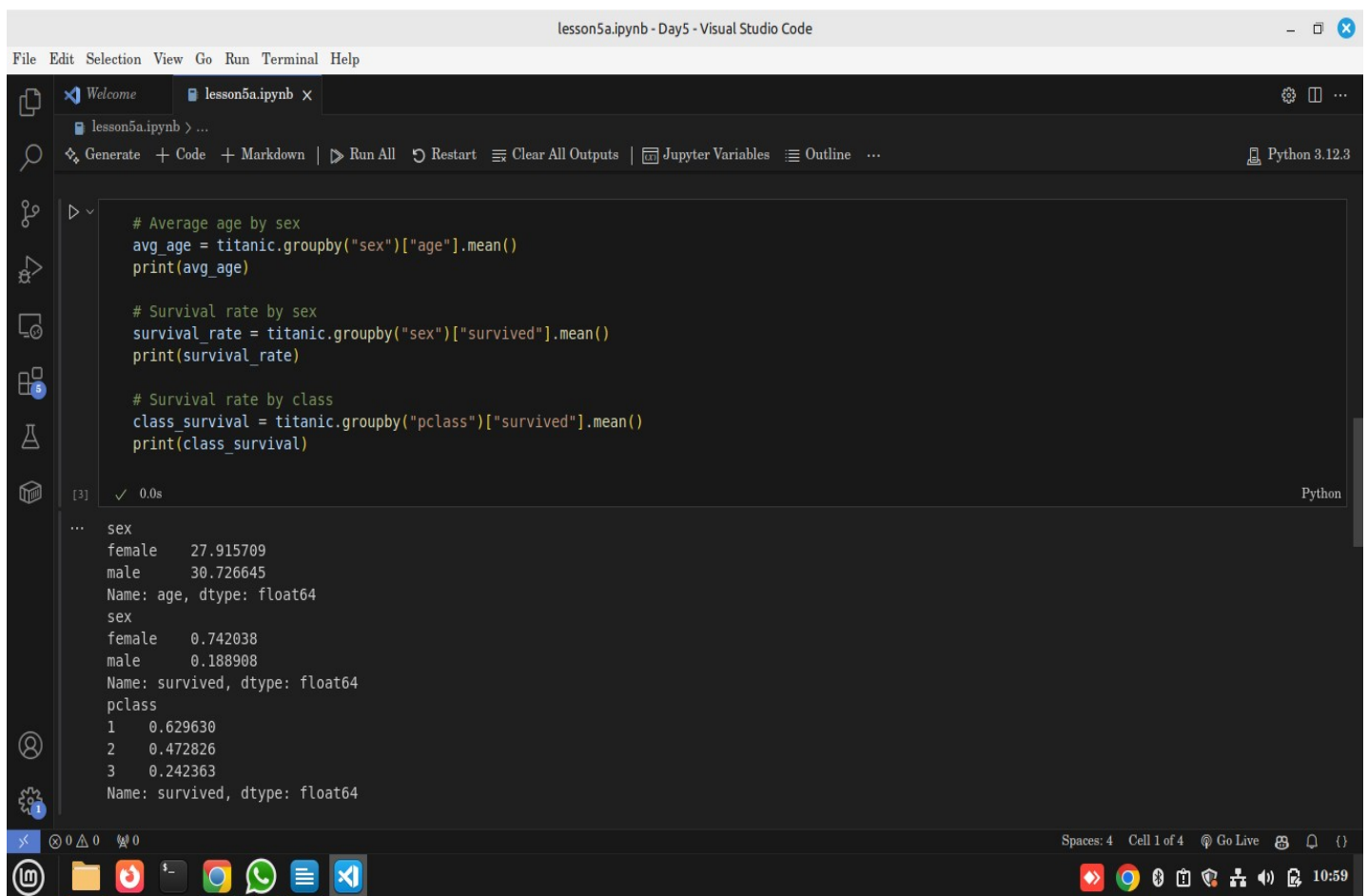
## Sample output

**Explanation**

- `df[condition]` → selects rows matching condition.

- `&` = AND, `|` = OR.

- Used for subsetting datasets.

**Use in Data Science** → Filtering helps analysts zoom in on specific groups (e.g., "Find customers over 40 in Nairobi with purchases > $500").

---

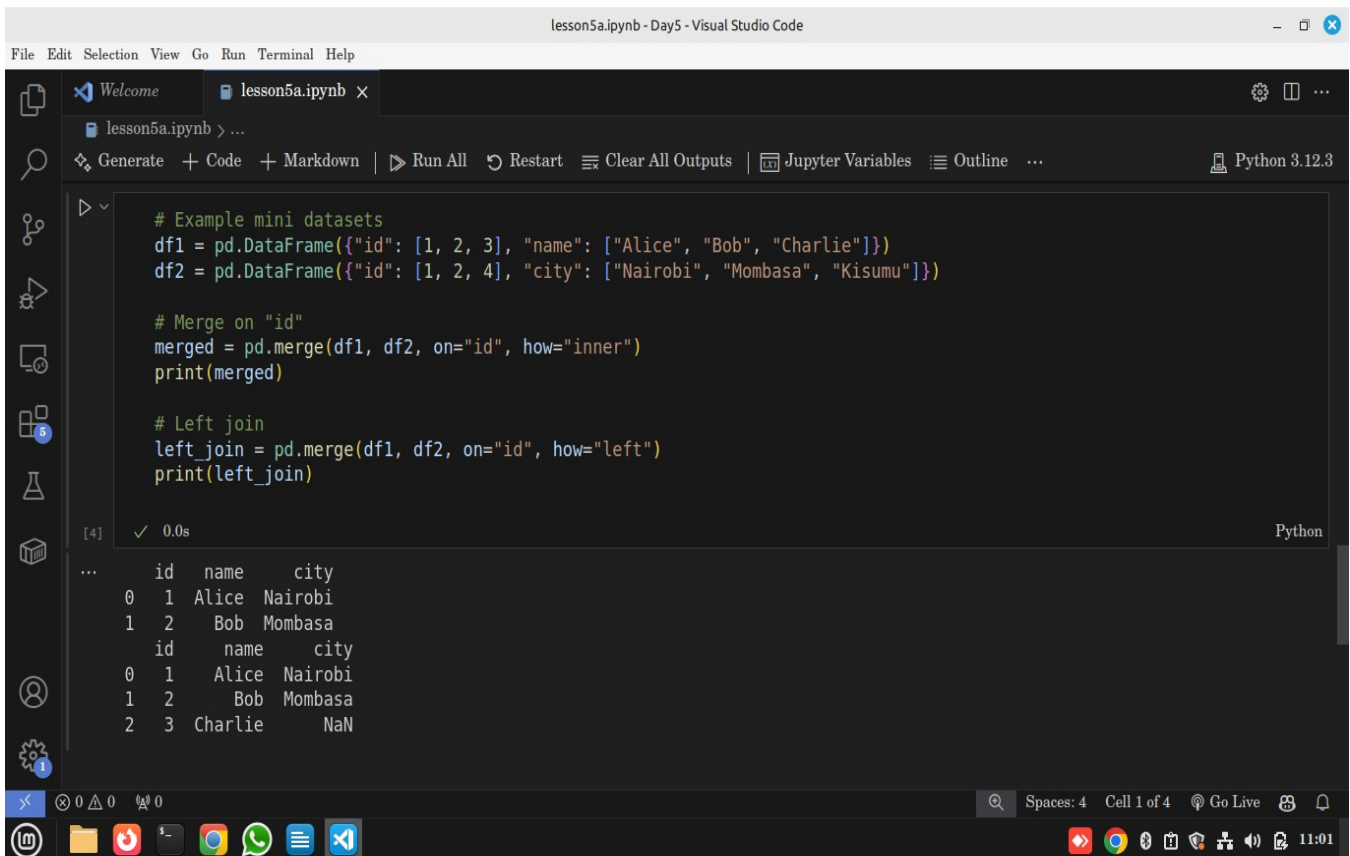## ◆ **GroupBy & Aggregations**

**Practical**



**Explanation**

- `groupby()` → splits data into groups.

- `.mean()`, `.sum()`, `.count()` → applies aggregation.

**Use in Data Science** → GroupBy helps compare segments (e.g., "average salary by department", "churn rate by subscription plan").

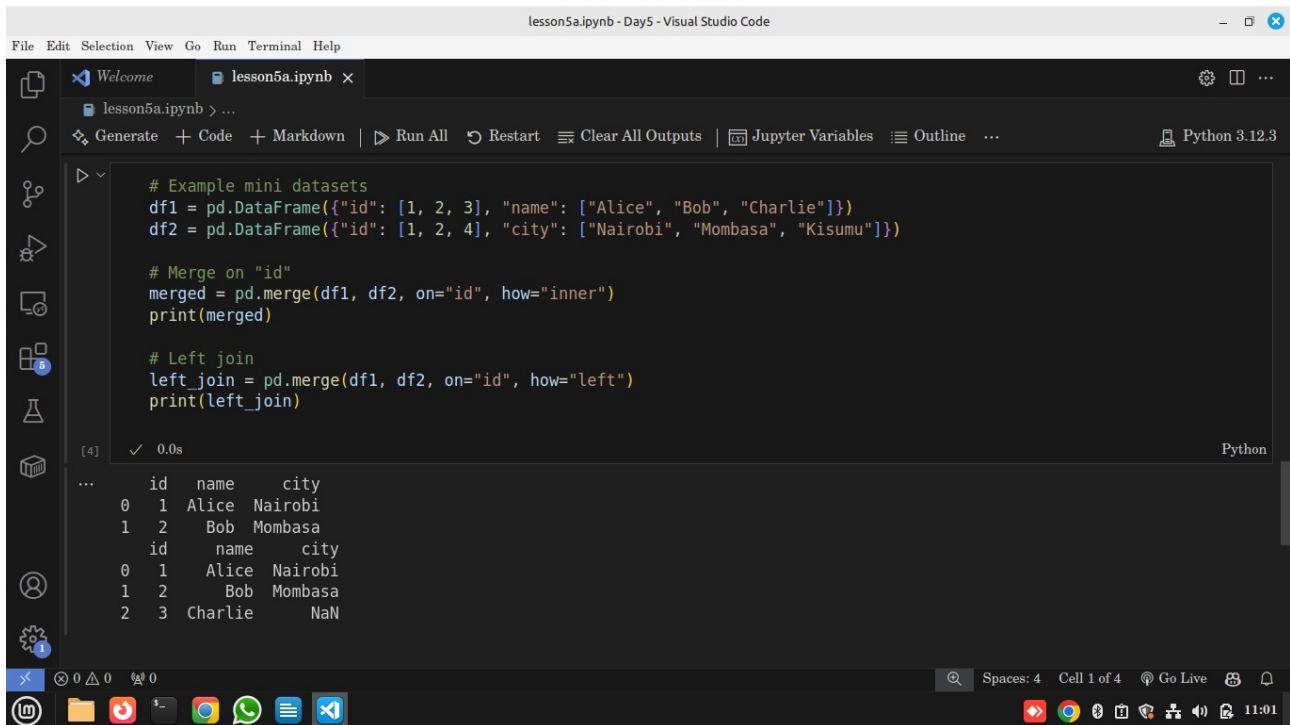---

◆ **Merge & Join**

# Practical



**Explanation**

- `pd.merge()` combines datasets by keys (like SQL joins).

- Types → `inner`, `left`, `right`, `outer`.

**Use in Data Science** → Merging combines info from multiple sources (e.g., transactions + customer demographics).

## 🔷 Pivot Tables

**Practical**



**Explanation**

- Pivot tables reorganize data to show summaries across **two dimensions**.

- Example → Survival rate by sex *and* passenger class.

**Use in Data Science** → Pivot tables simplify comparisons (like Excel pivot tables for quick insights).

---

## 🔷 Mini Challenge

👉 On Titanic dataset:

1. Filter passengers under 18.

2. Group by class and calculate survival rate.

3. Pivot survival rates by sex and class.

---

## 🔷 Reflection

- Wrangling transforms raw data into **analysis-ready datasets**.

- Without wrangling, you can't do visualization or ML.

- Which operation felt most powerful today? (Filtering, GroupBy, Merge, Pivot).