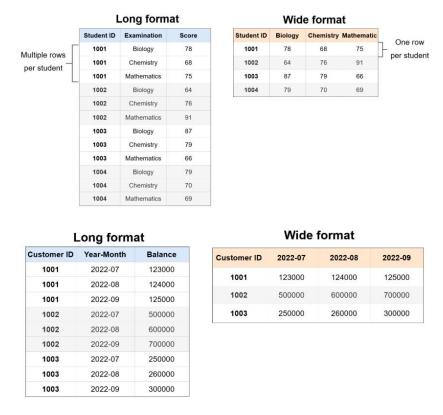# How to leverage pivot () to reshape data?

The pivot function in Pandas is a powerful tool used for **reshaping data** from a "long" format to a "wide" format. Think of it like creating a simple pivot table in a spreadsheet program, where you select specific columns to become new row labels, new column labels, and the values that fill the intersection of these new rows and columns.

### Long format

| Student ID | Examination | Score |
|---|---|---|
| 1001 | Biology | 78 |
| 1001 | Chemistry | 68 |
| 1001 | Mathematics | 75 |
| 1002 | Biology | 64 |
| 1002 | Chemistry | 76 |
| 1002 | Mathematics | 91 |
| 1003 | Biology | 87 |
| 1003 | Chemistry | 79 |
| 1003 | Mathematics | 66 |
| 1004 | Biology | 79 |
| 1004 | Chemistry | 70 |
| 1004 | Mathematics | 69 |

Multiple rows per student

### Wide format

| Student ID | Biology | Chemistry | Mathematic |
|---|---|---|---|
| 1001 | 78 | 68 | 75 |
| 1002 | 64 | 76 | 91 |
| 1003 | 87 | 79 | 66 |
| 1004 | 79 | 70 | 69 |

One row per student

### Long format

| Customer ID | Year-Month | Balance |
|---|---|---|
| 1001 | 2022-07 | 123000 |
| 1001 | 2022-08 | 124000 |
| 1001 | 2022-09 | 125000 |
| 1002 | 2022-07 | 500000 |
| 1002 | 2022-08 | 600000 |
| 1002 | 2022-09 | 700000 |
| 1003 | 2022-07 | 250000 |
| 1003 | 2022-08 | 260000 |
| 1003 | 2022-09 | 300000 |

### Wide format

| Customer ID | 2022-07 | 2022-08 | 2022-09 |
|---|---|---|---|
| 1001 | 123000 | 124000 | 125000 |
| 1002 | 500000 | 600000 | 700000 |
| 1003 | 250000 | 260000 | 300000 |

**Purpose of Pivot**

The primary **purpose** of pivot is to **rearrange data to make it more readable, summarize it in a cross-tabular format, or prepare it for specific types of analysis or visualization** that require a "wide" data structure. It transforms a DataFrame where certain values are repeated across rows into a more compact table where unique values from one column become new columns, and unique values from another become new rows.

**How Pivot Works and Why It Is Required**

pivot operates by taking three key pieces of information from your original DataFrame:

1. **index (required):**

- o **What it does:** This parameter specifies the column (or columns) from your original DataFrame whose unique values will become the **new row labels** (the index) of your reshaped DataFrame.

- o **Why it's required:** You need a clear way to identify the distinct entities or categories that will form the rows of your new, wider table. For example, if you have sales data and you want to see sales figures per month, 'Month' would be your index.

2. **columns (required):**

   - o **What it does:** This parameter specifies the column from your original DataFrame whose unique values will become the **new column headers** of your reshaped DataFrame.

   - o **Why it's required:** This defines the categories that will spread across the top of your new table. For instance, if you want to see sales figures per month *and* per product type, 'Product Type' would be your columns.

3. **values (optional, but common):**

   - o **What it does:** This parameter specifies the column (or columns) from your original DataFrame whose values will populate the cells at the intersection of the new rows and columns.

   - o **Why it's required:** This is the actual data you want to display in your reshaped table. If you're looking at sales, 'Sales Amount' would be your values. If omitted, Pandas will try to use all remaining columns as values, which can lead to a multi-level column index.

**The core mechanism:** pivot essentially takes a unique combination of an index value and a columns value and places the corresponding values in their new intersection.

**Important Constraint:** For pivot to work correctly, the combination of values in the index column(s) and columns column(s) must be **unique**. If there are duplicate combinations (e.g., two sales entries for 'January' and 'Product A'), pivot won't know which value to put in the cell and will raise an error. In such

cases, you would typically use pivot_table, which handles aggregation of duplicate entries.

**Why is Pivot Required?**

pivot is essential for several reasons in data analysis:

- **Improved Readability and Presentation**: It transforms data from a raw, transactional format into a more intuitive, summary-like table that is easier for humans to read and interpret. For example, seeing monthly sales for each product type in a single table is much clearer than scanning through many individual transactions.

- **Cross-Tabulation**: It allows for quick cross-tabulation of data, showing the relationship between two categorical variables (one as index, one as columns) and a numerical variable (as values).

- **Preparation for Specific Tools/Models**: Many statistical analysis tools, machine learning libraries, or visualization tools expect data in a "wide" format where each column represents a distinct feature or variable. pivot helps prepare data for these downstream processes.

- **Summarizing Data**: While it doesn't perform aggregations like groupby().sum(), it implicitly summarizes by bringing related values together into a single cell based on the unique index-column combination.

In summary, pivot is a powerful data reshaping operation that transforms "long" data into a more accessible and analytically useful "wide" format, making it easier to extract insights and prepare data for further steps in the data science workflow.