

# Week 4: Working with Data in Python

---

Dr Christian Engels  
*ce50@st-andrews.ac.uk*

FI5699 Dissertation Module  
Department of Finance  
University of St Andrews Business School



University of  
**St Andrews**

## Introduction

---

# What We Will Cover Today

This is a **hands-on session**. You are welcome to have your laptop open and follow along.

- 1 **Polars Data Structures** — DataFrames, Series, and data types
- 2 **Reading & Writing Data** — CSV and Excel I/O with Polars
- 3 **Expressions & Contexts** — the Polars way of querying data
- 4 **Working with Expressions** — arithmetic, comparisons, conditionals
- 5 **Working with Stock Returns** — downloading Yahoo Finance data
- 6 **Combining DataFrames** — joins and concatenation
- 7 **Reshaping Data** — pivot and unpivot
- 8 **Visualising Data** — building plots with `plotnine`

By the end of today, you will be able to load real financial data, query, transform, reshape, and visualise it in Python.



# Table of Contents

## Introduction

Roadmap

Recap

Getting Started

Motivation

## Polars Data Structures

The DataFrame

Data Types

Inspecting Data

Exercises

## Reading and Writing Data

CSV Files

Excel Files

Exercises

## Expressions and Contexts

The Core Idea

The Four Contexts

Exercises

## Working with Expressions

Arithmetic and Comparisons

Conditionals and Casting

Expression Expansion

Exercises

## Working with Stock Returns

Downloading Stock Prices

Computing Returns

Multiple Stocks

Exercises

## Combining DataFrames

Joins

Concatenation

Exercises

## Reshaping Data

Pivot and Unpivot

Exercises

## Visualising Data

The Grammar of Graphics

Building Plots

Exercises

## Summary

Key Takeaways

Self-Check

Next Steps



## Recap: What You Built in Week 2

Concept	What you learned
Positron	Console, Terminal, and Editor
uv	<code>uv init</code> , <code>uv add</code> to manage packages
Variables	<code>price = 142.50</code> , types ( <code>int</code> , <code>float</code> , <code>str</code> , <code>bool</code> )
Lists & Dicts	<code>["AAPL", "MSFT"]</code> and <code>{"ticker": "AAPL"}</code>
Loops	<code>for</code> and <code>while</code> — automating repetitive tasks

Today we move from Python basics to **working with real data** using **polars** and **plotnine**.



Every concept in this lecture has a runnable Python script in the **week-4-code** folder.

Script	Covers
<code>01_data_structures.py</code>	DataFrames, Series, types, nulls
<code>02_reading_writing_data.py</code>	CSV & Excel I/O
<code>03_expressions_contexts.py</code>	select, with_columns, filter, group_by
<code>04_working_with_expressions.py</code>	Arithmetic, conditionals, casting
<code>05_real_stock_data.py</code>	Yahoo Finance data & returns
<code>06_combining_dataframes.py</code>	Joins & concatenation
<code>07_reshaping_data.py</code>	Pivot & unpivot
<code>08_visualising_data.py</code>	plotnine charts

Open **week-4-code** in Positron and run each script to see the concepts in the `.py` files.



# Setting Up: Two Steps

## 1 Open the folder in Positron

File → Open Folder... → select `week-4-code`

## 2 Install dependencies

Open a Terminal in Positron and run:

```
1 uv sync
```

This reads `pyproject.toml` and installs Polars, plotnine, tidyfinance, and all other packages into a virtual environment.

## 3 Select the interpreter

Click **Select Interpreter** in the top-right corner of Positron and choose the `.venv` Python that `uv sync` just created. This ensures Positron uses the correct packages when you run scripts.

Once the interpreter is set, you are ready to go.  
Open any script and press **Ctrl+Shift+Enter** to run it.



In Week 2 you used lists and dictionaries to organise small amounts of data. But your dissertation will involve **thousands of rows** — stock prices, returns, firm characteristics.

	Lists / Dicts	Polars
<b>Data size</b>	A handful of values	Millions of rows
<b>Operations</b>	One value at a time	Entire columns at once
<b>Speed</b>	Pure Python (slow)	Written in Rust (fast)
<b>Structure</b>	Ad hoc	Typed columns, schemas

Think of Polars as a **supercharged spreadsheet** that you control with code.





## Polars Data Structures

---

A **DataFrame** is the core Polars data structure:

- A table of rows and columns — like a spreadsheet

Column (Series)

ticker	price	shares
<i>str</i>	<i>f64</i>	<i>i64</i>
AAPL	189.84	5 ← Row
MSFT	378.91	30
TSLA	248.42	20

Diagram illustrating a DataFrame structure. The table has three columns: 'ticker', 'price', and 'shares'. The 'price' column is labeled 'Column (Series)'. The 'shares' column has a value of 5 for the row with ticker 'AAPL', which is labeled 'Row'.



# A DataFrame Is a Table with Typed Columns

A **DataFrame** is the core Polars data structure:

- A table of rows and columns — like a spreadsheet
- Each column has a **name** and a **data type**

Column (Series)

ticker	price	shares
<i>str</i>	<i>f64</i>	<i>i64</i>
AAPL	189.84	5 ← Row
MSFT	378.91	30
TSLA	248.42	20

Diagram description: A table with 3 columns and 5 rows. The first row is a header with blue background and white text: 'ticker', 'price', 'shares'. The second row shows data types: 'str', 'f64', 'i64'. The third row shows 'AAPL', '189.84', '5'. The fourth row shows 'MSFT', '378.91', '30'. The fifth row shows 'TSLA', '248.42', '20'. A vertical line above the 'price' column is labeled 'Column (Series)'. A horizontal arrow points to the '5' in the 'AAPL' row, labeled 'Row'.



# A DataFrame Is a Table with Typed Columns

A **DataFrame** is the core Polars data structure:

- A table of rows and columns — like a spreadsheet
- Each column has a **name** and a **data type**
- Each column is internally a **Series**

Column (Series)

ticker	price	shares
<i>str</i>	<i>f64</i>	<i>i64</i>
AAPL	189.84	5 ← Row
MSFT	378.91	30
TSLA	248.42	20

Diagram illustrating a DataFrame structure. The columns are labeled 'ticker', 'price', and 'shares'. The rows are labeled 'AAPL', 'MSFT', and 'TSLA'. The data types for the columns are 'str', 'f64', and 'i64' respectively. An arrow points to the '5' value in the 'shares' column for the 'AAPL' row, labeled 'Row'.



# A DataFrame Is a Table with Typed Columns

A **DataFrame** is the core Polars data structure:

- A table of rows and columns — like a spreadsheet
- Each column has a **name** and a **data type**
- Each column is internally a **Series**
- All values in a column share the same type

Column (Series)

ticker	price	shares
<i>str</i>	<i>f64</i>	<i>i64</i>
AAPL	189.84	5 ← Row
MSFT	378.91	30
TSLA	248.42	20

Diagram illustrating a DataFrame structure. The columns are labeled 'ticker', 'price', and 'shares'. The rows represent individual data points. The 'price' column is labeled 'Column (Series)'. The 'shares' column value '5' is highlighted with a blue arrow pointing to it, labeled 'Row'.



# A DataFrame Is a Table with Typed Columns

A **DataFrame** is the core Polars data structure:

- A table of rows and columns — like a spreadsheet
- Each column has a **name** and a **data type**
- Each column is internally a **Series**
- All values in a column share the same type
- Different columns can have different types

Column (Series)

ticker	price	shares
<i>str</i>	<i>f64</i>	<i>i64</i>
AAPL	189.84	5
MSFT	378.91	30
TSLA	248.42	20

Row



# Creating Your First Polars DataFrame

```
1 import polars as pl
2
3 stocks = pl.DataFrame({
4     "ticker": ["AAPL", "MSFT", "TSLA", "AMZN", "GOOGL"],
5     "price": [189.84, 378.91, 248.42, 178.25, 141.80],
6     "shares": [50, 30, 20, 40, 35],
7     "sector": ["Tech", "Tech", "Auto", "Retail", "Tech"],
8 })
9 print(stocks)
```

shape: (5, 4)

ticker	price	shares	sector
str	f64	i64	str
AAPL	189.84	50	Tech
MSFT	378.91	30	Tech
TSLA	248.42	20	Auto
AMZN	178.25	40	Retail
GOOGL	141.8	35	Tech



# Every Column Has a Data Type

Polars *infers* types automatically, but it helps to know the common ones:

Type	Polars name	Example values
Whole number	<b>Int64</b>	Shares: 50, 30, 20
Decimal number	<b>Float64</b>	Price: 189.84, 378.91
Text	<b>String</b>	Ticker: "AAPL", "MSFT"
True / False	<b>Boolean</b>	Is profitable: true, false
Date	<b>Date</b>	2024-01-15
Date + time	<b>Datetime</b>	2024-01-15 09:30:00
Missing value	<b>Null</b>	Unknown or not applicable

You can see each column's type in the --- row of the printed output.





## A Series Is a Single Column

A **Series** is a one-dimensional array — one column of a DataFrame:

```
1 prices = pl.Series("price", [189.84, 378.91, 248.42])
2 print(prices)
3 print(f"Type: {prices.dtype}")
4 print(f"Mean: {prices.mean():.2f}")
```

```
shape: (3,)
Series: 'price' [f64]
[
  189.84
  378.91
  248.42
]
Type:   Float64
Mean:   272.39
```

You rarely create Series directly — they arise when you extract a column from a DataFrame.



# Five Ways to Look at Your Data

Method	What it shows	Example
<code>.head(n)</code>	First <code>n</code> rows (default 5)	<code>stocks.head(3)</code>
<code>.tail(n)</code>	Last <code>n</code> rows	<code>stocks.tail(3)</code>
<code>.shape</code>	(rows, columns) tuple	<code>stocks.shape</code>
<code>.columns</code>	List of column names	<code>stocks.columns</code>
<code>.schema</code>	Column names and types	<code>stocks.schema</code>

```
1 print(stocks.shape)
2 print(stocks.columns)
```

```
(5, 4)
['ticker', 'price', 'shares', 'sector']
```



## Summary Statistics with `describe()`

```
1 print(stocks.describe())
```

```
shape: (9, 5)
```

```
+-----+-----+-----+-----+
| statistic | ticker | price   | shares | sector |
| ---      | ---    | ---     | ---    | ---    |
| str      | str    | f64     | f64    | str    |
+-----+-----+-----+-----+
| count     | 5      | 5.0     | 5.0    | 5      |
| null_count| 0      | 0.0     | 0.0    | 0      |
| mean      | null   | 227.444 | 35.0   | null   |
| std       | null   | 92.939425 | 11.18034 | null   |
| min       | AAPL   | 141.8   | 20.0   | Auto   |
| 25%      | null   | 178.25  | 30.0   | null   |
| 50%      | null   | 189.84  | 35.0   | null   |
| 75%      | null   | 248.42  | 40.0   | null   |
| max       | TSLA   | 378.91  | 50.0   | Tech   |
+-----+-----+-----+-----+
```

One command gives you count, mean, standard deviation, min, max, and quartiles — the same summary statistics you would compute in Excel or Stata.



Polars uses `null` for missing data. Every data type supports it — unlike Excel's blank cells, a `null` is explicit and typed.

## Why it matters:

- Real-world data has gaps (no price on a holiday, missing filings)

## `null` vs `NaN`:

When you see `null` in your output, it means the data is **genuinely missing**, not zero.



Polars uses `null` for missing data. Every data type supports it — unlike Excel's blank cells, a `null` is explicit and typed.

## Why it matters:

- Real-world data has gaps (no price on a holiday, missing filings)
- Polars propagates `null` through calculations: `null + 5 = null`

## `null` vs NaN:

When you see `null` in your output, it means the data is **genuinely missing**, not zero.



Polars uses `null` for missing data. Every data type supports it — unlike Excel's blank cells, a `null` is explicit and typed.

## Why it matters:

- Real-world data has gaps (no price on a holiday, missing filings)
- Polars propagates `null` through calculations: `null + 5 = null`
- You can detect them with `.is_null()` or fill them with `.fill_null()`

## `null` vs NaN:

When you see `null` in your output, it means the data is **genuinely missing**, not zero.



Polars uses `null` for missing data. Every data type supports it — unlike Excel's blank cells, a `null` is explicit and typed.

## Why it matters:

- Real-world data has gaps (no price on a holiday, missing filings)
- Polars propagates `null` through calculations: `null + 5 = null`
- You can detect them with `.is_null()` or fill them with `.fill_null()`

## `null` vs NaN:

- `null` = “no value exists”

When you see `null` in your output, it means the data is **genuinely missing**, not zero.



Polars uses `null` for missing data. Every data type supports it — unlike Excel's blank cells, a `null` is explicit and typed.

## Why it matters:

- Real-world data has gaps (no price on a holiday, missing filings)
- Polars propagates `null` through calculations: `null + 5 = null`
- You can detect them with `.is_null()` or fill them with `.fill_null()`

## `null` vs `NaN`:

- `null` = “no value exists”
- `NaN` = “a float calculation failed” (e.g. `0/0`)

When you see `null` in your output, it means the data is **genuinely missing**, not zero.





Polars uses `null` for missing data. Every data type supports it — unlike Excel's blank cells, a `null` is explicit and typed.

## Why it matters:

- Real-world data has gaps (no price on a holiday, missing filings)
- Polars propagates `null` through calculations: `null + 5 = null`
- You can detect them with `.is_null()` or fill them with `.fill_null()`

## `null` vs `NaN`:

- `null` = “no value exists”
- `NaN` = “a float calculation failed” (e.g. `0/0`)
- In Polars, use `null` for missing data

When you see `null` in your output, it means the data is **genuinely missing**, not zero.



Polars uses `null` for missing data. Every data type supports it — unlike Excel's blank cells, a `null` is explicit and typed.

## Why it matters:

- Real-world data has gaps (no price on a holiday, missing filings)
- Polars propagates `null` through calculations: `null + 5 = null`
- You can detect them with `.is_null()` or fill them with `.fill_null()`

## `null` vs `NaN`:

- `null` = “no value exists”
- `NaN` = “a float calculation failed” (e.g. `0/0`)
- In Polars, use `null` for missing data
- `NaN` is only for `Float64` columns

When you see `null` in your output, it means the data is **genuinely missing**, not zero.



### Try It Yourself (5 minutes)

- 1 Create a new file `week4.py` and add `import polars as pl` at the top
- 2 Create a DataFrame called `portfolio` with columns:
  - `"ticker"`: at least 4 stocks you know
  - `"price"`: a price for each (make them up)
  - `"shares"`: how many shares of each
- 3 Print the DataFrame
- 4 Print `portfolio.shape` and `portfolio.describe()`
- 5 **Bonus:** Create a `pl.Series` of monthly returns and compute its `.mean()`



## Reading and Writing Data

---

Most financial data starts life as a CSV file. Polars reads them with `pl.read_csv()`:

```
1 import polars as pl
2
3 df = pl.read_csv("stock_prices.csv")
4 print(df.head())
```

Function	What it does
<code>pl.read_csv(path)</code>	Read an entire CSV into a DataFrame immediately
<code>pl.scan_csv(path)</code>	Return a <b>LazyFrame</b> — delays reading until you call <code>.collect()</code>

`scan_csv` is faster for large files because Polars can optimise the query before reading the data.



## Useful read\_csv Parameters

```
1 df = pl.read_csv(  
2     "stock_prices.csv",  
3     separator=";",           # default; use "\t" for TSV  
4     has_header=True,        # first row is column names  
5     skip_rows=0,            # skip N rows at the top  
6     n_rows=1000,            # read only the first 1000 rows  
7     null_values=["NA", ""], # treat these strings as null  
8     try_parse_dates=True,   # auto-detect date columns  
9 )
```

`try_parse_dates=True` is especially useful for financial data — it converts date strings into proper **Date** columns automatically.



After transforming your data, save it with `.write_csv()`:

```
1 stocks = pl.DataFrame({
2     "ticker": ["AAPL", "MSFT", "TSLA"],
3     "price": [189.84, 378.91, 248.42],
4     "shares": [50, 30, 20],
5 })
6
7 stocks.write_csv("my_portfolio.csv")
```

This writes a standard comma-separated file with a header row. You can customise:

Parameter	Effect
<code>separator="\t"</code>	Write a tab-separated file instead
<code>include_header=False</code>	Omit the header row



For large files, `pl.scan_csv()` returns a **LazyFrame** instead of loading everything:

```
1 lf = pl.scan_csv("large_dataset.csv")
2
3 # Build a query without reading data yet
4 result = (
5     lf
6     .filter(pl.col("ticker") == "AAPL")
7     .select("date", "close")
8     .collect() # NOW it reads and executes
9 )
```

## read\_csv

- Reads entire file immediately

## scan\_csv





For large files, `pl.scan_csv()` returns a **LazyFrame** instead of loading everything:

```
1 lf = pl.scan_csv("large_dataset.csv")
2
3 # Build a query without reading data yet
4 result = (
5     lf
6     .filter(pl.col("ticker") == "AAPL")
7     .select("date", "close")
8     .collect() # NOW it reads and executes
9 )
```

## read\_csv

- Reads entire file immediately
- Returns a **DataFrame**

## scan\_csv



For large files, `pl.scan_csv()` returns a **LazyFrame** instead of loading everything:

```
1 lf = pl.scan_csv("large_dataset.csv")
2
3 # Build a query without reading data yet
4 result = (
5     lf
6     .filter(pl.col("ticker") == "AAPL")
7     .select("date", "close")
8     .collect() # NOW it reads and executes
9 )
```

## read\_csv

- Reads entire file immediately
- Returns a **DataFrame**
- Simple and direct

## scan\_csv



For large files, `pl.scan_csv()` returns a **LazyFrame** instead of loading everything:

```
1 lf = pl.scan_csv("large_dataset.csv")
2
3 # Build a query without reading data yet
4 result = (
5     lf
6     .filter(pl.col("ticker") == "AAPL")
7     .select("date", "close")
8     .collect() # NOW it reads and executes
9 )
```

## read\_csv

- Reads entire file immediately
- Returns a **DataFrame**
- Simple and direct
- Fine for small-medium files

## scan\_csv



For large files, `pl.scan_csv()` returns a **LazyFrame** instead of loading everything:

```
1 lf = pl.scan_csv("large_dataset.csv")
2
3 # Build a query without reading data yet
4 result = (
5     lf
6     .filter(pl.col("ticker") == "AAPL")
7     .select("date", "close")
8     .collect() # NOW it reads and executes
9 )
```

## read\_csv

- Reads entire file immediately
- Returns a **DataFrame**
- Simple and direct
- Fine for small-medium files

## scan\_csv

- Defers reading until `.collect()`



For large files, `pl.scan_csv()` returns a **LazyFrame** instead of loading everything:

```
1 lf = pl.scan_csv("large_dataset.csv")
2
3 # Build a query without reading data yet
4 result = (
5     lf
6     .filter(pl.col("ticker") == "AAPL")
7     .select("date", "close")
8     .collect() # NOW it reads and executes
9 )
```

## read\_csv

- Reads entire file immediately
- Returns a **DataFrame**
- Simple and direct
- Fine for small-medium files

## scan\_csv

- Defers reading until `.collect()`
- Returns a **LazyFrame**



For large files, `pl.scan_csv()` returns a **LazyFrame** instead of loading everything:

```
1 lf = pl.scan_csv("large_dataset.csv")
2
3 # Build a query without reading data yet
4 result = (
5     lf
6     .filter(pl.col("ticker") == "AAPL")
7     .select("date", "close")
8     .collect() # NOW it reads and executes
9 )
```

## read\_csv

- Reads entire file immediately
- Returns a **DataFrame**
- Simple and direct
- Fine for small-medium files

## scan\_csv

- Defers reading until `.collect()`
- Returns a **LazyFrame**
- Polars optimises the query first



For large files, `pl.scan_csv()` returns a **LazyFrame** instead of loading everything:

```
1 lf = pl.scan_csv("large_dataset.csv")
2
3 # Build a query without reading data yet
4 result = (
5     lf
6     .filter(pl.col("ticker") == "AAPL")
7     .select("date", "close")
8     .collect() # NOW it reads and executes
9 )
```

## read\_csv

- Reads entire file immediately
- Returns a **DataFrame**
- Simple and direct
- Fine for small-medium files

## scan\_csv

- Defers reading until `.collect()`
- Returns a **LazyFrame**
- Polars optimises the query first
- Better for large files



Polars can also read `.xlsx` files directly:

```
1 df = pl.read_excel("portfolio.xlsx")
2
3 # Read a specific sheet
4 df = pl.read_excel("portfolio.xlsx", sheet_name="Returns")
```

**Required dependency:** install one of these Excel engines:

```
$ uv add fastexcel      # fastest (recommended, default)
$ uv add openpyxl      # slower but handles tricky files
```

Without a `sheet_name`, Polars reads the **first sheet** by default.





Save a DataFrame to Excel with `.write_excel()`:

```
1 stocks.write_excel("my_portfolio.xlsx")
2
3 # Specify a worksheet name
4 stocks.write_excel("my_portfolio.xlsx", worksheet="Holdings")
```

Required dependency:

```
$ uv addxlsxwriter
```

**Performance tip:** CSV and Parquet are much faster than Excel for large datasets. Use Excel when you need to share with non-programmers; use CSV or Parquet for your own analysis pipelines.



### Try It Yourself (5 minutes)

- 1 Create a DataFrame with at least 4 rows of stock data (ticker, price, shares)
- 2 Write it to a CSV file with `.write_csv("portfolio.csv")`
- 3 Read it back with `pl.read_csv("portfolio.csv")` and print the result — does it match?
- 4 **Bonus:** Write to Excel with `.write_excel("portfolio.xlsx")` and read it back
- 5 **Bonus:** Try `pl.scan_csv("portfolio.csv").collect()` — same result?



## Expressions and Contexts

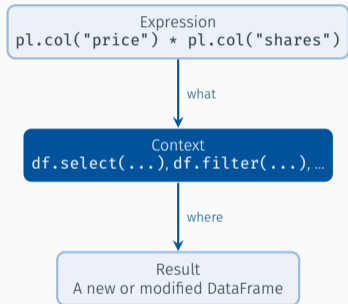
---

# Polars Separates *What* to Compute from *Where*

An **expression** is a recipe for a computation. It describes **what** to do, but does nothing until you place it in a **context**.

- Expression: “multiply price by shares”

This separation is what makes Polars fast and flexible.

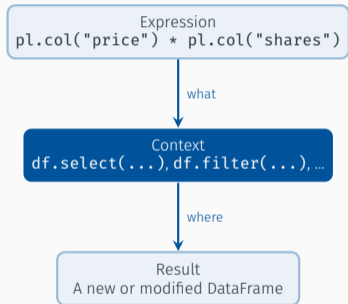


# Polars Separates *What* to Compute from *Where*

An **expression** is a recipe for a computation. It describes **what** to do, but does nothing until you place it in a **context**.

- Expression: “multiply price by shares”
- Context: “add it as a new column” or “replace all columns” or “use it to filter rows”

This separation is what makes Polars fast and flexible.



## Context 1: select — Pick and Transform Columns

select returns **only** the columns you ask for:

```
1 result = stocks.select(  
2     "ticker",  
3     (pl.col("price") * pl.col("shares")).alias("position_value"),  
4 )  
5 print(result)
```

shape: (5, 2)

ticker	position_value
str	f64
AAPL	9492.0
MSFT	11367.3
TSLA	4968.4
AMZN	7130.0
GOOGL	4963.0



## Context 2: with\_columns — Add New Columns

`with_columns` keeps **all existing** columns and adds new ones:

```
1 result = stocks.with_columns(  
2     position_value = pl.col("price") * pl.col("shares"),  
3 )  
4 print(result)
```

shape: (5, 5)

ticker	price	shares	sector	position_value
str	f64	i64	str	f64
AAPL	189.84	50	Tech	9492.0
MSFT	378.91	30	Tech	11367.3
TSLA	248.42	20	Auto	4968.4
AMZN	178.25	40	Retail	7130.0
GOOGL	141.8	35	Tech	4963.0



## Context 3: `filter` — Keep Matching Rows

`filter` keeps only the rows where a condition is `True`:

```
1 result = stocks.filter(pl.col("price") > 200)
2 print(result)
```

shape: (2, 4)

```
+-----+-----+-----+-----+
| ticker | price | shares | sector |
| ---   | ---   | ---   | ---   |
| str   | f64   | i64   | str   |
+=====+=====+=====+=====+
| MSFT  | 378.91 | 30    | Tech  |
| TSLA  | 248.42 | 20    | Auto  |
+-----+-----+-----+-----+
```

Think of `filter` as Excel's AutoFilter: only rows that match your condition survive.





## Context 4: `group_by` — Aggregate by Group

`group_by` splits the data into groups and applies aggregation functions:

```
1 result = stocks.group_by("sector").agg(  
2     pl.col("price").mean().alias("avg_price"),  
3     pl.len().alias("count"),  
4 )  
5 print(result.sort("sector"))
```

shape: (3, 3)

sector	avg_price	count
Auto	248.42	1
Retail	178.25	1
Tech	236.85	3

Like a pivot table in Excel: group rows by a category, then summarise each group.



## The Four Contexts at a Glance

Context	What it does	Excel analogy
<code>select</code>	Pick and transform columns	Choosing which columns to display
<code>with_columns</code>	Add or overwrite columns	Adding a formula column next to existing data
<b><code>filter</code></b>	Keep matching rows	AutoFilter on a column
<code>group_by</code>	Aggregate by group	Pivot table: group and summarise

Master these four and you can answer almost any data question. Everything else builds on them.



### Try It Yourself (5 minutes)

Using your `portfolio` DataFrame from Exercise 1:

- 1 Use `.select()` to show only the ticker and a new column `value` (`price × shares`)
- 2 Use `.with_columns()` to add a `value` column while keeping all originals
- 3 Use `.filter()` to show only stocks with price above 100
- 4 **Bonus:** Use `.group_by()` if you have a sector column — compute the mean price per sector

`.alias("name")` renames a computed column. Keyword syntax (`value = ...`) also works in `with_columns` and `select`.



## Working with Expressions

---

# Column Arithmetic Works Element-Wise

Polars applies operations to **entire columns at once** — no loops needed:

```
1 result = stocks.select(  
2     "ticker",  
3     (pl.col("price") + 10).alias("price_adj"),    # scalar  
4     (pl.col("price") * pl.col("shares")).alias("value"), # column  
5 )  
6 print(result)
```

shape: (5, 3)

ticker	price_adj	value
str	f64	f64
AAPL	199.84	9492.0
MSFT	388.91	11367.3
TSLA	258.42	4968.4
AMZN	188.25	7130.0
GOOGL	151.8	4963.0



## Comparisons Return Boolean Columns

```
1 result = stocks.select(  
2     "ticker", "price",  
3     is_expensive = pl.col("price") > 200,  
4     is_tech = pl.col("sector") == "Tech",  
5 )  
6 print(result)
```

shape: (5, 4)

```
+-----+-----+-----+  
| ticker | price | is_expensive | is_tech |  
| ---   | ---   | ---          | ---     |  
| str    | f64   | bool         | bool    |  
+=====+=====+=====+  
| AAPL   | 189.84 | false        | true    |  
| MSFT   | 378.91 | true         | true    |  
| TSLA   | 248.42 | true         | false   |  
| AMZN   | 178.25 | false        | false   |  
| GOOGL  | 141.8  | false        | true    |  
+-----+-----+-----+
```



## Combining Conditions with & and |

Use & (and) and | (or) to combine Boolean expressions. Parentheses are required.

```
1 # Tech stocks priced above 150
2 result = stocks.filter(
3     (pl.col("price") > 150) & (pl.col("sector") == "Tech")
4 )
5 print(result)
```

shape: (2, 4)

ticker	price	shares	sector
AAPL	189.84	50	Tech
MSFT	378.91	30	Tech

GOOGL (141.80) is Tech but not above 150, so it is excluded. Both conditions must be true.



## Conditionals: when / then / otherwise

Create new values based on conditions – like Excel's IF() but chainable:

```
1 result = stocks.select(  
2     "ticker", "price",  
3     size = pl.col("price") > 300).then(pl.lit("Large"))  
4         .when(pl.col("price") > 200).then(pl.lit("Mid"))  
5         .otherwise(pl.lit("Small"))  
6 )  
7 print(result)
```

shape: (5, 3)

ticker	price	size
str	f64	str
AAPL	189.84	Small
MSFT	378.91	Large
TSLA	248.42	Mid
AMZN	178.25	Small
GOOGL	141.8	Small





## Casting: Changing Data Types

Use `.cast()` to convert a column's type:

```
1 result = stocks.select(  
2     "ticker",  
3     pl.col("price").cast(pl.Int64).alias("price_int"),  
4     pl.col("shares").cast(pl.Float64).alias("shares_f"),  
5 )  
6 print(result)
```

```
shape: (5, 3)  
+-----+-----+  
| ticker | price_int | shares_f |  
| ---   | ---       | ---       |  
| str    | i64       | f64       |  
+-----+-----+  
| AAPL   | 189       | 50.0      |  
| MSFT   | 378       | 30.0      |  
| TSLA   | 248       | 20.0      |  
| AMZN   | 178       | 40.0      |  
| GOOGL  | 141       | 35.0      |  
+-----+-----+
```

**Note:** Floats are **truncated** (not rounded) when cast to integers. 189.84 becomes 189.



## Expression Expansion: Many Columns at Once

Instead of writing the same expression for each column, pass **multiple names** to `pl.col()`:

```
1 # Apply .mean() to both price and shares in one expression
2 result = stocks.select(
3     pl.col("price", "shares").mean(),
4 )
5 print(result)
```

```
shape: (1, 2)
+-----+-----+
| price | shares |
| ---  | ---  |
| f64   | f64   |
+=====+=====+
| 227.444 | 35.0 |
+-----+-----+
```

One expression, two columns. Polars “expands” the expression for each column name automatically.



## Naming Results with `alias`, `prefix`, and `suffix`

```
1 # .alias() renames a single result
2 pl.col("price").mean().alias("avg_price")
3
4 # .name.prefix() / .name.suffix() for batches
5 result = stocks.select(
6     pl.col("price", "shares").mean().name.prefix("avg_"),
7 )
8 print(result)
```

shape: (1, 2)

```
+-----+-----+
| avg_price | avg_shares |
| ---      | ---      |
| f64       | f64       |
+-----+-----+
| 227.444   | 35.0      |
+-----+-----+
```

Use `.alias()` for one column. Use `.name.prefix()` or `.name.suffix()` when expanding across several.



### Try It Yourself (5 minutes)

Using the `stocks` DataFrame:

- 1 Add a Boolean column `is_large` that is `True` when `price > 200`
- 2 Create a column `weight` as each stock's value divided by the total portfolio value:
  - Hint: `pl.col("value") / pl.col("value").sum()`
- 3 Use `when/then/otherwise` to label stocks as "Buy" (`price < 200`) or "Hold"
- 4 **Bonus:** Use expression expansion to compute `.max()` of both price and shares



## Working with Stock Returns

---

The `tidyfinance` package provides a clean interface to download stock prices:

```
$ uv add tidyfinance
```

```
1 import tidyfinance as tf
2
3 prices = tf.download_data(
4     domain="stock_prices",
5     symbols="AAPL",
6     start_date="2000-01-01",
7     end_date="2024-12-31",
8 )
9 print(prices.head())
```

This returns a DataFrame with columns: `symbol`, `date`, `volume`, `open`, `low`, `high`, `close`, and `adjusted_close`.



# Why Adjusted Close?

Stock prices are affected by events **after market close**:

- **Stock splits** — a 2-for-1 split halves the price but doubles shares

The **adjusted close** corrects for these events, giving you the “true” price trajectory for computing returns.



**Always** use adjusted close when computing returns.

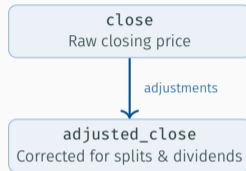


# Why Adjusted Close?

Stock prices are affected by events **after market close**:

- **Stock splits** — a 2-for-1 split halves the price but doubles shares
- **Dividends** — cash paid to shareholders reduces the stock price

The **adjusted close** corrects for these events, giving you the “true” price trajectory for computing returns.



**Always** use adjusted close when computing returns.





## Daily Returns: The Percentage Change

A **return** measures the percentage change in price from one day to the next:

$$r_t = \frac{p_t}{p_{t-1}} - 1$$

```
1 returns = (  
2     prices  
3     .sort("date")  
4     .with_columns(  
5         ret=pl.col("adjusted_close").pct_change()  
6     )  
7     .select("symbol", "date", "ret")  
8     .drop_nulls()  
9 )  
10 print(returns.head())
```

`.pct_change()` computes  $(x_t/x_{t-1}) - 1$  for each row. The first row is `null` because there is no previous price.



# Summary Statistics of Returns

```
1 print(returns.select("ret").describe())
```

```
shape: (9, 2)
```

```
+-----+-----+
| statistic | ret      |
| ---      | ---      |
| str       | f64      |
+=====+=====+
| count     | 6297.0   |
| null_count| 0.0      |
| mean      | 0.001    |
| std       | 0.02     |
| min       | -0.519   |
| 25%       | -0.008   |
| 50%       | 0.001    |
| 75%       | 0.012    |
| max       | 0.139    |
+-----+-----+
```

Summary statistics reveal the distribution of daily returns: mean, volatility (std), and the range of outcomes.



## Downloading Multiple Stocks at Once

Pass a list of symbols to download several stocks in one call:

```
1 symbols = ["AAPL", "MSFT", "TSLA", "AMZN", "GOOGL"]
2
3 prices_daily = tf.download_data(
4     domain="stock_prices",
5     symbols=symbols,
6     start_date="2020-01-01",
7     end_date="2024-12-31",
8 )
9 print(prices_daily.shape)
10 print(prices_daily.head())
```

The result is a **long-format** DataFrame — one row per stock per day. The `symbol` column identifies which stock each row belongs to.



## Computing Returns for Multiple Stocks

Use `.over("symbol")` to compute returns **within** each stock:

```
1 returns_daily = (  
2     prices_daily  
3     .sort("symbol", "date")  
4     .with_columns(  
5         ret=pl.col("adjusted_close")  
6             .pct_change()  
7             .over("symbol")  
8     )  
9     .select("symbol", "date", "ret")  
10    .drop_nulls()  
11 )  
12 print(returns_daily.head())
```

`.over("symbol")` is a **window function** — it applies `pct_change()` separately within each stock, preventing returns from crossing stock boundaries.



Compound daily returns within each month to get monthly returns:

```
1 returns_monthly = (  
2     returns_daily  
3     .with_columns(  
4         month=pl.col("date").dt.month_start()  
5     )  
6     .group_by("symbol", "month")  
7     .agg(  
8         ret_monthly=(  
9             (pl.col("ret") + 1).product() - 1  
10        )  
11    )  
12    .sort("symbol", "month")  
13 )  
14 print(returns_monthly.head())
```

The formula:  $(1 + r_1)(1 + r_2) \cdots (1 + r_n) - 1$  compounds daily returns into a single monthly return.



## Saving Your Data for Later

Once downloaded, save the data so you do not need to re-download it every time:

```
1 # Save to CSV
2 returns_daily.write_csv("returns_daily.csv")
3
4 # Read it back later
5 returns_daily = pl.read_csv(
6     "returns_daily.csv", try_parse_dates=True
7 )
8
9 # Or save to Excel for sharing
10 returns_monthly.write_excel(
11     "returns_monthly.xlsx", worksheet="Monthly Returns"
12 )
```

**Workflow tip:** Download once, save to CSV, and read from the file for all subsequent analysis. This is faster and avoids hitting API rate limits.



### Try It Yourself (10 minutes)

- 1 Install `tidyfinance`: `uv add tidyfinance`
- 2 Download daily prices for 3 stocks of your choice (2020–2024)
- 3 Compute daily returns using `.pct_change()` and `.over("symbol")`
- 4 Print `.describe()` for the returns column — which stock is most volatile?
- 5 Save your returns to a CSV file and read it back
- 6 **Bonus:** Compute monthly returns using `.group_by()` and `.product()`



## Combining DataFrames

---



# Your Data Often Lives in Multiple Tables

In practice, the information you need is rarely in one place:

Table	Contains
Stock prices	Ticker, date, price, volume
Company info	Ticker, sector, employees, founding year
Fund holdings	Fund name, ticker, weight

A **join** combines columns from two tables by matching on a shared key (like **ticker**).

If you've ever used **VLOOKUP** in Excel, a join does the same thing — but for entire tables at once.

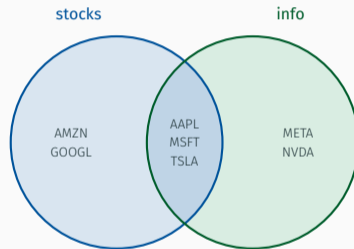


Imagine two tables sharing a **ticker** column:

- **stocks**: AAPL, MSFT, TSLA, AMZN, GOOGL

Some tickers appear in both (the overlap), some only on one side.

The **join type** determines which rows survive.

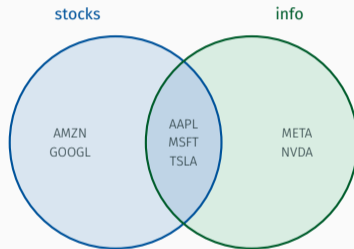


Imagine two tables sharing a **ticker** column:

- **stocks**: AAPL, MSFT, TSLA, AMZN, GOOGL
- **company\_info**: AAPL, MSFT, TSLA, META, NVDA

Some tickers appear in both (the overlap), some only on one side.

The **join type** determines which rows survive.



## Inner Join: Only Matching Rows Survive

```
1 company_info = pl.DataFrame({
2     "ticker": ["AAPL", "MSFT", "TSLA", "META", "NVDA"],
3     "employees": [164000, 221000, 128000, 67317, 29600],
4     "founded": [1976, 1975, 2003, 2004, 1993],
5 })
6
7 result = stocks.join(company_info, on="ticker", how="inner")
8 print(result)
```

shape: (3, 6)

ticker	price	shares	sector	employees	founded
AAPL	189.84	50	Tech	164000	1976
MSFT	378.91	30	Tech	221000	1975
TSLA	248.42	20	Auto	128000	2003

Only the 3 tickers appearing in **both** tables remain. AMZN, GOOGL, META, NVDA are dropped.



## Left Join: Keep All Left Rows

```
1 result = stocks.join(company_info, on="ticker", how="left")
2 print(result)
```

```
shape: (5, 6)
```

```
+-----+-----+-----+-----+-----+-----+
| ticker | price | shares | sector | employees | founded |
| ---   | ---   | ---   | ---   | ---   | ---   |
| str   | f64   | i64   | str   | i64   | i64   |
+-----+-----+-----+-----+-----+-----+
| AAPL  | 189.84 | 50    | Tech  | 164000 | 1976  |
| MSFT  | 378.91 | 30    | Tech  | 221000 | 1975  |
| TSLA  | 248.42 | 20    | Auto  | 128000 | 2003  |
| AMZN  | 178.25 | 40    | Retail | null   | null  |
| GOOGL | 141.8  | 35    | Tech  | null   | null  |
+-----+-----+-----+-----+-----+-----+
```

All 5 stocks from the left table survive. Where no match exists, Polars fills in `null`.



## Anti Join: Find What's Missing

```
1 result = stocks.join(company_info, on="ticker", how="anti")
2 print(result)
```

```
shape: (2, 4)
```

```
+-----+-----+-----+
| ticker | price | shares | sector |
| ---   | ---   | ---   | ---   |
| str    | f64   | i64   | str   |
+-----+-----+-----+
| AMZN   | 178.25 | 40    | Retail |
| GOOGL  | 141.8  | 35    | Tech   |
+-----+-----+-----+
```

An **anti join** returns left-table rows that have **no match** on the right. Useful for finding gaps in your data.



## Join Types at a Glance

Type	Rows kept	Use when...
<code>inner</code>	Only matching rows from both sides	You need complete records from both tables
<code>left</code>	All left rows, matches from right	You want to enrich a primary table
<code>full</code>	All rows from both sides	You need the full union of both datasets
<code>anti</code>	Left rows with <b>no</b> match	Finding gaps or missing data
<code>semi</code>	Left rows that have a match	Filtering by existence in another table

In finance, `left` joins are by far the most common — they let you enrich a primary dataset without losing any rows.



## Concatenation: Stacking DataFrames

When your data is split across files (e.g. one CSV per month), `pl.concat()` combines them:

```
1 jan = pl.DataFrame({"ticker": ["AAPL", "MSFT"], "return_pct": [5.2, 3.1]})
2 feb = pl.DataFrame({"ticker": ["AAPL", "MSFT"], "return_pct": [-1.4, 2.8]})
3
4 result = pl.concat([jan, feb], how="vertical")
5 print(result)
```

```
shape: (4, 2)
+-----+-----+
| ticker | return_pct |
| ---   | ---       |
| str    | f64       |
+-----+-----+
| AAPL   | 5.2       |
| MSFT   | 3.1       |
| AAPL   | -1.4      |
| MSFT   | 2.8       |
+-----+-----+
```

**vertical** = stack rows (same columns). **horizontal** = add columns side by side (same rows).





### Try It Yourself (5 minutes)

- 1 Create a second DataFrame `sectors` with columns `"ticker"` and `"sector"` for some (but not all) of your tickers
- 2 Use an **inner join** to combine `portfolio` and `sectors` — how many rows remain?
- 3 Use a **left join** instead — what happens to tickers without a sector?
- 4 Use an **anti join** to find which tickers are missing from `sectors`
- 5 **Bonus:** Use `pl.concat` to stack two DataFrames of returns (one for each month)



## Reshaping Data

---

# Wide vs Long: Two Ways to Organise the Same Data

## Wide format

Each variable gets its own column.

ticker	Jan	Feb
AAPL	5.2	-1.4
MSFT	3.1	2.8

Compact, human-readable. Good for summary tables.

## Long format

One row per observation.

ticker	month	return
AAPL	Jan	5.2
AAPL	Feb	-1.4
MSFT	Jan	3.1
MSFT	Feb	2.8

Tidy, easy to plot and analyse. Good for computation.



## Pivot: Reshape from Long to Wide

```
1 returns_long = pl.DataFrame({
2     "ticker":      ["AAPL", "AAPL", "MSFT", "MSFT"],
3     "month":       ["Jan",  "Feb", "Jan",  "Feb"],
4     "return_pct": [5.2,  -1.4, 3.1,   2.8],
5 })
6
7 result = returns_long.pivot(
8     index="ticker", on="month", values="return_pct"
9 )
10 print(result)
```

```
shape: (2, 3)
```

```
+-----+-----+-----+
| ticker | Jan | Feb |
| ---   | --- | --- |
| str    | f64 | f64 |
+=====+=====+=====+
| AAPL   | 5.2 | -1.4 |
| MSFT   | 3.1 | 2.8  |
+-----+-----+-----+
```

Each unique value in `month` becomes a new column. The table gets wider and shorter.



## Unpivot: Reshape from Wide to Long

```
1 returns_wide = pl.DataFrame({
2     "ticker": ["AAPL", "MSFT"],
3     "jan":    [5.2, 3.1],
4     "feb":    [-1.4, 2.8],
5 })
6
7 result = returns_wide.unpivot(index="ticker", on=["jan", "feb"])
8 print(result)
```

shape: (4, 3)

ticker	variable	value
str	str	f64
AAPL	jan	5.2
MSFT	jan	3.1
AAPL	feb	-1.4
MSFT	feb	2.8

Column names become values in a new **variable** column. The table gets taller and narrower.



## When to Use Which Format

	Wide	Long
Good for	Reading, summary tables, correlation matrices	Plotting, <code>group_by</code> , regression
Shape	Fewer rows, more columns	More rows, fewer columns
Reshape to	<code>.pivot()</code>	<code>.unpivot()</code>
Analogies	Excel pivot table output	Tidy data, “one row per observation”

**Rule of thumb:** If you need to `group_by` or plot, go long. If a human needs to read a table, go wide.



### Try It Yourself (5 minutes)

- 1 Create a long DataFrame with columns `ticker`, `month`, and `return_pct` for 3 tickers across 3 months (9 rows total)
- 2 Use `.pivot()` to reshape it to wide format (one column per month)
- 3 Take the wide result and use `.unpivot()` to go back to long — do you get the original back?
- 4 **Bonus:** Which format would you need to compute the mean return per ticker? Which to show a summary table to your supervisor?



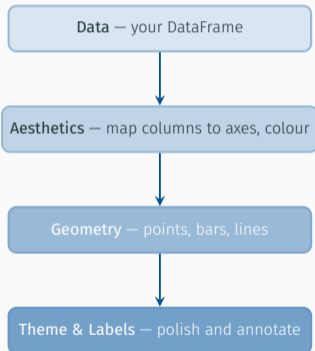
## Visualising Data

---



# Build Plots in Layers with `plotnine`

`plotnine` implements the **Grammar of Graphics** — every plot is built from independent layers:

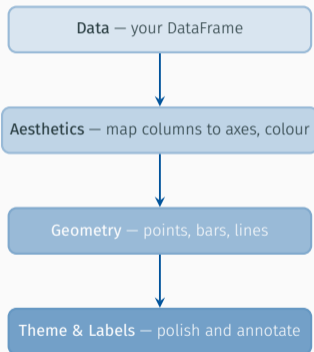


- Layers are combined with `+`



# Build Plots in Layers with `plotnine`

`plotnine` implements the **Grammar of Graphics** — every plot is built from independent layers:

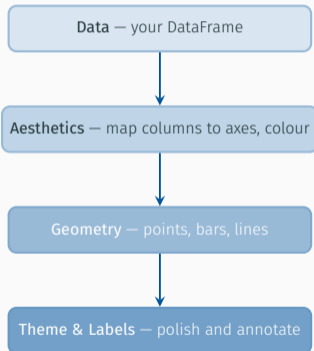


- Layers are combined with `+`
- Each layer is independent — swap `geom_col()` for `geom_point()` and you get a different chart from the same data



# Build Plots in Layers with `plotnine`

`plotnine` implements the **Grammar of Graphics** — every plot is built from independent layers:

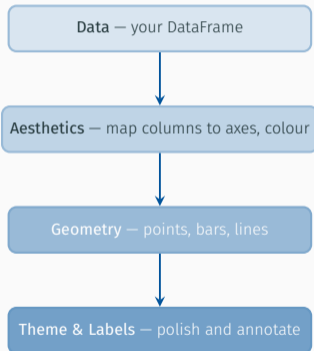


- Layers are combined with `+`
- Each layer is independent — swap `geom_col()` for `geom_point()` and you get a different chart from the same data
- If you know R's `ggplot2`, this is the same syntax



# Build Plots in Layers with `plotnine`

`plotnine` implements the **Grammar of Graphics** — every plot is built from independent layers:

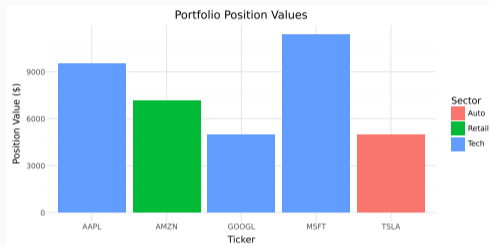


- Layers are combined with `+`
- Each layer is independent — swap `geom_col()` for `geom_point()` and you get a different chart from the same data
- If you know R's `ggplot2`, this is the same syntax
- `plotnine` works with pandas, so we convert: `stocks.to_pandas()`



# Plotting Stock Prices Over Time

```
1 from plotnine import *
2
3 # Convert Polars to pandas for
4 # plotnine
5 prices_pd = prices_daily.to_pandas()
6
7 (
8     ggplot(prices_pd, aes(
9         x="date",
10        y="adjusted_close",
11        color="symbol"))
12     + geom_line()
13     + labs(title="Stock Prices",
14           x="", y="Price ($)",
15           color="Symbol")
16     + theme_minimal()
17 )
```

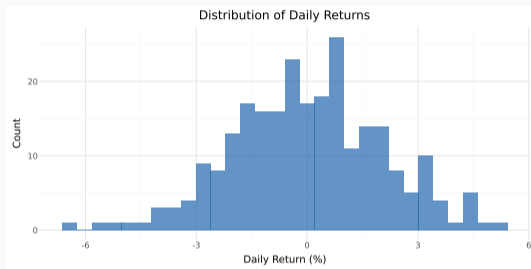


Adding `color="symbol"` in `aes()` automatically draws a separate line for each stock.



# Histogram: Distribution of Returns

```
1 aapl = (  
2   returns_daily  
3   .filter(  
4     pl.col("symbol") == "AAPL"  
5   )  
6   .to_pandas()  
7 )  
8 (  
9   ggplot(aapl, aes(x="ret"))  
10  + geom_histogram(  
11    bins=50,  
12    fill="#003366",  
13    alpha=0.7)  
14  + labs(  
15    title="AAPL Daily Returns",  
16    x="Return", y="Count")  
17  + theme_minimal()  
18 )  
19 )
```

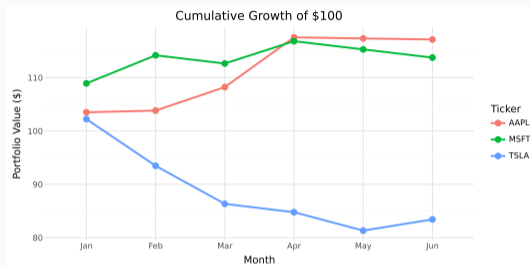


Returns are roughly bell-shaped but with **fat tails** – extreme returns occur more often than a normal distribution predicts.



# Cumulative Returns: Growth of \$1

```
1 cumulative = (  
2   returns_daily  
3   .sort("symbol", "date")  
4   .with_columns(  
5     growth=  
6       (1 + pl.col("ret"))  
7       .cum_prod()  
8       .over("symbol")  
9   )  
10  .to_pandas()  
11 )  
12  
13 (  
14   ggplot(cumulative, aes(  
15     x="date", y="growth",  
16     color="symbol"))  
17   + geom_line()  
18   + labs(title="Growth of $1",  
19         x="", y="Value ($)")  
20   + theme_minimal()  
21 )
```



`.cum_prod()` computes the running product:  $(1 + r_1)(1 + r_2) \cdots (1 + r_t)$ , showing how \$1 invested on day one grows over time.



## The Key `geom_*` Functions

Geom	Plot type	Use when...
<code>geom_col()</code>	Bar chart	Comparing categories
<code>geom_point()</code>	Scatter plot	Showing relationships between two variables
<code>geom_line()</code>	Line chart	Tracking values over time
<code>geom_histogram()</code>	Histogram	Showing the distribution of a variable
<code>geom_smooth()</code>	Trend line	Adding a fitted line to a scatter plot

The pattern: `ggplot(data, aes(...)) + geom_*() + labs() + theme_*()`.  
Swap the `geom` to change the chart type. Everything else stays the same.





### Try It Yourself (5 minutes)

Using the stock returns data from Exercise 5:

- 1 Plot the adjusted close prices over time for your stocks (use `geom_line()`)
- 2 Create a histogram of daily returns for one stock (use `geom_histogram()`)
- 3 Compute cumulative returns with `.cum_prod()` and plot the growth of \$1
- 4 **Bonus:** Map `color="symbol"` to compare multiple stocks on one chart

To save a plot: assign it to a variable and call `p.save("my_plot.pdf")`.



## Summary

---

## Data structures & I/O

- DataFrames, Series, and data types
- `read_csv` / `write_csv`
- `read_excel` / `write_excel`
- `scan_csv` for lazy reading

## Expressions & contexts

- `select`, `with_columns`, `filter`, `group_by`
- Arithmetic, comparisons, `when/then/otherwise`
- Casting and expression expansion

## Stock returns

- Downloading prices from Yahoo Finance
- Computing daily returns with `pct_change()`
- Window functions with `.over("symbol")`
- Compounding daily → monthly returns

## Transformations & visualisation

- Joins, concatenation, pivot/unpivot
- Grammar of Graphics with `plotnine`
- Price charts, return histograms, cumulative growth



Can you answer these in your own words?

- 1 What is the difference between `read_csv` and `scan_csv`?



Can you answer these in your own words?

- 1 What is the difference between `read_csv` and `scan_csv`?
- 2 What is the difference between `select` and `with_columns`?



Can you answer these in your own words?

- 1 What is the difference between `read_csv` and `scan_csv`?
- 2 What is the difference between `select` and `with_columns`?
- 3 Why do we use `adjusted_close` instead of `close` for returns?



Can you answer these in your own words?

- 1 What is the difference between `read_csv` and `scan_csv`?
- 2 What is the difference between `select` and `with_columns`?
- 3 Why do we use `adjusted_close` instead of `close` for returns?
- 4 What does `.pct_change().over("symbol")` do, and why is `.over()` needed?



Can you answer these in your own words?

- 1 What is the difference between `read_csv` and `scan_csv`?
- 2 What is the difference between `select` and `with_columns`?
- 3 Why do we use `adjusted_close` instead of `close` for returns?
- 4 What does `.pct_change().over("symbol")` do, and why is `.over()` needed?
- 5 When do you use a `left` join vs an `inner` join?





Can you answer these in your own words?

- 1 What is the difference between `read_csv` and `scan_csv`?
- 2 What is the difference between `select` and `with_columns`?
- 3 Why do we use `adjusted_close` instead of `close` for returns?
- 4 What does `.pct_change().over("symbol")` do, and why is `.over()` needed?
- 5 When do you use a `left` join vs an `inner` join?
- 6 In `plotnine`, what do `aes()` and `geom_line()` each do?



- Next week: **panel regression analysis** with financial datasets

If something breaks: read the error message, check your data types, try again.  
Still stuck? Email me at [ce50@st-andrews.ac.uk](mailto:ce50@st-andrews.ac.uk).



- Next week: **panel regression analysis** with financial datasets
- **Homework:** complete today's exercises — download stock data for your own tickers and explore the returns

If something breaks: read the error message, check your data types, try again.  
Still stuck? Email me at [ce50@st-andrews.ac.uk](mailto:ce50@st-andrews.ac.uk).



- Next week: **panel regression analysis** with financial datasets
- **Homework:** complete today's exercises — download stock data for your own tickers and explore the returns
- Resources:

If something breaks: read the error message, check your data types, try again.  
Still stuck? Email me at [ce50@st-andrews.ac.uk](mailto:ce50@st-andrews.ac.uk).



- Next week: **panel regression analysis** with financial datasets
- **Homework:** complete today's exercises — download stock data for your own tickers and explore the returns
- Resources:
  - Polars user guide: <https://docs.pola.rs/>

If something breaks: read the error message, check your data types, try again.  
Still stuck? Email me at [ce50@st-andrews.ac.uk](mailto:ce50@st-andrews.ac.uk).



- Next week: **panel regression analysis** with financial datasets
- **Homework:** complete today's exercises — download stock data for your own tickers and explore the returns
- Resources:
  - Polars user guide: <https://docs.pola.rs/>
  - Polars I/O guide: <https://docs.pola.rs/user-guide/io/>

If something breaks: read the error message, check your data types, try again.  
Still stuck? Email me at [ce50@st-andrews.ac.uk](mailto:ce50@st-andrews.ac.uk).



- Next week: **panel regression analysis** with financial datasets
- **Homework:** complete today's exercises — download stock data for your own tickers and explore the returns
- Resources:
  - Polars user guide: <https://docs.pola.rs/>
  - Polars I/O guide: <https://docs.pola.rs/user-guide/io/>
  - Tidy Finance with Python: <https://www.tidy-finance.org/python/>

If something breaks: read the error message, check your data types, try again.  
Still stuck? Email me at [ce50@st-andrews.ac.uk](mailto:ce50@st-andrews.ac.uk).



- Next week: **panel regression analysis** with financial datasets
- **Homework:** complete today's exercises — download stock data for your own tickers and explore the returns
- Resources:
  - Polars user guide: <https://docs.pola.rs/>
  - Polars I/O guide: <https://docs.pola.rs/user-guide/io/>
  - Tidy Finance with Python: <https://www.tidy-finance.org/python/>
  - plotnine gallery: <https://plotnine.org/>

If something breaks: read the error message, check your data types, try again.  
Still stuck? Email me at [ce50@st-andrews.ac.uk](mailto:ce50@st-andrews.ac.uk).

