# Data Curation & Feature Engineering

Session 2

# Content

- ➔ The Data Challenge in Investment Research
- ➔ Introduction to Data Curator
- ➔ Outlier Detection & Correction
- ➔ Feature Engineering
- ➔ Hands-On: Your First Data Pipeline

# The Data Challenge in Investment Research

Bottlenecks

# The Data Challenge

# 60–80%

of a quant researcher's time may be wasted on repetitive data wrangling tasks.

## Reproducibility Crisis

- Undocumented workflows prevent reliable validation of results.

## Scattered Sources

- Inconsistent provider formats require heavy normalization.

## Manual Errors

- Human intervention introduces silent data failures.

## Lost Velocity

- Time spent fixing pipelines reduces research productivity.
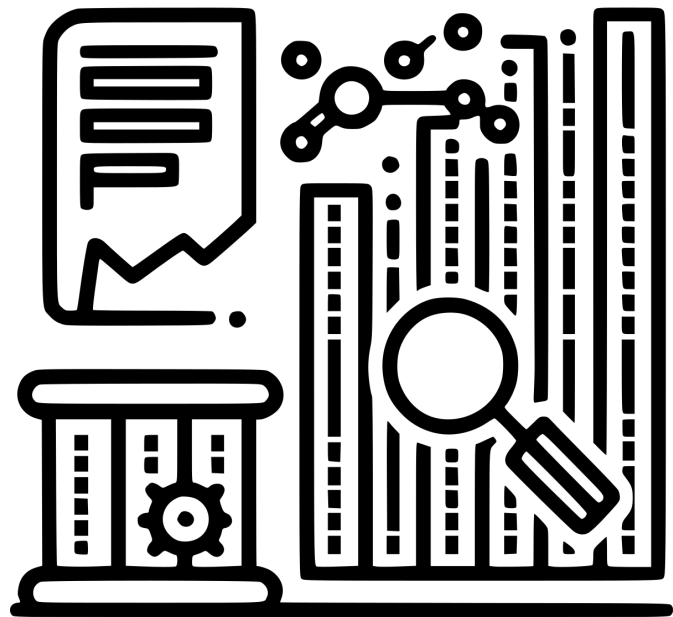
# The Messy Middle Problem

**Major bottlenecks**:

- Data quality
- Pipeline consistency
- Governance gaps
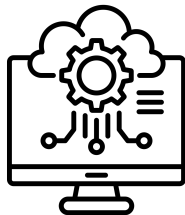
**In quant investing, this looks like**:

- 40 Jupyter notebooks
- 15 slightly different datasets
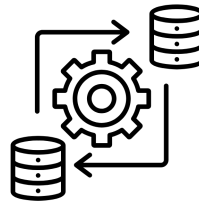- No reproducibility
- No audit trail

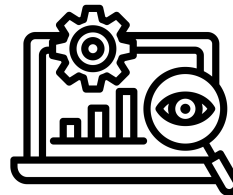Structured, versioned, reproducible data systems.

# The Tooling To Speed Up Feature Engineering

Data
Curator

Data
Weaver

Data
Analyzer

# Introduction to Data Curator

The Solution to Improve Data Wrangling Tasks

# Data Curator

A Python-based tool designed to **build structured databases** for market, fundamental, and alternative data, **transforming chaotic wrangling into a systematic workflow**.

| **Unified Data Access** | **Vendor Data Verification** | **Prebuilt Calculations** | **Custom Features** |
|---|---|---|---|
| Single interface for various data providers and custom in-house datasets, with homogenized tags. | Assess whether the provider's data is reliable or affected by serious issues (e.g., negative values, restated data). | Ready-to-use functions for outlier adjustments, ratios, momentum, trend, volatility, volume, etc. | Easily define your own custom feature functions to extend the capabilities of the Data Curator. |

**Excel-Based Configuration**: Specify tickers, dates, and features in a spreadsheet, **no complex coding required for setup**.

https://kaxanuk-data-curator.readthedocs.io
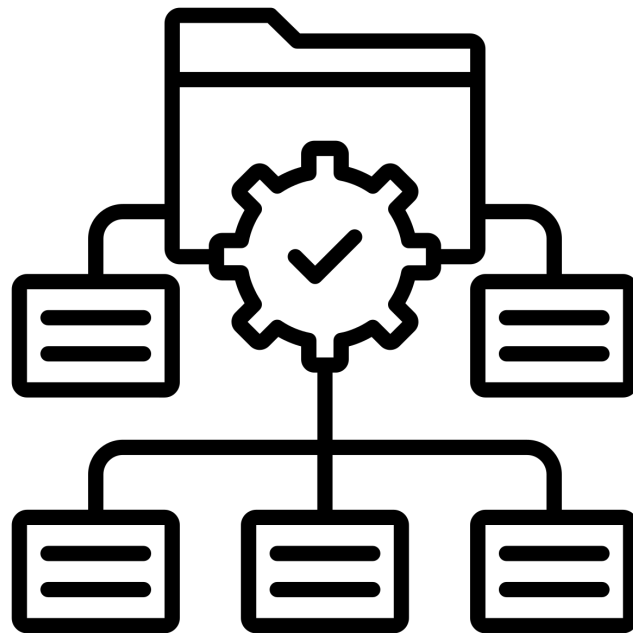
# Structured Financial Data

**Data is messy**:

- Missing earnings dates
- Inconsistent ticker mapping
- Outliers
- Corporate actions
- Different frequencies

**Structured data means**:

- Aligned timestamps
- Point-in-time data
- Cleaned values
- Standardized schema
- Version-controlled datasets

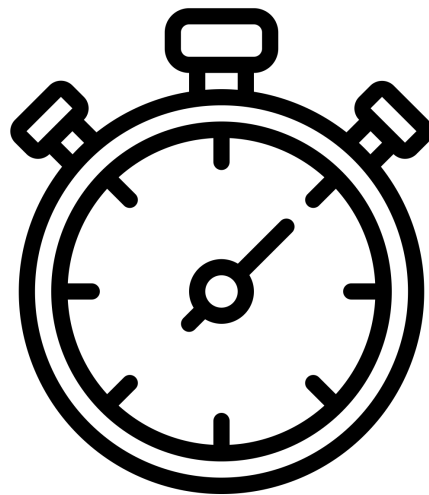# Time Matters: Avoiding Look-Ahead Bias

Every data point must respect:

- ❑ **Release date**
- ❑ **Availability date**
- ❑ **Revision history**

Example:

- Q1 earnings released on May 5
- You cannot use them in March backtests.

Feature Engineering must be:

**Chronologically honest.**

# Getting Structured Financial Data

## The 3-Step Workflow

### 01 Install

Set up the environment. Requires Python 3.12+.

**pip install kaxanuk.data_curator**

# Optional extension

**pip install kaxanuk.data_curator_extensions. yahoo_finance**

### 02 Configure

Initialize and edit settings in Excel.

**kaxanuk.data_curator init excel**

Config/parameters_datacurator.xlsx
Define Tickers, Dates, Providers

Config/.env Store API Keys securely

### 03 Execute

Run the pipeline to fetch and process data.

**kaxanuk.data_curator run**

Output/ Structured CSV/Parquet files ready for analysis

https://kaxanuk-data-curator.readthedocs.io/en/stable/user_guide/quick_start.html

# Before

**Fragmented infrastructure**:

- Multiple APIs with heterogeneous formats
- Inconsistent field naming across providers
- Missing or partially available data
- Manual joins and ad-hoc transformations
- Point-in-time misalignment
- Non-reproducible research notebooks
- Difficult feature sharing across teams
- High debugging overhead
- …

# Unified Time Series Output Across All Tools



Data Curator

# Outlier Detection & Correction

Ensuring Data Integrity

# Why Outliers Destroy Signals

**Outliers may come from**:

- Data entry errors
- API glitches
- Stock splits not adjusted
- Extreme corporate events

**Consequences**:

- Inflated z-scores
- False momentum signals
- Volatility distortion

# Data Entry Errors

Data
Analyzer

fis_weighted_average_basic_shares_outstanding

# Methods for Outlier Detection

**Some ideas we can test**:

- Rolling Z-Score Thresholds
- Median Absolute Deviation (MAD)
- Volatility-Scaled Return Filters
- Abnormal price gaps vs rolling volatility
- Structural break detection
- Stale value detection
- Near-zero denominator flags
- Corporate action validation (splits, dividends)
- Ratio sanity checks (negative or impossible values)
- Extreme growth without event justification
- Accounting restatement detection

# Methods for Outlier Correction

**We can implement**:

- Z-score filtering
- Rolling median deviation
- Winsorization
- Robust scaling
- Percentile clipping

Important:

**Outlier correction must be rule-based, reproducible, and grounded in financial logic (e.g., corporate actions and fundamental events).**

# Correction vs Removal

**Never "manually delete".**

**Choose**:

- Winsorize
- Cap values
- Replace with rolling median
- Flag for exclusion
- …

The pipeline must explain itself.

**Auditability > Convenience**.

# Outlier Detection and Correction

Data Curator enables **custom calculations** to diagnose, validate, and transform data into research-ready inputs.

**We define explicit rules.**

Those rules:

- Detect anomalies
- Correct distortions
- Standardize scaling
- Preserve chronology
- Remain reproducible

Outlier handling becomes part of the research system, **not a notebook hack**.

https://kaxanuk-data-curator.readthedocs.io/en/stable/api_reference/custom_calcs.html

## Custom Calculations

```python
def c_test(m_open, m_close):
    """
    Example features calculation function.

    Receives the market open and market close columns, and returns a column with their difference.

    For this function to generate an output column, you need to:
    1. Place it in the Config/custom_calculations.py file (if it doesn't exist you can copy this file there).
    2. Add c_test to the Output_Columns sheet in the Config/parameters_datacurator.xlsx file.

    Parameters
    ----------
    m_open : kaxanuk.data_curator.DataColumn
    m_close : kaxanuk.data_curator.DataColumn

    Returns
    -------
    kaxanuk.data_curator.DataColumn
    """
    # we're just doing a subtraction here, but you can implement any logic
    # just remember to return the same number of rows in a single column!
    return m_close - m_open
```

# Some Examples

## Outlier Adjusted Data

- **Price Spike Filter** — Flag returns exceeding a multiple of rolling volatility and cap or replace when no corporate action justifies the move.
- **Corporate Action Adjustment** — Detect split- or dividend-driven price jumps and adjust the historical series to preserve economic continuity.
- **Ratio Explosion Control** — Nullify or redesign ratios when denominators approach zero to prevent artificial signal distortions.
- **Abnormal Growth Check** — Identify extreme revenue or earnings changes without event support and winsorize or smooth accordingly.
- **Extreme Z-Score Cap** — Bound cross-sectional z-scores at predefined limits to stabilize rankings and reduce signal instability.
- **Missing Data Guardrail** — Apply exclusion or forward-fill only under predefined, financially justified rules.
- **Negative or Impossible Values Filter** — Flag and correct financially inconsistent observations (e.g., negative shares outstanding).
- **Stale Price Detection** — Detect abnormal price stagnation and classify as illiquidity or potential data error.

# Feature Engineering

From Data to Alpha Signals and Risk Factors

A **feature is a structured variable** derived from data that captures information relevant for **prediction** or **risk modeling**.

# Hypothesis-Driven Feature Engineering

**Avoid**: Data mining 200 features and picking the best.

**Instead**: Start from hypothesis.

Example

"**Markets underreact to earnings surprises.**"

**Then engineer**:

- Standardized earnings surprise
- Post-event drift
- Volume confirmation

**Feature engineering follows theory**.

# Types of Features

## Univariate (1×N) Features

- Derived from each asset's own history.
- Examples: Trend, rolling volatility, growth rates, historical z-scores.

Data Curator

## Multivariate (N×N) Features

- Relative measures across assets at a given time.
- Examples: Sector-neutral ranks, cross-sectional z-scores, relative strength.

Data Weaver

# Feature Analysis



Data
Analyzer

# Features Can Be Used as Alpha Signals or Risk Factors

## Alpha Signals

## Risk Factors

**Idiosyncratic predictive features** designed to capture pricing inefficiencies and generate **excess returns**.

**Systematic drivers of returns** that are common across assets and **explain risk premia**.

# Some Examples

## Alpha Signals

- **Earnings Surprise** — Abnormal returns following unexpected earnings results.
- **Short-Term Reversal** — Mean reversion after extreme short-term moves.
- **Residual Momentum** — Trend persistence unexplained by major risk factors.
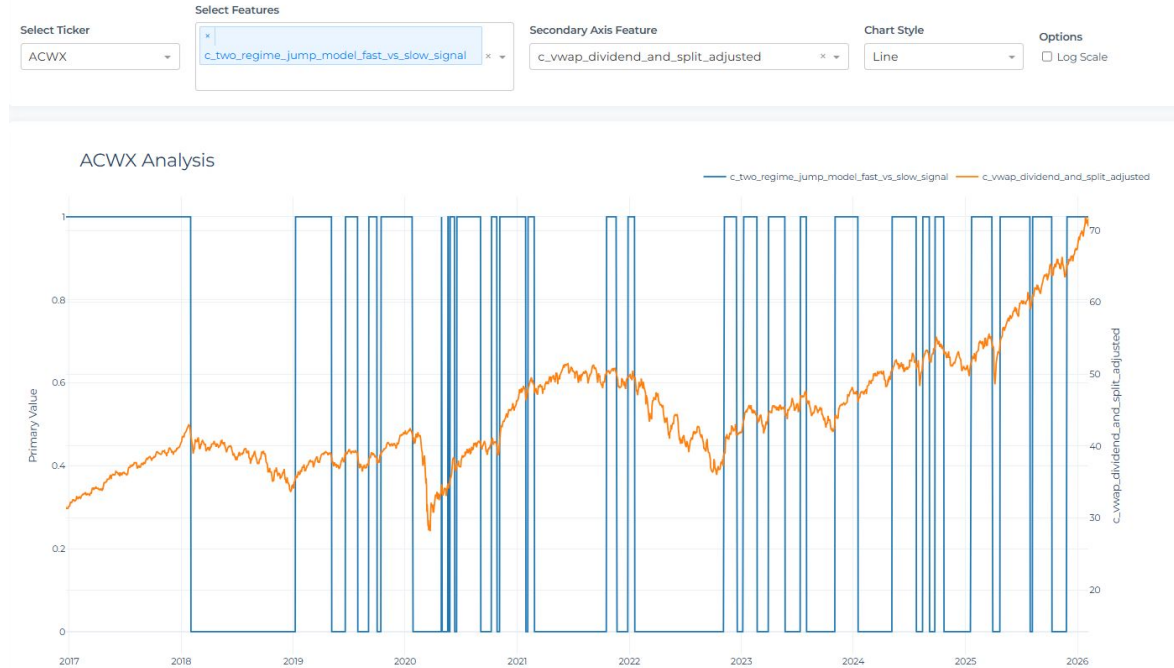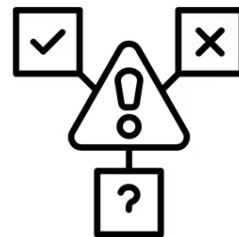- **Quality** — Excess returns linked to profitability and balance sheet strength.
- **Accruals** — Mispricing related to low earnings quality.
- **Sentiment** — Predictive signals derived from news or textual tone.
- **Dispersion** — Return opportunities arising from cross-sectional disagreement.
- **Event-Driven** — Temporary mispricing around corporate events (M&A, buybacks, guidance).

## Risk Factors

- **Market** — Exposure to overall asset class return movements.
- **Sectors** — Exposure to industry-specific economic dynamics.
- **Beta** — Sensitivity of an asset's returns to the market benchmark.
- **Size** — Exposure linked to company market capitalization.
- **Value** — Exposure to relatively undervalued companies.
- **Momentum** — Exposure to persistent price trends.
- **Volatility** — Exposure to return variability and risk level.
- **Growth** — Exposure to companies with strong fundamental expansion.

# Let's code!

Hands-On: Your First Data Pipeline

# Disclaimers

The content of this document is strictly informative and does not constitute an offer or recommendation of KaxaNuk S.C. to buy, sell or subscribe any kind of securities, or to perform specific transactions. KaxaNuk S.C. is not responsible for the interpretation given to the information and /or content of this document. KaxaNuk S.C. does not accept or will accept any liability for losses or damages resulting from investment decisions that would have been based on this document. The persons responsible for the preparation of this content certify that the opinions stated reflect their own point of view and do not represent the view of KaxaNuk S.C. nor of its officials. This document is based on publicly available information which are considered reliable, however KaxaNuk S.C. makes no warranty regarding its accuracy or completeness.

# Contact

[research@kaxanuk.mx](mailto:research@kaxanuk.mx)