

Python-Based Data Profilers

Here's a high-level overview of popular open-source data profilers available in the Python ecosystem, each designed for different use cases such as exploratory data analysis (EDA), data validation, unstructured data profiling, and ML monitoring.

1. [DataProfiler](#) (CapitalOne)

- **General-purpose, full-featured profiler**
- Developed by CapitalOne, **DataProfiler** supports structured and unstructured data including text and images. It automatically detects data types, computes statistics, identifies PII (Personally Identifiable Information), and detects data drift over time. It's ideal for use cases requiring comprehensive profiling of both tabular and unstructured datasets.

2. [Datamart-Profiler](#) (ViDA NYU)

- **Lightweight profiler focused on metadata generation**
- Part of the Auctus/Datamart ecosystem by NYU ViDA, this profiler detects data types, extracts statistical summaries, and metadata from structured data. It is lightweight and modular, making it suitable for integration into data discovery pipelines, especially for semantic search and dataset indexing.

3. [YData Profiling](#) (formerly Pandas Profiling)

- **Quick and automatic EDA for tabular data**
- YData Profiling is a go-to tool for rapid, visual profiling of pandas DataFrames. It generates an HTML report with statistics, correlations, missing value handling, and duplicate detection. While it doesn't handle unstructured data, it's highly effective for fast diagnostics of structured tabular datasets.

4. [WhyLogs / WhyLabs](#)

- **Statistical logging and monitoring for ML data**

- WhyLogs is a Python-native library for statistical data profiling, logging, and drift detection. Built for MLOps, it supports structured and unstructured data, and integrates with the WhyLabs platform for long-term monitoring. It's best suited for teams tracking model/data quality over time.

5. [Deequ / PyDeequ](#)

- **Data quality library for large-scale Spark pipelines**
- Developed by Amazon, Deequ is a Scala-based data quality tool, with PyDeequ offering a Python wrapper. It's tailored for **batch data validation on large datasets** in Spark, allowing custom constraints and profiling but lacks native support for unstructured data.

6. [OpenRefine](#) (with Python integration)

- **Interactive data wrangling and exploration tool**
- While not a profiler per se, OpenRefine allows semi-structured data cleaning and transformation through a GUI. It can be used alongside Python scripts via APIs, and is well-suited for **manual profiling, exploratory analysis, and data cleanup** tasks.

Python-Based Data Profilers – Feature Comparison

This table compares several Python-based data profilers based on core functionality for structured data.

Tool	Python Support	Open Source	Structured Data Summary Stats	Missing Values	Correlation	Duplicate Detection	Type Detection
DataProfiler	✓	✓ Apache 2	✓	✓	✓	✓	✓ (incl. unstructured)
Datamart-Profiler	✓	✓ MIT	✓	✓	✗	⚠ Basic support	✓ (custom text types)
YData Profiling	✓	✓ MIT	✓	✓	✓	✓	✓
WhyLogs / WhyLabs	✓	✓ Apache 2	✓ (rolling stats, sketches)	✓	⚠ Approximate	⚠ Approximate	✓ (typed features)
PyDeequ	✓ (via Py4J)	✓ Apache 2	✓ (through constraints)	✓	⚠ via checks	⚠ via constraints	✓ (schema-based)

Pandas `df.dtypes`

`df.dtypes` is not a profiler. It is a **Pandas attribute** that returns the **data type of each column** in a DataFrame. It's useful for quickly checking the schema of the dataset. Here's a recap of the **main inferred types** we'll encounter in `df.dtypes`:

- Numeric: `int64`, `float64`
- Text: `object`, `string`
- Date/Time: `datetime64[ns]`, `timedelta[ns]`
- Boolean: `bool`
- Categorical: `category`

Note: Pandas does not automatically infer `geometry`, `GeoJSON`, or WKT types — you must convert those explicitly (e.g., using `shapely` or `GeoPandas`).

Pandas `df.dtypes` vs. YData Profiling (formerly pandas-profiling)

Feature	<code>df.dtypes</code> (Pandas)	YData Profiling
Purpose	Shows column data types	Generates automated, detailed data reports
Scope	Only types	Types + statistics, distributions, correlations
Output	Simple text	Interactive HTML report or JSON
Ease of Insight	Manual inspection required	Immediate visual summary of data quality and content
Detects anomalies?	✗ No	✓ Yes (missing values, duplicates, type mismatches)
Visuals	✗ None	✓ Histograms, bar charts, heatmaps, etc.
Correlation checks?	✗ No	✓ Yes (Pearson, Spearman, Kendall)
Duplicate rows check?	✗ No	✓ Yes
Type detection & override?	Manual via <code>df.dtypes</code> and conversions	Automatic suggestion with override options

Installation needed?	✗ Built into Pandas	✓ <code>pip install ydata-profiling</code>
----------------------	---------------------	--

Strong Candidates

Feature	YData Profiling	DataProfiler	Datamart-Profiler
GUI HTML Reports	✓ Rich visuals	✗ Console/text output	✗ No built-in GUI; designed for backend metadata workflows. UI through the DataProfileVis tool.
Unstructured data support	✗ No	✓ Yes (text, logs, NLP)	✓ Limited
Auto ML Labeling	✗ No	✓ Optional (NER, classification)	✗ No
Extensibility	Medium	High (modular profiling engine)	Medium (custom types via lightweight plugin interface)
Data Type Inference	✓ Great	✓ Also great, more robust for mixed types	✓ Good for custom schema

Use Case Recommendation

- Use **datamart-profiler** for lightweight profiling focused on **metadata generation**. Ideal for scenarios with a **high emphasis on integration with other systems**.
- Use **YData Profiling** for **quick exploratory visualizations** and summary statistics.
- Use **DataProfiler** for **production-scale profiling**, especially with **mixed-type detection** and **non-tabular data** (e.g., logs, free text).