

Analytics Validation & Safe Insight Generation System

Yugaash Sridhar

Department of Computer Science and Engineering
IIITDM Kanchipuram, India

Abstract

In data-driven environments, analytics pipelines frequently operate directly on raw datasets without sufficient validation of data quality and structure. This practice often leads to misleading insights and unreliable decision-making. This paper presents a data-quality-first analytics validation system that evaluates the structural reliability of structured datasets before generating insights. The proposed approach introduces layered validation, interpretable health scoring, and safe, rule-based insight generation to prevent incorrect analytics. The system is domain-agnostic, scalable across diverse datasets, and designed using conservative validation logic.

Index Terms

Data Quality, Analytics Validation, Data Profiling, Health Score, Safe Analytics

I. INTRODUCTION

Analytics systems play a crucial role in modern decision-making across industries. However, most analytics tools assume that input datasets are analytically valid and directly generate summaries and visualizations. In real-world scenarios, datasets often contain identifiers, ambiguous numeric fields, unconverted date attributes, missing values, and constant columns. These structural issues can significantly distort analytical outcomes.

This work proposes a validation-first analytics framework that prioritizes correctness over aggressive insight generation.

II. PROPOSED METHODOLOGY

A. Data Ingestion

The system accepts structured CSV datasets without enforcing domain-specific assumptions.

B. Data Profiling

Profiling computes dataset structure, schema, and missing values.

C. Data Quality Validation

Rule-based checks detect identifiers, ambiguous identifiers, date-like columns, and constant attributes.

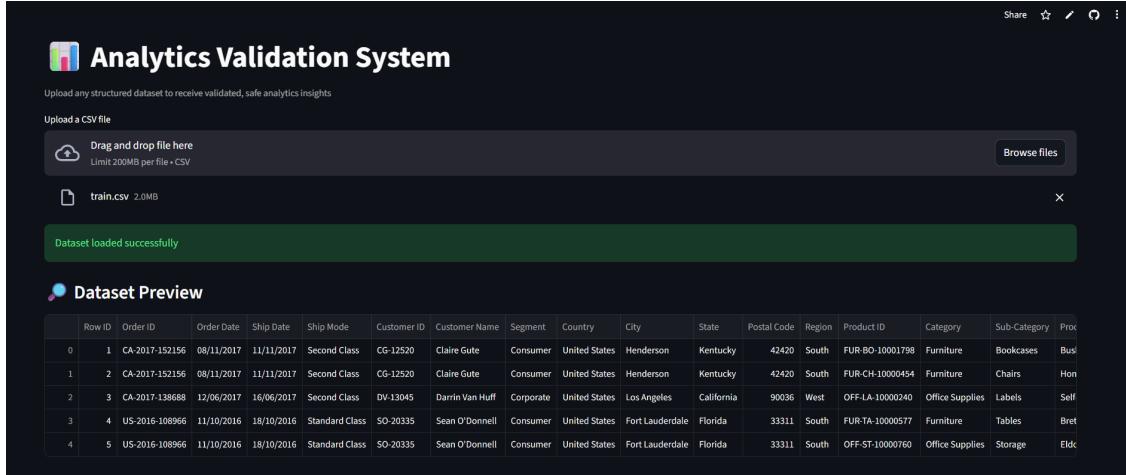


Fig. 1. Analytics Validation System – CSV upload interface

```

{
    "Row ID": "int64",
    "Order ID": "object",
    "Order Date": "object",
    "Ship Date": "object",
    "Ship Mode": "object",
    "Customer ID": "object",
    "Customer Name": "object",
    "Segment": "object",
    "Country": "object",
    "City": "object",
    "State": "object",
    "Postal Code": "float64",
    "Region": "object",
    "Product ID": "object",
    "Category": "object",
    "Sub-Category": "object",
    "Product Name": "object",
    "Sales": "float64"
}

```

Fig. 2. Dataset preview and automatically detected schema for train.csv

D. Data Health Scoring

A Data Health Score (0–100) quantifies dataset reliability.

E. Safe Insight Generation

Insights are generated only from validated attributes.

III. VISUALIZATION STRATEGY

Only safe insights are visualized.

IV. RESULTS AND DISCUSSION

Trend analysis is intentionally blocked when date columns are not validated.

```

Data Profiling Summary

{
  "rows": 9800,
  "columns": 18,
  "missing_values": [...],
  "numeric_columns": [...],
  "categorical_columns": [
    0: "Order ID",
    1: "Order Date",
    2: "Ship Date",
    3: "Ship Mode",
    4: "Customer ID",
    5: "Customer Name",
    6: "Segment",
    7: "Country",
    8: "City",
    9: "State",
    10: "Region",
    11: "Product ID",
    12: "Category",
    13: "Sub-Category",
    14: "Product Name"
  ]
}

```

Fig. 3. Data profiling summary showing rows, columns, and missing values

```

Data Quality Checks

{
  "date_like_columns": [
    0: "Order Date",
    1: "Ship Date"
  ],
  "ambiguous_date_columns": [],
  "identifier_columns": [
    0: "Row ID"
  ],
  "ambiguous_identifier_columns": [
    0: "Postal Code"
  ],
  "constant_columns": [
    0: "Country"
  ],
  "usable_numeric_columns": [
    0: "Sales"
  ],
  "usable_categorical_columns": [
    0: "Order ID",
    1: "Ship Mode",
    2: "Customer ID",
    3: "Customer Name",
    4: "Segment",
    5: "City",
    6: "State",
    7: "Region",
    8: "Product ID",
    9: "Category",
    10: "Sub-Category",
    11: "Product Name"
  ]
}

```

Fig. 4. Rule-based data quality checks identifying unsafe analytical columns



Fig. 5. Computed Data Health Score (77/100) with detected risk factors

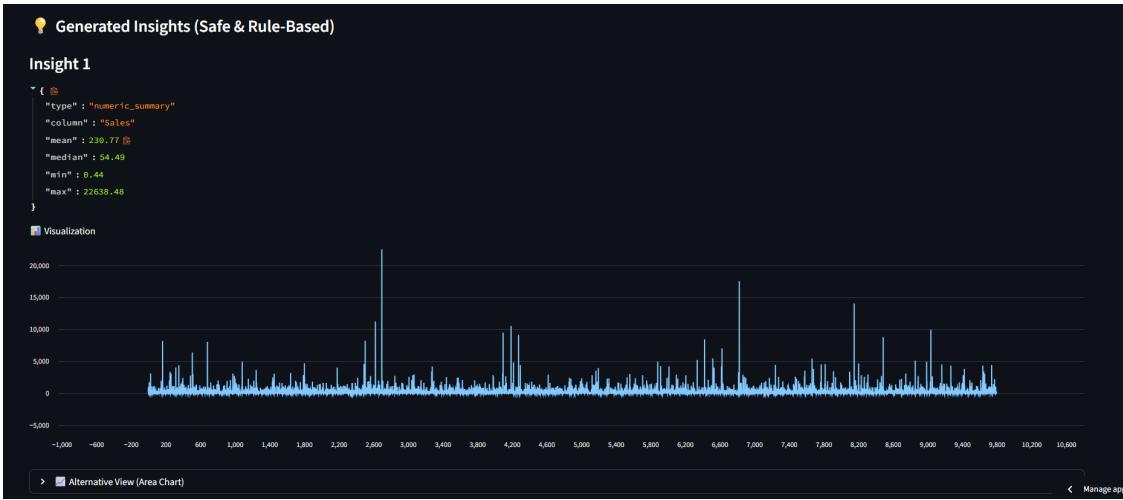


Fig. 6. Rule-based safe insights generated from validated columns

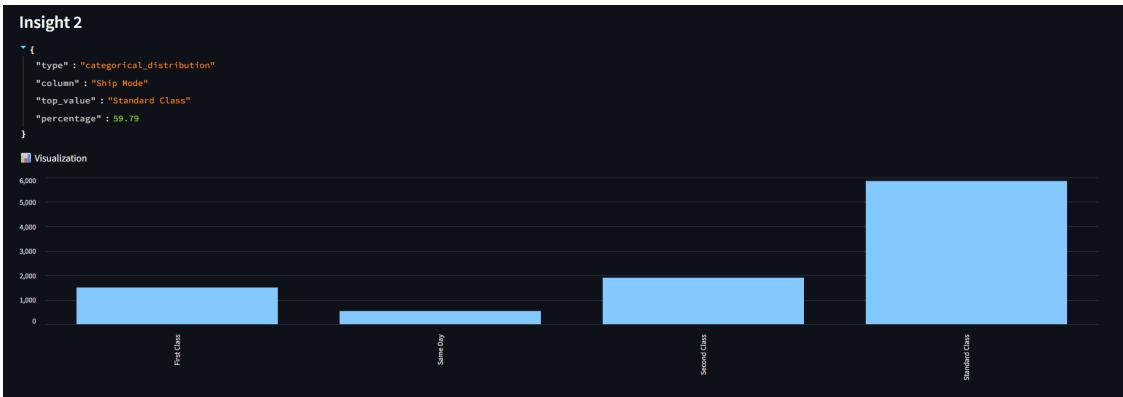


Fig. 7. Numeric sales insight visualized using line and area charts

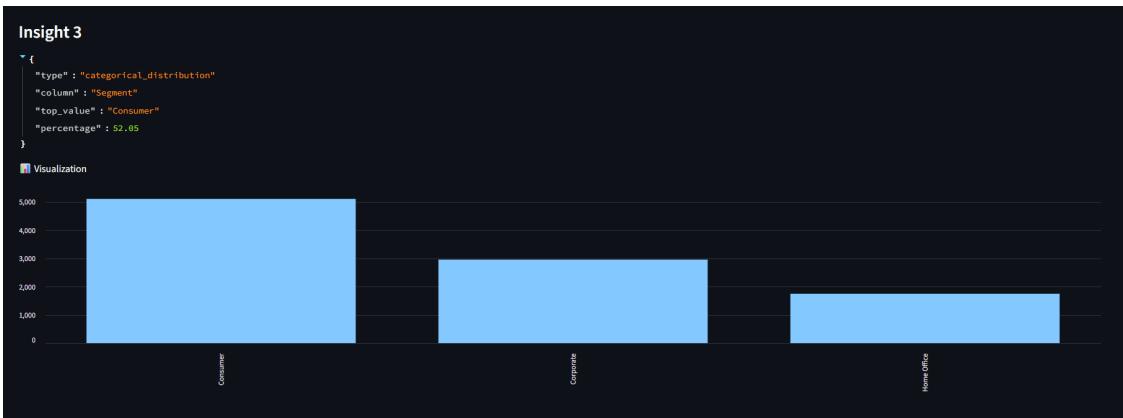


Fig. 8. Categorical distribution visualized using bar charts

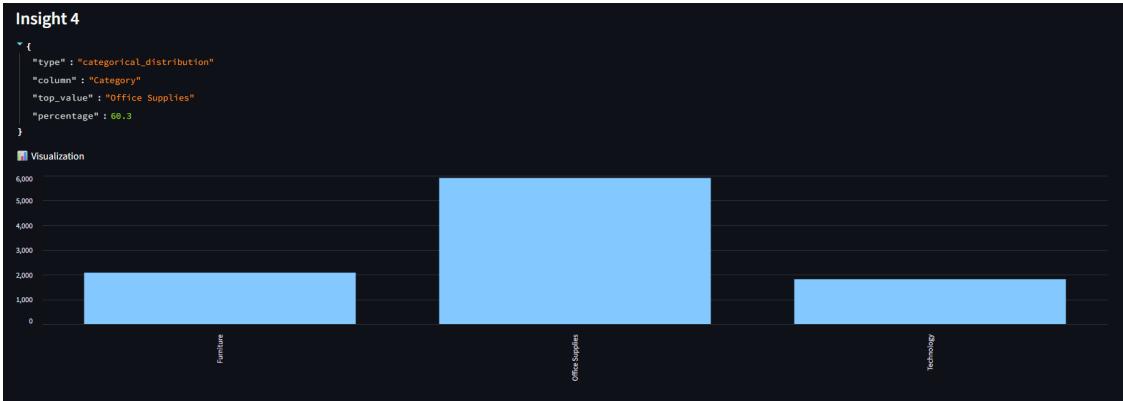


Fig. 9. Group comparison visualized using bar and radar charts

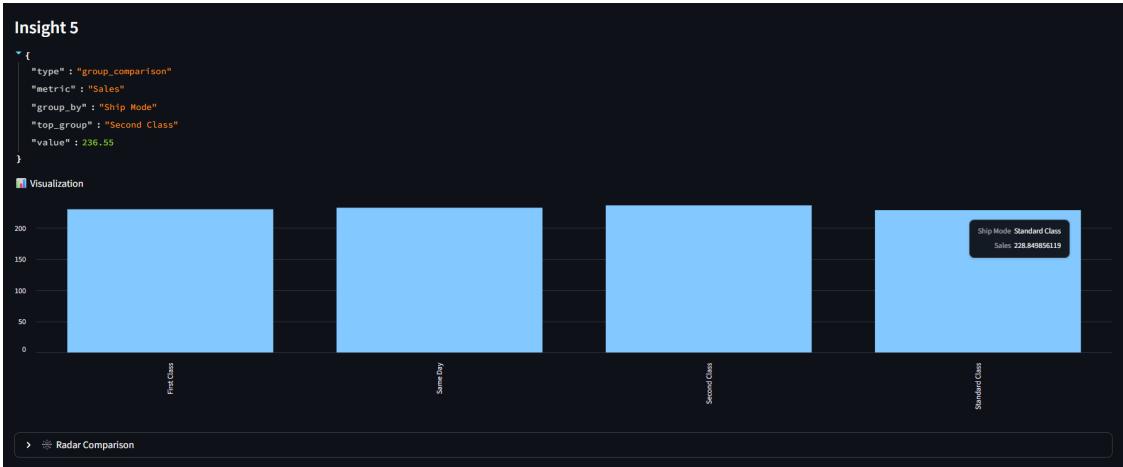


Fig. 10. Trend readiness insight indicating detected but unconverted date fields

The validation-first approach significantly improves trust in analytics outputs.

V. CONCLUSION

This paper presented a data-quality–first analytics validation system that prevents misleading analytics by validating datasets before insight generation. The approach improves correctness, transparency, and trustworthiness in data-driven decision-making.

REFERENCES

- [1] T. Redman, *Data Quality: The Accuracy Dimension*. Morgan Kaufmann, 2001.
- [2] ISO/IEC 25012:2008, *Software Engineering—Software Product Quality Requirements and Evaluation (SQuaRE)*.
- [3] A. Abedjan et al., “Detecting Data Errors: Where Are We and What Needs to Be Done?” *VLDB Endowment*, 2016.