

Verifying and Validating Data Sources

Edward Apostol, December 2023

In real-world scenarios, ensuring that data is accurate, reliable, and suitable for analysis and visualization involves several key standards, processes, and frameworks. These steps are critical to validate data as a trusted **source of truth**, ensuring that any insights or conclusions drawn from it are credible and actionable. Here's an overview of the key standards and processes used in industry:

1. Data Governance Frameworks

Data governance refers to the collection of practices, policies, and standards to manage the availability, usability, integrity, and security of data across an organization.

- **ISO/IEC 38505:** Part of the broader ISO 38500 standard on corporate IT governance, it establishes best practices for data governance in organizations.
- **DAMA-DMBOK (Data Management Body of Knowledge):** This framework outlines key data management functions, such as data quality, data architecture, data governance, and more. It's widely used as a guide for managing data as a key organizational asset.
- **GDPR (General Data Protection Regulation):** While focused on privacy and data protection, GDPR also sets standards for data accuracy and transparency, requiring organizations to ensure that personal data is kept accurate and up to date.
- **CMMI (Capability Maturity Model Integration):** A process improvement framework that helps organizations evaluate and improve their data management practices, from initial collection to processing and usage.

2. Data Quality Assessment and Validation Processes

Data quality is one of the most critical aspects of ensuring data can be trusted as a source of truth. Several standards and processes focus on validating and assessing data quality.

- **ISO 8000:** This standard provides guidelines on how to ensure data quality, focusing on elements like data accuracy, completeness, consistency, and timeliness.
- **Data Profiling:** The process of examining and understanding data at a granular level. Profiling tools (such as Talend, Informatica) help identify outliers, missing values, patterns,

and other quality issues.

- **Data Quality Dimensions:**
 - **Accuracy:** The degree to which data reflects the real-world objects it represents.
 - **Completeness:** Ensuring all required data is present.
 - **Consistency:** Data should be consistent across various systems and formats.
 - **Timeliness:** Data should be current and available when needed.
 - **Validity:** Data conforms to predefined formats or business rules.
 - **Uniqueness:** No duplicate records or entities in the dataset.

Organizations often deploy **data quality tools** (such as IBM InfoSphere QualityStage, Talend Data Quality) to automate this process.

3. Data Lineage and Traceability

Tracking the origin and transformation of data is crucial for establishing trust in the data source. **Data lineage** allows organizations to trace the flow of data from its source to its final destination (and all the steps in between).

- **Metadata Management:** Systems like **Collibra**, **Alation**, and **Apache Atlas** enable organizations to track data lineage by maintaining a comprehensive metadata catalog. This ensures visibility into where data came from, how it has been transformed, and who has accessed or modified it.
- **ETL Validation (Extract, Transform, Load):** Ensuring that data transformations are accurate and that the ETL processes do not introduce errors. ETL tools such as **Talend** and **Informatica** often have built-in validation and reconciliation features to ensure data integrity through each transformation stage.

4. Data Certification Processes

Data certification is a process where data is verified by a trusted authority or process to ensure that it is of sufficient quality and integrity to be used as a single source of truth.

- **Data Stewards and Data Certification Programs:** Many organizations designate data stewards or teams responsible for verifying datasets before they are approved for use. These teams might certify data based on quality checks, timeliness, accuracy, and compliance with business rules.
- **Certified Data Repositories:** Some industries (e.g., healthcare, finance) rely on certified data repositories, which are repositories that have undergone rigorous quality validation processes. Examples include the use of **FHIR (Fast Healthcare Interoperability**

Resources) standards in healthcare data.

5. Data Validation Techniques

Data validation ensures that the data is correct and meaningful before it is used for analysis or visualization. Some key techniques include:

- **Schema Validation:** Ensuring that data conforms to the expected structure, such as column types, constraints, or data formats. Tools like **JSON Schema** or **Apache Avro** are used to validate that data matches predefined schemas.
- **Business Rule Validation:** Ensuring that data meets predefined business rules or logic (e.g., the total revenue cannot be negative). Many organizations use rules engines (such as **Drools** or **Talend Rules Management**) to enforce such validations.
- **Anomaly Detection:** Identifying data points that don't fit within expected patterns or ranges. Statistical models or machine learning-based anomaly detection techniques (using tools like **DataRobot** or **H2O.ai**) can flag outliers or errors in data sets.

6. Data Auditing and Reporting

Auditing and tracking data-related processes help establish accountability and transparency, ensuring data integrity throughout its lifecycle.

- **Internal Data Audits:** Many organizations conduct regular data audits to review how data is collected, stored, and processed. This includes checking for data integrity, adherence to security policies, and ensuring compliance with regulations.
- **Automated Audits and Reports:** Tools like **Splunk**, **Logstash**, and **DataDog** can automatically track data logs, providing an audit trail that ensures data hasn't been tampered with or corrupted.

7. Standardization through Industry Standards

Different industries have established standards to ensure data reliability and interoperability:

- **FHIR (Fast Healthcare Interoperability Resources):** A standard for exchanging healthcare information electronically, ensuring that healthcare data is interoperable and accurate.
- **XBRL (eXtensible Business Reporting Language):** Used for financial reporting, XBRL ensures that financial data is standardized and can be used consistently across organizations.
- **ISO 27001:** An information security management standard that outlines best practices for

securing data, including policies to ensure data integrity and prevent unauthorized access or tampering.

8. Data Visualization Standards

Visualization standards ensure that data is communicated clearly, accurately, and without misrepresentation.

- **Data Visualization Best Practices:**
 - **Choosing the Right Visualization:** Ensure the data is presented in the correct chart type (bar charts for comparisons, line charts for trends, etc.). This helps avoid misleading or inaccurate interpretations.
 - **Consistency in Scales and Axes:** Inconsistent scales or truncated axes can distort the viewer's understanding of the data.
 - **Adherence to Accessibility Standards:** Ensure data visualizations are accessible to all users, including those with visual impairments (e.g., using colorblind-friendly palettes, providing alternative text for charts).
- **Tools for Data Visualization Quality Control:** Visualization tools like **Tableau**, **Power BI**, and **D3.js** offer built-in guidelines and warnings for ensuring data visualizations follow best practices and remain accurate.

9. Benchmarking and External Data Validation

Sometimes, internal data needs to be validated against external benchmarks or reference data to confirm its accuracy.

- **Benchmarking Against Industry Standards:** Many organizations compare their data with external benchmarks (e.g., industry performance metrics or government statistics) to validate their data's accuracy and relevance.
- **External Data Providers:** Some organizations acquire data from third-party providers (e.g., **Bloomberg**, **FactSet**, **Quandl**) which are known for having rigorous data validation and certification processes.

10. Regulatory Compliance

In some industries, compliance with regulatory standards for data handling is mandatory, and non-compliance can result in penalties or legal consequences.

- **Sarbanes-Oxley (SOX):** In the finance sector, SOX compliance mandates accurate

financial reporting, which includes verifying the integrity and validity of data.

- **HIPAA:** In healthcare, HIPAA regulations ensure that medical data is handled securely, and any patient information is accurate and protected.
- **Basel III:** In banking, data used for risk analysis must adhere to strict standards of validation and reporting as outlined by international banking regulations.

Summary of Key Practices:

- **Data Governance Frameworks** like ISO/IEC 38505, DAMA-DMBOK.
- **Data Quality Standards** using ISO 8000 and tools for data profiling.
- **Data Lineage and Traceability** through metadata management systems.
- **Data Certification** processes with data stewards.
- **Validation Techniques** such as schema validation and anomaly detection.
- **Auditing and Reporting** to ensure data integrity.
- **Industry-Specific Standards** like FHIR and XBRL.
- **Visualization Standards** to ensure clarity and accuracy in presentation.

These processes ensure that data is not only correct but also reliable, traceable, and meaningful, making it a valid source for analysis and visualization.