



(REVIEW ARTICLE)



## Data Quality as Code: Automating validation rules with declarative pipelines and CI/CD Integration

Raveendra Reddy Pasala \*

*Independent Researcher, Move.Inc, Enterprise Systems, Los Angeles, California, USA.*

International Journal of Science and Research Archive, 2024, 13(02), 4224-4233

Publication history: Received on 22 November 2024; revised on 26 December 2024 accepted on 29 December 2024

Article DOI: <https://doi.org/10.30574/ijrsra.2024.13.2.2665>

### Abstract

Data Quality as Code (DQaC) is an automated approach that embeds data validation rules into modern data pipelines, ensuring consistent and reliable data processing. By leveraging declarative pipelines and integrating with CI/CD frameworks, organizations can enforce quality checks at every stage of data transformation. This method enhances accuracy, consistency, and compliance, reducing risks associated with poor data quality. CI/CD tools such as GitHub Actions and Jenkins automate validation, preventing erroneous data from being deployed. Observability solutions like DataDog and Prometheus monitor data quality trends in real time. While DQaC improves scalability and governance, challenges include managing complex validation rules and minimizing performance overhead. Implementing DQaC enables organizations to establish a robust data quality framework, fostering trust in analytics and decision-making. As data ecosystems grow, AI-driven validation and real-time monitoring will further strengthen the role of DQaC in data governance.

Integrating DQaC within CI/CD pipelines ensures that data quality checks run automatically whenever a data pipeline is updated. This approach prevents the deployment of pipelines that introduce bad data into production.

**Keywords:** Data Quality; Validation Rules; Declarative Pipelines; CI/CD Integration; Automation; Governance

### 1. Introduction

Data quality is a fundamental aspect of modern data-driven organizations. Poor data quality can lead to incorrect analytics, flawed decision-making, and regulatory non-compliance. Traditional approaches to data quality management often involve manual processes, making them inefficient, error-prone, and difficult to scale. As data pipelines become more complex and dynamic, there is a need for an automated and integrated approach to ensure high-quality data.

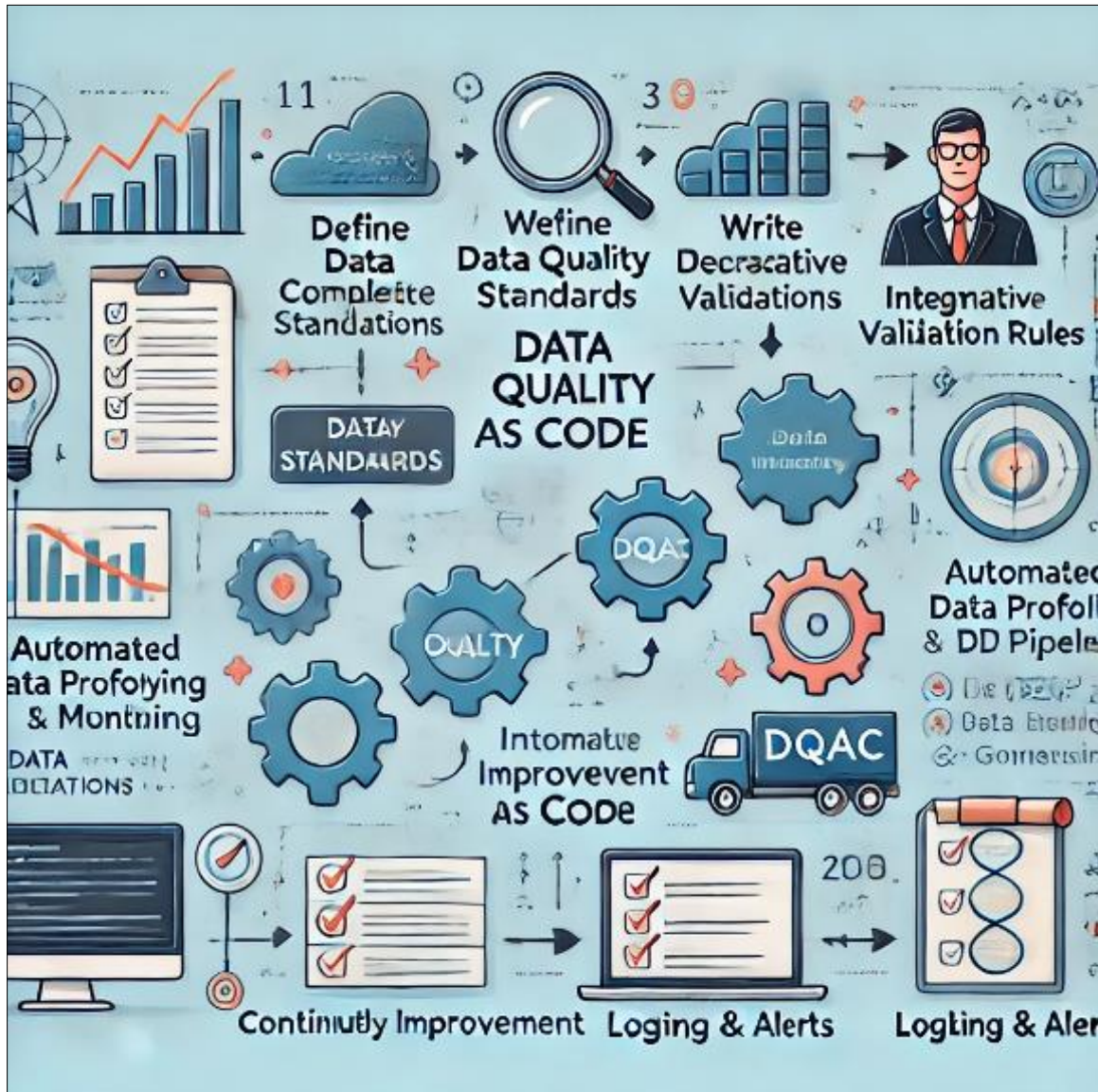
Data Quality as Code (DQaC) is a modern paradigm that treats data validation rules as code, enabling automated quality checks within data pipelines. By embedding validation rules into CI/CD processes, organizations can enforce data quality standards at every stage of the data lifecycle. This approach ensures that only clean and reliable data reaches analytical and operational systems.

The integration of DQaC with declarative pipelines allows teams to define validation criteria in a structured and reusable manner. This methodology improves consistency, reduces the risk of human errors, and enhances collaboration between data engineers, analysts, and quality assurance teams.

\* Corresponding author: Raveendra Reddy Pasala

This article explores the implementation of Data Quality as Code, highlighting its role in CI/CD integration, monitoring, and automation. It also discusses the key benefits, challenges, and best practices for adopting DQaC to improve data reliability and governance.

Integrating DQaC within CI/CD pipelines ensures that data quality checks run automatically whenever a data pipeline is updated. This approach prevents the deployment of pipelines that introduce bad data into production.



### 1.1. Implementing Data Quality Checks in CI/CD

CI/CD tools like GitHub Actions, Jenkins, or GitLab CI/CD can be configured to enforce data validation before merging changes. Implementing data quality checks in CI/CD ensures that every modification to a data pipeline undergoes rigorous validation, preventing bad data from entering production systems.

#### 1.1.1. Defining Data Validation Rules

Data validation rules define the expected structure, format, and integrity constraints of a dataset. These rules can be implemented using various frameworks such as Great Expectations, Deequ, or dbt tests.

For instance, in Great Expectations, a validation rule might check for:

- Completeness (no missing values in critical fields)

- Data type enforcement (e.g., ensuring a column contains only integer values)
- Range constraints (e.g., date fields falling within an expected range)

### 1.1.2. Integration with CI/CD Pipelines

To integrate data validation into a CI/CD pipeline, organizations need to define automated checks that trigger during different stages of the deployment lifecycle. Typically, this involves:

- Pre-Commit Hooks: Running lightweight checks locally before code is committed.
- Pre-Merge Validation: Ensuring data quality checks pass before merging changes into the main branch.
- Post-Deployment Validation: Running comprehensive tests after deployment to ensure continued data integrity.

### 1.1.3. Example: Data Quality Check in GitHub Actions

name: Data Quality CI

on: [push, pull\_request]

jobs:

validate:

runs-on: ubuntu-latest

steps:

- name: Checkout code

uses: actions/checkout@v3

- name: Install dependencies

run: pip install great\_expectations

- name: Run data validation

run: great\_expectations checkpoint run data\_quality\_check

This ensures that any commit to the repository triggers a data validation job, maintaining high data integrity.

### 1.1.4. Implementing Data Quality Checks in Jenkins

For organizations using Jenkins, data quality checks can be incorporated into a Jenkins pipeline as follows:

```
pipeline {
  agent any
  stages {
    stage('Checkout') {
      steps {
        git 'https://github.com/example/data-pipeline.git'
      }
    }
  }
}
```

```

stage ('Install Dependencies') {
  steps {
    sh 'pip install great_expectations'
  }
}

stage ('Run Data Quality Checks') {
  steps {
    sh 'great_expectations checkpoint run data_quality_check'
  }
}
}

```

This setup ensures that each stage of the pipeline enforces data quality before progressing further.

## 1.2. Automated Reporting and Alerts

Integrating data quality checks with CI/CD pipelines should also include robust reporting mechanisms to ensure visibility and proactive resolution of data quality issues. Automated reporting and alerting help teams stay informed about validation failures, anomalies, and trends in data quality. These mechanisms can be integrated into CI/CD workflows, allowing teams to respond quickly to potential data quality risks before they impact downstream applications.

### 1.2.1. Methods of Reporting and Alerting

There are several ways to implement automated reporting and alerting for data quality checks, including:

- **Email Notifications:** Sending automated emails whenever a data validation check fails, ensuring that stakeholders are promptly informed.
- **Slack/Teams Integration:** Using messaging platforms like Slack or Microsoft Teams to post real-time alerts in dedicated channels.
- **Incident Management Tools:** Connecting with tools like PagerDuty or Opsgenie to escalate critical data quality issues that require immediate attention.
- **Logging and Monitoring Dashboards:** Storing validation logs in centralized monitoring solutions like Prometheus, DataDog, or Splunk, enabling real-time tracking and historical analysis.
- **JIRA Ticketing System:** Automatically creating JIRA tickets for unresolved data quality failures, allowing teams to track and resolve issues systematically.

### 1.2.2. Example: Slack Alert for Data Quality Failures

A simple example of integrating Slack alerts into a CI/CD pipeline using GitHub Actions:

```
name: Data Quality Check
```

```
on: [push, pull_request]
```

```
jobs:
```

```
  validate:
```

runs-on: ubuntu-latest

steps:

- name: Checkout code

uses: actions/checkout@v3

- name: Install dependencies

run: pip install great\_expectations

- name: Run data validation

run: great\_expectations checkpoint run data\_quality\_check || echo "Validation failed" > failure.log

- name: Send Slack Alert

if: failure()

uses: rtCamp/action-slack-notify@v2

with:

webhook\_url: \${{ secrets.SLACK\_WEBHOOK\_URL }}

message: "Data quality validation failed. Check logs for details."

This workflow ensures that failed data quality checks trigger an immediate Slack notification, alerting the team in real-time.

### 1.2.3. Example: Prometheus and Grafana for Monitoring Data Quality Trends

To track data quality over time, teams can integrate their validation results with Prometheus and visualize trends using Grafana.

- **Store Validation Metrics:** Modify data quality tests to log validation results as Prometheus metrics.
- **Configure Prometheus Scraping:** Set up Prometheus to scrape these metrics periodically.
- **Visualize in Grafana:** Create dashboards displaying key data quality indicators such as schema drift, missing values, and validation pass rates.

This approach enables proactive monitoring and helps teams identify trends in data quality before issues escalate.

### 1.2.4. Benefits of Automated Reporting and Alerts

- **Faster Issue Resolution:** Automated alerts help teams quickly identify and resolve data quality issues before they affect production systems. Real-time notifications reduce the time spent manually diagnosing problems.
- **Improved Visibility and Transparency:** Automated reporting ensures that stakeholders across teams are informed about the current state of data quality. This fosters collaboration between data engineers, analysts, and business teams.
- **Reduced Operational Risks:** Proactive monitoring and alerts prevent poor-quality data from propagating through data pipelines, reducing the risk of erroneous insights and business disruptions.
- **Compliance and Auditability:** Automated reporting provides a detailed log of validation checks and data quality trends, helping organizations meet regulatory requirements and maintain data governance.
- **Scalability and Efficiency:** Automating alerts allows organizations to scale data quality monitoring across large datasets and complex pipelines without increasing manual workload.
- **Historical Analysis and Trend Detection:** Centralized logging and monitoring solutions enable teams to analyze long-term data quality trends, helping to identify recurring issues and improve validation rules over time.

- **Minimized Business Impact:** Early detection of data issues prevents incorrect data from affecting key business decisions, ensuring data reliability in analytical and operational workflows.
- **Customizable Alerting Mechanisms:** Organizations can tailor alerting mechanisms based on severity levels, ensuring that critical failures receive immediate attention while minor warnings are logged for review.
- **Seamless Integration with DevOps and DataOps Practices:** Automated reporting aligns with DevOps and DataOps principles, embedding data quality monitoring into continuous delivery pipelines for enhanced agility.
- **Reduced Maintenance Overhead:** Automation minimizes the need for manual monitoring and intervention, freeing up resources for more strategic data engineering initiatives.

#### 1.2.5. Example: Slack Alert for Data Quality Failures

A simple example of integrating Slack alerts into a CI/CD pipeline using GitHub Actions:

```
name: Data Quality Check
```

```
on: [push, pull_request]
```

```
jobs:
```

```
  validate:
```

```
    runs-on: ubuntu-latest
```

```
    steps:
```

```
      - name: Checkout code
```

```
        uses: actions/checkout@v3
```

```
      - name: Install dependencies
```

```
        run: pip install great_expectations
```

```
      - name: Run data validation
```

```
        run: great_expectations checkpoint run data_quality_check || echo "Validation failed" > failure.log
```

```
      - name: Send Slack Alert
```

```
        if: failure ()
```

```
        uses: rtCamp/action-slack-notify@v2
```

```
        with:
```

```
          webhook_url: ${{secrets.SLACK_WEBHOOK_URL}}
```

```
          message: "Data quality validation failed. Check logs for details."
```

This workflow ensures that failed data quality checks trigger an immediate Slack notification, alerting the team in real-time.

---

## 2. Observability and Monitoring

Observability and monitoring are critical components of Data Quality as Code, ensuring that data integrity is maintained in real-time and issues are promptly addressed. By implementing advanced monitoring solutions, organizations can gain deep insights into data pipeline health, detect anomalies, and take corrective actions proactively.

### 2.1.1. Key Components of Observability

- **Data Quality Metrics:** Establishing key performance indicators (KPIs) such as data completeness, accuracy, consistency, and timeliness to monitor the health of data pipelines.
- **Real-time Monitoring:** Leveraging tools like Prometheus, Grafana, and DataDog to visualize data quality trends and detect deviations in real-time.
- **Logging and Tracing:** Capturing detailed logs of data processing stages to track errors, failures, and inconsistencies, enabling root cause analysis.
- **Automated Anomaly Detection:** Using machine learning models to identify patterns and detect outliers that may indicate potential data quality issues.
- **Alerting and Incident Management:** Integrating monitoring tools with alerting platforms to notify teams when thresholds are breached, ensuring rapid response to issues.

### 2.1.2. Benefits of Observability in Data Quality

- **Proactive Issue Detection:** Continuous monitoring ensures that data anomalies are identified before they affect business decisions.
- **Enhanced Data Pipeline Reliability:** Observability helps maintain stable and efficient data pipelines by preventing data quality issues from escalating.
- **Faster Root Cause Analysis:** Centralized logging and tracing enable teams to pinpoint the source of data quality problems quickly.
- **Compliance and Governance:** Organizations can ensure compliance with industry regulations by maintaining thorough records of data quality metrics and validation checks.
- **Optimized Performance:** Monitoring resource utilization and pipeline efficiency helps optimize data processing workflows, reducing overhead costs.

By integrating observability and monitoring into Data Quality as Code, organizations can create a robust framework for maintaining high-quality data and minimizing disruptions in their data ecosystems.

### 2.1.3. Best Practices for Adopting Data Quality as Code (DQaC)

Adopting Data Quality as Code (DQaC) successfully requires a structured approach to defining, implementing, and maintaining data quality checks within data pipelines. Below are the key best practices organizations should follow to ensure a smooth and effective adoption of DQaC:

---

## 3. Define Clear Data Quality Standards

Before implementing DQaC, organizations must establish well-defined data quality standards. These should align with business goals and regulatory requirements.

- **Identify Key Data Quality Metrics:** Ensure validation rules address completeness, consistency, accuracy, validity, timeliness, and uniqueness.
- **Establish Governance Guidelines:** Define policies on how data should be structured, validated, and monitored.
- **Document Standards:** Maintain a central repository of data quality definitions for easy reference.

---

## 4. Use Declarative Validation Rules

Instead of writing custom scripts for every data validation check, leverage declarative frameworks that allow defining validation rules in a structured format.

- Frameworks like Great Expectations, Deequ, and dbt tests allow teams to define data expectations using YAML or JSON configurations.
- **Benefits of declarative validation:**
  - Standardizes quality checks across teams.
  - Enables easy updates and modifications.
  - Improves collaboration between data engineers, analysts, and governance teams.

## 5. Integrate with CI/CD Pipelines

Embedding data validation into Continuous Integration/Continuous Deployment (CI/CD) ensures that data quality is continuously enforced.

- Implement pre-commit hooks: Run lightweight validation checks before committing code.
- Pre-merge validations: Automate checks before merging data pipeline updates.
- Post-deployment tests: Run comprehensive validation to confirm data quality in production.

Example CI/CD Integration in GitHub Actions

yaml

CopyEdit

name: Data Quality Check

on: [push, pull\_request]

jobs:

validate:

runs-on: ubuntu-latest

steps:

- name: Checkout code

uses: actions/checkout@v3

- name: Install dependencies

run: pip install great\_expectations

- name: Run data validation

run: great\_expectations checkpoint run data\_quality\_check

---

## 6. Implement Robust Logging and Monitoring

Observability is key to ensuring real-time data quality management.

- Use monitoring tools like Prometheus, DataDog, or Grafana to track data quality trends.
- Enable automated logging and tracing for debugging failed validation checks.
- Set up real-time alerts through Slack, Microsoft Teams, or PagerDuty to notify teams of anomalies.

Example: Sending a Slack Alert for Validation Failures

yaml

CopyEdit

- name: Send Slack Alert

if: failure ()

uses: rtCamp/action-slack-notify@v2



with:

```
webhook_url: ${{secrets.SLACK_WEBHOOK_URL}}
```

```
message: "Data quality validation failed. Check logs for details."
```

---

## 7. Automate Data Profiling and Anomaly Detection

Automating data profiling helps in proactively identifying potential quality issues before they become critical.

- Profile datasets regularly using tools like pandas-profiling or DataProfiler.
  - Use AI-driven anomaly detection to spot trends and detect deviations.
  - Compare historical data trends to prevent schema drifts and missing values.
- 

## 8. Establish Data Lineage and Impact Analysis

Understanding how data flows across systems and the impact of changes is crucial.

- Track data lineage using tools like Apache Atlas or OpenLineage.
  - Perform impact analysis before implementing changes to ensure downstream data remains intact.
  - Document dependencies between datasets, transformations, and validation rules.
- 

## 9. Ensure Data Security and Compliance

Data validation should also enforce security and compliance to avoid regulatory risks.

- Mask sensitive data using encryption techniques.
  - Ensure compliance with GDPR, HIPAA, and industry regulations.
  - Audit all data validation checks to maintain transparency and accountability.
- 

## 10. Foster Collaboration Between Teams

Since data quality is a shared responsibility, fostering collaboration across teams is critical.

- Define clear roles and responsibilities between data engineers, analysts, and governance teams.
  - Use a central repository for validation rules to promote transparency.
  - Encourage feedback loops by allowing stakeholders to report data issues and suggest improvements.
- 

## 11. Continuously Improve Validation Rules

Data quality requirements evolve over time, so continuous improvement is necessary.

- Conduct regular reviews of validation rules and CI/CD pipelines.
  - Gather feedback from end-users to refine checks.
  - Incorporate new technologies such as AI-driven validation methods.
- 

## 12. Sort Small and Scale Gradually

Rather than implementing Dak across the entire data ecosystem at once:

- Begin with a pilot project in a specific data pipeline.

- Gradually expand based on lessons learned.
  - Standardize best practices as adoption scales.
- 

### **13. Conclusion**

Data Quality as Code transforms data validation from a manual, ad-hoc process into an automated, scalable, and repeatable practice. By leveraging declarative pipelines, CI/CD integration, and observability tools, organizations can ensure their data remains accurate and reliable. As data ecosystems grow in complexity, adopting DQaC will be essential in maintaining trust in data-driven decision-making.

By implementing the strategies outlined in this article, businesses can enhance their data quality frameworks, reduce operational risks, and unlock more reliable insights from their data assets. Future advancements in AI-driven data validation and real-time monitoring will further solidify DQaC as a critical component of data governance and engineering workflows.

---

### **Compliance with ethical standards**

#### *Disclosure of conflict of interest*

No conflict of interest to be disclosed.

---

### **References**

- [1] DataOps Manifesto. Principles of Data Quality Automation (DataOps Manifesto)
- [2] TDWI. Best Practices in Automated Data Quality (TDWI Research)
- [3] Snowflake. Ensuring Data Quality at Scale (Snowflake Whitepapers)