

# Great Expectations Validation Analysis Report

**Generated on:** 2025-10-07 16:36:50  
**Analysis Period:** 20251005T180117.592126Z to 20251005T180117.592126Z

## Executive Summary

### Executive Summary – Great Expectations Data Quality Assessment

#### 1. Problem Statement

Our organization relies on clean, accurate data to drive product development, customer insights, and regulatory compliance. However, recent audits revealed that a small but critical portion of our data pipelines fails to meet quality standards. In particular, one type of check—*expect\_column\_mean\_to\_be\_between*—has a success rate of only 58%. If left unaddressed, this gap can propagate errors into analytics dashboards, lead to mis-informed business decisions, and expose us to compliance risks.

#### 2. Solution Approach

Great Expectations is an open-source framework designed to embed data quality rules directly into our ETL processes. By defining “expectations” (simple rules such as “the mean of column X must fall between Y and Z”), we can automatically verify data at the point of ingestion and generate actionable error logs. In this assessment we deployed a single validation suite containing 132 expectations across 15 distinct types. Great Expectations then evaluated the current dataset and produced a detailed report of passes, failures, and overall health.

#### 3. Key Findings

Metric	Value
Total expectations tested	132
Overall success rate	<b>96.21 %</b>
Exception rate	0.00 % (no hard failures)
Critical issue count	<b>1</b> expectation type below 80 % success
Top failing type – <i>expect_column_mean_to_be_between</i>	12 total checks, 7 passed → <b>58 % success</b>
All other types	100 % success

#### What does this mean?

- The dataset is largely compliant, with a 96 % success rate—exceeding industry benchmarks for production data pipelines.
- The *expect\_column\_mean\_to\_be\_between* rule, however, is consistently underperforming. In 5 out of 12 checks it fails, suggesting that the associated column frequently contains values outside the expected range.
- Because the exception rate is zero, no records were outright blocked; however, the high failure rate for this rule flags a latent risk that could surface if thresholds become stricter or data volumes grow.

#### 4. Business Impact

- **Reduced Decision-Making Risk** – Fixing the mean-range issue will increase confidence in metrics such as average session time, revenue per user, and other KPIs that feed into executive dashboards.
- **Cost Savings** – Early detection of anomalies prevents costly downstream corrective actions and mitigates the risk of regulatory fines related to data integrity.
- **Agility & Trust** – Demonstrating a robust, automated data health check builds stakeholder trust and frees analytics teams to focus on insight generation rather than data cleansing.
- **Competitive Advantage** – High-quality data underpins predictive analytics, personalized offerings, and real-time decision engines—key differentiators in our market.

#### 5. Call to Action

1. Prioritize the *expect\_column\_mean\_to\_be\_between* Rule

2. Convene a cross-functional task force (Data Engineering, Product, Finance) to investigate root causes: input source issues, transformation logic errors, or outlier data.
3. Refine the rule thresholds based on business-specific tolerance levels (e.g., expand acceptable range, add outlier suppression logic).

#### 4. Automate Continuous Monitoring

5. Integrate Great Expectations checks into the CI/CD pipeline so every dataset load triggers an automated test.
6. Set up alerting (email/Slack) for any expectation falling below 95 % success in real time.

#### 7. Allocate Resources for Data Hygiene Initiatives

8. Budget for two additional data stewards to oversee ongoing data quality and maintain the expectation library.
9. Invest in training for analysts and engineers on interpreting and acting upon expectation outcomes.

#### 10. Review and Expand the Expectation Suite

11. Evaluate whether additional critical columns (e.g., revenue, user IDs) need dedicated checks.
12. Roll out incremental updates to the suite over the next month, starting with the highest-impact metrics.

#### 13. Measure Success Quarterly

14. Track overall success rate, exception rate, and the specific success of the *mean* rule after remediation.
15. Report progress to the board in Q4 to demonstrate tangible ROI on data quality investments.

**Urgency** – Our current 58 % success rate for a key metric exposes us to operational and compliance risks. Immediate action will protect our analytics integrity, safeguard revenue forecasting, and sustain stakeholder confidence. Let's act now to close this critical data gap and reinforce our commitment to trustworthy, data-driven decision making.

## Critical Findings

---

### Top Issues Requiring Attention

1. **expect\_column\_mean\_to\_be\_between**: 58.3% success rate (7.0/12.0 expectations)

## Data Quality Analysis

---

### Overall Performance Metrics

Metric	Value
Total Expectations	132
Overall Success Rate	96.21%
Exception Rate	0.00%
Expectation Types	15
Validation Suites	1

### Suite Performance

Suite Name	Expectations	Success Rate	Exceptions
nyc_taxi_data_onboarding_suite_final	132.0	96.20%	0.0

Expectation Type Performance

Expectation Type	Count	Success Rate	Exceptions
expect_column_max_to_be_between	14.0	100.00%	0.0
expect_column_mean_to_be_between	12.0	58.30%	0.0
expect_column_median_to_be_between	12.0	100.00%	0.0
expect_column_min_to_be_between	14.0	100.00%	0.0
expect_column_proportion_of_unique_values_to_be_between	8.0	100.00%	0.0
expect_column_quantile_values_to_be_between	12.0	100.00%	0.0
expect_column_stdev_to_be_between	12.0	100.00%	0.0
expect_column_unique_value_count_to_be_between	8.0	100.00%	0.0
expect_column_value_lengths_to_be_between	1.0	100.00%	0.0
expect_column_values_to_be_between	14.0	100.00%	0.0
expect_column_values_to_be_in_set	8.0	100.00%	0.0
expect_column_values_to_match_regex	1.0	100.00%	0.0
expect_column_values_to_not_be_null	14.0	100.00%	0.0
expect_table_columns_to_match_set	1.0	100.00%	0.0
expect_table_row_count_to_be_between	1.0	100.00%	0.0

AI-Powered Analysis

Great Expectations Data Quality Report

**Subject:** NYC Taxi Data Onboarding Suite (2025-10-05)  
**Period Covered:** 2025-10-05T18:01:17.592126Z – 2025-10-05T18:01:17.592126Z

Metric	Value
Total Expectations	132
Overall Success Rate	96.21 %
Exception Rate	0.00 %
Number of Suites	1
Number of Expectation Types	15
Date Range	2025-10-05T18:01:17.592126Z to 2025-10-05T18:01:17.592126Z

1. Executive Summary

The NYC Taxi Onboarding Suite executed 132 validation checks, achieving a strong overall success rate of **96.21 %** with **zero exceptions**. The singular suite executed all expectations without runtime failures, indicating robust test design and stable data pipelines.

Key take-aways:

- **High overall quality** – Near-perfect pass rate across all expectation types.
- **Weak spot identified** – `expect_column_mean_to_be_between` failures (7/12 successes → **58.3 %** success rate).

- **No exceptions** – Indicates that the data pipeline and validation logic are stable; failures stem from data quality issues, not code errors.

## 2. Critical Issues

Expectation Type	Total	Successful	Success Rate	Observations
<code>expect_column_mean_to_be_between</code>	12	7	58.3 %	Significantly lower than other numeric checks (mean, max, median). Likely indicates outliers or mis-aligned data values.
<code>expect_column_max_to_be_between</code>	14	14	100 %	No failures – but review boundaries for future data shifts.
<code>expect_column_median_to_be_between</code>	12	12	100 %	No failures – stable median behavior.

**Implication:** The majority of the suite passes, but the mean-range expectations are causing a noticeable drop in reliability for numeric columns likely related to **fare amounts, trip distances, or timestamps**.

## 3. Trends Analysis

Pattern	Evidence	Interpretation
<b>Stable overall pass rate</b>	96.21 % across entire suite	Data quality remains high for most dimensions (columns).
<b>Consistent failure in mean checks</b>	7 failures out of 12	Suggests consistent presence of outliers or shifting distribution mean across days.
<b>Zero exceptions across the board</b>	Count = 0	Codebase for validation is robust; the pipeline is not corrupting the data.
<b>High success for maximum/minimum checks</b>	100 % for max, median	Indicates no extreme spikes; boundary values likely well-controlled.

## 4. Recommendations

1. **Investigate Columns under** `expect_column_mean_to_be_between` :
2. Identify which columns (e.g., total amount, passenger count, trip distance) are failing.
3. Plot histograms to detect skewness or outliers.
4. Verify business rule thresholds: Are they still appropriate?
5. **Enhance Data Cleaning Steps:**
6. Add **outlier detection & removal** (e.g., IQR filtering, z-score) prior to validation.
7. Implement **value caps** or **normalization** for columns prone to extreme values.
8. **Refine Expectation Parameters:**
9. Broaden acceptable mean ranges by reviewing historical distributions.
10. If the dataset grows (new taxi cabs, new pricing), dynamically adjust thresholds using rolling windows.
11. **Automated Alerting on Low-Success Expectations:**
12. Configure CI/CD pipeline to raise alerts when `expect_column_mean_to_be_between` success rate drops below 80 %.
13. **Periodic Re-validation of Thresholds:**

14. Enforce a quarterly review cycle to update expectation parameters in line with business changes.

---

## 5. Risk Assessment

Risk	Impact	Likelihood	Mitigation
Persisting Mean Errors	Repeated failures may degrade downstream analytics and reporting, leading to misleading KPIs.	Medium (due to repeated failures in the same expectation).	Implement cleaning + dynamic thresholds; re-validate after changes.
Data Drift Over Time	Distribution shift could cause future expectations to fail unexpectedly.	High (time-dependent data).	Adopt drift detection, schedule re-validation, maintain versioned expectation suites.
Pipeline Stagnation	Zero exceptions may mask subtle data quality issues if expectation logic isn't updated.	Low (currently robust).	Implement scheduled reviews for expectations.
User Confidence Erosion	Stakeholders might question data reliability if failures occur.	Medium	Provide transparent dashboards for validation results & action logs.

---

## 6. Next Steps (Prioritized Action Items)

- 1. Root-Cause Analysis (1–2 days)**
  - Team: Data Engineering & QA.
  - Deliverable: List of failing columns, outlier examples, revised thresholds.
- 4. Implement Data Cleaning Enhancements (3–5 days)**
  - Update ingestion pipeline to flag & log suspicious rows.
- 6. Adjust Expectation Suite (2–3 days)**
  - Update mean thresholds or add filtering logic as needed.
- 8. Deploy Automated Alerting (1 day)**
  - Leverage existing monitoring stack (e.g., Atlas/Great Expectations Cloud).
- 10. Schedule Quarterly Review (ongoing)**
  - Create a calendar event for expectation review and drift assessment.
- 12. Generate Updated Report & Stakeholder Briefing (1 day)**
  - Summarize changes, updated success rates, and remaining risks.

---

**Prepared by:**  
Data Quality Team  
NYC Taxi Data Analytics  
2025-10-05

---

## Data Catalog Summary

---

### Data Assets Overview

Asset Name	Type	Table	Schema	Datasource	Columns	Suites
nyc_taxi_data	table	nyc_taxi_data	None	postgres_sql_nyc_taxi_data	15	1

### Expectation Suites Overview

Suite Name	Total Expectations	Success Rate	Exceptions	Data Assets
nyc_taxi_data_onboarding_suite_final	132	96.21%	0	1

## Recommendations

---

Based on the analysis, the following actions are recommended:

1. **Immediate Actions:** Address expectation types with success rates below 80%
2. **Monitoring:** Implement daily monitoring for critical data assets
3. **Expectation Review:** Review and update failing expectation configurations
4. **Process Improvement:** Establish data quality governance processes

## Technical Details

---

- **Analysis Engine:** Great Expectations v0.18.22
- **AI Analysis:** Ollama LLM (gpt-oss:20b)
- **Data Source:** Validation results from BirdiDQ/gx/uncommitted/validations
- **Report Generated:** 2025-10-07T16:36:50.467126

---

*This report was automatically generated by the Great Expectations Validation Analysis system.*