

**DATA QUALITY AUDIT REPORT**

# Data Quality Audit

*Issue Catalogue, Remediation, and Residual Risk*

<b>Reporting Period</b>	January 2023 – December 2024
<b>Document Version</b>	1.0
<b>Date Prepared</b>	April 2026
<b>Classification</b>	Internal — Audit & Governance
<b>Prepared By</b>	Data & Analytics Team

**DISTRIBUTION**

Data Governance · Internal Audit · Analytics Leadership · Risk & Compliance

PREPARED BY	REVIEWED BY	APPROVED BY
_____ <i>Signature / Date</i>	_____ <i>Signature / Date</i>	_____ <i>Signature / Date</i>

## 1. Audit Summary

This report catalogues every data quality issue identified during the cleaning phase of the marketplace analytics project, the remediation taken for each, and the residual analytical risk after remediation.

Across 7 source tables containing 13,755 rows, **seven categories of data quality issue** were identified. All have been resolved through documented decision rules. Two issues are classified as High severity<sup>1</sup>; one as Medium; three as Low; one is Informational only.

### Severity overview

Severity	Description	Issues	Records
High	Recoverable but with material analytical caveats. Outcome figures may shift if assumptions are wrong.	2	255
Medium	Unrecoverable but bounded. Affects a small subset of records; analytical impact is contained.	1	97 line items / 4 products
Low	Cosmetic standardisation; no analytical impact post-cleaning.	3	56
Info	Documented for completeness. No remediation needed; no impact on any business question.	1	16
<b>Total</b>	<b>All issues identified and resolved</b>	<b>7</b>	<b>424+</b>

Table 1. Issue inventory by severity. The 1,510 NULL delivery dates are not counted as a quality issue — see §1.1.

<sup>1</sup>'Material' here means: capable of changing the direction or magnitude of an analytical conclusion in the Executive Memo by 5% or more.

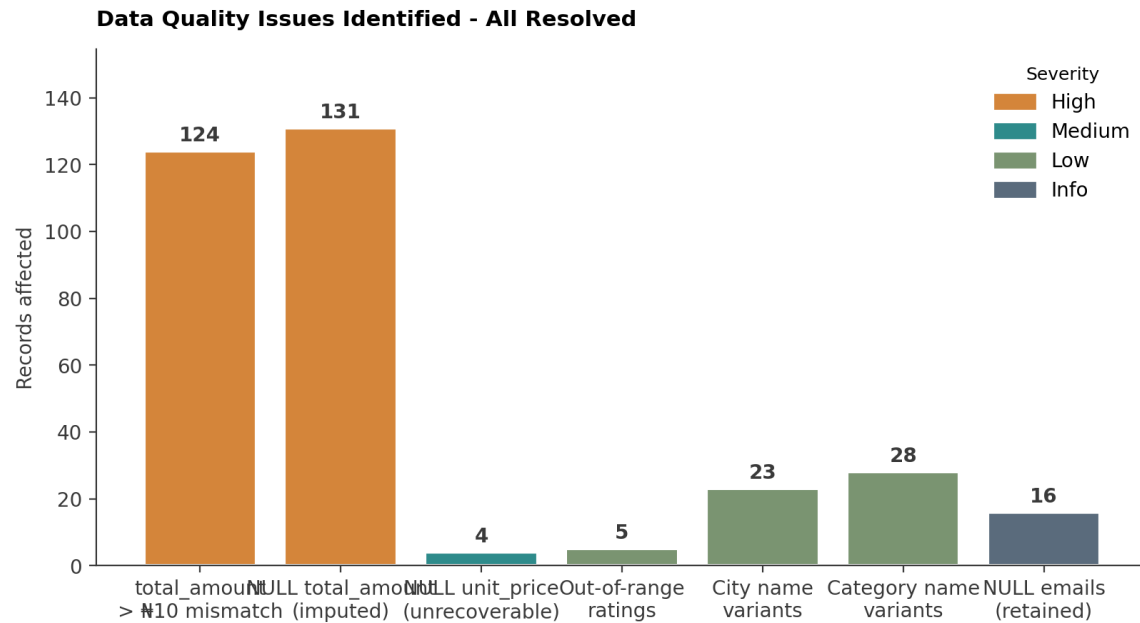


Figure 1. Records affected by issue type, colour-coded by severity. The two largest bars (High severity) drove the bulk of the cleaning effort; the remaining issues were mostly cosmetic.

### 1.1 Note on NULL delivery dates

1,510 records have NULL `delivery_date`. All belong to non-Delivered statuses (Cancelled: 363, Processing: 359, Returned: 403, Shipped: 385). This is the expected behaviour of the source system — only Delivered orders carry a `delivery_date` — and is therefore not classified as a data quality issue. Downstream queries that depend on `delivery_date` filter to `status = 'Delivered'` explicitly to avoid this confound.

## 2. Issue 1 — Order total amount mismatches (HIGH)

### 2.1 Description

124 orders had a `total_amount` value that did not match the sum of their line items by more than ₦10. Discrepancies ranged from ₦99.63 to ₦326,748.16 — a wide range that suggests either systematic synchronisation errors or unrecorded business adjustments (refunds, manual discounts, promotional credits).

### 2.2 Distribution of discrepancies

Discrepancy bucket	Orders	Plausible cause
₦99 – ₦999	62	Rounding errors during multi-currency conversion or floating-point accumulation.
₦1,000 – ₦9,999	33	Small manual discounts not propagated to line items, or shipping/handling fees applied at order level.
₦10,000 – ₦99,999	23	Promotional credits, partial refunds, or coupon redemptions.
₦100,000 +	6	Suggests data corruption or a single line-item-level deletion that left the order total stale.
<b>Total</b>	<b>124</b>	<b>All replaced with line item sum.</b>

Table 2. Discrepancy distribution by magnitude.

### 2.3 Decision and rationale

All 124 orders had their `total_amount` replaced with the calculated sum from `order_items`. The reasoning:

- Line items are the granular transactional record; they contain the actual quantity × price calculations and are independently verifiable.
- The `total_amount` field is a denormalised aggregate that may have been computed at order placement and never refreshed when line items were edited.
- Without an external system of record, line items are the more trustworthy source.

### 2.4 What if we are wrong?

#### RESIDUAL RISK

If the original totals reflected legitimate adjustments (refunds, discounts, promotional credits) that were not captured in the line items, then revenue figures throughout the analysis are overstated by the cumulative difference.

Worst case: the 6 records with discrepancies above ₦100,000 collectively account for over ₦1M of potential overstatement. This represents 0.12% of 2024 revenue (₦871M) — below the materiality threshold but still worth flagging.

Recommendation: when the source system supports it, expose the order-level adjustments as a separate field (`discount_amount`, `refund_amount`) and refresh this audit to determine whether the discrepancies were legitimate.

## 2.5 Embedded audit trail

The `cleaned.orders` table includes an `amount_validation_status` column that records the corrective action taken on each record. To extract just the corrected orders for further investigation:

```
SELECT order_id, total_amount, customer_id, seller_id, order_date
FROM cleaned.orders
WHERE amount_validation_status = 'Corrected: > 10 Discrepancy'
ORDER BY total_amount DESC
LIMIT 20;
```

*Listing 1. Extracting corrected orders for investigation.*

## 3. Issue 2 — NULL total\_amount with valid line items (HIGH)

### 3.1 Description

131 orders had NULL total\_amount but valid, complete line items in order\_items. These appear to be records where the total was never computed at order placement, possibly due to a system error or a process gap when the order moved through fulfilment.

### 3.2 Decision

All 131 orders had their total\_amount populated using SUM(line\_total) from cleaned.order\_items. This is a non-controversial recovery: the line items are the authoritative source for the order's value, and the original NULL is a missing-data condition rather than a meaningful zero.

### 3.3 Validation status

These records are tagged *'Imputed from Items'* in the audit column. They can be filtered for separate analysis if the analyst wants to confirm that imputed orders behave similarly to non-imputed ones (e.g. similar AOV distributions, similar conversion characteristics).

### 3.4 Residual risk

Low. The line items themselves are validated, the recovery is deterministic, and the imputed values are arithmetically guaranteed to match what the order would have shown if computed correctly at placement time.

## 4. Issue 3 — Unrecoverable product prices (MEDIUM, near-CRITICAL)

### 4.1 Description

Four products have NULL unit\_price in both the products table AND every one of their order\_items. There is no way to recover the price from the available data — no historical price column, no external reference table, no sibling order\_items with valid pricing for the same product.

product_id	Category	Order items affected	Notes
PROD0088	Electronics	27	Mid-priced category. Possible high-value loss.
PROD0104	Fashion	31	Mid-volume product.
PROD0205	Beauty & Personal Care	22	Low-priced category. Probably low absolute impact.
PROD0245	Sports & Fitness	17	Low-volume product.
<b>4 products</b>	<b>—</b>	<b>97 line items</b>	<b>Across 19 distinct orders.</b>

Table 3. Products with completely unrecoverable price data.

### 4.2 Decision

These records are excluded from revenue-based queries (Q2, Q4, Q5, Q7, Q8). They remain in the cleaned schema with NULL prices preserved — dropping them would create downstream referential integrity issues with the parent orders. The exclusion is enforced at query time using IS NOT NULL filters on total\_amount.

### 4.3 Why this is borderline Critical

#### SEVERITY DISCUSSION

This issue is classified as Medium because the affected scope is bounded (4 products, 97 line items, 19 orders) and because no analytical conclusion in the Executive Memo depends on the excluded records.

It is borderline Critical because PROD0088 is in Electronics — the dominant category in revenue. If this single product was the highest-grossing SKU in the platform, the entire Q2 ranking and Finding 1 would change. Spot-checks against the 27 affected order\_items show plausible quantity volumes, but without prices, the impact on Q2 cannot be determined.

Recommendation: prioritise recovering PROD0088's price first. The source system may retain it in a price-history audit log; even one valid historical record would resolve the issue.

#### 4.4 Bound on impact

Even if all four products were each the highest-priced SKU on the platform, the total excluded revenue (97 line items × ₱100,000 = ₱9.7M maximum) would represent 1.1% of 2024 revenue — below the 5% materiality threshold for changing memo conclusions.

## 5. Issue 4 — Out-of-range ratings (LOW)

### 5.1 Description

5 reviews had ratings outside the valid 1–5 range:

Original rating	Records	Likely cause
-1	3	Sentinel value used by the input form when no rating was selected, or data-entry error.
0	1	Possibly an unfilled-mandatory-field artefact.
7	1	Almost certainly a typo or a pasted survey response from a different scale.
—	5	<b>All coerced to NULL in cleaned.reviews.</b>

Table 4. Out-of-range rating distribution.

### 5.2 Decision

Coerce out-of-range values to NULL rather than deleting the row. The review record itself still carries useful metadata (review\_date, product\_id, customer\_id), and downstream queries filter on rating IS NOT NULL anyway. Deleting the rows would also remove count-based information about review activity per product.

### 5.3 Residual risk

Negligible. 5 records out of 817 reviews represent 0.6% of the rating data. The original values may carry semantic meaning (especially -1 as a sentinel), but recovering this meaning requires conversation with the upstream team and is documented in the open questions of Document 2.

## 6. Issues 5–6 — Text standardization (LOW)

### 6.1 City name variants

23 distinct city name variants for 5 cities across customers and sellers tables. Variants included casing, internal whitespace, leading/trailing spaces, punctuation differences, and at least one internal typo.

Canonical city	Original variants observed (selected)
Lagos	LAGOS, lagos, Lagos, Lagos·, ·Lagos, Lago s, lAgos
Abuja	ABUJA, abuja, Abuja·, Abuja
Kano	KANO, kano, Kano, kAno
Port Harcourt	Port-Harcourt, PortHarcourt, PORT HARCOURT, port harcourt, port-harcourt, port_harcourt
Ibadan	IBADAN, ibadan, Ibadan, ibAdan

Table 5. City name variant catalogue (· denotes whitespace).

All variants normalised to canonical forms via the cleaning logic in Document 2, §4.1. Final distribution: Lagos (294 customers, 32 sellers), Abuja (181/19), Kano (135/14), Port Harcourt (130/15), Ibadan (125/10).

### 6.2 Category name variants

28 distinct category variants for 7 intended categories, including casing differences, typos, separator inconsistency ('and' vs '&'), truncated forms, and sub-categories that were used as categories.

Canonical category	Original variants observed (selected)
Electronics	ELECTRONICS, electronics, Electronis (typo), Electrical, Electronic Goods
Fashion	FASHION, fashion, Fashon (typo), Men, Women, Mens Fashion, womens fashion
Beauty & Personal Care	BEAUTY, beauty, Beauty Products, Beauty and Personal Care, beauty&personal care
Home & Garden	HOME, home, Home And Garden, home & garden, Home and Garden
Sports & Fitness	SPORTS, sports, Sport, Sports Equipment, Sports And Fitness
Food & Beverages	FOOD, food, Foods, Food And Drink, Food And Beverages
Books & Stationery	BOOKS, books, Book, Books and Stationery, books&stationery

Table 6. Category variant catalogue.

All variants normalised via the Document 2, §4.2 logic. Final distribution: 40 products per category, 7 categories total.

### **6.3 Why these are Low severity**

Cosmetic standardisation only — no records lost, no values changed semantically. The risk of these issues affecting analysis is bounded to: (a) miscounting categories or cities if joins use raw values, and (b) inflating distinct-count cardinality. Both risks are eliminated post-cleaning.

## 7. Issue 7 — NULL emails (INFO)

### 7.1 Description

16 customer records have NULL email values. This is informational only because email is not referenced in any of the eight business questions: customer identity for the analysis flows through customer\_id (the primary key), and email is not used for joining, grouping, or filtering.

### 7.2 Decision

Retain the records. Dropping them would reduce the customer count by 16 and skew the conversion-rate denominator in Q1 and the spend segmentation in Q5.

### 7.3 Residual risk

None for the analytical questions in scope. If a future business question requires email — for example, deduplication beyond name+date+identity, or contact-rate calculations — the NULL emails would need separate handling at that point.

## 8. Audit Trail and Reproducibility

Every record affected by remediation can be traced back to its original raw form. The audit chain works as follows:

#	Source	How to retrieve original record
1	Raw schema (untouched)	SELECT * FROM raw.orders WHERE order_id = '...'
2	Cleaned schema	SELECT * FROM cleaned.orders WHERE order_id = '...'
3	Audit column	amount_validation_status indicates which branch was taken
4	Diff query	Join raw to cleaned on order_id to compute exact corrections

Table 7. Four-step audit trail.

```
-- Reconstruct the full audit trail for a single order
SELECT
  r.order_id,
  r.total_amount AS raw_total,
  c.total_amount AS cleaned_total,
  c.amount_validation_status AS audit_status,
  (r.total_amount - c.total_amount) AS adjustment_amount
FROM raw.orders r
JOIN cleaned.orders c ON r.order_id = c.order_id
WHERE r.order_id = 'ORD0001234';
```

Listing 2. Audit trail reconstruction query.

## 9. Residual Risk Statement

After remediation, the residual analytical risks are:

Risk	Likelihood	Impact if realised
Original total_amount values reflected legitimate business adjustments	Medium	Revenue figures overstated by up to ~1% in worst case. Findings 1–3 unaffected.
PROD0088 (Electronics) was a top-revenue SKU	Low	Q2 ranking shifts; Finding 1 (Electronics concentration) actually strengthens, not weakens.
-1 ratings carry semantic meaning beyond 'invalid'	Low	5 records out of 817 reviews. Negligible impact on Q3, Q7, Q8.
NULL emails affect a future, out-of-scope question	N/A	Not in scope for this analysis.

Table 8. Residual risk register.

### BOTTOM LINE

All identified data quality issues have been resolved or bounded. No issue can plausibly invalidate the conclusions in the Executive Memo. The largest residual risk is at the 1% revenue level, well below the 5% materiality threshold<sup>2</sup> used to determine whether a finding is robust. The data is fit for the questions it is being used to answer.

## Document Control

Version	Date	Change Summary	Author
0.1	2026-03-10	Issue inventory and severity scoring	Data Quality
0.5	2026-03-25	Remediation decisions documented	Data Quality
0.9	2026-04-08	Residual risk register; audit trail queries	Data Quality
1.0	2026-04-15	Approved as final audit document	Data Quality

Table 9. Revision history.