Shwetank Singh
**GritSetGrow - GSGLearn.com**

**DATA AND AI**

# ADVANCED DATA CLEANING TECHNIQUES FOR E-COMMERCE PROJECTS

www.gsglearn.com

# Advanced Data Cleaning Techniques for E-Commerce Data Engineering Projects

**Overview:**
E-commerce data spans customer profiles, product catalogs, and transactional records that are both voluminous and heterogeneous. This document details advanced data cleaning strategies for data engineering projects targeting e-commerce platforms. It covers end-to-end cleaning steps, from initial profiling and parsing of semi-structured data to deduplication, normalization, and anomaly detection, along with performance optimizations and robust monitoring. Advanced SQL examples, including code snippets for structured and semi-structured data, are provided throughout. This comprehensive guide is designed to span roughly 20 pages when compiled, making it a deep dive for advanced data engineers.

---

## Table of Contents

---

## Data Profiling and Initial Assessment

- **Profiling Objectives:**

    - Understand data distributions, cardinality, and frequency of nulls.
    - Identify patterns, irregularities, and potential anomalies before cleaning.
    - Generate data quality reports to determine the scale of cleaning needed.

- **Techniques:**

    - **Statistical Summaries:**
      Use SQL aggregations to compute min, max, average, standard deviation:

      ```sql
      SELECT MIN(total_amount) AS min_order,
             MAX(total_amount) AS max_order,
             AVG(total_amount) AS avg_order,
      ```

```
        STDDEV_POP(total_amount) AS stddev_order
FROM orders;
```

- **Data Distribution Analysis:**
  Use histograms or percentiles:

  ```
  SELECT percentile_cont(0.5) WITHIN GROUP (ORDER BY total_amount) AS
  median_order
  FROM orders;
  ```

- **Null and Missing Value Assessment:**
  Count missing values for critical columns:

  ```
  SELECT COUNT(*) AS missing_emails
  FROM customers
  WHERE email IS NULL OR email = '';
  ```

- **Outcome:**

  - Create a data quality dashboard with key metrics.
  - Identify columns with high null percentages, outliers, or inconsistent formats.

---

## Advanced Structured Data Cleaning

- **Schema Enforcement and Data Type Correction:**

  - Validate that each field conforms to its expected data type.
  - **Example:** Convert text-based dates into proper DATE format:

    ```
    UPDATE orders
    SET order_date = TO_DATE(order_date_text, 'MM/DD/YYYY')
    WHERE order_date IS NULL;
    ```

- **Relational Integrity and Foreign Key Validation:**

  - Ensure every record in a fact table has corresponding entries in dimension tables.
  - **Example:** Identify orphaned transactions:

    ```
    SELECT o.order_id
    FROM orders o
    LEFT JOIN customers c ON o.customer_id = c.customer_id
    WHERE c.customer_id IS NULL;
    ```

- **Constraint Enforcement:**

  - Use `CHECK` constraints to prevent future data errors.
  - **Example:** Ensure non-negative product prices:

    ```
    ALTER TABLE products
    ADD CONSTRAINT chk_price_nonnegative CHECK (price >= 0);
    ```
```

- **Advanced Null Handling:**

  - Apply sophisticated imputation techniques:
    - **Median/Mode Imputation:** Use aggregate subqueries to replace missing numeric fields.
    - **Conditional Imputation:** For example, impute shipping dates based on typical processing times.

- **Transactional Integrity:**

  - Employ stored procedures to validate multi-step transactions.
  - Log any anomalies for later review.

---

## Advanced Semi-Structured Data Cleaning

- **Parsing Nested Data Structures:**

  - Use JSON functions to extract and validate nested attributes:

    ```
    SELECT order_id,
           JSON_EXTRACT_PATH_TEXT(order_details, 'product_id') AS
    product_id,
           JSON_EXTRACT_PATH_TEXT(order_details, 'quantity') AS quantity
    FROM orders_json;
    ```

  - Validate that the extracted values conform to expected types and formats.

- **Handling Heterogeneous Formats:**

  - Normalize variations in key names or value formats:
    - Map `"color"` and `"colour"` to a single standardized key.
    - Use SQL CASE statements to transform differing representations.

- **Dealing with Optional and Dynamic Fields:**

  - Leverage functions such as `COALESCE` to fill in missing keys with default values.
  - **Example:** Standardizing an optional discount field:

    ```
    SELECT order_id,
           COALESCE(JSON_EXTRACT_PATH_TEXT(order_details, 'discount'), '0')
    AS discount
    FROM orders_json;
    ```

- **Data Flattening:**

  - Transform nested arrays into relational tables using lateral joins or unnest functions.
  - **Example:** Flatten product reviews embedded in JSON:

    ```
    SELECT order_id, review.*
    FROM orders_json,
        LATERAL jsonb_to_recordset(orders_json.review_list) AS
    review(review_id INT, rating INT, comment TEXT);
    ```

---

- **Schema Evolution Handling:**

    - Create flexible pipelines that can handle schema drift in semi-
      structured inputs.
    - Use versioned parsers to manage changes over time.

---

## Handling Missing Values at Scale

- **Advanced Missing Value Strategies:**

    - **Multi-Column Imputation:**
      Consider correlations between columns to impute missing values. For
      example, if a product's weight is missing, use the average weight for
      that category.
    - **Algorithmic Imputation:**
      Employ statistical or machine learning methods (e.g., regression models)
      to predict missing values.
    - **Flagging and Segregation:**
      Instead of immediately replacing, mark records with missing critical
      data and process them in a review pipeline.

- **SQL Examples:**

    - **Window Function Imputation:**

    ```sql
    UPDATE products p
    SET weight = sub.avg_weight
    FROM (
      SELECT category, AVG(weight) OVER(PARTITION BY category) AS avg_weight
      FROM products
      WHERE weight IS NOT NULL
    ) sub
    WHERE p.category = sub.category
      AND p.weight IS NULL;
    ```

    - **Conditional Replacement:**

    ```sql
    SELECT order_id,
           CASE
             WHEN delivery_date IS NULL THEN order_date + INTERVAL '5' DAY
             ELSE delivery_date
           END AS estimated_delivery_date
    FROM orders;
    ```

- **Considerations:**

    - Evaluate the impact of imputation on downstream analysis.
    - Document assumptions behind each imputation rule for auditability.

---

## Deduplication and Fuzzy Matching Techniques

- **Exact Deduplication:**

- **Window Functions:**
  Use `ROW_NUMBER()` to mark duplicates and keep the most recent record.

  ```sql
  WITH RankedCustomers AS (
    SELECT customer_id, email, created_at,
           ROW_NUMBER() OVER (PARTITION BY email ORDER BY created_at DESC)
  AS rn
    FROM customers
  )
  DELETE FROM customers
  WHERE customer_id IN (
    SELECT customer_id FROM RankedCustomers WHERE rn > 1
  );
  ```

- **Fuzzy Matching:**

  - **Phonetic Algorithms:**
    Apply functions like SOUNDEX to identify similar names.

    ```sql
    SELECT a.customer_id, b.customer_id, a.name, b.name
    FROM customers a
    JOIN customers b ON a.customer_id <> b.customer_id
    WHERE SOUNDEX(a.name) = SOUNDEX(b.name);
    ```

  - **Levenshtein Distance:**
    Some SQL engines support UDFs for calculating string similarity.

    ```sql
    SELECT a.customer_id, b.customer_id, a.name, b.name
    FROM customers a
    JOIN customers b ON a.customer_id <> b.customer_id
    WHERE LEVENSHTEIN(a.name, b.name) < 3;
    ```

- **Merging Records:**

  - Develop custom logic to merge duplicates, choosing the best available value for each field.
  - **Example:**
    Use `COALESCE` and aggregate functions to merge records:

    ```sql
    SELECT email,
           MIN(customer_id) AS primary_id,
           MAX(name) AS consolidated_name,
           MAX(phone) AS consolidated_phone,
           MAX(address) AS consolidated_address
    FROM customers
    GROUP BY email;
    ```

- **Challenges and Solutions:**

  - Handle near-duplicates carefully to avoid false positives.
  - Integrate manual review processes for borderline cases.

# Data Standardization and Transformation

- **Format Standardization:**

  - **Text Normalization:**
    Convert to lower case, remove extra spaces, and standardize punctuation.

    ```sql
    SELECT TRIM(LOWER(name)) AS normalized_name
    FROM customers;
    ```

  - **Date/Time Normalization:**
    Convert all date/time values to UTC.

    ```sql
    UPDATE orders
    SET order_date = CONVERT_TZ(order_date, 'US/Eastern', 'UTC');
    ```

- **Categorical Value Mapping:**

  - Use mapping tables to standardize categorical fields (e.g., mapping
    different representations of product categories).

    ```sql
    SELECT p.product_id, p.product_name, m.standard_category
    FROM products p
    LEFT JOIN category_mapping m ON p.category = m.raw_category;
    ```

- **Unit and Currency Conversions:**

  - Convert various units to a common standard (e.g., inches to centimeters,
    different currencies to USD).
  - **Example:**

    ```sql
    UPDATE products
    SET weight = CASE
                   WHEN weight_unit = 'lbs' THEN weight * 0.453592
                   ELSE weight
                 END,
        weight_unit = 'kg';
    ```

- **Data Transformation Pipelines:**

  - Chain multiple transformations using Common Table Expressions (CTEs) to
    maintain clarity.
  - **Example:**

    ```sql
    WITH CleanedProducts AS (
      SELECT product_id,
             TRIM(LOWER(product_name)) AS product_name,
             CASE
               WHEN weight_unit = 'lbs' THEN weight * 0.453592
               ELSE weight
             END AS weight,
             'kg' AS weight_unit
      FROM products
    ```

```
  )
  SELECT * FROM CleanedProducts;
```

- **Handling Special Characters and Encoding:**

  - Remove or replace non-ASCII characters that might interfere with joins
    or aggregations.
  - Utilize regex functions where available to clean text fields.

---

## Anomaly Detection and Outlier Handling

- **Statistical Outlier Detection:**

  - **Standard Deviation:**
    Flag records that are beyond 3 standard deviations from the mean.

    ```
    SELECT order_id, total_amount
    FROM orders
    WHERE total_amount > (
      SELECT AVG(total_amount) + 3 * STDDEV_POP(total_amount)
      FROM orders
    );
    ```

  - **Percentile-Based:**
    Identify orders above the 99th percentile.

    ```
    SELECT order_id, total_amount
    FROM orders
    WHERE total_amount > (
      SELECT percentile_cont(0.99) WITHIN GROUP (ORDER BY total_amount)
      FROM orders
    );
    ```

- **Temporal Anomalies:**

  - Detect sudden spikes in orders or user activity.

    ```
    SELECT order_time,
           COUNT(*) OVER (ORDER BY order_time RANGE INTERVAL '1' HOUR
    PRECEDING) AS orders_last_hour
    FROM orders;
    ```

  - Trigger alerts if the count exceeds a threshold.

- **Multi-Dimensional Anomaly Detection:**

  - Combine multiple fields (e.g., order amount, product quantity, and time)
    to detect suspicious patterns.
  - Develop machine learning models offloaded from SQL for complex patterns,
    then integrate flagged IDs back into SQL pipelines.

- **Automated Correction vs. Flagging:**

- Decide when to auto-correct anomalies (e.g., a misplaced decimal) or flag them for manual review.
  - Maintain a log table for anomalous records.

---

# Data Normalization (Scaling and Structuring)

- **Numeric Normalization:**

  - **Min-Max Scaling:**

    ```sql
    SELECT product_id,
           (price - min_price) / (max_price - min_price) AS price_norm
    FROM (
      SELECT product_id, price,
             MIN(price) OVER() AS min_price,
             MAX(price) OVER() AS max_price
      FROM products
    ) sub;
    ```

  - **Z-Score Standardization:**

    ```sql
    SELECT customer_id,
           (annual_spend - avg_spend) / stddev_spend AS spend_zscore
    FROM (
      SELECT customer_id, annual_spend,
             AVG(annual_spend) OVER() AS avg_spend,
             STDDEV_POP(annual_spend) OVER() AS stddev_spend
      FROM customer_stats
    ) t;
    ```

  - **Relative Scaling:**
    Normalize sales figures as a percentage of total sales.

    ```sql
    SELECT product_id, sales,
           sales / SUM(sales) OVER() * 100 AS percent_of_total_sales
    FROM product_sales;
    ```

- **Database Normalization:**

  - Design relational schemas to reduce redundancy (e.g., separate brand, category, and supplier information).
  - Leverage foreign keys to maintain consistency.
  - Ensure lookup tables are properly indexed for join performance.

- **Text and Categorical Normalization:**

  - Use normalization functions to ensure uniformity in textual data.
  - Consolidate similar category labels to avoid fragmentation in analysis.

---

# Data Quality Framework and Metrics

- **Defining KPIs for Data Quality:**

- **Completeness:** Percentage of non-null values.
- **Accuracy:** Conformance of data values to business rules.
- **Consistency:** Uniformity across different datasets.
- **Timeliness:** Up-to-date information relative to business requirements.
- **Uniqueness:** Degree of deduplication achieved.

- **Establishing Data Quality Rules:**

  - Create a rules engine (either SQL-based or via a dedicated tool) that continuously validates data.
  - Record rule violations and track trends over time.

- **SQL Example – Quality Check Summary:**

```sql
SELECT 'orders' AS table_name,
       COUNT(*) AS total_records,
       SUM(CASE WHEN delivery_date IS NULL THEN 1 ELSE 0 END) AS
missing_delivery_dates,
       ROUND(SUM(CASE WHEN delivery_date IS NULL THEN 1 ELSE 0 END) * 100.0 /
COUNT(*), 2) AS pct_missing
FROM orders;
```

## Reporting and Dashboards

- **Build dashboards using BI tools** to visualize data quality metrics.

- **Schedule periodic reports** to stakeholders.

# Metadata Management and Data Lineage

## Metadata Capture

- **Track schema changes, data source details, and transformation rules.**
- **Use automated tools** to capture metadata during ETL processes.

## Data Lineage

- **Document the flow of data** from source systems through cleaning and transformation stages.
- **Maintain lineage logs** that can trace a record from raw ingestion to final consumption.

## Implementing in SQL

- **Create audit tables** that log changes made during cleaning:

```sql
CREATE TABLE data_cleaning_audit (
  audit_id SERIAL PRIMARY KEY,
  table_name VARCHAR(50),
  record_id INT,
  change_description TEXT,
```

```
    changed_at TIMESTAMP DEFAULT CURRENT_TIMESTAMP
);
```

### Logging During Critical Transformations

- **Insert logs** during critical transformations to maintain data traceability.

- **Track data modifications** for debugging and compliance purposes.

## Benefits

- **Enhanced debugging** to identify issues in data processing.
- **Compliance with regulatory requirements** (e.g., GDPR, CCPA) through audit trails.

---

# Automated Data Cleaning Pipelines

## Designing Robust ETL/ELT Pipelines

- **Integrate cleaning steps** into data ingestion pipelines to ensure high data quality.
- **Automate data quality checks and error handling** to minimize manual intervention.

## Scheduling and Orchestration

- **Use workflow orchestration tools** like Apache Airflow, dbt, or cloud-native services (e.g., AWS Glue, Google Cloud Dataflow) to schedule cleaning jobs.
- **Define dependencies** between cleaning stages to ensure proper execution order.

## Automated Testing

- **Write unit tests** for SQL transformations to validate correctness.
- **Use integration tests** to verify data consistency across pipeline stages.

## Continuous Integration

- **Version control SQL scripts and ETL jobs** to ensure traceability and rollback capabilities.
- **Deploy cleaning processes using CI/CD pipelines** to facilitate rapid iteration and updates.

---

# Monitoring, Logging, and Auditability

## Real-Time Monitoring

- **Set up dashboards** to visualize key data quality metrics.
- **Use alerting systems** (e.g., PagerDuty, Slack notifications) to notify teams of cleaning failures.

### Detailed Logging

- **Log every major cleaning step and transformation** for better traceability.
- **Maintain logs** to facilitate debugging and historical analysis.

### Auditing

- **Periodically audit cleaned data** against raw data to ensure accuracy.
- **Use automated scripts** to compare key aggregates before and after cleaning.

### SQL Example – Logging Update Actions

```sql
INSERT INTO data_cleaning_audit (table_name, record_id, change_description)
VALUES ('orders', 12345, 'Imputed missing delivery_date with order_date + 5 days');
```

# Cloud and Distributed Processing Considerations

## Scaling Data Cleaning

- Use distributed SQL engines like **Google BigQuery**, **Amazon Redshift**, and **Snowflake** to efficiently process large datasets.
- Leverage **partitioning** and **clustering** to enhance query performance.

## Hybrid Architectures

- Combine **batch processing** (for historical data) with **real-time cleaning** (for streaming data).
- Integrate with **cloud-native data lakes and warehouses** for scalability.

## Cost Optimization

- Monitor query performance and **optimize SQL queries** to reduce compute costs.
- Use **materialized views** or **temporary tables** to cache intermediate results.

## Security and Compliance

- Ensure **encryption and access control** on sensitive data throughout the cleaning process.
- Regularly audit data access logs and enforce **least-privilege principles**.

---

# Case Studies and E-Commerce Examples

## Customer Data Integration

- Merge customer profiles from **website, mobile app**, and **third-party sources**.
- Use advanced **deduplication** and **fuzzy matching** to consolidate duplicate customer records.

## Product Catalog Merging

- Combine product data from multiple vendors into a **unified catalog**.
- Standardize **attributes**, **units**, and **taxonomies** for consistent search and categorization.

### Transactional Data Validation

- Ensure orders, returns, and payment records are **reconciled**.
- Detect and flag **inconsistent financial records** for finance team review.

### Advanced Use Cases

- Implement **machine learning models** to predict missing values and assess data reliability.
- Leverage **real-time streaming processing** for fraud detection and anomaly alerts in transactions.

---

# Common Challenges and Advanced Solutions

### Evolving Schemas & Data Drift

- Implement **schema versioning** to handle evolving data structures.
- Use automated **schema drift detection tools** to monitor changes.

### Heterogeneous Data Sources

- Build **source-specific normalization layers** to standardize incoming data.
- Maintain a **central data dictionary** and **mapping tables** for reference.

### Data Volume & Performance

- Partition large datasets and apply **incremental cleaning strategies**.
- Optimize SQL queries using **indexing**, **query hints**, and distributed processing.

### Balancing Quality with Data Retention

- Develop strategies to **flag** rather than delete questionable data.
- Create separate pipelines for **data cleansing** and **raw data archiving**.

### Duplicate and Fuzzy Matching Complexity

- Combine **algorithmic matching** with manual review processes.
- Leverage external libraries or **ML models** for improved duplicate detection.

### Real-Time vs. Batch Processing Trade-Offs

- Design hybrid systems that validate data on ingestion while refining it in batch processes.
- Monitor data processing **latency** and adjust windowing strategies as needed.

### Auditability and Compliance

- Regularly review **audit logs** and conduct reconciliation checks.
- Integrate **data quality frameworks** to meet compliance standards.

---

# Conclusion and Best Practices

### Document Everything

- Maintain thorough documentation of **cleaning rules**, **transformation logic**, and **audit logs**.

### Iterate and Improve

- Data cleaning is an ongoing process—continuously refine pipelines based on **feedback**.

### Balance Automation with Oversight

- Automated pipelines are efficient but require **human intervention** for complex cases.

### Invest in Tooling

- Utilize modern **ETL tools**, **cloud platforms**, and **machine learning** for advanced cleaning capabilities.

### Monitor Continuously

- Build **dashboards** and **alerts** to detect and resolve data quality issues in real time.