# Data Quality Exploratory Data Analysis (EDA) Report by Team 2

## Dataset Overview

This report reviews key data quality and structural issues identified across datasets related to learners, cohorts, marketing campaigns, and user demographics. The objective is to provide clear observations, risks, and recommended mitigation strategies to improve data integrity and analytic reliability.

Datasets Covered:

- Learner Dataset
- Content and Tracking Data
- Cohort Data
- Marketing Campaign Data
- Data Import Logs
- User Demographic Data

## Summary Stats

- The status field in the Learner dataset contains undocumented numeric codes.
- Duplicate learner_id and apply_date combinations detected.
- Multiple cohort assignments found per learner.
- Tracking data shows many NULL entries for questions and inconsistent naming.
- Cohort sizes range from unrealistically large to zero-day durations.
- Marketing campaigns have duplicate names and inactive campaigns with spend.
- Data import suffers from malformed rows and inconsistent ID formats.
- User demographics include missing birthdates, inconsistent gender values, and invalid emails.

## Missing & Duplicates

- Duplicate learner_id/apply_date pairs could bias enrollment counts.
- Missing critical fields such as tracking questions and demographic attributes.
- NULL and literal "NULL" values present in NOT NULL fields.
- Duplicate campaign names (e.g., "Copy 3", "Copy 4") affect marketing analysis.
- Malformed rows cause import failures and incomplete records.

## Outliers

- Epoch timestamps in UNIX format require conversion.
- Cohort sizes up to 100,000 and durations spanning years flagged.
- Marketing cost per result varies significantly; some campaigns show low engagement despite high reach.
- User demographic zip codes show inconsistent formats.
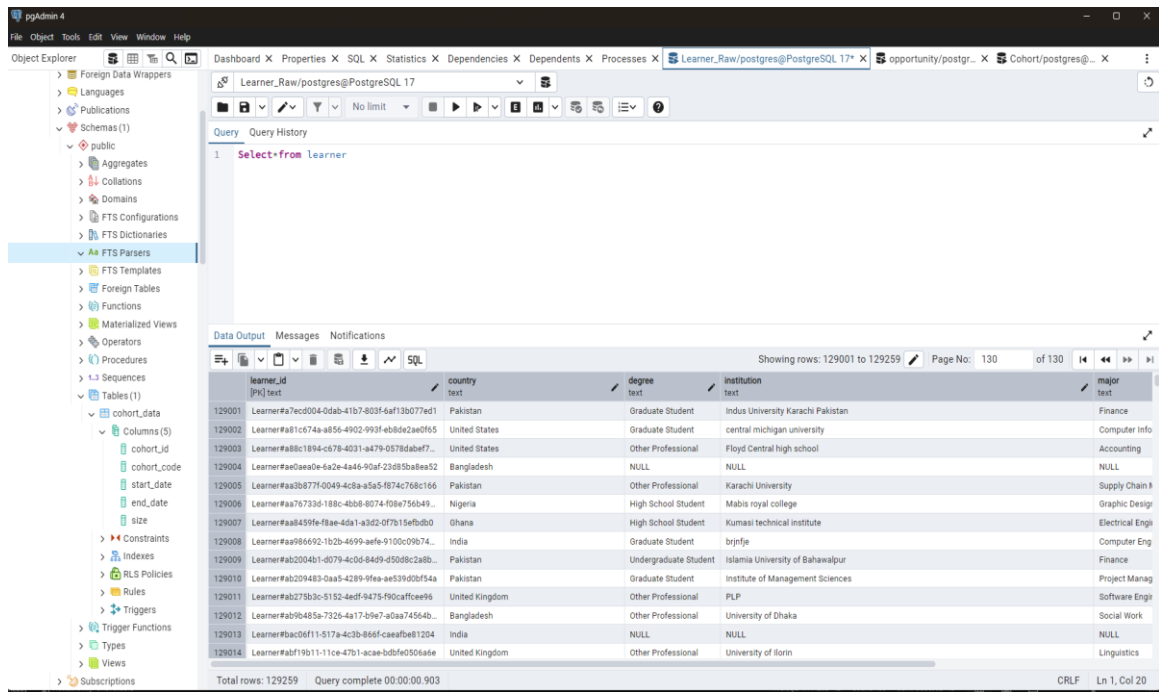- Outlier cohort durations and size discrepancies need validation.

## Trend Visuals

- Learner enrollment over time, highlighting duplicates.
- Cohort size and duration distribution.
- Marketing campaign spend vs results over periods.
- Demographic data completeness trends.
- Status code frequency and anomaly detection.

## Key Insights

- Enforce data validation rules and create data dictionaries for undocumented codes.
- Implement deduplication and uniqueness constraints on learner and cohort data.
- Standardize naming conventions and clean tracking data entries.
- Set thresholds and validation rules for cohort sizes and durations.
- Optimize marketing campaign lifecycle management to avoid wasted spend.
- Improve import procedures with schema validation and pre-processing.
- Normalize and enrich demographic data to improve analytic accuracy.
- Adopt consistent timestamp and ID formatting standards across datasets.

# 1. Image of Learner_Raw_Data successfully imported in PostgreSQL



# 2. Image of Opportunity_Data successfully Imported.

## 3. Image of Cohort_Data successfully imported



## 4. Image of Marketing_Data successfully imported

## 5. Image of LearnerOpportunity _Raw_Data successfully imported



## 6. Image of Cognito_Data successfully imported

# Trend Visuals Using Power BI



# Summarize Data Sample:

```sql
-- Replace NULL and blank assigned_cohort
UPDATE learneropportunity
SET assigned_cohort = (
    SELECT assigned_cohort
    FROM learneropportunity
    WHERE assigned_cohort IS NOT NULL AND TRIM(assigned_cohort) <> ''
    GROUP BY assigned_cohort
    ORDER BY COUNT(*) DESC
    LIMIT 1
)
WHERE assigned_cohort IS NULL OR TRIM(assigned_cohort) = '';

-- Replace NULL status
UPDATE learneropportunity
SET status = (
    SELECT status
    FROM learneropportunity
    WHERE status IS NOT NULL
    GROUP BY status
    ORDER BY COUNT(*) DESC
    LIMIT 1
)
WHERE status IS NULL;
```

```sql
SELECT
    q1,
    median_size,
    q3,
    (q3 - q1) AS iqr,
    (q1 - 1.5 * (q3 - q1)) AS lower_bound,
    (q3 + 1.5 * (q3 - q1)) AS upper_bound,
    min_size,
    max_size,
    avg_size,
    outlier_count,
    mode_cte.size AS mode_value
FROM size_stats, outliers, mode_cte;
```

```sql
--fill the null values with corresponding mode value
UPDATE "Cognito_Raw2"
SET
    gender = COALESCE(gender, 'Female'),
    city = COALESCE(city, 'Queens Village'),
    zip = COALESCE(zip, '11428'),
    state = COALESCE(state, 'NY')
WHERE
    gender IS NULL OR city IS NULL OR zip IS NULL OR state IS NULL;
```

Data Output   Messages   Notifications

UPDATE 42870

Object   SQL ✕   Statistics ✕   Dependencies ✕   Dependents ✕   Processes ✕   Execelerate/postgr... ✕   Execelerate/postgres@PostgreSQL 17* ✕

Execelerate/postgres@PostgreSQL 17

No limit

Query   Query History

```
1    --count the no. of duplicate records
2  ∨ SELECT COUNT(*) AS total_duplicate_records
3    FROM (
4      SELECT "user_id", "email", "gender", "UserCreateDate", "UserLastModifiedDate", "birthdate", "city", "zip",
5      "state",
6      COUNT(*)
7      AS duplicate_count
8      FROM "Cognito_Raw2"
9      GROUP BY "user_id", "email", "gender", "UserCreateDate", "UserLastModifiedDate", "birthdate", "city", "zip",
10     "state"
11     HAVING COUNT(*) > 1
12   ) AS duplicates;
13
```

Data Output   Messages   Notifications

Showing rows: 1 to 1   Page No: 1   of 1

| total_duplicate_records |
| bigint |
| 1 | 0 |

---

Cognito_Raw2
- Columns (9)
  - user_id
  - email
  - gender
  - UserCreateDate
  - UserLastModifiedDate
  - birthdate
  - city
  - zip
  - state
- Constraints
- Indexes
- RLS Policies
- Rules
- Triggers
- cohortraw
- learneropportunity
- learners
- marketing_campaigns
  - Columns
  - Constraints
  - Indexes
  - RLS Policies
  - Rules

Execelerate/postgres@PostgreSQL 17

Query   Query History                                          Scratch Pad ✕

```
1    --count the no. of null values in each column
2  ∨ SELECT
3      COUNT(*) FILTER (WHERE enrollment_id IS NULL) AS null_er
4      COUNT(*) FILTER (WHERE learner_id IS NULL) AS null_learn
5      COUNT(*) FILTER (WHERE assigned_cohort IS NULL OR TRIM(a
6      COUNT(*) FILTER (WHERE apply_date IS NULL) AS null_apply
7      COUNT(*) FILTER (WHERE status IS NULL) AS null_status
8    FROM learneropportunity;
9
```

Data Output   Messages   Notifications

Showing rows: 1 to 1   Page No: 1   of 1

| null_enrollment_id | null_learner_id | null_or_blank_assigned_cohort | null_apply_date | null_status |
| bigint | bigint | bigint | bigint | bigint |
| 1 | 0 | 0 | 13318 | 188 | 186 |