

WEEK 2

DATA TRANSFORMATION & MASTER TABLE CREATION

By

(SLU 0704 DVA | TEAM 15)



DATA VISUALIZATION ASSOCIATE EARLY INTERNSHIP

DURATION - 1 MONTH

APRIL - 2025

USER DATA REPORT

ETL REPORT

Dataset Name - Learner_Raw.csv

Records - 129,259

Columns - 5

- learner_id (Primary Key)
- country
- degree
- institution
- major

Observations -

- High percentage of missing values in degree, institution, and major.
- Case inconsistencies in institution (e.g., "Saint Louis University" vs "saint louis university").
- Duplicate records possible.

Table Design -

Table Name - learner_raw

Column Name	Data Type	Constraints	Description
learner_id	TEXT	PRIMARY KEY	Unique identifier for each learner
country	TEXT	CHECK (country <> "")	Standardized country name
degree	TEXT		Education status or degree
institution	TEXT		Cleaned institution name
major	TEXT		Cleaned major/field of study

Indexes -

- idx_country, idx_degree, idx_institution for faster queries.

Purpose - This table holds cleaned, standardized, and deduplicated learner records.

Table Creation SQL Script -

```

Query Query History
1 v CREATE TABLE learner_raw (
2   learner_id TEXT,
3   country TEXT,
4   degree TEXT,
5   institution TEXT,
6   major TEXT
7 );
8 v COPY learner_raw (learner_id, country, degree, institution, major)
9 FROM 'D:\Learner_Raw.csv'
10 DELIMITER ','
11 CSV HEADER;
12 SELECT * FROM learner_raw
13

```

Data Output Messages Notifications

learner_id text	country text	degree text	institution text	major text
Learner#00004f18-8b86-4fe4-ad7e-6c8d988f5335	Nigeria	Undergraduate Student	Federal University of Technology Owerri	Civil Engineering
Learner#00006478-745f-49bf-b126-02584e830720	Nigeria	NULL	NULL	NULL
Learner#00010567-1336-433c-a941-a612b3d2fb...	Kenya	Graduate Student	UNICAF UNIVERSITY	Environmental Sustain...
Learner#00011c80-0c5c-4601-9696-b2ca787e264f	Bangladesh	NULL	NULL	NULL
Learner#000141a7-4c82-401fa2e6-dd12b4b260...	Nigeria	NULL	NULL	NULL
Learner#0acf3501-57a8-4585-91e3-6bbb89d3c8...	Ghana	NULL	NULL	NULL
Learner#0001ca2c-7bec-4a33-833c-b844a29f4dea	Nigeria	Graduate Student	Nasarawa State University, Keffi	Accounting
Learner#0003bed9-d9d9-49a7-a755-a9562aaa0d...	Pakistan	Graduate Student	CTTI College KP Campus	Part Time

Total rows: 129259 Query complete 00:00:00.197 CRLF Ln 12, Col 1

Stored Procedure SQL Script -

```
Query Query History
1 CREATE OR REPLACE PROCEDURE populate_learner_master()
2 LANGUAGE plpgsql
3 AS $$*
4 BEGIN
5     DELETE FROM learner_master;
6     INSERT INTO learner_master (learner_id, country, degree, institution, major)
7     SELECT DISTINCT
8         learner_id,
9         INITCAP(TRIM(country)),
10        INITCAP(TRIM(degree)),
11        INITCAP(TRIM(institution)),
12        INITCAP(TRIM(major))
13    FROM learner_raw
14    WHERE learner_id IS NOT NULL;
15 END;
16 $$;

Data Output Messages Notifications
CREATE PROCEDURE

Query returned successfully in 69 msec.
```

Execution Steps -

1. Open pgAdmin → Query Tool.
2. Run the table creation script.
3. Run the stored procedure creation script.
4. Execute the procedure - **CALL populate_learner_master();**
5. Verify with - **SELECT * FROM learner_master LIMIT 10;**

DATA QUALITY REPORT

1. Issues Detected -

- Missing Values:
 - country - 2,275 missing
 - degree - 52,693 missing
 - institution - 53,073 missing
 - major - 52,871 missing
- Inconsistent Text Formats:
 - e.g., “Saint Louis University” vs. “saint louis university”
- Duplicate Records:
 - Not quantified directly but DISTINCT used to prevent them.

2. Cleaning Logic -

- Applied TRIM() to remove extra spaces.
- Used INITCAP() to standardize casing.
- Used DISTINCT to eliminate duplicates.
- Preserved NULLs where appropriate.

3. Validation Checks -

The screenshot shows a SQL query editor interface. The top section displays a multi-line SQL script with numbered lines (42 to 56) for validating learner data. The script includes calculations for null values and counts of learners per ID, followed by a distinct selection of institutions. The bottom section shows the resulting data output in a table format, with columns for row number, institution name, and a lock icon. The table contains 17 rows of data. A status bar at the bottom indicates 27186 total rows and a query completion time of 00:00:00.302.

	initcap text
9	@Pir Abdul Qadir Shah Jilani Institute Of Medical Science Gambat
10	+2 B M High School Pinjrawan
11	© 2023 Tomtom, © Openstreetmap Sri Venkateswara College Of Engineering
12	02 Academy
13	1
14	1. Hussaini Adamu Federal Polytechnic Kazaure. 2. Collage Of Health Sciences And Technology Tsafe.
15	10
16	100
17	10alvtics Academv.Uk

Total rows: 27186 Query complete 00:00:00.302

Challenges & Resolutions -

- Challenge - Text inconsistencies
Solution - Case normalization using INITCAP()
- Challenge - High null values
Solution - Preserved for data integrity, potential for enrichment later
- Challenge - Duplicate records
Solution - Removed via DISTINCT

Conclusion -

- The ETL process successfully cleaned and integrated raw learner data.
- The stored procedure ensures repeatability and scalability.
- Data is now ready for use in analysis or reporting.

OPPORTUNITY DATA REPORT

ETL REPORT

Dataset Name - Opportunity_Raw.csv

- Records - 187
- Columns - 5
 - opportunity_id (Primary Key)
 - opportunity_name
 - category
 - opportunity_code
 - tracking_questions

Missing Data -

- tracking_questions - 69 nulls (36.9%)
- All other fields: No missing values

Table Design -

Table Name - opportunity_raw

Column Name	Data Type	Constraints
opportunity_id	TEXT	PRIMARY KEY
opportunity_name	TEXT	NOT NULL
category	TEXT	NOT NULL
opportunity_code	TEXT	UNIQUE NOT NULL
tracking_questions	TEXT	

Table Creation SQL Script -

```

Query History
CREATE TABLE opportunity_raw (
    opportunity_id TEXT,
    opportunity_name TEXT,
    category TEXT,
    opportunity_code TEXT,
    tracking_questions TEXT
);

COPY opportunity_raw (opportunity_id, opportunity_name, category, opportunity_code, tracking_questions)
FROM 'D:\Opportunity_Raw.csv'
DELIMITER ','
CSV HEADER;

SELECT * FROM opportunity_raw

Data Output  Messages  Notifications
Showing rows: 1 to 187 | Page: 1 of 1
opportunity_id | opportunity_name | category | opportunity_code | tracking_questions
text          | text           | text     | text            | text
Opportunity#00000000G127E8VYE08TXBT6X | Choosing and Planning for Your Major | Event | E501873 | NULL
Opportunity#00000000G2PB6VB4ANR28CV2P | The Financial: Article Writing Competition Test | Competition | M523594 | NULL
Opportunity#00000000G4AM4J9BMPK3TJ... | Entrepreneurship and Innovation | Internship | I289641 | NULL
Opportunity#00000000G4F19XBEXPWK38F3N | Statement of Purpose (SOP) Writing Workshop | Event | E258709 | NULL
Opportunity#00000000G4KVEP36NNJR5YWTJ | Project Management | Internship | I584159 | NULL
Opportunity#00000000G8BW90E86ARRKM3... | Cybersecurity: Defensive Hacking | Internship | I155449 | NULL
Opportunity#00000000G8JG2FEA12SVNXXEN | Esports and Game Design | Internship | I860340 | NULL
Opportunity#00000000G95BD07NB0181K0XD | Data Visualization | Internship | I660879 | NULL

Total rows: 187 | Query complete 00:00:00.064

```

Stored Procedure SQL Script -

```
Query Query History
1 ✓ CREATE OR REPLACE PROCEDURE populate_opportunity_master()
2 LANGUAGE plpgsql
3 AS $$
4 BEGIN
5     DELETE FROM opportunity_master;
6     INSERT INTO opportunity_master (opportunity_id, opportunity_name, category, opportunity_code, tracking_questions)
7     SELECT DISTINCT
8         opportunity_id,
9         INITCAP(TRIM(opportunity_name)),
10        INITCAP(TRIM(category)),
11        UPPER(TRIM(opportunity_code)),
12        TRIM(tracking_questions)
13    FROM opportunity_raw
14    WHERE opportunity_id IS NOT NULL;
Data Output Messages Notifications
CREATE PROCEDURE
Query returned successfully in 46 msec.
```

Execution Steps -

- Open pgAdmin → Query Tool.
- Create the opportunity_master table.
- Create the stored procedure.
- Execute the procedure - **CALL populate_opportunity_master();**
- Verify with - **SELECT * FROM opportunity_master;**

DATA QUALITY REPORT

1. Issues Detected -

- tracking_questions has 36.9% missing values.
- Some opportunity names follow repetitive templates (e.g., “Entrepreneurship...”).
- No duplicates found in opportunity_id or opportunity_code.

2. Cleaning Logic -

- Trimmed all text fields using TRIM().
- Converted opportunity_name and category to Title Case using INITCAP().
- Standardized opportunity_code to uppercase using UPPER().

3. Validation Checks -

```
Query Query History
CALL populate_opportunity_master();
SELECT * FROM opportunity_master;
SELECT COUNT(*) FROM opportunity_raw;
SELECT COUNT(*) FROM opportunity_master;
SELECT COUNT(*) FROM opportunity_master WHERE tracking_questions IS NULL;
SELECT opportunity_id, COUNT(*) FROM opportunity_master GROUP BY opportunity_id HAVING COUNT(*) > 1;

Data Output Messages Notifications
Showing rows: 1 to 187 | Page No: 1 of 1 | << <> >> >>
opportunity_id [PK] opportunity_name category opportunity_code tracking_questions
text          text      text      text      text
Opportunity#00000000107DT690M8FCRX5A3S Graphic Design Associate Career ATK71AK {"serial_number": "1", "is_required_for_badge_award": "false", "code": "Q546OEQ", "question": "What is your favorite design software?"}
Opportunity#0000000010JB6ZKSEENQP5KCS Financial Controller Career A1K1ITG {"serial_number": "1", "is_required_for_badge_award": "false", "code": "QB3K3XR", "question": "What is your favorite financial management tool?"}
Opportunity#00000000109PRD2VKY0B574SMF Customer Acquisition And Retention Virtual Internship Internship ION83IM {"serial_number": "1", "is_required_for_badge_award": "true", "code": "OD4KHSI", "question": "What is your strategy for customer acquisition?"}
Opportunity#00000000102S9XBS7SM4KB0Q57 Secrets To Operational Excellence Course Course URMUHH7 {"serial_number": "1", "is_required_for_badge_award": "true", "code": "OQPPQPN", "question": "What are the secrets to operational excellence?"}
Opportunity#000000000G127EBVVE0BTXBT6X Choosing And Planning For Your Major Event E501873 NULL
Opportunity#0000000010Y6FBMAKXMP4PBW... Final Video Reflection - Your Journey, Captured Engageme... NSNCSUW {"serial_number": "1", "is_required_for_badge_award": "true", "code": "QNWK3PE", "question": "What was your final video reflection about your journey?"}
Opportunity#0000000010KMTAJ2JZCPRAHGW... The Happiness Project - Exploring The Science And Philosophy Of Well-Being Event EF8URPJ {"serial_number": "1", "is_required_for_badge_award": "true", "code": "QFPEBF", "question": "What did you learn from The Happiness Project?"}
Opportunity#0000000010FM369Z61B4WD1F4H The Creator%27s Journey: Turning Content Into Revenue Course U385067 {"serial_number": "1", "is_required_for_badge_award": "true", "code": "QKZ36J0", "question": "What is the key takeaway from The Creator's Journey?"}

Total rows: 187   Query complete 00:00:00.072
CRLF Ln 46, Col 1
```

Challenges & Resolutions -

- Challenge - High null values in tracking_questions
Solution - Preserved, assuming optional field
- Challenge - Text formatting inconsistencies
Solution - Used TRIM() and INITCAP()
- Challenge - Standardizing codes
Solution - Applied UPPER() for consistency

Conclusion -

- Data was successfully cleaned and loaded using PostgreSQL stored procedure.
- Results are standardized, validated, and ready for analysis or production usage.
- Procedure ensures consistency and reusability for future ETL runs.

COHORT DATA REPORT

ETL REPORT

Dataset Name - CohortRaw.csv

- Records - 639
- Columns - 5
 - cohort_id
 - cohort_code [Primary Key]
 - start_date
 - end_date
 - size

Observations -

- No missing values
- High variability in cohort size and duration
- Some cohorts last 0 days (likely same start & end date)
- Outliers in size and duration

Table Design -

Table Name - cohort

Column Name	Data Type	Constraints	Notes
cohort_id	TEXT	NOT NULL	Cohort reference ID
cohort_code	TEXT	PRIMARY KEY	Unique identifier
start_date	DATE	NOT NULL	Cohort start
end_date	DATE	NOT NULL	Cohort end
size	INTEGER	CHECK (size >= 0)	Number of participants
duration	INTEGER	GENERATED	Derived as end_date - start_date

Table Creation SQL Script -

Query Query History

```

1 ✓ CREATE TABLE cohort (
2   cohort_id TEXT,
3   cohort_code TEXT,
4   start_date TEXT,
5   end_date TEXT,
6   size TEXT
7 );
8 ✓ COPY cohort (cohort_id, cohort_code, start_date, size)
9 FROM 'D:\CohortRaw.csv'
10 DELIMITER ','
11 CSV HEADER;
12
13 SELECT * FROM cohort

```

Data Output Messages Notifications

cohort_id	cohort_code	start_date	end_date	size
1	B456514	1.6805E+12	1.68296E+12	1500
2	Cohort# B328821	1.67385E+12	1.67691E+12	1000
3	Cohort# B289256	1.6684E+12	1.67108E+12	100000
4	Cohort# B0VCBOF	1.66389E+12	1.66389E+12	40
5	Cohort# B908347	1.67324E+12	1.67631E+12	100000
6	Cohort# B306047	1.67324E+12	1.67631E+12	10000
7	Cohort# B883644	1.65665E+12	1.65924E+12	1500
8	Cohort# B280844	1.6684E+12	1.67108E+12	100000

Total rows: 639 Query complete 00:00:00.090

Stored Procedure SQL Script -

```
Query Query History
1 ✓ CREATE OR REPLACE PROCEDURE populate_cohort_master()
2 LANGUAGE plpgsql
3 AS $$
4 BEGIN
5     DELETE FROM cohort_master;
6     INSERT INTO cohort_master (cohort_id, cohort_code, start_date, end_date, size, duration)
7     SELECT
8         TRIM(cohort_id),
9         TRIM(cohort_code),
10        TO_DATE(start_date, 'YYYY-MM-DD'),
11        TO_DATE(end_date, 'YYYY-MM-DD'),
12        size,
13        (TO_DATE(end_date, 'YYYY-MM-DD') - TO_DATE(start_date, 'YYYY-MM-DD'))
14    FROM cohart
Data Output Messages Notifications
CREATE PROCEDURE
Query returned successfully in 46 msec.
```

Execution Steps -

- Create the cohort_master table.
- Create the stored procedure.
- Execute the procedure - **CALL populate_cohort_master();**
- Verify with - **SELECT * FROM cohort_master;**

DATA QUALITY REPORT

1. Issues Detected -

- No missing values
- Some durations = 0 days (possibly one-day sessions or data entry issues)
- Very large cohort sizes (up to 100,000) – possible outliers

2. Cleaning Logic -

- Trimmed IDs and codes to remove extra spaces
- Converted string dates to DATE using TO_DATE()
- Computed duration as end_date - start_date
- Preserved NULLs for fields where applicable (none in this case)

3. Validation Checks -

```
Query Query History
8
9
10
11
12    SELECT COUNT(*) FROM cohart;
13    SELECT COUNT(*) FROM cohort_master;
14
15
16    SELECT * FROM cohort_master WHERE duration < 0;
17    SELECT * FROM cohort_master WHERE duration = 0;
18
19    -- Size outlier check
20    SELECT * FROM cohort_master WHERE size > 10000;
21
Data Output Messages Notifications
count
bigint
1      639
```

Query Query History

```

1  SELECT *
2  FROM cohart
3  WHERE NOT (
4      start_date ~ '^\d{4}-\d{2}-\d{2}$' AND
5      end_date ~ '^\d{4}-\d{2}-\d{2}$'
6  );
7  SELECT COUNT(*) FROM cohart;
8  SELECT COUNT(*) FROM cohort_master;
9
10 SELECT * FROM cohort_master WHERE duration < 0;
11 SELECT * FROM cohort_master WHERE duration = 0;
12 SELECT * FROM cohort_master WHERE size > 10000;
13
14
15

```

Data Output Messages Notifications



	cohort_id text	cohort_code [PK] text	start_date text	end_date text	size integer
1	Cohort#	B456514	1.6805E+12	1.68296E+12	1500
2	Cohort#	B328821	1.67385E+12	1.67691E+12	1000
3	Cohort#	B289256	1.6684E+12	1.67108E+12	100000
4	Cohort#	B0VCB0F	1.66389E+12	1.66389E+12	40
5	Cohort#	B908347	1.67324E+12	1.67631E+12	100000
6	Cohort#	B306047	1.67324E+12	1.67631E+12	10000
7	Cohort#	B883644	1.65665E+12	1.65924E+12	1500
8	Cohort#	B280844	1.6684E+12	1.67108E+12	100000
9	Cohort#	B466039	1.66477E+12	1.66745E+12	100000
Total rows: 639		Query complete 00:00:00.096			

Challenges & Resolutions -

- Challenge - Duration = 0 days
Solution - Preserved; could be valid short programs
- Challenge - Large cohort sizes
Solution - Preserved; flagged as possible outliers
- Challenge - Format consistency in dates
Solution - Enforced using TO_DATE()

Conclusion -

- Data has been successfully cleaned, standardized, and loaded.
- ETL procedure is reusable and handles date and duration derivation.
- Dataset is now analytics-ready and can be used in dashboards or reporting.

LEARNER OPPORTUNITY DATA REPORT

ETL REPORT

Dataset Name - LearnerOpportunity_Raw.csv

- Records - 113,602
 - Columns - 5
- 1.enrollment_id (Primary Key)
 - 2.learner_id
 - 3.assigned_cohort
 - 4.apply_date
 - 5.status

Missing Values -

- assigned_cohort - 13,318 (11.72%)
- apply_date - 188 (0.17%)
- status - 186 (0.16%)

Uniqueness -

- Unique enrollments - 57,966
- Unique learners - 187 → ~310 enrollments per learner on average

Table Design -

Table Name - learner_opportunity

Column Name	Data Type	Constraints
enrollment_id	TEXT	PRIMARY KEY
learner_id	TEXT	NOT NULL
assigned_cohort	TEXT	
apply_date	DATE	
status	TEXT	

Table Creation SQL Script -

```

1 ✓ CREATE TABLE IF NOT EXISTS learner_opportunity (
2     enrollment_id TEXT,
3     learner_id TEXT,
4     assigned_cohort TEXT,
5     apply_date TEXT,
6     status TEXT
7 );
8 ✓ COPY learner_opportunity (enrollment_id, learner_id, assigned_cohort, apply_date, status)
9 FROM 'D:\LearnerOpportunity_Raw.csv'
10 DELIMITER ','
11 CSV HEADER;
12
13 SELECT * FROM learner_opportunity
14

```

Data Output Messages Notifications

enrollment_id	learner_id	assigned_cohort	apply_date	status
Leamer#4e6f78a9-f9b2-4352-ad22-d43dc46f5ff7	Opportunity#000000010WCBS50CYGDX97ES4	BAM6HBR	2024-04-10T06:28:31.902Z	1070
Leamer#4e79d245-3436-4fec-9906-901a03639a...	Opportunity#000000010WCBS50CYGDX97ES4	BAM6HBR	2023-11-15T03:08:17.442Z	1120
Leamer#4e9f9fc6b5-0576-4dbc-b7f5-1faef29b2df	Opportunity#000000010WCBS50CYGDX97ES4	BAM6HBR	2024-04-06T14:07:01.322Z	1070
Leamer#4ea61aa9-17da-4b60-9872-359b8e1e16...	Opportunity#000000010WCBS50CYGDX97ES4	BT4YTCR	2024-04-11T22:01:13.548Z	1070
Leamer#4eb218c7-467a-470a-9e3e-a2b7bc649e...	Opportunity#000000010WCBS50CYGDX97ES4	BT4YTCR	2024-10-22T15:44:13.402Z	1120
Leamer#4ec728db-7d09-4a8e-b1ab-dab3011a55...	Opportunity#000000010WCBS50CYGDX97ES4	BT4YTCR	2024-04-12T07:00:26.574Z	1070
Leamer#4edc150c-ea73-4144-993f-77ca8124de...	Opportunity#000000010WCBS50CYGDX97ES4	BT4YTCR	2024-11-24T11:07:11.742Z	1070
Leamer#4ef3715a-e420-4296-908c-eab9e5b797...	Opportunity#000000010WCBS50CYGDX97ES4	BC69M2K	2025-02-18T12:41:51.125Z	1070

Total rows: 113602 Query complete: 00:00:00.300

Stored Procedure SQL Script -

The screenshot shows a database interface with a query history tab and a SQL editor tab. The SQL editor contains the following code:

```
1 v CREATE OR REPLACE PROCEDURE populate_learner_opportunity_master()
2 LANGUAGE plpgsql
3 AS $$
4 BEGIN
5   DELETE FROM learner_opportunity_master;
6   INSERT INTO learner_opportunity_master (enrollment_id, learner_id, assigned_cohort, apply_date, status)
7   SELECT DISTINCT
8     TRIM(enrollment_id),
9     TRIM(learner_id),
10    TRIM(assigned_cohort),
11    CASE
12      WHEN apply_date ~ '^\d{4}-\d{2}-\d{2}$' THEN TO_DATE(apply_date, 'YYYY-MM-DD')
13      ELSE NULL
14    END,
15  END;
```

The Data Output tab shows the results of the query, which is an empty table with columns: enrollment_id, learner_id, assigned_cohort, apply_date, and status. The table has 8 rows, each with null values.

Execution Steps -

- Create the learner_opportunity_master table.
- Create the stored procedure.
- Execute the procedure - **CALL populate_learner_opportunity_master();**
- Verify with - **SELECT * FROM learner_opportunity_master LIMIT 10;**

DATA QUALITY REPORT

1. Issues Detected -

- 11.7% missing assigned_cohort
- Minor nulls in apply_date and status
- Some potential non-date strings in apply_date field
- Very high number of enrollments per learner — ~310 (expected behavior per system design)

2. Cleaning Logic -

- Trimmed all text fields
- Filtered valid date formats using regex before converting to DATE
- Removed any malformed apply_date entries
- Used SELECT DISTINCT to eliminate exact duplicates

3. Validation Checks -

The screenshot shows a database interface with a query history tab and a SQL editor tab. The SQL editor contains the following validation queries:

```
SELECT apply_date FROM learner_opportunity WHERE apply_date !~ '^\d{4}-\d{2}-\d{2}$';
SELECT
  COUNT(*) FILTER (WHERE assigned_cohort IS NULL) AS null_cohort,
  COUNT(*) FILTER (WHERE apply_date IS NULL) AS null_date,
  COUNT(*) FILTER (WHERE status IS NULL) AS null_status
FROM learner_opportunity_master;

SELECT enrollment_id, COUNT(*)
FROM learner_opportunity_master
GROUP BY enrollment_id
HAVING COUNT(*) > 1;
```

The Data Output tab shows the results of the validation queries. It includes two tables: one for null counts and one for learners with multiple enrollments.

enrollment_id	count
Learner#2444d3b7-3204-4b66-a1e2-72172db26...	3
Learner#d7b8c0bd-8fc7-4a9c-a617-369e3fc6ed...	2
Learner#6b9871b2-7830-4866-81ad-8674120eb...	2
Learner#725951b3-90ed-4494-9bb7-573721fcce...	2
Learner#0b248ca5-aa1e-49c6-b447-4b85701d4...	4

Challenges & Resolutions -

- Challenge - Inconsistent date formats
Solution - Used regex filter + TO_DATE() with fallback to NULL
- Challenge - High number of enrollments
Solution - Kept as-is; expected behavior from system design
- Challenge - Minor missing values
Solution - Preserved to reflect accurate source data

Conclusion -

- The ETL process successfully cleaned, standardized, and loaded learner opportunity data into the master table.
- The stored procedure ensures consistency, efficiency, and reusability.
- Cleaned dataset is ready for analysis, reporting, or integration into applications.

COGNITO DATA REPORT

ETL REPORT

Dataset Name - Cognito_Raw2.csv

- Records - 129,178
- Columns - 9
- 1. user_id [Primary Key]
- 2. email
- 3. gender
- 4. UserCreateDate
- 5. UserLastModifiedDate
- 6. birthdate
- 7. city
- 8. zip
- 9. state

Missing Values -

- state - 42,937 (~33%)
- zip - 42,869 (~33%)
- city - 42,866 (~33%)
- gender - 42,862 (~33%)
- birthdate - 42,862 (~33%)

Approximately 1 in 3 users lack complete demographic or geographic information.

Table Design -

Table Name - cognito

Column Name	Data Type	Constraints
user_id	TEXT	PRIMARY KEY
email	TEXT	NOT NULL UNIQUE
gender	TEXT	
birthdate	DATE	
age	INTEGER	
city	TEXT	
zip	TEXT	
state	TEXT	
created_at	TIMESTAMP	NOT NULL
last_modified_at	TIMESTAMP	

Table Creation SQL Script -

```
Query Query History
1 ✓ CREATE TABLE IF NOT EXISTS cognito (
2     user_id TEXT PRIMARY KEY,
3     email TEXT,
4     gender TEXT,
5     UserCreateDate TEXT,
6     UserLastModifiedDate TEXT,
7     birthdate TEXT,
8     city TEXT,
9     zip TEXT,
10    state TEXT
11 );
12 ✓ COPY cognito (user_id, email, gender, UserCreateDate, UserLastModifiedDate, birthdate, city, zip, state)
13 FROM 'D:\Cognito_Raw2.csv'
14 DELIMITER ','
15 CSV HEADER;

Data Output Messages Notifications
COPY 129178

Query returned successfully in 1 secs 558 msec.
```

Stored Procedure SQL Script -

```
Query Query History
1 ✓ CREATE OR REPLACE PROCEDURE populate_cognito_master()
2 LANGUAGE plpgsql
3 AS $$$
4 BEGIN
5     DELETE FROM cognito_master;
6     INSERT INTO cognito_master (
7         user_id, email, gender, birthdate, age, city, zip, state, created_at, last_modified_at
8     )
9     SELECT DISTINCT
10        TRIM(user_id),
11        TRIM(email),
12        INITCAP(TRIM(gender)),
13        CASE
14            WHEN birthdate ~ '^\d{4}-\d{2}-\d{2}$' THEN TO_DATE(birthdate, 'YYYY-MM-DD')
15            ELSE NULL
16        END,
17        CASE
18            WHEN birthdate ~ '^\d{4}-\d{2}-\d{2}$' THEN DATE_PART('year', AGE(TO_DATE(birthdate, 'YYYY-MM-DD')))
19            ELSE NULL
20        END,
21        INITCAP(TRIM(city)),
22        TRIM(zip),
23        INITCAP(TRIM(state)),
24        TO_TIMESTAMP(UserCreateDate, 'YYYY-MM-DD"T"HH24:MI:SS'),
25        TO_TIMESTAMP(UserLastModifiedDate, 'YYYY-MM-DD"T"HH24:MI:SS')
26     FROM cognito
27     WHERE user_id IS NOT NULL;
28 END;
$$;
```

Execution Steps -

- Run the table creation script.
- Create the stored procedure.
- Execute the procedure - **CALL populate_cognito_master();**
- Verify with - **SELECT * FROM cognito_master LIMIT 10;**

DATA QUALITY REPORT

1. Issues Detected -

- ~33% missing demographic/location fields
- Some invalid or malformed birthdates
- Age range includes outliers (e.g., 3 years, 101 years)
- Variations in case (e.g., “female”, “FEMALE”, “Female”)

2. Cleaning Logic -

- Trimmed and case-normalized all text fields
- Validated date formats using regex
- Computed age using birthdate → AGE()
- Converted ISO-like strings in date columns to proper TIMESTAMP
- Removed exact duplicates using DISTINCT

3. Validation Checks -

The screenshot shows a SQL query editor interface. The top section is labeled "Query History" and contains the following SQL code:

```
1  SELECT
2      COUNT(*) FILTER (WHERE gender IS NULL) AS null_gender,
3      COUNT(*) FILTER (WHERE city IS NULL) AS null_city,
4      COUNT(*) FILTER (WHERE birthdate IS NULL) AS null_birthdate
5  FROM cognito_master;
6
7  SELECT MIN(age), MAX(age), AVG(age) FROM cognito_master;
8  SELECT * FROM cognito_master WHERE age < 5 OR age > 90;
9
10 SELECT COUNT(*) FROM cognito;
11 SELECT COUNT(*) FROM cognito_master;
```

The bottom section is labeled "Data Output" and displays a single row of results:

	count
1	129178

Challenges & Resolutions -

- Challenge - Missing ~33% of demographic data
Solution - Preserved as-is; used NULLs for analysis integrity
- Challenge - Age range has outliers
Solution - Outliers flagged in analysis, not removed for now
- Challenge - Inconsistent date formats
Solution - Handled using TO_TIMESTAMP() and regex validation
- Challenge - Gender and location casing
Solution - Normalized using INITCAP()

Conclusion -

- Data successfully transformed, standardized, and loaded.
- Stored procedure ensures repeatable ETL pipeline.
- Cleaned dataset is analytics-ready with additional age, created_at, and last_modified_at fields.

MARKETING DATA REPORT

ETL REPORT

Dataset Name - Marketing Campaign Data All Accounts (2023-2024)(Detail1).csv

Records - 155

Columns - 13

- Ad Account Name
- Campaign name
- Delivery status
- Delivery level
- Reach
- Outbound clicks
- Landing page views
- Result type
- Results
- Cost per result
- Amount spent (AED)
- CPC (cost per link click)
- Reporting starts

Table Design -

Table Name - marketing

Column Name	Data Type	Description
ad_account_name	TEXT	Name of the ad account
campaign_name	TEXT	Campaign name
delivery_status	TEXT	Status of delivery
delivery_level	TEXT	Ad level (e.g., ad set, campaign)
reach	INTEGER	Total users reached
outbound_clicks	INTEGER	Number of outbound clicks
landing_page_views	INTEGER	Number of users who viewed the landing page
result_type	TEXT	Objective of the ad campaign
results	INTEGER	Measured results based on result type
cost_per_result	NUMERIC	Cost per result (calculated or provided)
amount_spent_aed	NUMERIC	Total ad spend in AED
cpc	NUMERIC	Cost per link click (CPC)
reporting_starts	DATE	Start date of reporting

Table Creation SQL Script -

```
Query Query History
1 CREATE TABLE IF NOT EXISTS marketing (
2     ad_account_name TEXT,
3     campaign_name TEXT,
4     delivery_status TEXT,
5     delivery_level TEXT,
6     reach TEXT,
7     outbound_clicks TEXT,
8     landing_page_views TEXT,
9     result_type TEXT,
10    results TEXT,
11    cost_per_result TEXT,
12    amount_spent_aed TEXT,
13    cpc TEXT,
14    reporting_starts TEXT
15 );
16

Data Output Messages Notifications
CREATE TABLE

Query returned successfully in 37 msec.
```

Stored Procedure SQL Script -

```
Query History
CREATE OR REPLACE PROCEDURE populate_marketing_master()
LANGUAGE plpgsql
AS $$

BEGIN
    DELETE FROM marketing_master;
    INSERT INTO marketing_master (
        ad_account_name, campaign_name, delivery_status, delivery_level,
        reach, outbound_clicks, landing_page_views, result_type,
        results, cost_per_result, amount_spent_aed, cpc, reporting_starts
    )
    SELECT DISTINCT
        TRIM("Ad Account Name"),
        TRIM("Campaign name"),
        INITCAP(TRIM("Delivery status")),
        INITCAP(TRIM("Delivery level")),
        NULLIF("Reach", '')::INTEGER,
        NULLIF("Outbound clicks", '')::INTEGER,
        NULLIF("Landing page views", '')::INTEGER,
        INITCAP(TRIM("Result type")),
        NULLIF("Results", '')::INTEGER,
        REGEXP_REPLACE("Cost per result", '[^0-9.]', '', 'g')::NUMERIC,
        REGEXP_REPLACE("Amount spent (AED)", '[^0-9.]', '', 'g')::NUMERIC,
        REGEXP_REPLACE("CPC (cost per link click)", '[^0-9.]', '', 'g')::NUMERIC,
        TO_DATE("Reporting starts", 'YYYY-MM-DD')
    FROM marketing;
END;
$$;
```

Execution Steps -

- Run the table creation script.
- Create the stored procedure.
- Execute the procedure - **CALL populate_marketing_master();**
- Verify with - **SELECT * FROM marketing_master LIMIT 10;**

DATA QUALITY REPORT

1. Missing Values -

- Ad Account Name - 7
- Campaign Name - 8
- Delivery status - 7
- Delivery level - 7
- Reach - 7
- Outbound clicks - 9
- Landing page views - 9
- Result type - 7
- Results - 6
- Cost per result - 7
- Amount spent (AED) - 6
- CPC - 8
- Reporting starts -7

~5–10% of rows have at least one missing metric, primarily in CPC or landing views.

2. Cleaning Logic -

- Trimmed and case-normalized text columns (TRIM, INITCAP)
- Converted currency fields like AED and CPC from strings with symbols (e.g., “AED 1,200”) to NUMERIC
- Handled blanks by converting them to NULL
- Removed or ignored rows with invalid or malformed numeric formats
- Preserved NULLs to reflect incomplete submissions

3. Validation Checks -

```
Query Query History
1 CALL populate_marketing_master();
2 SELECT * FROM marketing_master;
3
4 SELECT COUNT(*) FROM marketing;
5 SELECT COUNT(*) FROM marketing_master;
6
7 ▾ SELECT
8   COUNT(*) FILTER (WHERE campaign_name IS NULL) AS null_campaigns,
9   COUNT(*) FILTER (WHERE amount_spent_aed IS NULL) AS null_spend,
10  COUNT(*) FILTER (WHERE reporting_starts IS NULL) AS null_reporting_start
11 FROM marketing_master;
12
13 ▾ SELECT campaign_name, amount_spent_aed, cpc, cost_per_result
14   FROM marketing_master
15   ORDER BY amount_spent_aed DESC
16   LIMIT 10;
17
18 ▾ SELECT campaign_name, COUNT(*)
19   FROM marketing_master
20   GROUP BY campaign_name
21   HAVING COUNT(*) > 1;
```

```
Query Query History
1 SELECT cost_per_result, amount_spent_aed, cpc FROM marketing_master LIMIT 10;
2 SELECT * FROM marketing_master WHERE reach > 1000000 OR amount_spent_aed > 100000;
3 SELECT COUNT(*) FROM marketing_master WHERE landing_page_views IS NULL;
4 SELECT campaign_name, cpc FROM marketing_master ORDER BY cpc DESC LIMIT 10;
```

Data Output Messages Notifications

count	bigint
1	0

Challenges & Resolutions -

- Challenge - No Primary Key in Raw Data

Solution - Used DISTINCT to avoid duplicates; assumed combination of fields (e.g., campaign name + start date) as surrogate key for analysis.

- Challenge - Missing or Incomplete Records

Solution - Retained incomplete records by preserving NULLs for accurate reporting; flagged for potential future enrichment.

- Challenge - Currency Formatting Inconsistencies

Solution - Used REGEXP_REPLACE to strip currency symbols (e.g., AED, commas) and cast to NUMERIC.

- Challenge - Text-Numeric Mixing in Metrics Columns

Solution - Applied conversion logic with NULLIF and ::INTEGER/NUMERIC to handle blanks or malformed text entries.

- Challenge - Case Inconsistency in Text Fields (e.g., delivery status)

Solution - Normalized with INITCAP() to ensure consistency in reporting and grouping.

- Challenge - Date Format Validation

Solution - Used TO_DATE() with fallback to NULL for invalid or missing entries in reporting_starts.

- Challenge - Skewed Distribution in Spend/Reach

Solution - Identified outliers using histograms and range filters; preserved extreme values to maintain data integrity.

- Challenge - Ambiguity in Field Granularity (e.g., same campaign name repeated)

Solution - Treated repeated names as valid where other metrics differed (e.g., different spend or delivery level).

Conclusion -

- The ETL successfully cleaned and standardized the marketing performance dataset.

- Currency fields, numeric metrics, and date parsing were properly handled.

- The resulting table marketing_master is now clean, structured, and ready for marketing efficiency analysis.

MASTER TABLE REPORT

DATASETS OVERVIEW -

Dataset Name	PostgreSQL Table	Description	Primary Key
Learner_Raw.csv	learner_raw	Learner education details (degree, institution, major)	learner_id
Cognito_Raw2.csv	cognito	Learner demographic data (email, gender, age, location)	user_id
LearnerOpportunity_Raw.csv	learner_opportunity	Mapping of learner to cohort + application status	enrollment_id
CohortRaw.csv	cohort	Details about cohorts (batch, duration, size)	cohort_code
Opportunity_Raw.csv	opportunity_raw	Opportunity data (type, name, tracking questions)	opportunity_code
Marketing Campaign Data.csv	marketing	Ad performance metrics	<i>No clear primary key</i>

🛠 STEP-BY-STEP: CREATING THE MASTER TABLE

✓ Step 1 - Create Required Tables

All raw CSV files are first imported into PostgreSQL as individual tables (learner_raw, cognito, etc.)

✓ Step 2 - Create the Master Table Using JOINs

```

Query  Query History
1  ✓ CREATE TABLE master_table AS
2   SELECT
3     lo.enrollment_id,
4     lr.learner_id,
5     lr.degree,
6     lr.institution,
7     lr.major,
8     lo.assigned_cohort,
9     ch.start_date AS cohort_start_date,
10    ch.end_date AS cohort_end_date,
11    ch.size AS cohort_size,
12    lo.apply_date,
13    lo.status,
14    cg.email,
15    cg.gender,
16    cg.UserCreateDate,
17    cg.UserLastModifiedDate,
18    cg.birthdate,
19    cg.city,
20    cg.zip,
21    cg.state,
22    op.opportunity_id,
23    op.opportunity_name,
24    op.category,
25    op.opportunity_code,
26    op.tracking_questions
27   FROM learner_opportunity lo
28   LEFT JOIN learner_raw lr ON lo.learner_id = lr.learner_id
29   LEFT JOIN cohort ch ON lo.assigned_cohort = ch.cohort_code
30   LEFT JOIN cognito cg ON lr.learner_id = cg.user_id
31   LEFT JOIN opportunity_raw op ON lo.learner_id = op.opportunity_id;

```

✓ Step 3 - Applying Primary & Foreign Keys to the Master Table

```
Query Query History
1 ✓ CREATE TABLE final_master_table (
2   ---cognito_raw
3     user_id TEXT NOT NULL,
4     email TEXT,
5     gender TEXT,
6     birthdate TEXT,
7     age TEXT,
8     city TEXT,
9     zip TEXT,
10    ----learner_raw
11      learner_id TEXT NOT NULL,
12      country TEXT,
13      degree TEXT,
14      institution TEXT,
15      major TEXT,
16      ----Opportunity_raw
17      opportunity_id TEXT NOT NULL,
18      opportunity_name TEXT,
19      category TEXT,
20      opportunity_code TEXT,
21      ----cohort_raw
22      cohort_code TEXT NOT NULL,
23      start_date TEXT,
24      end_date TEXT,
25      cohort_size TEXT,
26      ----learneropportunity_raw
27      enrollment_id TEXT NOT NULL,
28      opportunity_lp_id TEXT NOT NULL,
29      assigned_cohort TEXT,
30      status TEXT,
31      PRIMARY KEY (user_id, learner_id, opportunity_id, cohort_code),
32      FOREIGN KEY (assigned_cohort) REFERENCES cohort_raw_staging(cohort_code),
33      FOREIGN KEY (enrollment_id) REFERENCES learner_raw_staging(learner_id_ref),
34      FOREIGN KEY (opportunity_lp_id) REFERENCES opportunity_raw_staging(opportunity_id_ref)
35    );
```

⚠ WHY MARKETING DATASET CANNOT BE JOINED

Problem -

- This dataset contains aggregated ad performance (e.g., CPC, results, cost per result) with no learner identifiers.

No Join Key -

None of these columns relate to learner IDs, enrollment, or opportunity data.

Recommendation -

Use it for separate campaign dashboards — not as part of master_table.

FIELDS USED

- enrollment_id, assigned_cohort, apply_date, status (learner_opportunity)
- learner_id, degree, institution, major, country (learner_raw)
- start_date, end_date, size (cohort)
- email, gender, birthday, city, zip, state, etc. (cognito)
- opportunity_id, opportunity_name, category, tracking_questions (opportunity_raw)

STORED PROCEDURE SQL SCRIPT

Query Query History Data Output Messages Notifications

```
38 ----- Inserting data into master table
39 v CREATE OR REPLACE PROCEDURE sp_create_final_master_table()
40 LANGUAGE plpgsql
41 AS $$
42 BEGIN
43     -- Step 1: Clear the table (if it exists)
44     DELETE FROM final_master_table;
45
46 v     INSERT INTO final_master_table
47     SELECT DISTINCT
48         TRIM(c.user_id) AS user_id,
49         LOWER(TRIM(c.gender)) AS gender,
50         c.birthdate,
51         EXTRACT(YEAR FROM AGE(CURRENT_DATE, CAST(birthdate AS DATE))) AS age,
52         LOWER(TRIM(c.city)) AS city,
53         LOWER(TRIM(c.zip)) AS zip,
54
55         TRIM(l.learner_id_ref) AS learner_id,
56         LOWER(TRIM(l.country)) AS country,
57         LOWER(TRIM(l.degree)) AS degree,
58         LOWER(TRIM(l.institution)) AS institution,
59         LOWER(TRIM(l.major)) AS major,
60
61         TRIM(o.opportunity_id_) AS opportunity_id,
62         LOWER(TRIM(o.opportunity_name)) AS opportunity_name,
63         LOWER(TRIM(o.category)) AS category,
64         LOWER(TRIM(o.opportunity_code)) AS opportunity_code,
65
66         LOWER(TRIM(ch.cohort_code)) AS cohort_code,
67         TO_TIMESTAMP(CAST(ch.start_date AS DOUBLE PRECISION) / 1000.0, 0) AS start_date,
68         TO_TIMESTAMP(CAST(ch.end_date AS DOUBLE PRECISION) / 1000.0, 0) AS end_date,
69         ch.size AS cohort_size,
```

```
77     FROM cognito_raw_staging c
78     inner join learner_raw_staging l on c.user_id= l.learner_id_ref
79     inner join learneropportunity_raw_staging lo on lo.learner_enrollment_id=l.learner_id_ref
80     inner join opportunity_raw_staging o on lo.opportunity_id = o.opportunity_id_ref
81     inner join cohort_raw_staging ch on ch.cohort_code=lo.assigned_cohort;
82
83 v     WHERE
84         c.user_id <> 'NULL' AND c.email <> 'NULL' AND l.learner_id <> 'NULL' AND o.opportunity_id <> 'NULL'
85         AND lo.enrollment_id <> 'NULL' AND ch.cohort_code <> 'NULL' AND lo.status <> 'NULL'
86         AND ch.size <> 'NULL' AND ch.start_date <> 'NULL' AND ch.end_date <> 'NULL'
87         AND c.birthdate <> 'NULL';
88
89     -- Create indexes for optimization
90     CREATE INDEX idx_user_id ON final_master_table(user_id);
91     CREATE INDEX idx_learner_id ON final_master_table(learner_id);
92     CREATE INDEX idx_opportunity_id ON final_master_table(opportunity_id);
93     CREATE INDEX idx_cohort_code ON final_master_table(cohort_code);
94     CREATE INDEX idx_enrollment_id ON final_master_table(enrollment_id);
95 END;
96 $$;
97
98 CALL sp_create_final_master_table();
```

EXECUTION STEPS -

- Run the table creation script.
- Create the stored procedure.
- Execute the procedure - **CALL sp_create_final_master_table();**
- Verify with - **SELECT * FROM final_master_table LIMIT 10;**

DATA QUALITY REPORT

1. Completeness Checks

Table	Column	% Missing
learner_raw	degree	40.7%
learner_raw	institution	41.0%
learner_raw	major	40.9%
cognito	gender, city, state, zip	~33%
learner_opportunity	assigned_cohort	11.7%
learner_opportunity	apply_date	0.17%

2. Consistency Checks

Check	Query	Result
Unique learner IDs	<code>SELECT COUNT(DISTINCT learner_id) FROM learner_raw;</code>	✓ Passed
Unique enrollment IDs	<code>SELECT COUNT(DISTINCT enrollment_id) FROM learner_opportunity;</code>	✓ Passed
Duplicate check in master table	<code>SELECT enrollment_id, COUNT(*) ... HAVING COUNT(*) > 1;</code>	✓ No duplicates

3. Referential Integrity

All JOINS in the master table were verified using LEFT JOIN and null-checks -

```
SELECT COUNT(*) FROM final_master_table WHERE learner_id IS NULL;  
SELECT COUNT(*) FROM final_master_table WHERE email IS NULL;
```

- Learner-Cognito match rate: ~66%
- Learner-Cohort match rate: ~88%

4. Validation

📌 Date Range Checks -

```
SELECT MIN(apply_date), MAX(apply_date) FROM final_master_table;  
SELECT MIN(birthdate), MAX(birthdate) FROM final_master_table;
```

- Apply dates range from June 2022 to late 2024
- Birthdates range from 1922 to 2021 (some outliers exist: ages < 5 or > 100)

Gender Distribution -

SELECT gender, COUNT(*) FROM final_master_table GROUP BY gender;

- Female: Majority
- Male: Present
- Null: ~33%

Status Distribution -

SELECT status, COUNT(*) FROM final_master_table GROUP BY status;

- Status 1070: Dominates (~67%)
- Others: 1030, 1055, 1120 are secondary

5. Marketing Data

Due to absence of foreign keys or joinable columns -

- marketing dataset was excluded from the final_master_table.
- This dataset is analyzed separately using stand-alone dashboards and summaries.

6. Validation Checks

- Check if the Table Exists
- Preview Sample Data
- Count Total Records
- Check for Nulls in Key Fields
- Check for Duplicate Enrollments
- Verify Data Types
- Status Distribution check

Data Output																	
Showing rows: 1 to 1000																	
Page No: 1 of 100																	
user_id	gender	birthdate	age	city	zip	learner_id	country	degree	institution	major	opportunity_id	text	text	text	text	text	text
00004f18-8086-4fe4-ad7e-6c8d988f53...	male	6/23/2001	23	owerri	460103	00004f18-8b86-4fe4-ad7e-6c8d988f53...	nigeria	undergraduate student	federal university of technology owerri	civil engineering	000000000GWQAXC5X						
00010567-1336-433c-9a41-e612bd2d7...	female	5/4/1996	28	naivasha	20117	00010567-1336-433c-9a41-e612bd2d7...	kenya	graduate student	unicaf university	environmental sustainability	00000000100lPM3AD0						
0001ca2c-7bec-4a33-833c-b844a2914...	male	6/6/1986	38	abuja	900211	0001ca2c-7bec-4a33-833c-b844a2914...	nigeria	graduate student	nasarawa state university, keffi	accounting	000000000GN2AOY7X						
0001ca2c-7bec-4a33-833c-b844a2914...	male	6/6/1986	38	abuja	900211	0001ca2c-7bec-4a33-833c-b844a2914...	nigeria	graduate student	nasarawa state university, keffi	accounting	0000000010BDV2YMK						
0001ca2c-7bec-4a33-833c-b844a2914...	male	6/6/1986	38	abuja	900211	0001ca2c-7bec-4a33-833c-b844a2914...	nigeria	graduate student	nasarawa state university, keffi	accounting	0000000010GJ1ZMA						
0003bed9-d9d9-49a7-a755-a9562aaa0...	male	4/12/1999	26	khanur	64100	0003bed9-d9d9-49a7-a755-a9562aaa0...	pakistan	graduate student	ctti college kp campus	part time	0000000010SAZKDAE						
0004295c-717e-4953-b3bf-f22daff9e903	male	7/8/2002	22	ikorodu	101242	0004295c-717e-4953-b3bf-f22daff9e903	nigeria	undergraduate student	caleb university	computer science	0000000010RQJKA9N						
00049a81-94a9-4b25-92ed-62d017f3b...	female	8/11/1988	36	antipolo	1870	00049a81-94a9-4b25-92ed-62d017f3b...	philippines	undergraduate student	our lady of fatima university	nursing	0000000010GJ1ZMA						
00048381-3917-4dd3-9639-a501aa08...	female	9/1/2007	17	barwani	451551	00048381-3917-4dd3-9639-a501aa08...	india	high school student	saket international school	medicine	0000000010WCBSS0C						
000926e9-66af-4e40-a788-538e74b9a...	male	7/10/1999	25	wah cantt	47100	000926e9-66af-4e40-a788-538e74b9a...	pakistan	graduate student	air university	marketing	0000000010AWJ1XAB						
0009e8ae-95a0-4b58-b5dd-6bb90d914...	female	8/30/2000	24	dhangadhi	10901	0009e8ae-95a0-4b58-b5dd-6bb90d914...	nepal	undergraduate student	global college international	travel and tourism	0000000010GH4NB8Q						
000bce06-dbee-4d9f-98bc-26ed26ab9...	male	5/18/1999	25	kottayam	686501	000bce06-dbee-4d9f-98bc-26ed26ab9...	india	undergraduate student	mg university	cyber security	0000000010ZKNWJ4						
000eb279-b2a9-4891-ad00-e284d52a0...	male	7/13/2005	19	islamabad	44100	000eb279-b2a9-4891-ad00-e284d52a0...	pakistan	undergraduate student	fast nuces, islamabad	computer science	0000000010GGH4NB8Q						
000e3b4-383d-4412-a5a2-e80766179...	male	7/22/2000	24	future city	4911010	000e3b4-383d-4412-a5a2-e80766179...	egypt	undergraduate student	seif hitham nabil salah	undergraduate	0000000010GNTFT74M						
000ff2e3-c653-4153-b6b9-3b4daee7f1c...	male	6/1/1996	28	surulere	101283	000ff2e3-c653-4153-b6b9-3b4daee7f1c...	nigeria	undergraduate student	university of lagos	philosophy	0000000010GJ1ZMA						
00108b20-0435-4ad7-9d27-8878f635d...	male	2/19/2003	22	ghaziabad	201009	00108b20-0435-4ad7-9d27-8878f635d...	india	undergraduate student	abes engineering college	computer engineering	0000000010WQAXC5X						
0012caf1-83ee-4a25-9f78-71dfcb18e...	female	9/16/2005	19	agra	282001	0012caf1-83ee-4a25-9f78-71dfcb18e...	india	undergraduate student	baikunthi devi kanya mahavidyalaya	business administration	0000000010GH4NB8Q						
0012caf1-83ee-4a25-9f78-71dfcb18e...	female	9/16/2005	19	agra	282001	0012caf1-83ee-4a25-9f78-71dfcb18e...	india	undergraduate student	baikunthi devi kanya mahavidyalaya	business administration	00000000106FMJNZK						
0012caf1-83ee-4a25-9f78-71dfcb18e...	female	9/16/2005	19	agra	282001	0012caf1-83ee-4a25-9f78-71dfcb18e...	india	undergraduate student	baikunthi devi kanya mahavidyalaya	business administration	0000000010106DCAK3F						
0012caf1-83ee-4a25-9f78-71dfcb18e...	female	9/16/2005	19	agra	282001	0012caf1-83ee-4a25-9f78-71dfcb18e...	india	undergraduate student	baikunthi devi kanya mahavidyalaya	business administration	0000000010SAZKDAE						
0012caf1-83ee-4a25-9f78-71dfcb18e...	female	9/16/2005	19	agra	282001	0012caf1-83ee-4a25-9f78-71dfcb18e...	india	undergraduate student	baikunthi devi kanya mahavidyalaya	business administration	0000000010YMF7EFAP						
00140b73-0f1a-4d50-a8c5-c969bf1d5...	male	4/26/2003	21	delhi	110053	00140b73-0f1a-4d50-a8c5-c969bf1d5...	india	undergraduate student	manav rachna international institute of resea...	nutrition and dietetics	0000000010WCBS50C						
0015182a-3e78-4986-9c9f-6e1058887f...	female	9/14/2002	22	iligan city	9200	0015182a-3e78-4986-9c9f-6e1058887f...	philippines	undergraduate student	msu-it	nursing	000000001045Z1BFR6						
0015182a-3e78-4986-9c9f-6e1058887f...	female	9/14/2002	22	iligan city	9200	0015182a-3e78-4986-9c9f-6e1058887f...	philippines	undergraduate student	msu-it	nursing	0000000010WCBS50C						
00158775-51aa-4848-923c-8c8cf3cb9...	female	8/18/2002	22	lagos	104102	00158775-51aa-4848-923c-8c8cf3cb9...	nigeria	undergraduate student	university of lagos	health information manageme...	000000000GH4NB8Q						
00159faf-a845-48ea-95e1-ee479421cd...	female	7/4/2005	19	lahore	54000	00159faf-a845-48ea-95e1-ee479421cd...	pakistan	high school student	superior group of colleges	medicine	0000000010WCBS50C						

Query History

```

1  SELECT COUNT(*) FROM master_table;
2  ✓ SELECT
3      COUNT(*) FILTER (WHERE learner_id IS NULL) AS null_learners,
4      COUNT(*) FILTER (WHERE assigned_cohort IS NULL) AS null_cohorts,
5      COUNT(*) FILTER (WHERE email IS NULL) AS null_emails
6  FROM master_table;
7  ✓ SELECT enrollment_id, COUNT(*)
8  FROM master_table
9  GROUP BY enrollment_id
10 HAVING COUNT(*) > 1;

```

Data Output Messages Notifications

Showing rows: 1 to 100

	enrollment_id	count
	text	bigint
25	Learner#34a0a304-ffc1-4ea4-8414-4682119a4f...	2
26	Learner#b9b208bc-ee2c-42dd-a5f5-1f517da31a...	2
27	Learner#6c8afa3e-ef03-4ae6-a688-8376e4fb6fe9	2
28	Learner#78dc9e53-1641-4051-975b-acd2c596d...	6
29	Learner#432c7d97-3707-4bfa-871c-116a75d7e...	2
30	Learner#164e08fa-a725-4aaa-94c9-15a58d8dda...	2
31	Learner#b8bdac15-50de-4a70-8dc9-9d5fc75339...	4
32	Learner#1a123198-726f-490f-8490-1a22c092cd...	2
33	Learner#6a9b7adb-f436-4409-a883-2d40a85f17...	15
34	Learner#3a9a8617-1d1c-40c7-b4a3-ffd27c517f...	3
35	Learner#6be4c7d4-0d2c-40d5-802c-5d745a0d8...	2
36	Learner#fe242c76-03fb-42cf-b7d6-e266481b3b...	2

Total rows: 20572 Query complete 00:00:00.132

Query History Data Output Messages

Execute script F5

	user_id	learner_id_ref	opportunity_id_ref	learner_enrollment_id	opportunity_id
	text	text	text	text	text
1	54898df1-8a76-4d64-91ae-46df03211d9e	54898df1-8a76-4d64-91ae-46df03211d9e	0000000010WCBS50CYGDX97ES4	54898df1-8a76-4d64-91ae-46df03211d9e	0000000010WCBS50CYGDX97ES4
2	5499900e-9030-49b8-a81d-1294aa61fd...	5499900e-9030-49b8-a81d-1294aa61fd...	0000000010WCBS50CYGDX97ES4	5499900e-9030-49b8-a81d-1294aa61fd...	0000000010WCBS50CYGDX97ES4
3	54b07361-3334-46e8-b880-9eceaa3cb54...	54b07361-3334-46e8-b880-9eceaa3cb54...	0000000010WCBS50CYGDX97ES4	54b07361-3334-46e8-b880-9eceaa3cb54...	0000000010WCBS50CYGDX97ES4
4	135b4498-9bf3-416f-a357-be4753888e14	135b4498-9bf3-416f-a357-be4753888e14	000000000G4AM4J9NBMPK3T...	135b4498-9bf3-416f-a357-be4753888e14	000000000G4AM4J9NBMPK3T...
5	87e80e08-09ba-42c8-98de-59c309a71e...	87e80e08-09ba-42c8-98de-59c309a71e...	0000000010WCBS50CYGDX97ES4	87e80e08-09ba-42c8-98de-59c309a71e...	0000000010WCBS50CYGDX97ES4
6	54f29f7e-7d08-4164-9643-2455ce3085ce	54f29f7e-7d08-4164-9643-2455ce3085ce	0000000010WCBS50CYGDX97ES4	54f29f7e-7d08-4164-9643-2455ce3085ce	0000000010WCBS50CYGDX97ES4
7	551efa40-c68e-4ca3-a2ac-c2efb5f37d6c	551efa40-c68e-4ca3-a2ac-c2efb5f37d6c	0000000010WCBS50CYGDX97ES4	551efa40-c68e-4ca3-a2ac-c2efb5f37d6c	0000000010WCBS50CYGDX97ES4
8	5536dc67-8182-4e02-9ee8-1b2c609df2...	5536dc67-8182-4e02-9ee8-1b2c609df2...	0000000010WCBS50CYGDX97ES4	5536dc67-8182-4e02-9ee8-1b2c609df2...	0000000010WCBS50CYGDX97ES4
9	5537fb70-e218-432b-a606-f68d05e47609	5537fb70-e218-432b-a606-f68d05e47609	0000000010WCBS50CYGDX97ES4	5537fb70-e218-432b-a606-f68d05e47609	0000000010WCBS50CYGDX97ES4
10	55449ca1-d46a-4679-9e1c-20de13a3fce	55449ca1-d46a-4679-9e1c-20de13a3fce	0000000010WCBS50CYGDX97ES4	55449ca1-d46a-4679-9e1c-20de13a3fce	0000000010WCBS50CYGDX97ES4
11	12dd9c42-ebe2-4ded-b84d-9e43d91ba8...	12dd9c42-ebe2-4ded-b84d-9e43d91ba8...	000000000GWQAXC5X45C2MH...	12dd9c42-ebe2-4ded-b84d-9e43d91ba8...	000000000GWQAXC5X45C2MH...
12	557127c7-a35e-4a57-bf5a-ebf6d1e79de4	557127c7-a35e-4a57-bf5a-ebf6d1e79de4	0000000010WCBS50CYGDX97ES4	557127c7-a35e-4a57-bf5a-ebf6d1e79de4	0000000010WCBS50CYGDX97ES4
13	2591fb35-fc49-4a62-927a-606be516a8b0	2591fb35-fc49-4a62-927a-606be516a8b0	000000000GN2A0AY7XK8C5FZPP	2591fb35-fc49-4a62-927a-606be516a8b0	000000000GN2A0AY7XK8C5FZPP
14	55808aca8-b461-4bad-8ebd-cb7e937cc6...	55808aca8-b461-4bad-8ebd-cb7e937cc6...	0000000010WCBS50CYGDX97ES4	55808aca8-b461-4bad-8ebd-cb7e937cc6...	0000000010WCBS50CYGDX97ES4
15	558abf2a-6297-458c-90d9-47bfdeba57fe	558abf2a-6297-458c-90d9-47bfdeba57fe	0000000010WCBS50CYGDX97ES4	558abf2a-6297-458c-90d9-47bfdeba57fe	0000000010WCBS50CYGDX97ES4
16	559a4f8f-716d-4c4c-bb16-0530eef9cd72	559a4f8f-716d-4c4c-bb16-0530eef9cd72	0000000010WCBS50CYGDX97ES4	559a4f8f-716d-4c4c-bb16-0530eef9cd72	0000000010WCBS50CYGDX97ES4
17	559b141c-7d4f-4e88-8339-a2f3e026f5a5	559b141c-7d4f-4e88-8339-a2f3e026f5a5	0000000010WCBS50CYGDX97ES4	559b141c-7d4f-4e88-8339-a2f3e026f5a5	0000000010WCBS50CYGDX97ES4
18	55a600a4-bdf4-4ac5-96a7-a79b5f315642	55a600a4-bdf4-4ac5-96a7-a79b5f315642	0000000010WCBS50CYGDX97ES4	55a600a4-bdf4-4ac5-96a7-a79b5f315642	0000000010WCBS50CYGDX97ES4
19	55a72611-9632-403c-8804-83e6bfe7e9...	55a72611-9632-403c-8804-83e6bfe7e9...	0000000010WCBS50CYGDX97ES4	55a72611-9632-403c-8804-83e6bfe7e9...	0000000010WCBS50CYGDX97ES4
20	55b133e7-bbc8-4ccc-b133-89c4c54e14...	55b133e7-bbc8-4ccc-b133-89c4c54e14...	0000000010WCBS50CYGDX97ES4	55b133e7-bbc8-4ccc-b133-89c4c54e14...	0000000010WCBS50CYGDX97ES4
21	55c4b07d-92b7-4a7c-a0ab-2be8e2a989...	55c4b07d-92b7-4a7c-a0ab-2be8e2a989...	0000000010WCBS50CYGDX97ES4	55c4b07d-92b7-4a7c-a0ab-2be8e2a989...	0000000010WCBS50CYGDX97ES4
22	55f556b3-dcda-4ee1-af9f-4cef1f5f9805	55f556b3-dcda-4ee1-af9f-4cef1f5f9805	0000000010WCBS50CYGDX97ES4	55f556b3-dcda-4ee1-af9f-4cef1f5f9805	0000000010WCBS50CYGDX97ES4
23	56090ae5-4559-4a24-a40a-4283b26521...	56090ae5-4559-4a24-a40a-4283b26521...	0000000010WCBS50CYGDX97ES4	56090ae5-4559-4a24-a40a-4283b26521...	0000000010WCBS50CYGDX97ES4
24	566ae0a7-82db-4a8a-ab31-58e2482704...	566ae0a7-82db-4a8a-ab31-58e2482704...	0000000010WCBS50CYGDX97ES4	566ae0a7-82db-4a8a-ab31-58e2482704...	0000000010WCBS50CYGDX97ES4
25	566c7f0c-6fb4-4cc2-b660-23217372f93c	566c7f0c-6fb4-4cc2-b660-23217372f93c	0000000010WCBS50CYGDX97ES4	566c7f0c-6fb4-4cc2-b660-23217372f93c	0000000010WCBS50CYGDX97ES4
26	568f977b-2c04-4efb-a85a-63ebefc3aafa	568f977b-2c04-4efb-a85a-63ebefc3aafa	0000000010WCBS50CYGDX97ES4	568f977b-2c04-4efb-a85a-63ebefc3aafa	0000000010WCBS50CYGDX97ES4
27	569c2f8a-4aa2-4f02-a890-8dd642beb421	569c2f8a-4aa2-4f02-a890-8dd642beb421	0000000010WCBS50CYGDX97ES4	569c2f8a-4aa2-4f02-a890-8dd642beb421	0000000010WCBS50CYGDX97ES4

Total rows: 99875 Query complete 00:00:00.698

A screenshot of a database query results window. The table has three columns: opportunity_id_ref (text), opportunity_id (text), and count (bigint). The data shows various IDs and their counts, such as 000000000G127E8VYE08TGBT6X with a count of 13, and 000000000G8BW90E86ARRKM3... with a count of 427.

	opportunity_id_ref text	opportunity_id text	count bigint
1	000000000G127E8VYE08TGBT6X	000000000G127E8VYE08TGBT6X	13
2	000000000G2PB6VB4ANR28CV2P	000000000G2PB6VB4ANR28CV2P	1
3	000000000G4AM4J9NBMPK3TJ...	000000000G4AM4J9NBMPK3TJ...	1045
4	000000000G4F19XBEXPWKS8F3N	000000000G4F19XBEXPWKS8F3N	6
5	000000000G4KVEP36NNJR5YWTJ	000000000G4KVEP36NNJR5YWTJ	68
6	000000000G8BW90E86ARRKM3...	000000000G8BW90E86ARRKM3...	427
7	000000000G8JG2FEA12SVNXXEN	000000000G8JG2FEA12SVNXXEN	438
8	000000000G95BD07NB0181K0XD	000000000G95BD07NB0181K0XD	105
9	000000000GBZ5VRTC3YS9T716N	000000000GBZ5VRTC3YS9T716N	489
10	000000000GCKFV5K6Q8FWGFH...	000000000GCKFV5K6Q8FWGFH...	30
11	000000000GCTJ4F7QXJWWMBD...	000000000GCTJ4F7QXJWWMBD...	60
12	000000000GDD59YDSJCXA2H46X	000000000GDD59YDSJCXA2H46X	19
13	000000000GEHAHYHGRSY59TR1D	000000000GEHAHYHGRSY59TR1D	25
14	000000000GG3B9VDBKAQM1D9...	000000000GG3B9VDBKAQM1D9...	1
15	000000000GGJG260YZ8XEWZV...	000000000GGJG260YZ8XEWZV...	20
16	000000000GH4BAHF58NTZC0489	000000000GH4BAHF58NTZC0489	8
17	000000000GHB4N83QX9KJM48K2	000000000GHB4N83QX9KJM48K2	9025

CHALLENGES & RESOLUTIONS

- Challenge - High missing demographic data
Solution - Improve Cognito capture process
- Challenge - Duplicate institutions (e.g., case mismatch)
Solution - Normalize data during frontend entry
- Challenge - No opportunity-opportunity_code mapping
Solution - Add an opportunity_enrollment mapping table
- Challenge - Marketing data unlinked
Solution - Use only for campaign performance dashboards

✓ SUMMARY TABLE

Check	Query	Expected Outcome
Table exists	information_schema.tables	✓ Found
Row count	SELECT COUNT(*)	≈ learner_opportunity rows
Nulls in key columns	COUNT(*) FILTER (WHERE ...)	Minimal
Duplicates	GROUP BY enrollment_id	0
Join validation	email, cohort date checks	Mostly matched
Status distribution	GROUP BY status	1070 dominant
Age validity	AGE(birthdate)	5–100 typical
Country/institution	Top 10 counts	Balanced spread

- Master table integrates 5 out of 6 datasets
- Marketing Dataset is excluded due to structural limitations
- Data quality is acceptable for analysis; improvements noted
- The pipeline is automated using a stored procedure and ready for dashboard use

VISUALISATIONS

