

UNIVERSIDAD TECNOLÓGICA DE CHIHUAHUA

TECNOLOGÍAS DE LA INFORMACIÓN



Extracción de Conocimiento en Bases de Datos

**II.3. Reporte de solución de caso de estudio de técnicas de
limpieza de datos**

IDGS91N

PRESENTA:

JORGE ALEJANDRO HERNANDEZ CONTRERAS

DOCENTE:

Enrique Mascote

Chihuahua, Chih., 12 Octubre de 2025

Introducción.....	3
Limpieza de datos.....	3
Eliminación de duplicados.....	3
Normalización de texto.....	3
Formato de fecha.....	4
Revisión de datos numéricos.....	4
Valores únicos en status.....	4
Hechos y dimensiones.....	4
Tabla de hechos.....	4
Tablas de dimensiones.....	4
Modelo relacional.....	6
Conclusiones.....	7
Referencias.....	8

Introducción

El archivo que se trabajó tiene información sobre personas migrantes, como su país de ciudadanía, tipo de visa, país de residencia, tipo de pasajero y fechas. El propósito de este trabajo fue limpiar los datos, identificar tablas de hechos y dimensiones para un modelo de datos y crear un diseño relacional normalizado. Todo esto sirve para que la información se pueda analizar correctamente y se pueda usar en un data warehouse sin errores.

Limpieza de datos

Al revisar los datos, se encontró que no había valores nulos en las columnas principales. Aun así, se analizaron otros posibles problemas que afectan la calidad de la información. Estas son las acciones aplicadas:

Eliminación de duplicados

Se buscó si existían filas repetidas completamente. Para eso se usó quitar duplicados en Excel. Las filas idénticas fueron eliminadas para evitar que los conteos se dupliquen.

Normalización de texto

Se detectó que algunas columnas como citizenship, visa, country_of_residence y passenger_type podían tener inconsistencias como:

- Espacios de más
- Mayúsculas y minúsculas mezcladas
- Variaciones en nombres

Se aplicaron funciones como:

- =ESPACIOS()
- =NOMPROPIO()
- =MAYUSC()

Esto se hizo para que todo quedara estandarizado.

Formato de fecha

La columna year_month estaba como texto, así que se separó en:

- Año (year)
- Mes (month_num)

Revisión de datos numéricos

Las columnas estimate y standard_error se convirtieron a número en caso de que se detectaran como texto.

Valores únicos en status

Se revisó que los valores fueran consistentes, corrigiendo posibles variantes como minúsculas o errores.

Todo cambio realizado se registró en una hoja de auditoría con la siguiente información:

- Valor original
- Valor corregido
- Motivo

Hechos y dimensiones

Después de limpiar los datos, se identificaron las tablas necesarias para un modelo tipo data warehouse.

Tabla de hechos

Fact_Migration contiene los valores medibles:

- estimate
- standard_error
- Llaves de todas las dimensiones relacionadas

Tablas de dimensiones

Dim_Date

- Año, mes, nombre del mes, trimestre.

Dim_PassengerType

- Tipo de pasajero (ej. Long-term migrant).

Dim_Direction

- Arrivals o Departures.

Dim_Citizenship

- País de ciudadanía.

Dim_Visa

- Tipo de visa.

Dim_CountryResidence

- País de residencia anterior.

Dim_Status

- Estado del dato (Provisional o Final).

Modelo relacional

```
2  CREATE TABLE dim_date(
3      date_key DATE PRIMARY KEY,
4      year INTEGER NOT NULL,
5      month INTEGER NOT NULL,
6      month_name VARCHAR(20),
7      quarter INTEGER
8  );
9
10 CREATE TABLE dim_passenger_type(
11     passenger_type_key SERIAL PRIMARY KEY,
12     passenger_type VARCHAR(100) UNIQUE NOT NULL
13 );
14
15 CREATE TABLE dim_direction(
16     direction_key SERIAL PRIMARY KEY,
17     direction VARCHAR(50) UNIQUE NOT NULL
18 );
19
20 CREATE TABLE dim_citizenship(
21     citizenship_key SERIAL PRIMARY KEY,
22     citizenship VARCHAR(150) UNIQUE NOT NULL,
23     iso_code VARCHAR(10)
24 );
25
26 CREATE TABLE dim_visa(
27     visa_key SERIAL PRIMARY KEY,
28     visa_type VARCHAR(150) UNIQUE NOT NULL
29 );
30
31 CREATE TABLE dim_country_residence(
32     country_residence_key SERIAL PRIMARY KEY,
33     country_name VARCHAR(150) UNIQUE NOT NULL,
34     iso_code VARCHAR(10),
35     region VARCHAR(100)
36 );
37
38 CREATE TABLE dim_status(
39     status_key SERIAL PRIMARY KEY,
40     status VARCHAR(50) UNIQUE NOT NULL
41 );
42
43 --Tabla de hechos
44 CREATE TABLE fact_migration(
45     fact_id SERIAL PRIMARY KEY,
46     date_key DATE NOT NULL,
47     passenger_type_key INTEGER NOT NULL,
48     direction_key INTEGER NOT NULL,
49     citizenship_key INTEGER NOT NULL,
50     visa_key INTEGER NOT NULL,
51     country_residence_key INTEGER NOT NULL,
52     status_key INTEGER NOT NULL,
53     estimate BIGINT NOT NULL,
54     standard_error BIGINT,
55     CONSTRAINT fk_date FOREIGN KEY(date_key) REFERENCES dim_date(date_key),
56     CONSTRAINT fk_passenger FOREIGN KEY(passenger_type_key) REFERENCES dim_passenger_type(passenger_type_key),
57     CONSTRAINT fk_direction FOREIGN KEY(direction_key) REFERENCES dim_direction(direction_key),
58     CONSTRAINT fk_citizenship FOREIGN KEY(citizenship_key) REFERENCES dim_citizenship(citizenship_key),
59     CONSTRAINT fk_visa FOREIGN KEY(visa_key) REFERENCES dim_visa(visa_key),
60     CONSTRAINT fk_countryres FOREIGN KEY(country_residence_key) REFERENCES dim_country_residence
61         (country_residence_key),
61     CONSTRAINT fk_status FOREIGN KEY(status_key) REFERENCES dim_status(status_key)
```

Conclusiones

Este trabajo ayudó a entender cómo limpiar y organizar un conjunto grande de datos para que se pueda usar en un modelo relacional y después en un data warehouse. Aunque el archivo no tenía valores faltantes, sí era importante estandarizar el texto, revisar duplicados y convertir los formatos de fecha.

También se aprendió cómo separar la información en hechos y dimensiones, lo que facilita el análisis y evita datos repetidos. El modelo normalizado permite mantener integridad y trabajar de forma más profesional con la información.

Referencias

Microsoft Support. (2024). *Limpiar datos en Excel*. support.microsoft.com.

<https://support.microsoft.com/es-es/office/limpiar-datos-en-excel>

Talend. (2024). *Introducción a la normalización de bases de datos*. talend.com.

<https://www.talend.com/es/resources/database-normalization/>