



UNIVERSIDAD TÉCNICA PARTICULAR DE LOJA
La Universidad Católica de Loja

Informe Final
Fundamentos de Bases de
Datos

Autor: Castillo Paladines Byron Josué

Octubre 2022 – Febrero 2023

Contenido

Introducción	3
Modelos	4
Modelo Conceptual.....	4
Modelo Lógico	5
Dependencias Funcionales – Tabla Universal	6
Modelo Físico	7
Normalización	7
First Normal Form (1NF).....	7
Second Normal Form (2NF)	8
Third Normal Form (3NF)	8
Pasos para Realizar el Proyecto	15
Creación de un “Schema”	15
Conexión del “Schema” a la Base De Datos	17
Importación del CSV.....	17
Creación de funciones que Permitan Extraer y Limpiar los Datos.....	20
LIMPIEZA COLUMNA CREW	33
Consultas.....	38
Conclusiones	40

Introducción

Las bases de datos han demostrado a lo largo del tiempo que son esenciales en muchos sistemas informáticos ya que nos da la opción de ingresar, manipular, almacenar y/o eliminar grandes cantidades de información, este presente informe tiene como objetivo aplicar todos los conocimientos obtenidos en todo este tiempo de estudio con base a la materia de Base de Datos, aplicando la importación del CSV desde una base de datos usando el sistema de gestión de base de datos MySQL, además de pasar por distintas fases, las cuales incluyen la inserción de datos, la limpieza de dichos datos, la carga y sobre todo la explotación de datos. Igualmente, podemos diferenciar los métodos empleados al realizar todas las fases propuestas en este trabajo.

Este informe busca producir resultados orientados a la gestión y uso eficiente de la información relevante y de interés, a partir del conjunto de datos de trabajo. Además, nos permitió recordar varios temas estudiados en la materia tales como diagramas de Entidad/Relación, modelos conceptuales, lógicos y físicos, la creación de instructivos SQL, entre muchos más temas que se muestran a lo largo del informe.

Modelos

Modelo Conceptual

Después de organizar los datos en una tabla universal, los normalizamos haciendo uso del Modelado Entidad/Relación, además nos permite representar gráficamente la estructura de los datos y las relaciones entre ellos tomando en cuenta en como estarían organizados los datos del archivo CSV y las características de cada columna, identificando los atributos correspondientes a cada relación.

La normalización implica separar los datos repetidos y asegurar su integridad mediante creación de tablas individuales para cada entidad y relación.

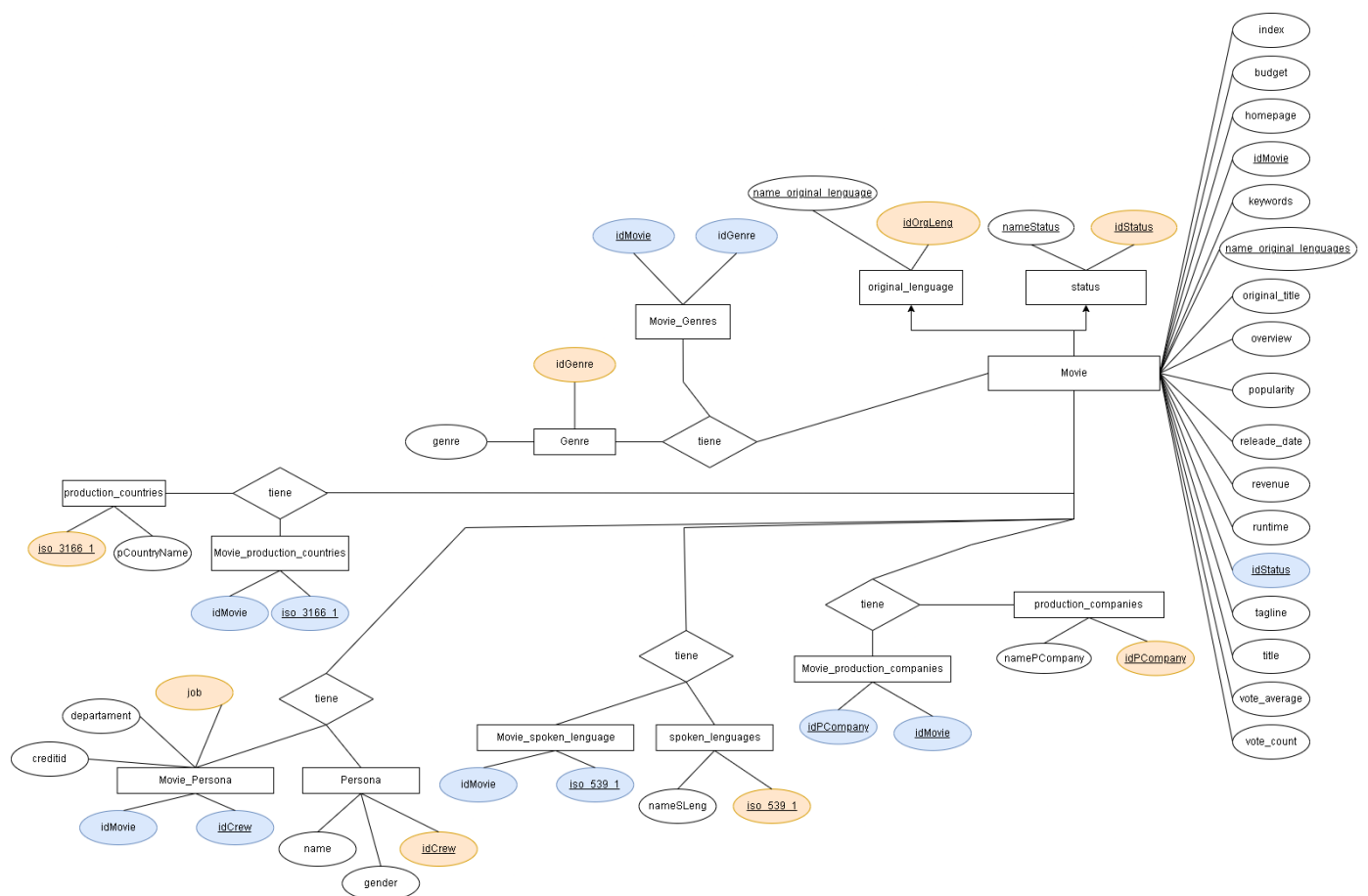


Figura1: Representación del Modelo Conceptual

Modelo Lógico

Permite especificar las relaciones entre las tablas mediante claves foráneas y claves primarias. Una clave primaria es una columna o conjunto de columnas que identifica de manera única a cada registro en una tabla. Por otro lado, una clave foránea se refiere a una columna en una tabla que hace referencia a una clave primaria en otra tabla. De esta forma, el Modelo Lógico SQL permite representar las relaciones entre las entidades de manera clara y estructurada.

Con el modelo conceptual ya realizado, se seleccionan los atributos apropiados para cada tabla y se crea un prototipo de estas basado en las entidades y relaciones identificadas. En cada tabla, se declara una clave primaria y, en su caso, una clave foránea, esto con el objetivo de tener una base sólida para desarrollar el esquema de la base de datos.

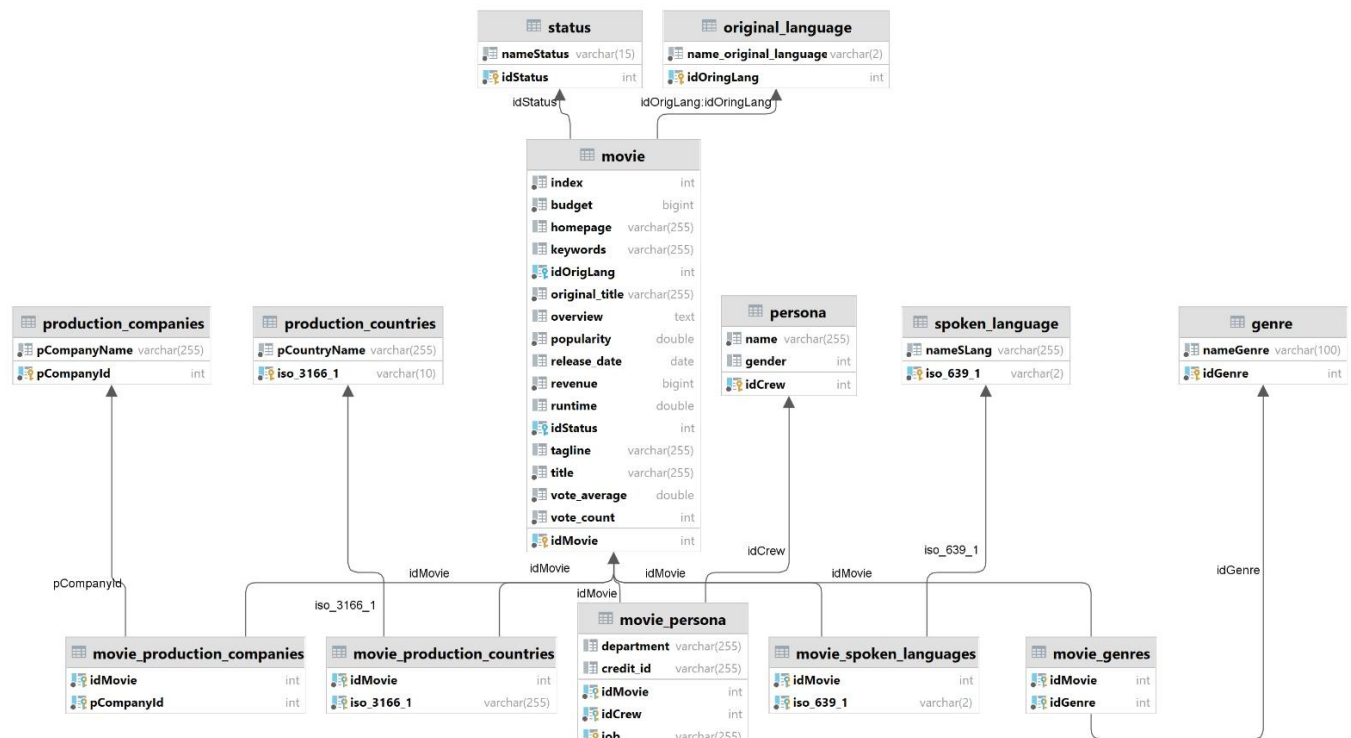
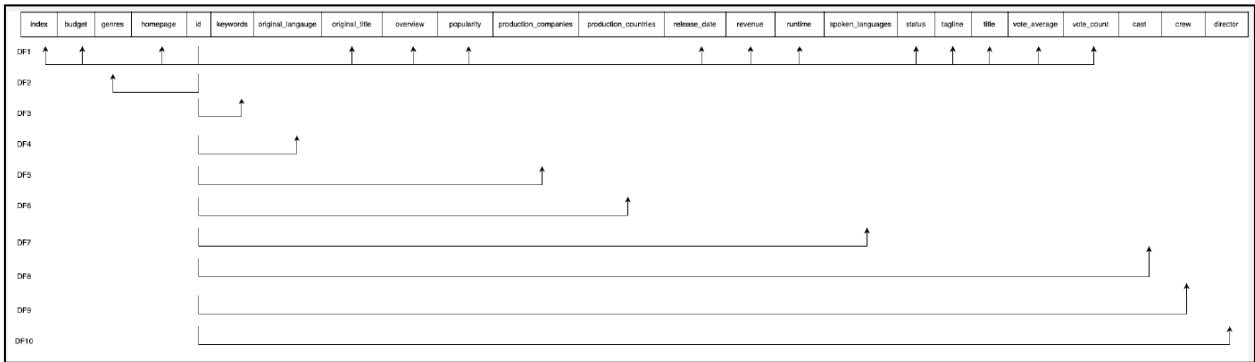


Figura 2: Representación gráfica del Modelo Lógico

Dependencias Funcionales – Tabla Universal



Se refiere a las relaciones que existen entre los atributos de una entidad y como un cambio en un atributo puede afectar a otros atributos, es importante para garantizar la integridad de los datos en una base de datos y para optimizar el rendimiento de las consultas. Una vez establecido el modelo lógico y conceptual se creó las dependencias funcionales existentes en la tabla universal.

Figura 3: La imagen muestra una descripción gráfica de las dependencias existentes en las columnas del dataset

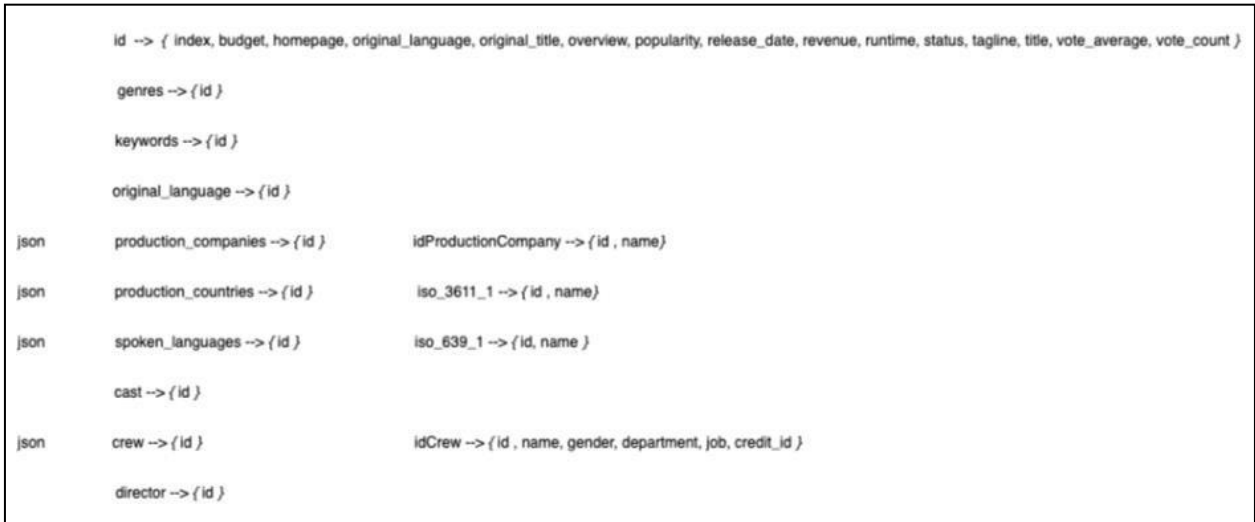
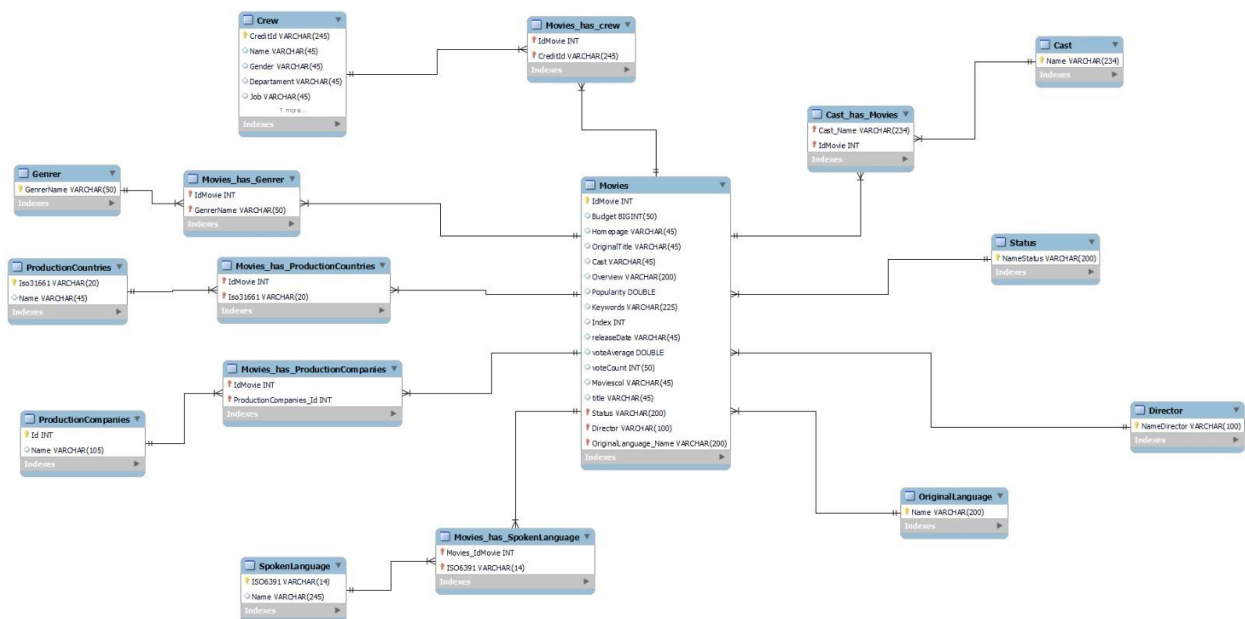


Figura 4: Dependencias Funcionales textualizadas para un mejor entendimiento

Modelo Físico

Este modelo describe cómo se almacenan los datos en el sistema de almacenamiento y cómo se accede a ellos, con este modelo podemos pasar a generar las respectivas sentencias DDL, incluye información detallada sobre la implementación de las tablas, columnas, índices, claves primarias y foráneas, y otras características de la base de datos, tales como la configuración de la memoria caché, el almacenamiento de datos, los tamaños de página, etc.



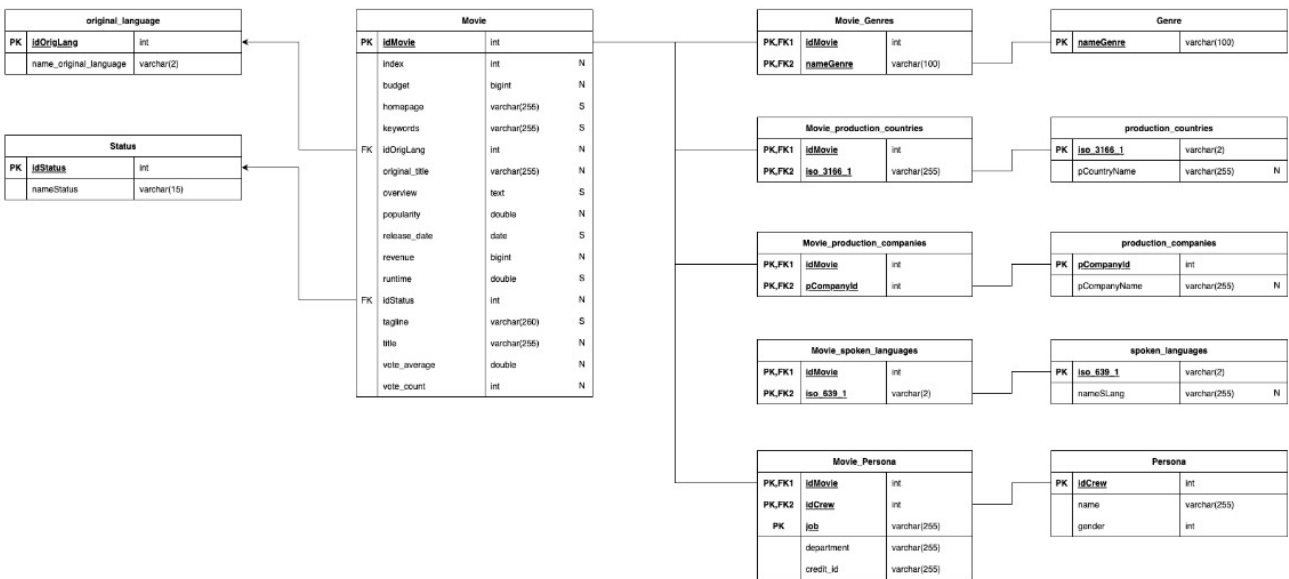


Figura 5: Representación gráfica del Modelo Físico

Normalización

La normalización se la puede describir como un proceso en el cual se organizan los datos por medio de creación de elementos como tablas y relaciones las cuales pueden ser modelos de bases relacionales, tomando como sentido el funcionamiento de las bases de datos relacionales la normalización busca una mejor organización de los datos y a la vez eliminar la redundancia

First Normal Form (1NF)

- Una relación en donde la intersección de cada fila y columna contiene un solo valor.
- No existen filas repetidas

- La fase antes de 1NF es Unnormalized Form (UNF) la cual es una tabla que contiene uno más grupos repetidos.
- Para transformar la tabla no normalizada a la primera forma normal, identificamos y eliminamos los grupos que se repiten dentro de la tabla.
- Todos los atributos son atómicos, simples e individuales.
- Un grupo repetitivo es un atributo, o grupo de atributos, dentro de una tabla que ocurre con múltiples valores para una sola ocurrencia de los atributos clave designados para esa tabla.
- Identificar cada conjunto de relacionados con la clave principal.
- Elimina grupos repetidos en tablas individuales.

Second Normal Form (2NF)

- Una relación que está en la primera forma normal y cada atributo que no es de clave principal depende funcionalmente de la clave principal.
- Crear tablas separadas para el conjunto de valores y los registros múltiples, estas tablas se deben relacionar con una clave externa.
- Los registros no deben depender de otra cosa que la clave principal de la tabla, incluida la clave compuesta si es necesario.
- La normalización de las relaciones 1NF a 2NF implica la eliminación de dependencias parciales. Si existe una dependencia parcial, eliminamos los atributos parcialmente dependientes de la relación colocándolos en una nueva relación junto con una copia de su determinante.

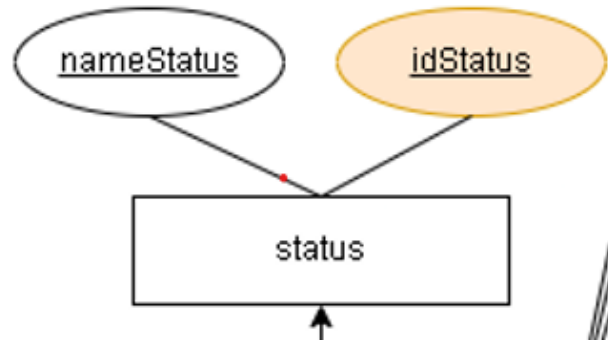
Third Normal Form (3NF)

- Comprueba las dependencias transitivas, eliminando campos que no dependen de la clave principal.
- Una relación que está en primera y segunda forma normal y en la que ningún atributo que no sea de clave principal depende transitivamente de la clave principal.
- Los valores que no dependen de la clave principal no pertenecen a la tabla.
- Los campos que no pertenecen a la clave principal colóquelos en una tabla aparte y relacionen ambas tablas por medio de una clave externa.
- La normalización de las relaciones 2NF a 3NF implica la eliminación de las dependencias transitivas.
- Los campos que no pertenecen a la clave principal colóquelos en una tabla aparte y relacionen ambas tablas por medio de una clave externa.

- Si existe una dependencia transitiva, eliminamos los atributos transitivamente dependientes de la relación colocando los atributos en una nueva relación junto con una copia del determinante.

Primero identificamos y eliminamos grupos repetitivos dentro de la tabla. Un grupo repetitivo es un atributo, o grupo de atributos, dentro de una tabla que ocurre con múltiples valores para una sola ocurrencia de los atributos clave designados para esa tabla.

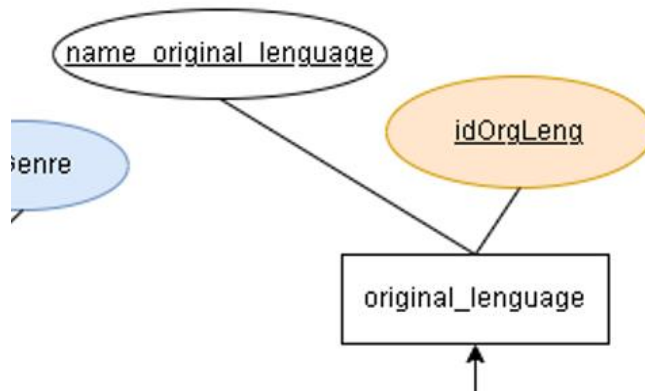
Tabla Uno a Muchos (1 : N)



Status es una columna con tres posibles valores (released, rumored, post-production)

- Una película puede tener un Status, pero un Status puede representar a muchas películas.
- Para solucionar, se crea una tabla separada que se llame Status.
- La llave primaria es nameStatus.
- Como solución, utilizamos el nameStatus como llave foránea en

Original_language es una columna con el nombre del lenguaje original de la Movie.



- Una película tiene un original_language, pero un original_language puede tener muchas películas.

- Para solucionar, se crea una tabla separada que se llame original_language.

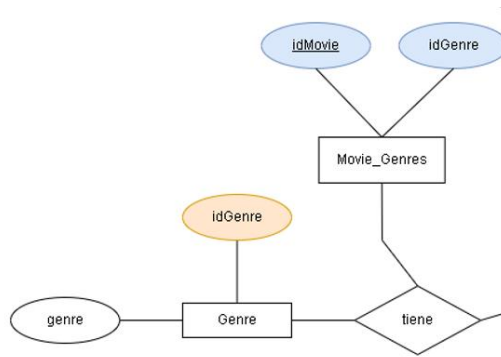
La llave primaria es idOrgLeng.

- Como solución, utilizamos el idOrgLeng como llave foránea en original_language.

Tabla Muchos a Muchos (N : N)

En nuestro caso, ninguna de las 4803 entradas se repite, cada película es diferente (única).

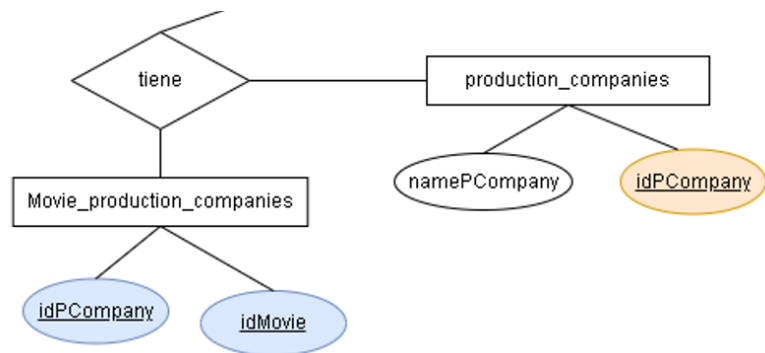
Cada Movie tiene un index y un id diferente. Como llave primaria utilizamos id el cual nombramos idMovie para no confundirlo con otros id que existan en otras columnas.



- Esta columna contiene un String de géneros que puede contener 0 a n géneros.
 - Para solucionar, cada columna debe contener un solo valor.
 - En este caso tenemos una lista de géneros. Existen múltiples valores.
 - La solución es crear una tabla separada que se llame Genres.

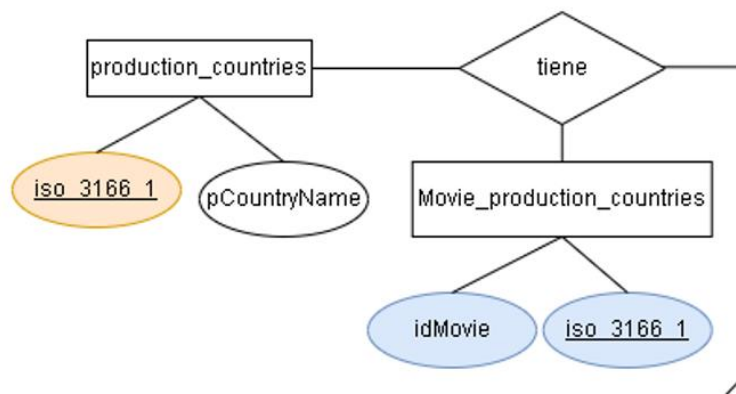
La llave primaria es idGenre

- Una Movie puede tener muchos géneros, y un género puede estar en muchas Movie.
- La relación (muchos a muchos) se soluciona con una tabla llamada Movie_Genres utilizando la llave primaria de cada uno.



- Production_companies contiene un String de JSON que contiene un name y un id que puede contener 0 a n production_companies.
 - Para solucionar, cada columna debe contener un solo valor.
 - Las columnas de “Movie_production_companies” y “production_countries” son String de JSON.

- La solución es crear una tabla separada que se llame `production_companies`.
- Tiene dos atributos, `namePCompany` e `idPCompany` el cual es la primary key.
- Una `Movie` puede tener muchos `production_companies`, y un `production_companies` puede aparecer en muchas `Movie`.
- La relación (muchos a muchos) se soluciona con una tabla llamada `Movie_Production_companies` utilizando la llave primaria de cada uno.



- `production_countries` contiene un String de JSON que contiene un `iso_3166_1` y un `pCountryName`
 - Para solucionar, cada columna debe contener un solo valor.
 - La solución es crear una tabla separada que se llame `Movie_production_countries`.
 - Tiene dos atributos, `iso_3166_1` y `idMovie`.
 - Una `Movie` puede tener muchos `production_countries`, y un `production_countries` puede aparecer en muchas `Movie`.

- spoken_languages contiene un String de JSON que contiene un iso_639_1 y un name que puede contener 0 a n spoken_languages.
 - Para solucionar, cada columna debe contener un solo valor.
 - En este caso tenemos una lista de spoken_languages. Existen múltiples valores.
 - La solución es crear una tabla separada que se llame spoken_languages.
 - Tiene dos atributos, iso_639_1 y name los cuales los nombramos, iso_639_1 el cual es la primary key y spokLangName.
 - Una Movie puede tener muchos spoken_languages, y un spoken_languages puede aparecer en muchas Movie.
 - La relación (muchos a muchos) se soluciona con una tabla llamada Movie_Spoken_languages utilizando la llave primaria de cada uno.

- Cast contiene un String que contiene los nombres de los actores de una película que puede contener 0 a n Cast member.
 - Para solucionar, cada columna debe contener un solo valor.
 - En este caso tenemos una lista de cast. Existen múltiples valores.
 - La solución es crear una tabla separada que se llame Cast.
 - La primary key es CastName.

- Una Movie puede tener muchos Cast member, y un Cast member puede aparecer en muchas Movie.
 - La relación (muchos a muchos) se soluciona con una tabla llamada Movie_Cast utilizando la llave primaria de cada uno.
-
- Crew contiene un String de JSON que contiene un name, gender, department, job, credit_id, id que puede contener 0 a n Crew members.
 - Para solucionar, cada columna debe contener un solo valor.
 - En este caso tenemos una lista de Crew members. Existen múltiples valores.
 - La solución es crear una tabla separada que se llame Crew.
 - Tiene los atributos name, gender, department, job, credit_id, id los cuales nombraremos crewName, gender, department, job, credit_id, idCrew el cual es la primary key.
 - Una Movie puede tener muchos Crew members, y un Crew members puede aparecer en muchas Movie.
 - La relación (muchos a muchos) se soluciona con una tabla llamada Movie_Crew utilizando la llave primaria de cada uno.
-
- Cada Movie tiene uno de los siguientes:
 - Production_countries
 - Production_companies

- Spoken_languages
- Crew

Pasos para Realizar el Proyecto

Creación de “Schema”.

MEDIANTE LA CREACIÓN AUTOMÁTICA

Para crear automáticamente el esquema en la base de datos de MySql, primero debes acceder a la parte superior izquierda de la pantalla donde se encuentran los íconos. Luego, debes seleccionar el ícono destinado para la creación de un esquema, como se puede ver en el anexo proporcionado.:

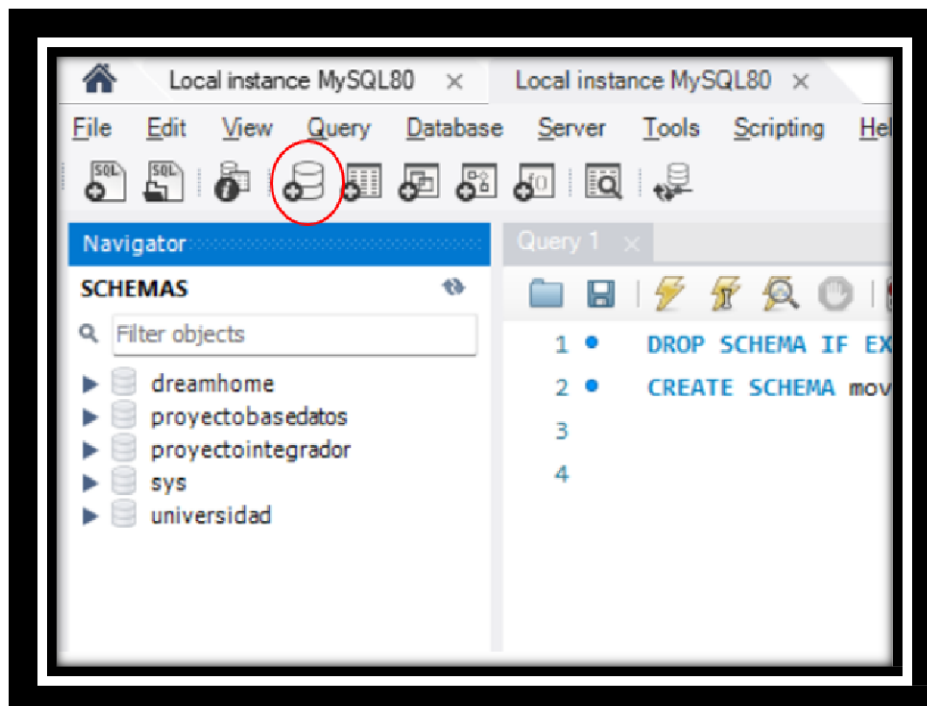


Figura 6: Representación de la barra de opciones perteneciente a Workbench

Después de seleccionar el ícono para crear un esquema, se le dará un nombre, en este caso "movie_dataset". En este esquema, se pueden realizar varias modificaciones, como la selección de un juego de caracteres, en este caso se ha utilizado "utf8":

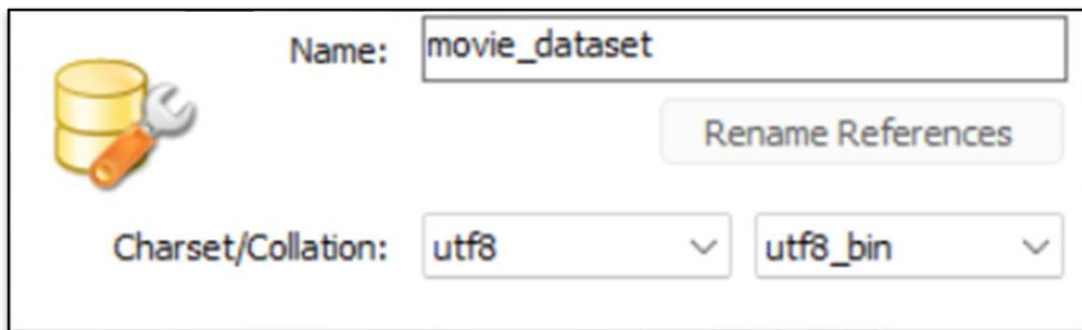


Figura 7: Representación de la funcionalidad de creación de Schema automáticamente

MEDIANTE SENTENCIAS SQL

Para la creación del “Schema” o base de datos usando sentencias SQL, es necesario hacer lo siguiente:

```
CREATE DATABASE IF NOT EXISTS `movie_dataset` DEFAULT CHARACTER SET utf8 ;
```

Es necesario especificarle la codificación de caracteres utf8 ya que en el archivo CSV se usa caracteres especiales y esta misma codificación los transforma.

Conexión del “Schema” a la Base De Datos.

Para poderse conectar, por así decirlo, es necesario hacerlo mediante sentencias SQL, la sentencia es la siguiente:

```
USE `Movie_Dataset` ;
```

Importación del CSV.

DataGrip es un entorno de desarrollo integrado (IDE) desarrollado por JetBrains para la gestión de bases de datos. Este contiene una herramienta que facilita la importación de diferentes tipos de archivos a las tablas de MySQL. Los siguientes son los pasos que se describen en la imagen proporcionada:

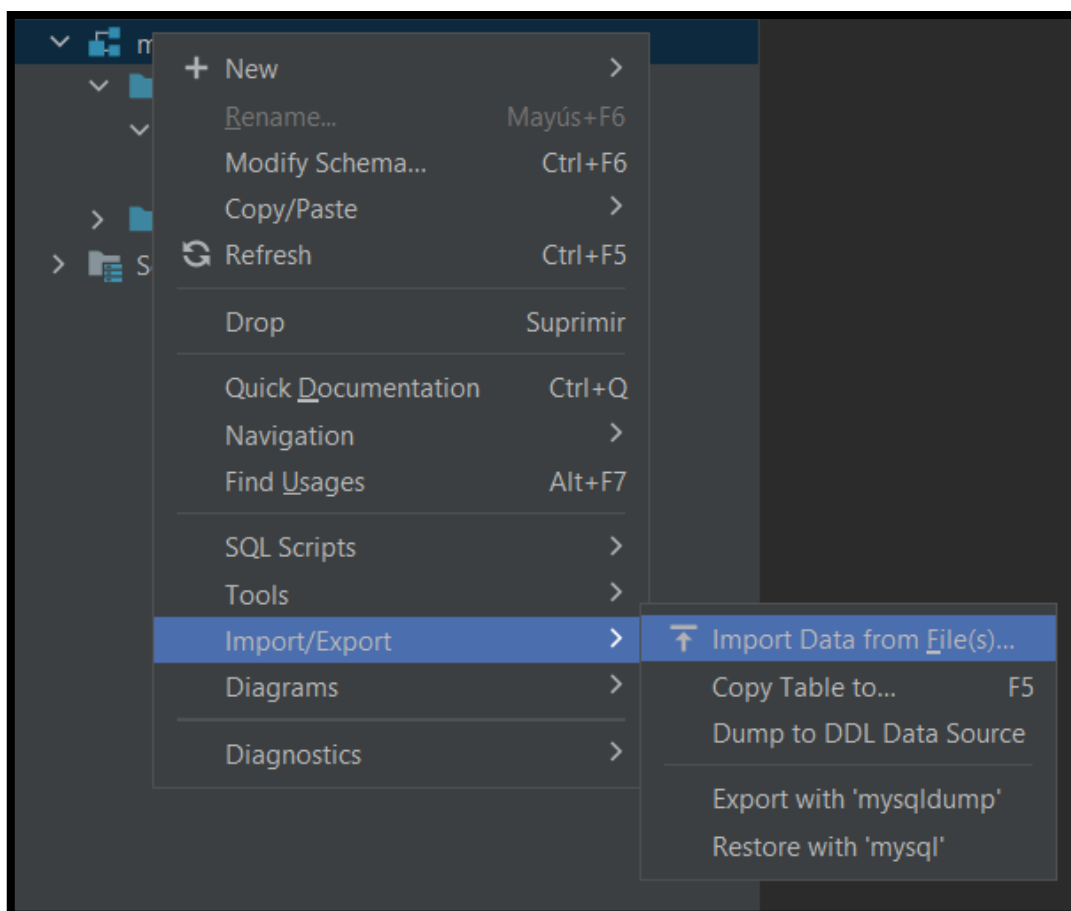


Figura 8: Representación de la funcionalidad de importación de archivos en DataGrip

Una vez entrado a este apartado, se habilita una interfaz de importación en donde especificamos que el separador es la coma, y que la primera fila son los títulos de las columnas

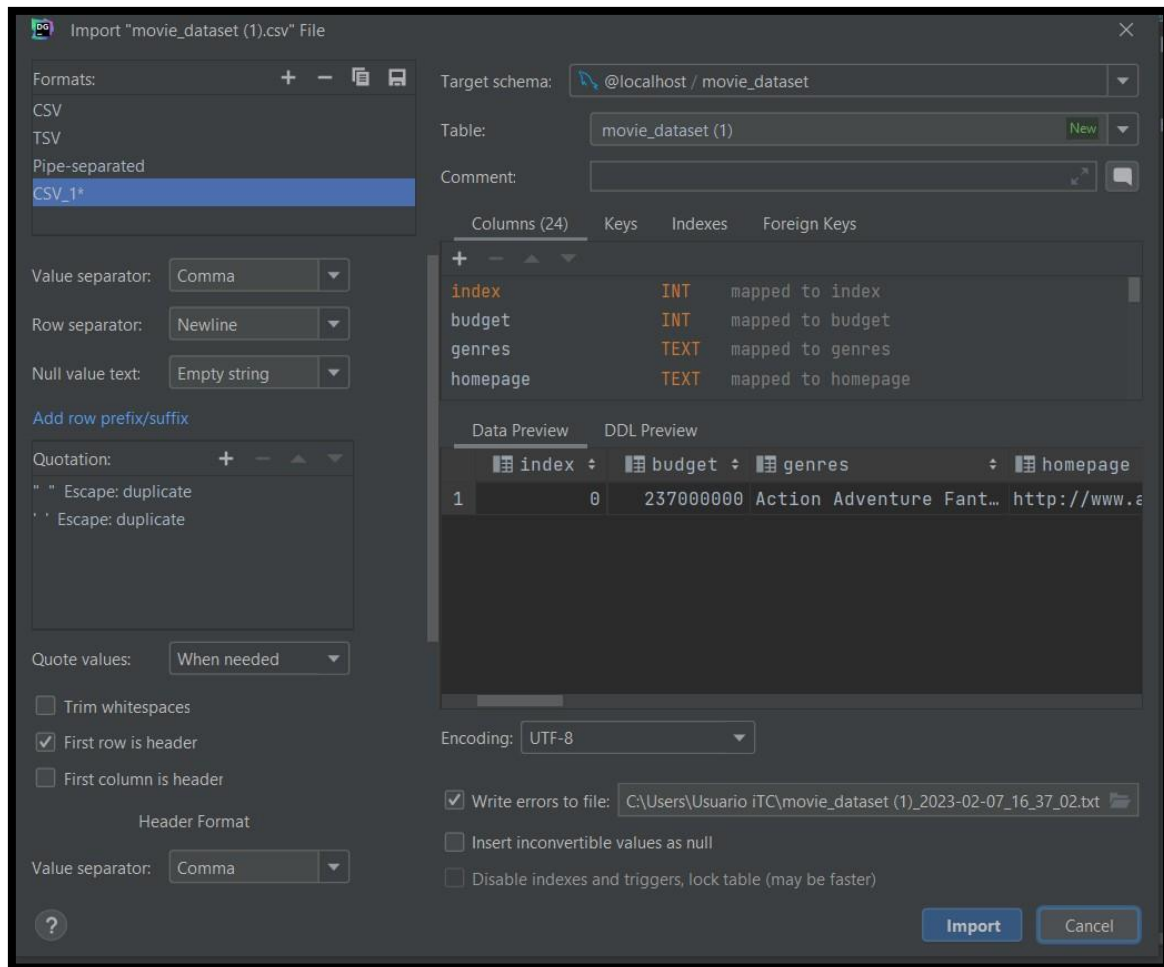


Figura 8: Representación de la interfaz de importación de archivos en DataGrip

Creación de funciones que Permitan Extraer y Limpiar los Datos.

En el procedimiento almacenado de SQL, se comienza verificando si existe un procedimiento con el nombre "TablaDirector". Si existe, se elimina. Luego, se declaran dos variables: la primera es "done" con un valor predeterminado de "FALSE", y la segunda es "nameDirector", que es de tipo VARCHAR y tiene una longitud máxima de 100 caracteres.

A continuación, se declara un cursor llamado "CursorDirector", que selecciona los nombres únicos de los directores en la tabla "movie_dataset". Se establece un manejador de continuación para detectar cuándo se ha llegado al final del cursor. Luego, se abre el cursor y se inicia un ciclo que recorre cada uno de los nombres de los directores en el cursor.

En cada iteración del ciclo, se asigna el nombre del director a la variable "nameDirector". Si se ha alcanzado el final del cursor, se sale del ciclo. Si el nombre de director es nulo, se asigna un valor vacío. Además, se tienen casos especiales en los nombres que pueden ser solucionados con la función "Replace" en la misma sentencia de código.

Luego, se crea una consulta dinámica que inserta el nombre del director en la tabla "director". Se prepara y ejecuta la consulta dinámica y se libera la memoria de la consulta preparada. Se repite este proceso para cada nombre de director en el cursor.

Finalmente, se cierra el cursor y se finaliza el procedimiento almacenado. En resumen, este procedimiento permite crear y llenar una tabla con los nombres únicos de los directores en una tabla de origen. Este código nos servirá para poblar estas tres tablas.

Población Tabla “Director”

```
1 • USE Movies2023;
2
3 • DROP PROCEDURE IF EXISTS TablaDirector;
4
5 DELIMITER $$
6 • CREATE PROCEDURE TablaDirector()
7   BEGIN
8
9     DECLARE done INT DEFAULT FALSE;
10    DECLARE nameDirector VARCHAR(100);
11
12    -- Declarar el cursor
13    DECLARE CursorDirector CURSOR FOR
14      SELECT DISTINCT CONVERT(director USING UTF8MB4) AS names from movie_dataset;
15
16    -- Declarar el handler para NOT FOUND (esto es marcar cuando el cursor ha llegado a su fin)
17    DECLARE CONTINUE HANDLER FOR NOT FOUND SET done = TRUE;
18
19    -- Abrir el cursor
20    OPEN CursorDirector;
21    CursorDirector_loop: LOOP
22      FETCH CursorDirector INTO nameDirector;
23
24      -- Si alcanzo el final del cursor entonces salir del ciclo repetitivo
25    IF done THEN
```



```

        LEAVE CursorDirector_loop;
27     END IF;
28     IF nameDirector IS NULL THEN
29         SET nameDirector = '';
30     END IF;
31     SET @_oStatement = CONCAT('INSERT INTO directorCURSOR (name) VALUES (\'',
32     REPLACE(REPLACE(nameDirector, '\\', '\\\\'), '\\u', '\\\\u')
33     ,'\');');
34     PREPARE sent1 FROM @_oStatement;
35     EXECUTE sent1;
36     DEALLOCATE PREPARE sent1;
37
38 END LOOP;
39 CLOSE CursorDirector;
40 END $$
41 DELIMITER ;
42
43 CALL TablaDirector();
44
45 DROP TABLE IF EXISTS directorCURSOR;
46
47 CREATE TABLE directorCURSOR (
48     name varchar(255) PRIMARY KEY
49 );

```

Población Tabla “Status”

```

1  USE Movies2023;
2
3  DROP PROCEDURE IF EXISTS TablaStatus;
4
5  DELIMITER $$
6  CREATE PROCEDURE TablaStatus()
7  BEGIN
8
9      DECLARE done INT DEFAULT FALSE;
10     DECLARE nameStatus VARCHAR(100);
11
12     -- Declarar el cursor
13     DECLARE CursorStatus CURSOR FOR
14         SELECT DISTINCT CONVERT(status USING UTF8MB4) AS names from movie_dataset;
15
16     -- Declarar el handler para NOT FOUND (esto es marcar cuando el cursor ha llegado a su fin)
17     DECLARE CONTINUE HANDLER FOR NOT FOUND SET done = TRUE;
18
19     -- Abrir el cursor
20     OPEN CursorStatus;
21     CursorStatus_loop: LOOP
22         FETCH CursorStatus INTO nameStatus;
23
24         -- Si alcanzo el final del cursor entonces salir del ciclo repetitivo
25         IF done THEN

```

```

        LEAVE CursorStatus_loop;
27     END IF;
28     IF nameStatus IS NULL THEN
29         SET nameStatus = '';
30     END IF;
31     SET @_oStatement = CONCAT('INSERT INTO statusCURSOR (name) VALUES (\'',
32         nameStatus, '\');');
33     PREPARE sent1 FROM @_oStatement;
34     EXECUTE sent1;
35     DEALLOCATE PREPARE sent1;
36
37 END LOOP;
38 CLOSE CursorStatus;
39 END $$
40 DELIMITER ;
41
42 CALL TablaStatus();
43
44 DROP TABLE IF EXISTS statusCURSOR;
45
46 CREATE TABLE statusCURSOR (
47     name varchar(255) PRIMARY KEY
48 );
49
50 SELECT * FROM statusCURSOR;

```

Población Tabla “Original_language”

```

1  USE Movies2023;
2
3  DROP PROCEDURE IF EXISTS TablaOriginalLanguage;
4
5  DELIMITER $$
6  CREATE PROCEDURE TablaOriginalLanguage()
7  BEGIN
8
9      DECLARE done INT DEFAULT FALSE;
10     DECLARE nameOL VARCHAR(100);
11
12     -- Declarar el cursor
13     DECLARE CursorOL CURSOR FOR
14         SELECT DISTINCT original_language AS names from movie_dataset;
15
16     -- Declarar el handler para NOT FOUND (esto es marcar cuando el cursor ha llegado a su fin)
17     DECLARE CONTINUE HANDLER FOR NOT FOUND SET done = TRUE;
18
19     -- Abrir el cursor
20     OPEN CursorOL;
21     CursorOL_loop: LOOP
22         FETCH CursorOL INTO nameOL;
23
24         -- Si alcanzo el final del cursor entonces salir del ciclo repetitivo
25     IF done THEN

```

```

26      LEAVE Cursor0L_loop;
27  END IF;
28  IF name0L IS NULL THEN
29      SET name0L = '';
30  END IF;
31  SET @_oStatement = CONCAT('INSERT INTO original_languageCURSOR (name) VALUES (\',
32      name0L, '\');');
33  PREPARE sent1 FROM @_oStatement;
34  EXECUTE sent1;
35  DEALLOCATE PREPARE sent1;
36
37  END LOOP;
38  CLOSE Cursor0L;
39  END $$
DELIMITER ;

42  CALL TablaOriginalLanguage();
43
44  DROP TABLE IF EXISTS original_languageCURSOR;
45
46  CREATE TABLE original_languageCURSOR (
47      name varchar(255) PRIMARY KEY
48  );

50  SELECT * FROM original_languageCURSOR;

```

Debido a que la tabla "Movie" hace uso de claves externas, era necesario crear las tablas "Director", "Status" y "original_language" para utilizar sus claves primarias como claves externas en la tabla "Movie" (REFERENCIAS).

Población Tabla “Movie”

El siguiente código SQL es un procedimiento almacenado que crea una tabla llamada "películas" en la base de datos. La tabla se crea a partir de los datos de un conjunto de datos llamado "movie_dataset". El procedimiento almacenado comienza eliminando todos los procedimientos con el mismo nombre. Luego, se declaran varias variables para almacenar los valores de cada columna en la tabla "movie_dataset".

A continuación, se crea un cursor llamado "CursorMovie" que selecciona los datos de la tabla "movie_dataset". Los cursores se utilizan para recorrer cada fila de una tabla y extraer valores de cada columna. Se declara un "gestor" para manejar el caso en el que se llegue al final del cursor.

El cursor se abre y comienza a recorrer cada fila de la tabla "movie_dataset" en un ciclo llamado "CursorMovie_loop". Dentro del ciclo, se utiliza la declaración FETCH para obtener valores de la fila actual y almacenarlos en variables previamente declaradas.

Si se llega al final del cursor, el ciclo termina. En caso contrario, se verifica si el nombre del director está vacío. Si es así, se le asigna un valor de cero. Se ejecuta una consulta para obtener la identificación del director a partir de su nombre y se guarda en la variable "Director_idDirector".

Finalmente, se inserta una nueva fila en la tabla "películas" con los valores obtenidos de la tabla "movie_dataset" y el ID del director. Cuando se completa el ciclo, el cursor se cierra y el procedimiento almacenado concluye.

```

1 DROP PROCEDURE IF EXISTS TablaMovie;
2
3 DELIMITER $$
4 CREATE PROCEDURE TablaMovie()
5 BEGIN
6
7     DECLARE done INT DEFAULT FALSE;
8     DECLARE Movindex INT;
9     DECLARE Movbudget BIGINT;
10    DECLARE Movhomepage VARCHAR(1000);
11    DECLARE MovidMovie INT;
12    DECLARE Movkeywords TEXT;
13    DECLARE Movoriginal_language VARCHAR(255);
14    DECLARE Movoriginal_title VARCHAR(255) ;
15    DECLARE Movoverview TEXT;
16    DECLARE Movpopularity DOUBLE;
17    DECLARE Movrelease_date VARCHAR(255);
18    DECLARE Movrevenue BIGINT;
19    DECLARE Movruntime DOUBLE;
20    DECLARE Movstatus VARCHAR(255);
21    DECLARE Movtagline VARCHAR(255);
22    DECLARE Movtitle VARCHAR(255);
23    DECLARE Movvote_average DOUBLE;
24    DECLARE Movvote_count INT;
25    DECLARE nameDirector VARCHAR(255);
26
27    DECLARE Director_nameDirector varchar(255);
28    DECLARE Director_nameStatus varchar(255);
29    DECLARE Director_nameOriginal_language varchar(255);
30
31    -- Declarar el cursor
32    DECLARE CursorMovie CURSOR FOR
33        SELECT `index`,budget,homepage,id,keywords,original_language,original_title,overview,popularity,release_date,revenue,
34            tagline,title,vote_average,vote_count,director FROM movie_dataset;
35
36    -- Declarar el handler para NOT FOUND (esto es marcar cuando el cursor ha llegado a su fin)
37    DECLARE CONTINUE HANDLER FOR NOT FOUND SET done = TRUE;
38
39    -- Abrir el cursor
40    OPEN CursorMovie;
41
42    CursorMovie_loop: LOOP
43        FETCH CursorMovie INTO Movindex,Movbudget,Movhomepage,MovidMovie,Movkeywords,Movoriginal_language,Movoriginal_title,M
44        Movpopularity,Movrelease_date,Movrevenue,Movruntime,Movstatus,Movtagline,Movtitle,Movvote_average,Movvote_count,nameD
45
46        -- Si alcanzo el final del cursor entonces salir del ciclo repetitivo
47        IF done THEN
48            LEAVE CursorMovie_loop;
49        END IF;
50        IF nameDirector IS NULL THEN
51            SET nameDirector = '';

```

```

51     END IF;
52
53     SELECT `name` INTO Director_nameDirector FROM directorCURSOR WHERE directorCURSOR.name = nameDirector;
54     SELECT `name` INTO Director_nameStatus FROM statusCURSOR WHERE statusCURSOR.name = Movstatus;
55     SELECT `name` INTO Director_nameOriginal_language FROM original_languageCURSOR WHERE original_languageCURSOR.name = M
56
57     INSERT INTO MovieCURSOR (`index`,budget,homepage,id,keywords,original_language,original_title,overview,popularity,rel
58     tagline,title,vote_average,vote_count,director)
59     VALUES (Movindex,Movbudget,Movhomepage,MovidMovie,Movkeywords,Director_nameOriginal_language,Movoriginal_title,Movove
60     Movpopularity,Movrelease_date,Movrevenue,Movruntime,Director_nameStatus,Movtagline,Movtitle,Movvote_average,Movvote_c
61
62     END LOOP;
63     CLOSE CursorMovie;
64     END $$
65     DELIMITER ;
66
67     CALL TablaMovie ();
68
69     DROP TABLE IF EXISTS MovieCURSOR;
70
71     CREATE TABLE MovieCURSOR (
72         `index` int,
73         budget bigint,
74         homepage varchar(1000),
75         id int PRIMARY KEY,
76         keywords TEXT,
77         original_language varchar(255),
78         original_title varchar(255),
79         overview TEXT,
80         popularity double,
81         release_date varchar(255),
82         revenue bigint,
83         runtime double,
84         `status` varchar(255),
85         tagline varchar(255),
86         title varchar(255),
87         vote_average double,
88         vote_count int,
89         director varchar(255),
90         FOREIGN KEY (original_language) REFERENCES original_languageCURSOR(name),
91         FOREIGN KEY (status) REFERENCES statusCURSOR(name),
92         FOREIGN KEY (director) REFERENCES directorCURSOR(name)
93     );
94
95     DROP TABLE IF EXISTS MovieCURSOR;
96
97     SELECT COUNT(*) FROM MovieCursor;
98     SELECT * FROM MovieCursor;

```

Tenemos tres columnas con valores JSON, para las cuales seguiremos el mismo procedimiento para las tres que son `spoken_languages`, `production_companies` y `production_countries`

Este código es un procedimiento que se encarga de extraer información de una tabla llamada `"movie_dataset"`, en particular la columna `"production_companies"`, que contiene información en formato JSON sobre las compañías productoras de películas.

El procedimiento comienza declarando variables, un cursor que se encarga de recorrer las filas de la tabla `"movie_dataset"` y un controlador de eventos `"CONTINUE HANDLER"` para determinar cuándo se ha alcanzado el final de los datos. Luego, se crea una tabla temporal llamada `"production_companieTem"` para almacenar los datos extraídos del formato JSON.

El cursor se usa para recorrer las filas de la tabla `"movie_dataset"` y, para cada fila, se extrae la información de las compañías productoras de las películas utilizando la función `"JSON_EXTRACT"`. Los datos extraídos se insertan en la tabla `"production_companieTem"`.

Después de que se han recorrido todas las filas, se seleccionan los datos únicos de la tabla `"production_companieTem"` y se insertan en la tabla `"production_companie"`. Finalmente, se cierra el cursor y se elimina la tabla temporal `"production_companieTem"`.

Población Tabla “Production_companies”

```
1 USE Movies2023;
2
3 DROP PROCEDURE IF EXISTS TablaProduction_companies;
4
5 DELIMITER $$
6 CREATE PROCEDURE TablaProduction_companies ()
7
8 BEGIN
9
10 DECLARE done INT DEFAULT FALSE ;
11 DECLARE jsonData json ;
12 DECLARE jsonId varchar(250) ;
13 DECLARE jsonLabel varchar(250) ;
14 DECLARE resultSTR LONGTEXT DEFAULT '';
15 DECLARE i INT;
16
17 -- Declarar el cursor
18 DECLARE myCursor
19 CURSOR FOR
20 SELECT JSON_EXTRACT(CONVERT(production_companies USING UTF8MB4), '$[*]') FROM movie_dataset ;
21
22 -- Declarar el handler para NOT FOUND (esto es marcar cuando el cursor ha llegado a su fin)
23 DECLARE CONTINUE HANDLER
24 FOR NOT FOUND SET done = TRUE ;
25
26 -- Abrir el cursor
27 OPEN myCursor ;
28 drop table if exists production_companietem;
29 SET @sql_text = 'CREATE TABLE production_companietem ( id int, nameCom VARCHAR(100));';
30 PREPARE stmt FROM @sql_text;
31 EXECUTE stmt;
32 DEALLOCATE PREPARE stmt;
33 cursorLoop: LOOP
34 FETCH myCursor INTO jsonData;
35
36 -- Controlador para buscar cada uno de los arrays
37 SET i = 0;
38
39 -- Si alcanzo el final del cursor entonces salir del ciclo repetitivo
40 IF done THEN
41 LEAVE cursorLoop ;
42 END IF ;
43
44 WHILE(JSON_EXTRACT(jsonData, CONCAT('$[', i, ']')) IS NOT NULL) DO
45 SET jsonId = IFNULL(JSON_EXTRACT(jsonData, CONCAT('$[', i, '].id')), '');
46 SET jsonLabel = IFNULL(JSON_EXTRACT(jsonData, CONCAT('$[', i, '].name')), '');
47 SET i = i + 1;
48
49 SET @sql_text = CONCAT('INSERT INTO production_companietem VALUES (', REPLACE(jsonId, '\'', ''), ', ', jsonLabel, '); ');
50 PREPARE stmt FROM @sql_text;
```



```

51      EXECUTE stmt;
52      DEALLOCATE PREPARE stmt;
53
54  END WHILE;
55
56  END LOOP ;
57
58  select distinct * from production_companieTem;
59  INSERT INTO production_companiesCURSOR
60  SELECT DISTINCT id, nameCom
61  FROM production_companieTem;
62  drop table if exists production_companieTem;
63  CLOSE myCursor ;
64
65  END$$
66  DELIMITER ;
67
68  call TablaProduction_companies();
69
70  CREATE TABLE production_companiesCURSOR (
71      id INT PRIMARY KEY,
72      name varchar(100)
73  );
74
75  DROP TABLE production_companiesCURSOR;

```

Población Tabla “Production_countries”

```

1  USE Movies2023;
2
3  DROP PROCEDURE IF EXISTS TablaProduction_countries;
4
5  DELIMITER $$
6  CREATE PROCEDURE TablaProduction_countries ()
7
8  BEGIN
9
10     DECLARE done INT DEFAULT FALSE ;
11     DECLARE jsonData json ;
12     DECLARE jsonId varchar(250) ;
13     DECLARE jsonLabel varchar(250) ;
14     DECLARE resultSTR LONGTEXT DEFAULT '';
15     DECLARE i INT;
16
17     -- Declarar el cursor
18     DECLARE myCursor
19     CURSOR FOR
20         SELECT JSON_EXTRACT(CONVERT(production_countries USING UTF8MB4), '$[*]') FROM movie_dataset ;
21
22     -- Declarar el handler para NOT FOUND (esto es marcar cuando el cursor ha llegado a su fin)
23     DECLARE CONTINUE HANDLER
24     FOR NOT FOUND SET done = TRUE ;
25

```

```

26 -- Abrir el cursor
27 OPEN myCursor ;
28 drop table if exists production_countriesTem;
29 SET @sql_text = 'CREATE TABLE production_countriesTem ( iso_3166_1 varchar(2), nameCountry VARCHAR(100));';
30 PREPARE stmt FROM @sql_text;
31 EXECUTE stmt;
32 DEALLOCATE PREPARE stmt;
33 cursorLoop: LOOP
34     FETCH myCursor INTO jsonData;
35
36     -- Controlador para buscar cada uno de los arrays
37     SET i = 0;
38
39     -- Si alcanzo el final del cursor entonces salir del ciclo repetitivo
40     IF done THEN
41         LEAVE cursorLoop ;
42     END IF ;
43
44     WHILE(JSON_EXTRACT(jsonData, CONCAT('$[', i, ']')) IS NOT NULL) DO
45         SET jsonId = IFNULL(JSON_EXTRACT(jsonData, CONCAT('$[', i, '].iso_3166_1')), '');
46         SET jsonLabel = IFNULL(JSON_EXTRACT(jsonData, CONCAT('$[', i, '].name')), '');
47         SET i = i + 1;
48
49         SET @sql_text = CONCAT('INSERT INTO production_countriesTem VALUES (', REPLACE(jsonId, '\\', '\\\\'), ', ', jsonLabel, '); ');
50         PREPARE stmt FROM @sql_text;
51         EXECUTE stmt;
52         DEALLOCATE PREPARE stmt;
53
54     END WHILE;
55
56 END LOOP ;
57
58 select distinct * from production_countriesTem;
59 INSERT INTO production_countriesCURSOR
60     SELECT DISTINCT iso_3166_1, nameCountry
61     FROM production_countriesTem;
62 drop table if exists production_countriesTem;
63 CLOSE myCursor ;
64
65 END$$
66 DELIMITER ;
67
68 call TablaProduction_countries();
69
70 SELECT * FROM production_countriesCURSOR;
71
72 CREATE TABLE production_countriesCURSOR (
73     iso_3166_1 varchar(2) PRIMARY KEY,
74     name varchar(100)
75 );

```

Población Tabla “Spoken_languages”

```
1 USE Movies2023;
2
3 DROP PROCEDURE IF EXISTS TablaSpokenLanguages;
4
5 DELIMITER $$
6 CREATE PROCEDURE TablaSpokenLanguages ()
7
8 BEGIN
9
10 DECLARE done INT DEFAULT FALSE ;
11 DECLARE jsonData json ;
12 DECLARE jsonId varchar(250) ;
13 DECLARE jsonLabel varchar(250) ;
14 DECLARE resultSTR LONGTEXT DEFAULT '';
15 DECLARE i INT;
16
17 -- Declarar el cursor
18 DECLARE myCursor
19 CURSOR FOR
20 SELECT JSON_EXTRACT(CONVERT(spoken_languages USING UTF8MB4), '$[*]') FROM movie_dataset ;
21
22 -- Declarar el handler para NOT FOUND (esto es marcar cuando el cursor ha llegado a su fin)
23 DECLARE CONTINUE HANDLER
24 FOR NOT FOUND SET done = TRUE ;
25
26 -- Abrir el cursor
27 OPEN myCursor ;
28 drop table if exists spokenLanguagesTem;
29 SET @sql_text = 'CREATE TABLE spokenLanguagesTem ( iso_639_1 varchar(2), nameLang VARCHAR(100));';
30 PREPARE stmt FROM @sql_text;
31 EXECUTE stmt;
32 DEALLOCATE PREPARE stmt;
33 cursorLoop: LOOP
34 FETCH myCursor INTO jsonData;
35
36 -- Controlador para buscar cada uno de los arrays
37 SET i = 0;
38
39 -- Si alcanzo el final del cursor entonces salir del ciclo repetitivo
40 IF done THEN
41 LEAVE cursorLoop ;
42 END IF ;
43
44 WHILE (JSON_EXTRACT(jsonData, CONCAT('$[', i, ']')) IS NOT NULL) DO
45 SET jsonId = IFNULL(JSON_EXTRACT(jsonData, CONCAT('$[', i, '].iso_639_1')), '');
46 SET jsonLabel = IFNULL(JSON_EXTRACT(jsonData, CONCAT('$[', i, '].name')), '');
47 SET i = i + 1;
48
49 SET @sql_text = CONCAT('INSERT INTO spokenLanguagesTem VALUES (', REPLACE(jsonId, '\'', ''), ', ', jsonLabel, '); ');
50 PREPARE stmt FROM @sql_text;
```

```

51      EXECUTE stmt;
52      DEALLOCATE PREPARE stmt;
53
54  END WHILE;
55
56  END LOOP ;
57
58  select distinct * from spokenLanguagesTem;
59  INSERT INTO spoken_languagesCURSOR
60  SELECT DISTINCT iso_639_1, nameLang
61  FROM spokenLanguagesTem;
62  drop table if exists spokenLanguagesTem;
63  CLOSE myCursor ;
64
65  END$$
66  DELIMITER ;
67
68  call TablaSpokenLanguages();
69
70  SELECT * FROM spoken_languagesCURSOR;
71
72  CREATE TABLE spoken_languagesCURSOR (
73      iso_639_1 varchar(2) PRIMARY KEY,
74      name varchar(100)
75  );

```

LIMPIEZA COLUMNA CREW

El procedimiento comienza declarando algunas variables: "idMovie", "idCrew", "idJSON", "jobJSON", "departmentJSON" y "credit_idJSON", todas ellas son utilizadas para almacenar valores temporales durante la ejecución del procedimiento.

Luego, se define un cursor llamado "myCursor" que selecciona todos los registros de la tabla "movie_dataset" y convierte una columna "crew" en un formato que sea compatible con JSON.

Se define un handler llamado "CONTINUE HANDLER" que se ejecutará cuando el cursor llegue al final de los registros y marcará la variable "done" como "TRUE".

El cursor se abre y se inicia un loop infinito que se repetirá hasta que se alcance el final del cursor o se salga explícitamente del loop. Dentro del loop, se usa la función "FETCH" para

obtener un registro a la vez y se guardan los valores en las variables correspondientes.

Se usa un segundo loop "WHILE" para procesar los datos en formato JSON en el

campo "idCrew". Se extraen los valores de "job", "id", "department" y "credit_id" utilizando la función "JSON_EXTRACT" y se guardan en las variables correspondientes.

Finalmente, se construye una sentencia SQL dinámica usando la función "CONCAT" y

se ejecuta usando las funciones "PREPARE", "EXECUTE" y "DEALLOCATE PREPARE".

Una vez que se han procesado todos los registros, se cierra el cursor y se finaliza el

procedimiento. Para ejecutar el procedimiento, se puede llamar "call TablaCrew();".

```
-- Tabla Crew--
DROP PROCEDURE IF EXISTS TablaCrew;
DELIMITER $$
CREATE PROCEDURE TablaCrew ()
BEGIN
    DECLARE done INT DEFAULT FALSE;
    DECLARE idMovie int;
    DECLARE idCrew text;
    DECLARE idJSON text;
    DECLARE jobJSON text;
    DECLARE departmentJSON text;
    DECLARE credit_idJSON text;
    DECLARE i INT;
    -- Declarar el cursor
    DECLARE myCursor
    CURSOR FOR
    SELECT id, CONVERT(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE
```

```

    (REPLACE(crew, '\'', '\'), '{\'', '{')
    '\': '\', ': ': ', '\', '\', '\', '\', '\': ', ': ', '\', '\', '\')
    USING UTF8mb4 ) FROM movie_dataset;
-- Declarar el handler para NOT FOUND (esto es marcar cuando el cursor ha llegado a su fin)
DECLARE CONTINUE HANDLER
FOR NOT FOUND SET done = TRUE ;
-- Abrir el cursor
OPEN myCursor ;
cursorLoop: LOOP
    FETCH myCursor INTO idMovie, idCrew;
    -- Controlador para buscar cada uno de los arrays
    SET i = 0;
    -- Si alcanzo el final del cursor entonces salir del ciclo repetitivo
    IF done THEN
        LEAVE cursorLoop ;
    END IF ;
    WHILE(JSON_EXTRACT(idCrew, CONCAT('$[', i, '].id')) IS NOT NULL) DO
        SET jobJSON = JSON_EXTRACT(idCrew, CONCAT('$[', i, '].job')) ;
        SET idJSON = JSON_EXTRACT(idCrew, CONCAT('$[', i, '].id')) ;
        SET departmentJSON = JSON_EXTRACT(idCrew, CONCAT('$[', i, '].department')) ;

        SET credit_idJSON = JSON_EXTRACT(idCrew, CONCAT('$[', i, '].credit_id')) ;
        SET i = i + 1;
        SET @sql_text = CONCAT('INSERT Ignore INTO Crew VALUES (', idMovie, ', ',
            REPLACE(idJSON, '\'', ''), ', ', REPLACE(idJSON, '\'', ''), ', ',
            REPLACE(departmentJSON, '\'', ''), ', ', REPLACE(credit_idJSON, '\'', ''), '); ');
        PREPARE stmt FROM @sql_text;
        EXECUTE stmt;
        DEALLOCATE PREPARE stmt;
    END WHILE;
END LOOP ;
CLOSE myCursor ;
END$$
DELIMITER ;
call TablaCrew(); 5 m 36 s

```

```

34  crew,
35  ': 0,', ': \'0\','), ': 1,', ': \'1\','), ': 2,', ': \'2\','), '\id\': ', '\id\': \')',
36
37
38  # Felipe
39  'Louis D\'Esposito', 'Louis D\\\'Esposito'), 'Anthony D\'Agostino', 'Anthony D\\\'Agostino'),
40  'Robin L. D\'Arcy', 'Robin L. D\\\'Arcy'), 'Greg D\'Auria', 'Greg D\\\'Auria'), 'Julie D\'
41  'Johnny D\'Mara', 'Johnny D\\\'Mara'), 'Angela Frances D\'Anna', 'Angela Frances D\\\'Anna'),
42  'Gloria D\'Alessandro', 'Gloria D\\\'Alessandro'), 'James D. D\'Amico', 'James D. D\\\'Amico'),
43  'Pamella D\'Pella', 'Pamella D\\\'Pella'), 'Andrea D\'Amico', 'Andrea D\\\'Amico'), 'Vince
90  # Isaac
91  'Kelly L\'Estrange', 'Kelly L\\\'Estrange '), 'Martin L\'Heureux', 'Martin L\\\'Heureux'),
92  'Derek L\'Estrange', 'Derek L\\\'Estrange'), 'Jennifer L\'Huillier', 'Jennifer L\\\'Huillier'),
93  'Marie-Claude L\'Heureux', 'Marie-Claude L\\\'Heureux'), 'Louis L\'Amour', 'Louis L\\\'Amour'),
94  'Sloane U\'Ren', 'Sloane U\\\'Ren'), 'Muhammed M\'Barek', 'Muhammed M\\\'Barek '), 'James
95  'Hannah d\'Angerio', 'Hannah d\\\'Angerio '), 'Susan d\'Arcy', 'Susan d\\\'Arcy '), 'Alex
96  'Eric d\'Arbeloff', 'Eric d\\\'Arbeloff'), 'Brian Van\'t Hul', 'Brian Van\\\'t Hul'), 'K
97  'Ken\'ichi Yoshida', 'Ken\\\'ichi Yoshida')
124  # Carlos Montero
125  'Andrew O\'Connor', 'Andrew O\\\'Connor'), 'James O\'Brien', 'James O\\\'Brien'), 'David (
126  'Thaddeus O\'Sullivan', 'Thaddeus O\\\'Sullivan'), 'Mark O\'Rowe', 'Mark O\\\'Rowe'), 'Jir
127  'Lorenzo O\'Brien', 'Lorenzo O\\\'Brien'), 'Aoife O\'Sullivan', 'Aoife O\\\'Sullivan'), 'I
128  'Bill O\'Brien', 'Bill O\\\'Brien'), 'Steven O\'Meagher', 'Steven O\\\'Meagher'), 'Edwin (
129  'Liam O\'Brien', 'Liam O\\\'Brien'), 'Gerardine O\'Flynn', 'Gerardine O\\\'Flynn'), 'Rache
130  'Damien O\'Donnell', 'Damien O\\\'Donnell'), 'Rockne S. O\'Bannon', 'Rockne S. O\\\'Bannon'),
131  'Dan O\'Meara', 'Dan O\\\'Meara'), 'Maureen O\'Connell', 'Maureen O\\\'Connell'), 'Johnny
132  'Anthony O\'Brien', 'Anthony O\\\'Brien'), 'Martha O\'Neill', 'Martha O\\\'Neill'), 'Stewa
133  'Maura O\'Connell', 'Maura O\\\'Connell'), 'Sean O\'Donnell', 'Sean O\\\'Donnell'), 'Ken (
134  'Christine O\'Malley', 'Christine O\\\'Malley'), 'Mike O\'Connell', 'Mike O\\\'Connell'),

```

Conclusiones

Este proyecto lo llevé a cabo siguiendo unas normas de modelado de base de datos, como normalización, diseño conceptual, lógico, etc. Después de mucho esfuerzo por mi parte y de todos los demás miembros del grupo, logramos cumplir con todos los objetivos establecidos al inicio del desarrollo. Debemos destacar que enfrentamos varios desafíos a medida que avanzábamos en el proyecto, pero gracias a las directrices de nuestros profesores y a nuestra capacidad de resolver problemas, pudimos satisfacer los requisitos dados.

Una vez que completé todos los pasos necesarios para finalizar el proyecto, pude ver un gran progreso en mi capacidad de tomar decisiones sobre el modelado y la población de una base de datos, además de aprender nuevas estrategias de programación aplicables a situaciones de la vida real. El desarrollo de este proyecto sin duda me acerca más a ser un profesional seguro a la hora de trabajar en proyectos futuros.

Aprendí que los datos deben ser procesados cuidadosamente antes de su tratamiento para evitar errores y consistencia. También comprendí la importancia de las herramientas al momento de interactuar con la base de datos. Puedo concluir que todo lo que estudié este ciclo fue muy importante, ya que lo utilicé en este proyecto integrador. Debemos destacar que este proyecto integrador fue un desafío para mí como estudiante, ya que es la primera vez que trabajo con tanta cantidad de datos

