UNIVERSIDAD TÉCNICA PARTICULAR DE LOJA



TITULACIÓN DE INGENIERÍA EN CIENCIAS DE LA COMPUTACIÓN

TEMA:

Ciclo de vida de bases de datos relacionales normalizada

Tercer semestre.

Alumno: Denis Alexander Cuenca Buele

Tutor: Ing. Nelson Piedra

SEGUNDO BIMESTRE 2021-2022

1. Introducción:

En el presente proyecto se aspira explicar el proceso de normalización de una base de datos "Movie dataset" la cual presenta e1errores de diseño los cuales imposibilitan el uso de la misma, motivo por el cual se aplicaran varias estrategias para la modificación de esta base de datos, que nos permitan volverla un modelo que responda de forma eficiente y estable reglas de normalización y su uso en futuros procesos de implementación. En este documento se encuentra detallado el proceso para construir un modelo estable, detallando los campos existentes en las tablas, su errores, recomendaciones de diseño y aplicación de estos desde un enfoque práctico.

2. Desarrollo:

A continuación se nos propone un modelo de base de datos "movie dataset" la misma que contiene información sobre varias películas como su nombre, presupuesto, cast, entre otras que se detallaran más adelante; la base de datos cuenta con una tabla única "Movie" con un total de 24 atributos entre los cuales resaltan: una clave primaria, claves candidatas, atributos multivaluados, atributos con formato JSON; todo esto en un mismo archivo csv.

A fin de crear un modelo sostenible y consistente se pretende la aplicación de reglas de normalización sobre este primer estado de la base de datos, para lo cual se debe comprender a fondo el modelo con el que se trabajará, así pues se comienza por describir la única entidad que contiene el modelo, la tabla "Movie".

La entidad "Movie" contiene los siguientes atributos:

MOVIE:

- index
- -budget
- -homepage
- -id
- -original language
- -original title
- -overview
- -popularity
- -realese date
- -revenue
- -runtime
- -status
- -tagline
- -title
- -vote average
- -vote count
- genres
- -keywords
- -production companies
- -production countries
- -spoken languages
- -cast
- crew

2.1. Antes de aplicar Normalización:

En la entidad "Movie" existen varios campos que no son compatibles con el concepto de modelo sostenible, según la normalización, por ejemplo existen atributos multivaluados ademas de una dependencia funcional parcial en toda la tabla, lo cual se entra en conflicto con la segunda y tercera forma normal.

Los campos que presentan estos problemas en cuestión son los siguientes:

2.2. Campos Multivaluados (con sus respectivos atributos):

```
- genres: id_movie, genre
```

-production companies: id_movie, id, company

-production countries: id_movie, iso_3166_1, country

-spoken languages: Iso_6391, id_movie, name

-cast: name, id, id_movie

- crew: id_movie, id_crew, job, name, gender, credit_id, department

-Director: name, id, id_movie

Los campos aquí expuestos representan un problema de primera forma normal, al ser campos que almacenan valores multivaluados, por lo cual se los separa en tablas diferentes con esperando así generar un mejor diseño, el resultado final se presenta a continuación donde se describe cada tabla con sus atributos:

2.3. Dependencias funcionales de las nuevas tablas:

A continuación se describe las dependencias funcionales de cada una de las nuevas tablas:

```
    - genres:
        id_movie → genre
    -keywords:
        id_movie, id → keyword
    -movie_production companies:
        id_movie, id → id_movie, id
    -production companies:
        id → company name
    -movie_production countries:
```

```
id movie, iso 3166 1 -> id movie, id
```

-production countries :

 $iso_3166_1 \rightarrow country$

-movie_spoken languages :

Iso_639_1, id_movie → Iso_639_1, id_movie

-spoken languages:

 $Iso_639_1 \rightarrow name$

- crew:

id_movie, credit_id --> job, department, id_crew

- person:

credit_id → name, gender,

-movie_director:

id_movie, credit_id --> id_movie, credit_id

Siendo todas las tablas, a excepción de crew dependencias funcionales completas las cuales responden al modelo que nos propone la segunda forma normal.

ENTIDAD MOVIE:

| Columna | Dominio | Opcion al | Multivaluad o | Proposito | Tipo | comentario |
|----------------------|-------------|--------------|------------------|---------------------------------------------------------------|-------------------|------------------------------------|
| id | int(10) | FALSE | FALSE | Identificar las pelicula de forma única | Clave primaria | |
| Vote Average | double | FALSE | FALSE | Promedio de los votos recolectados | | |
| Vote Count | int(8) | FALSE | FALSE | Total de votos recolectados | | |
| home Page | varchar(20) | FALSE | FALSE | Página oficial de la película | | El valor introducido es una url |
| revenue | double | TRUE | FALSE | Ganancia recolectada por la película | | |
| budget | double | TRUE | FALSE | Presupuesto destinado para la producción de la película | | |
| status | varchar(10) | FALSE | FALSE | Estatus de producción de la película | | |
| tagline | varchar(20) | TRUE | FALSE | Lema de la película | | |
| original language | varchar(20) | FALSE | FALSE | Lenguaje original en el que la película fue grabada | | |
| original title | varchar(20) | FALSE | FALSE | Titulo seleccionado para la película | | |

| runtime | double | TRUE | FALSE | Duración de la proyección de la película | |
|-----------------|---------------|-------|-------|--------------------------------------------------|-------------------------------------------|
| release date | date | FALSE | FALSE | Fecha de salida del juego | Valor introducido en el formato: D/M/A |
| overview | Varchar(50 0) | TRUE | FALSE | Sinopsis de la trama principal de la película | |
| popuarity | Int(15) | TRUE | FALSE | Popularidad en puntos que obtuvo la película | |
| title | Varchar(30 | TRUE | FALSE | Título de la película | |
| Caste | Varchar(30 0) | TRUE | FALSE | Nombre de la trabajadores | |

Entidad Production_companies:

| Linuada Froductio | | | | | | |
|-------------------|--------------|----------|------------------|--------------------------------------------------------|-----------------------|------------|
| Columna | tipo de dato | Opcional | Multivaluad o | Propósito | tipo | comentario |
| id | Int(10) | FALSE | FALSE | Identificar de forma única la columna | Clave primari a | |
| name | Varchar(20) | FALSE | FALSE | Nombre de la compañía productora del película | | |

Entidad movie_production_companies:

| Columna | tipo de dato | Opcional | Multivaluad o | Propósito | tipo | comentario |
|-----------|--------------|----------|------------------|--------------------------------------------------|---------------------------|------------|
| id | Int(10) | FALSE | FALSE | Identificar de forma única la columna | Clave primari a | |
| Id(Movie) | int(10) | FALSE | FALSE | Identificar las pelicula de forma única | Clave primari a y foránea | |

Entidad Production_countries:

| Columna | tipo de dato | Opcional | Multivaluado | Propósito | tipo | comentario |
|------------|--------------|----------|--------------|----------------|--------|------------|
| | | FALSE | FALSE | Nombre del | | |
| | Warehar(1E) | | | país donde se | | |
| name | Varchar(15) | | | produjo la | | |
| | | | | película. | | |
| | | | | Código de | Clave | |
| iso_3166_1 | Varchar(2) | FALSE | FALSE | abreviación de | primar | |
| | | | | un país | ia | |

Entidad Production_countries:

| Columna | tipo de dato | Opcional | Multivaluado | Propósito | tipo | comentario |
|------------|--------------|----------|--------------|-----------------|--------|------------|
| | _ | | | Código de | Clave | |
| iso_3166_1 | Varchar(2) | FALSE | FALSE | abreviación de | primar | |
| | | | | un país | ia | |
| | | | | | Clave | |
| | | | | Identificar las | primar | |
| Id(Movie) | int(10) | FALSE | FALSE | pelicula de | ia y | |
| | | | | forma única | forane | |
| | | | | | a | |

Entidad Spoken_languages:

| Columna | tipo de dato | Opcional | Multivaluado | Propósito | tipo | comentario |
|-------------|--------------|----------|--------------|-------------|----------|------------|
| nama | | FALSE | FALSE | Nombre del | | |
| | Varchar(20) | | | lenguaje | | |
| name | Varchar(20) | | | usado en la | | |
| | | | | película. | | |
| | | FALSE | FALSE | Código de | | |
| Inc. 620, 1 | Varabar(15) | | | abreviación | Clave | |
| Iso_639_1 | Varchar(15) | | | de un | primaria | |
| | | | | lenguaje | | |

Entidad movie spoken languages:

| Entitude movie_spoken_tanguages. | | | | | | | | |
|----------------------------------|--------------|----------|--------------|--------------|----------|------------|--|--|
| Columna | tipo de dato | Opcional | Multivaluado | Propósito | tipo | comentario | | |
| Iso_639_1 | | FALSE | FALSE | Código de | | | | |
| | Varabar(15) | | | abreviación | Clave | | | |
| | Varchar(15) | | | de un | primaria | | | |
| | | | | lenguaje | | | | |
| | | FALSE | FALSE | Identificar | Clave | | | |
| Id(Morrio) | ;n+(10) | | | las película | primaria | | | |
| Id(Movie) | int(10) | | | de forma | y | | | |
| | | | | única | foranea | | | |

Entidad crew:

| Columna | tipo de dato | Opcional | Multivaluado | Propósito | tipo | comentario |
|-----------|--------------|----------|--------------|------------------------------------------------------|-----------------------|--------------------------------------------------------|
| name | Varchar(40) | FALSE | FALSE | Nombre del trabajador | | |
| Gender | Int(1) | FALSE | FALSE | Genero del trabajador | | Atributo identificado con 2 para hombre y 1 para mujer |
| Credit_id | Varchar(20) | FALSE | FALSE | Identificador del crédito del trabajador | | |
| Id | Int(5) | FALSE | FALSE | Identifica de forma única a cada trabajador | Clave primari a | |
| Id(Movie) | int(10) | FALSE | FALSE | Identificar las pelicula de forma única | Clave primari a | |

movie_director:

| Columna | tipo de dato | Opcional | Multivaluado | Propósito | tipo | comentario |
|------------|--------------|----------|--------------|---------------|----------|------------|
| Credit_id | | FALSE | FALSE | Identificado | | |
| | Varchar(20) | | | r del crédito | | |
| | | | | del | | |
| | | | | trabajador | | |
| | int(10) | FALSE | | Identificar | | |
| Id(Morrio) | | | FALSE | las pelicula | Clave | |
| Id(Movie) | | | | de forma | primaria | |
| | | | | única | | |

Genres:

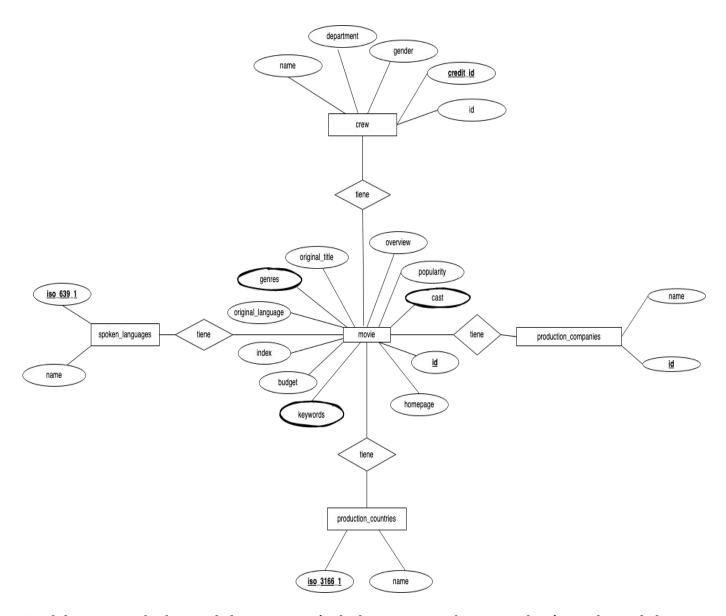
| Columna | tipo de dato | Opcional | Multivaluado | Propósito | tipo | comentario |
|------------|--------------|----------|--------------|--------------|----------|------------|
| Camus | 17 1 (20) | EAT CE | FALSE FALSE | Nombre del | | |
| Genre | Varchar(20) | FALSE | | género | | |
| | int(10) | FALSE | FALSE | Identificar | | |
| Id(Morrio) | | | | las pelicula | Clave | |
| Id(Movie) | | | | de forma | primaria | |
| | | | | única | | |

person:

| Columna | tipo de dato | Opcional | Multivaluado | Propósito | tipo | comentario |
|-----------|--------------|----------|--------------|---------------|----------|------------|
| namo | Varchar(20) | FALSE | FALSE | Nombre del | | |
| name | vaichai(20) | FALSE | FALSE | trabajador | | |
| | | | | Identificado | | |
| Credit id | Manchan(20) | EALCE | FALSE | r del crédito | Clave | |
| Credit_id | Varchar(20) | FALSE | FALSE | del | primaria | |
| | | | | trabajador | | |
| gondon | ;n+(10) | FALSE | FALSE | Genero del | | |
| gender | int(10) | FALSE | FALSE | trabajador | | |

2.4 Diagrama Entidad Relación:

A continuación se presenta el resultado de separar cada una de las tablas y se presenta la relación existente entre ellas, nótese que por el estilo de flecha la carnalidad es de muchos a muchos (m:n) entre la mayoría de las entidades



En el diagrama, todas las entidades, a excepción de department y job, tienen relación con la entidad Movie, pues es esa la entidad de donde surgen, la relación de movie con keywords es de muchos a

muchos, con ello se podrá guardar por una película varias palabras clave, y una palabra clave será capaz de describir a varias películas; similar es el caso para las entidades genres, director, cast y spoken languages. Production companies y movie comparten la misma cardinalidad muchos a muchos en donde una compañía será capaz de producir varias películas y una película producida por muchas compañías; lo mismo sucede con la entidad production cuontrie, La columna crew conserva un formato JSON en donde se detallan a forma de clave-valor los campos que ha de contenerse en una tabla de forma individual separada, en esta tabla las persona pertenecientes al crew son capaces de pertenecer a uno o más departamentos, y estas mismas personas podrán ejercer una o mas funciones de trabajo en la película, por lo que estas son entidades que comparten una cardinalidad de muchos a muchos con la entidad crew.

2.5. Tercera forma normal:

La tercera forma normal nos indica que no debe existir dependencias funcionales transitivas, ademas se añade que sin una tabla se encuentra en primera y segunda forma normal, entonces la tercera norma ya está aplicada por lo que se puede decir que nuestro esquema se encuentra normalizado.

2.6. Implementación en sql:

Estrategia de creación tablas y carga de datos:

Para empezar s e realiza la carga directa de los datos del archivo CSV en nuestra base de datos, previo a realizar la carga en tablas temporales se realiza la limpieza de los atributos con formato JSON:

```
(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLACE(REPLA
```

Luego se da inicio a la creación de procedimientos para extraer los datos cargados en la tabla movie_dataset:

```
96 - WE STATE THE VATIS DECAME BE TO THE STATE OF THE STA
```

```
simple_loop: LOOP
                                                                                                           id Movie,
124
                                                                            JSON_EXTRACT(CONVERT(production_companies | JSON_EXTRACT(production_companies | JSON_E
                                                                                                                                                                                                                                                                     ing utf8mb4), CONCAT("$[",a,"].id")) AS id_production_company,
ing utf8mb4), CONCAT("$[",a,"].name")) AS id_production_compan
                                                                                                                                                                                                                                                                                                                                                                                                                                            id_production_company
                                                                                                          tmp_production_countries (id_movie, iso_3166_1, country)
128
129
                                                                           JSON_EXTRACT(production_countries, CONCAT('$[',a,'].iso_3166_1')) AS iso
JSON_EXTRACT(production_countries , CONCAT('$[',a,'].name')) AS country
                                                                                                                                                                                                                                                                                                                                                                        iso_3166_1,
                                                                                                         tmp_spoken_languages (id_movie, iso_639_1, `language`)
134
                                                                           JSON_EXTRACT(spoken_languages , CONCAT('$[',a,'].iso_639_1')) AS iso JSON_EXTRACT(spoken_languages , CONCAT('$[',a,'].name')) AS language
                                                                                                                                                                                                                                                                                                                                                        iso 639 1.
                                                                            JSON_EXTRACT(C
JSON_EXTRACT(C
140
                                                                                                                                                                                                                                                                                 T("$[",a,"].job")) AS job,
T("$[",a,"].name")) AS name
                                                                                                                                                                                                                 utf8mb4),
                                                                                                                                                                                                                                                                                 T("$[",a,"].name")) AS name,
T("$[",a,"].gender")) AS gender,
T("$[",a,"].credit_id")) AS cred:
T("$[",a,"].department")) AS dep.
                                                                              JSON_EXTRACT(
                                                                                                                                                                                                                utf8mb4),
                                                                              JSON_EXTRACT(
                                                                              JSON_EXTRACT(
                                                                                                                                                                                                                  utf8mb4),
                                                                                                                                                                                                                                                                                                                                                                                                  credit_id,
                                                                             JSON_EXTRACT(
                                                                                                                                                                                                                                                                                                                                                                                                      department
                                                                              movie dataset cleaned m;
```

Una vez se crean y cargan los datos dentro de las tablas temporales se da inicio a la creación de las tablas finales, donde se almacenaran los datos de nuestro modelo ya normalizado.

Primero se borran las tablas implementadas, en caso de existir

```
16
               ΙF
                          spoken_languages;
17
               ΙF
                          production countries;
18
               ΙF
                         production companies;
19
               ΙF
                         worker:
20
                         credit;
21
               ΙF
                         movies_crew;
                         movies_production_companies;
               ΙF
                         movies_production_countries;
23
               ΙF
24
               ΙF
                         movies_spoken_languages;
25
               ΙF
                         movie;
```

```
TOLANT TOBEL MOVIE AS
SELECT DISTRIBUT 'index', budget, genres, homepage, id, keywords, original_language, original_title, overview, popularity, release_date, revenue, runtime, `status`,
tagline, title, vote_average, vote_count, cast, director
1809 movie_dataset_cleaned;
ALIGN LANGE ONLE AND PRIMARY KEY (id);
```

```
iso_639_1 VARCHAR(5) NOT NULL,
ilanguage` VARCHAR(100),
PRIMARY KEY(iso_639_1)

; );

CREATE TABLE movie_genre(
    id_movie INT not null,
    genre VARCHAR(100),
    PRIMARY KEY(iso_639_1)

CONSTRAINT FK_id_movie5 FOREIGN KEY (id_movie) REFERENCES movie(id)

);
```

```
CREATE TABLE production_countries(
    iso_3166_1 VARCHAR(5) NOT NULL,
    country VARCHAR (100),
    PRIMARY KEY(iso_3166_1)
);
```

```
CREATE TABLE production_companies(
    id_production_company INT NOT NULL,
    name_production_company VARCHAR (100),
    PRIMARY KEY(id_production_company)
);
```

```
CREATE TABLE person(
credit_id CHAR(100) NOT NULL,
name` VARCHAR(500),
gender INT,
PRIMARY KEY (credit_id)
);
```

```
CREATE TABLE movie_production_companies (
   id_movie INT MOT NULL,
   id_production_company INT MOT NULL,
   PRIMARY KEY(id_Movie, id_production_company),
   CONSTRAINT FK_id_movie FOREIGN KEY(id_movie) REFERENCES movie(id),
   CONSTRAINT FK_id_production_company FOREIGN KEY(id_production_company)
   | REFERENCES production_companies(id_production_company)
);
```

```
CREATE TABLE movie_production_countries (
   id_movie INT NOT NULL,
   iso_3166_1 VARCHAR (5) NOT NULL,
   PRIMARY KEY(id_movie, iso_3166_1),
   CONSTRAINT FK_id_movie2 FOREIGN KEY(id_movie)
   REFERENCES movie(id),
   CONSTRAINT FK_iso_3166_1 FOREIGN KEY(iso_3166_1)
   REFERENCES production_countries(iso_3166_1)
);
```

```
CREATE TABLE movie_spoken_languages (
    id_movie INT NOT NULL,
    iso_639_1 VARCHAR (5) NOT NULL,
    PRIMARY KEY(id_movie, iso_639_1),
    CONSTRAINT FK_id_movie3 FOREIGN KEY(id_movie)
    REFERENCES movie(id),
    CONSTRAINT FK_iso_639_1 FOREIGN KEY(iso_639_1)
    REFERENCES spoken_languages(iso_639_1)
);
```

```
CREATE TABLE movie_crew (
   id_movie INT NOT NULL,
   credit_id CHAR(100) NOT NULL,
   id_crew INT,
   job VARCHAR(100),
   departament VARCHAR(100),
   PRIMARY KEY(id_movie, credit_id),
   CONSTRAINT FK_id_movie4 FOREIGN KEY (id_movie)
   REFERENCES movie(id),
   CONSTRAINT FK_credit_id FOREIGN KEY (credit_id)
   REFERENCES person(credit_id)
);
```

```
CREATE TABLE movie_keyword(
   id_movie INT met null,
   keyword VARCHAR(100) NOT NULL ,
   PRIMARY KEY (id_movie, keyword),
   CONSTRAINT FK_id_movie6 FOREIGN KEY (id_movie)
   REFERENCES movie(id)
);
```

```
CREATE TABLE movie_cast(
   id_movie INT NUT NULL,
   cast VARCHAR(100),
   PRIMARY KEY (id_movie, cast),
   CONSTRAINT FK_id_movie7 FOREIGN KEY (id_movie)
   REFERENCES movie(id)
);
```

```
CREATE TABLE movie_director(
   id_movie INTEGER NOT NULL,
      credit_id CHAR(100),
    PRIMARY KEY (id_movie, credit_id),
    CONSTRAINT FK_id_movie8 FOREIGN KEY (id_movie)
    REFERENCES movie(id),
    CONSTRAINT FK_credit_id2 FOREIGN KEY (credit_id)
    REFERENCES person(credit_id)
);
```

Con las tablas finales creadas se procede a cargar los datos dentro de las mismas.

```
INSERT INTO production_countries(iso_3166_1, country)
SELECT DISTINCT iso_3166_1, country
FROM tmp_production_countries;

-- Se insertan los datos de la tabla movie_production_countries:
INSERT INTO movie_production_countries(id_movie, iso_3166_1)
SELECT DISTINCT id_movie, iso_3166_1
FROM tmp_production_countries;
-- Se insertan los datos de la tabla production_companies:
INSERT INTO production_companies(id_production_company, name_production_company)
SELECT DISTINCT id_production_company, name_production_company
FROM tmp_production_companies;
```

```
INSERT INTO movie_production_companies(id_movie, id_production_company)
SELECT DISFINCT id_movie, id_production_company
FROM tmp_production_companies;

-- Se insertan los datos de la tabla spoken_languages:
INSERT INTO spoken_languages(iso_639_1, `language`)
SELECT DISFINCT iso_639_1, `language`
FROM tmp_spoken_languages;

-- Se insertan los datos de la tabla movie_spoken_languages:
INSERT INTO movie_spoken_languages(id_movie, iso_639_1)
SELECT DISTINCT id_movie, iso_639_1
FROM tmp_spoken_languages;

-- Se insertan los datos de la tabla person:
INSERT INTO person(credit_id, name, gender)
SELECT DISTINCT credit_id, name, gender
FROM tmp_crew;
```

```
INSERT INTO movie_crew(id_movie, credit_id, id_crew, job, departament)
SELECT DISTINCT id_movie, credit_id, id_crew, job, department
FROM tmp_crew;

-- Se insertan los datos de la tabla movie_crew:
INSERT INTO movie_cast(id_movie, cast)
SELECT DISTINCT id_movie, cast
FROM tmp_cast;

-- Se insertan los datos de la tabla movie_genre:
INSERT INTO movie_genre(id_movie, genre)
SELECT DISTINCT id_movie, genre
FROM tmp_genres;

-- Se insertan los datos de la tabla movie_keyword:
INSERT INTO movie_keyword(id_movie, keyword)
SELECT DISTINCT id_movie, keyword
FROM tmp_movie_keyword;
```

```
INSERT INTO movie_director(id_movie, credit_id)

SELECT DISTINCT m.id, c.credit_id

FROM movie_dataset m, person c

WHERE m.director = REPLACE(c.name, '"', '');
```

De esta manera se obtiene nuestro modelo normalizado y con todos sus datos cargados, podemos entones proceder a eliminar las tablas temporales.

```
tmp_cast;
ΙF
ΙF
          tmp_genres;
ΙF
          tmp_production_countries;
          tmp_production_companies;
ΙF
ΙF
          tmp_spoken_languages;
ΙF
          tmp_crew;
          movie_dataset_cleaned;
ΙF
          tmp_spoken_languages;
ΙF
          tmp movie keyword;
ΙF
```

3. Conclusiones:

- El uso de las formas normales es fundamental en la creación de una base de datos, donde la información está siempre expuesta a errores de almacenamiento, y se corre el riesgo de perderla o hacer una mal uso de esta por malas practicas.
- La base de datos movie dataset tenia errores de diseño los cuales imposibilitan el uso correcto de los datos almacenados, volviéndola para algunos fines una base de datos inútil, al menos como se presento en una primera etapa, sin embargo tras la aplicación de formas normales y otras modificaciones que nos permitan su mejor uso movie dataset se funcional y estable.

Universidad Técnica Particular de Loja

- •Ingeniería en Ciencias de la Computación
- •Materia: Fundamentos de Base de Datos Octubre 2021 Febrero 2022
- •Proyecto Final Ciclo de vida de bases de datos relacionales normalizada
- •Estudiante: Denis Alexander Cuenca Buele | [@]utpl.edu.ec
- •Link del proyecto en Github: https://github.com/DenisCuenca/Proyecto_integrador.git

Profesor: Nelson Piedra | http://investigacion.utpl.edu.ec/nopiedra E

•Fecha: Loja, 8 de febrero del 2022

l a

C

e s

a

u

n

S

t

0

e

X

t

e r

n

0

.