

# PR\_06.2 Dani Gayol Rodríguez

|   |   |
|---|---|
| PR_06.2 Dani Gayol Rodríguez.....                           | 1 |
| Gestion de tablas y formatos.....                           | 1 |
| 1. Carga de Datos Delimitados (Básico): .....               | 2 |
| 2. Definición de Esquema y Comentarios: .....               | 2 |
| 3. Simulación de Pérdida de Datos (DROP TABLE): .....       | 3 |
| Tipos de Datos Complejos .....                              | 4 |
| 1. Uso de STRUCT (Información Personal): .....              | 5 |
| 2. Uso de ARRAY (Géneros de Película): .....                | 5 |
| 3. Uso de MAP (Puntuación de Usuarios por Ocupación): ..... | 5 |
| Consultas Básicas y JOINS.....                              | 6 |
| 1. Agregación Básica: .....                                 | 7 |
| 2. Filtrado y Agrupación: .....                             | 7 |
| 3. INNER JOIN: .....  | 8 |
| 4. LEFT OUTER JOIN: .....                                   | 9 |
| 5. LEFT SEMI JOIN (Simulación de IN/EXISTS): .....          | 9 |

# Gestion de tablas y formatos

## 1. Carga de Datos Delimitados (Básico):

```
1 CREATE EXTERNAL TABLE ml_ratings_ext (  
2     user_id INT,  
3     movie_id INT,  
4     rating INT,  
5     time_stamp BIGINT  
6 )  
7 ROW FORMAT DELIMITED  
8 FIELDS TERMINATED BY '\t'  
9 LINES TERMINATED BY '\n'  
10 LOCATION '/user/hive/movielens/u.data';
```

```
1 LOAD DATA INPATH '/user/maria_dev/movielens/u.data' INTO TABLE ml_ratings_ext;
```

| ml_ratings_ext.user_id | ml_ratings_ext.movie_id | ml_ratings_ext.rating | ml_ratings_ext.time_stamp |
|------------------------|-------------------------|-----------------------|---------------------------|
| 196                    | 242                     | 3                     | 881250949                 |
| 186                    | 302                     | 3                     | 891717742                 |
| 22                     | 377                     | 1                     | 878887116                 |
| 244                    | 51                      | 2                     | 880606923                 |
| 166                    | 346                     | 1                     | 886397596                 |
| 298                    | 474                     | 4                     | 884182806                 |

## 2. Definición de Esquema y Comentarios:

```

1 CREATE TABLE ml_users_managed (
2     user_id INT,
3     age INT,
4     gender STRING,
5     occupation STRING COMMENT 'Tipo de ocupación puesta por el usuario',
6     zipcode STRING
7 )
8 COMMENT 'Tabla gestionada con información de usuarios'
9 ROW FORMAT DELIMITED
10 FIELDS TERMINATED BY '|'
11 STORED AS TEXTFILE;

```

```

1 LOAD DATA INPATH '/user/maria_dev/movielens/u.user' INTO TABLE ml_users_managed;

```

| ml_users_managed.user_id | ml_users_managed.age | ml_users_managed.gender | ml_users_managed.occupation | ml_users_managed.zipcode |
|--------------------------|----------------------|-------------------------|-----------------------------|--------------------------|
| 1                        | 24                   | M                       | technician                  | 85711                    |
| 2                        | 53                   | F                       | other                       | 94043                    |
| 3                        | 23                   | M                       | writer                      | 32067                    |
| 4                        | 24                   | M                       | technician                  | 43537                    |
| 5                        | 33                   | F                       | other                       | 15213                    |








### 3. Simulación de Pérdida de Datos (DROP TABLE):

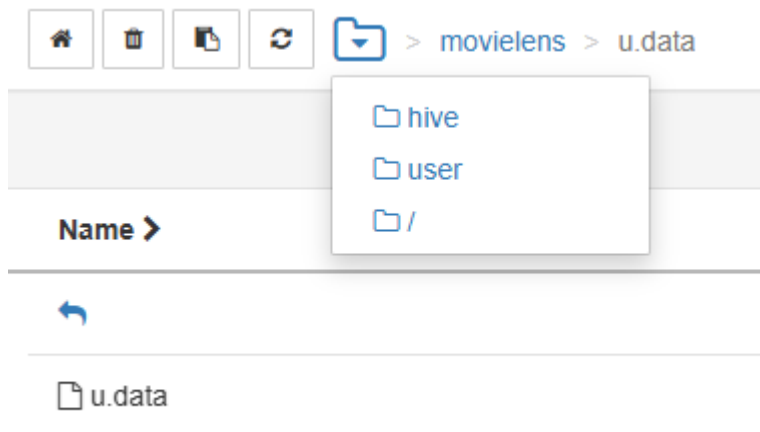
Para la tabla “ml\_ratings\_ext”, al ser “external”, Hive borra el metadata pero no elimina los archivos del directorio asociado

```

1 DROP TABLE ml_ratings_ext;

```

| TABLES   4  |                  |   |  |
|---|------------------|---|---|
| Search  |                  |  |   |
|  | ml_users_managed |   |   |
|  | peliculas        |   |   |
|  | usuarios         |   |   |
|  | votos            |   |   |



Para la tabla “ml\_users\_managed”, al ser “managed”, Hive controla el directorio y borra la carpeta y su contenido

```
1 DROP TABLE ml_users_managed;
```

| TABLES   3  |           |     |  |
|---|-----------|---|---|
| Search  |           |  |   |
|  | películas |   |   |
|  | usuarios  |   |   |
|  | votos     |   |   |

# Tipos de Datos Complejos

## 1. Uso de STRUCT (Información Personal):

```
1 CREATE TABLE ml_user_info (  
2     user_id INT,  
3     personal_data STRUCT<  
4         age: INT,  
5         gender: STRING  
6     >,  
7     occupation STRING,  
8     zipcode STRING  
9 )  
10 ROW FORMAT DELIMITED  
11 FIELDS TERMINATED BY '|';
```

```
1 SELECT user_id, personal_data.gender  
2 FROM ml_user_info  
3 WHERE personal_data.gender = 'F';
```

## 2. Uso de ARRAY (Géneros de Película):

Primera consulta

```
1 SELECT movie_id, size(genres_list) AS num_genres  
2 FROM ml_movies_array;
```

Segunda consulta

```
1 SELECT genres_list[1]  
2 FROM ml_movies_array  
3 WHERE movie_id = 50;
```

## 3. Uso de MAP (Puntuación de Usuarios por Ocupación):

```
1 SELECT map_keys(avg_rating_by_occupation)  
2 FROM ml_occupation_ratings;
```



# Consultas Básicas y JOINS

## 1. Agregación Básica:

Como eliminamos esta tabla anteriormente, solo basta con crearla otra vez ya que los datos siguen guardados

```
1 CREATE EXTERNAL TABLE ml_ratings_ext (  
2     user_id INT,  
3     movie_id INT,  
4     rating INT,  
5     time_stamp BIGINT  
6 )  
7 ROW FORMAT DELIMITED  
8 FIELDS TERMINATED BY '\t'  
9 LINES TERMINATED BY '\n'  
10 LOCATION '/user/hive/movielens/u.data';
```

Ahora hacemos la consulta

```
1 SELECT AVG(rating) AS avg_rating  
2 FROM ml_ratings_ext;
```

avg\_rating

3.52986

## 2. Filtrado y Agrupación:

Como eliminamos esta tabla anteriormente, aquí si hay que crear la tabla y insertarle los datos otra vez

```

1 CREATE TABLE ml_users_managed (
2     user_id INT,
3     age INT,
4     gender STRING,
5     occupation STRING COMMENT 'Tipo de ocupación declarada por el usuario',
6     zipcode STRING
7 )
8 COMMENT 'Tabla gestionada con información de usuarios MovieLens'
9 ROW FORMAT DELIMITED
10 FIELDS TERMINATED BY '|'
11 STORED AS TEXTFILE;

```

```

1 LOAD DATA INPATH '/user/maria_dev/movielens/u.user' INTO TABLE ml_users_managed;

```

Ahora ya podemos hacer la consulta

```

1 SELECT gender, occupation, COUNT(*) AS total_users
2 FROM ml_users_managed
3 GROUP BY gender, occupation
4 ORDER BY total_users DESC;

```

| gender | occupation | total_users |
|--------|------------|-------------|
| M      | student    | 136         |
| M      | other      | 69          |
| M      | educator   | 69          |
| M      | engineer   | 65          |
| F      | student    | 60          |

### 3. INNER JOIN:

Si quieres que la media no salga con muchas decimales puedes usar “ROUND” para redondear al número de decimales que quieras

```

1 SELECT u.occupation, AVG(r.rating) AS avg_rating
2 FROM ml_ratings_ext r
3 INNER JOIN ml_users_managed u
4     ON r.user_id = u.user_id
5 GROUP BY u.occupation
6 ORDER BY avg_rating DESC;

```



| u.occupation | avg_rating         |
|--------------|--------------------|
| none         | 3.779134295227525  |
| lawyer       | 3.7353159851301116 |
| doctor       | 3.688888888888889  |
| educator     | 3.6706206312221985 |
| artist       | 3.653379549393414  |

#### 4. LEFT OUTER JOIN:

```
1 CREATE TABLE unrated_users (user_id INT);
2
3 INSERT INTO unrated_users VALUES (1), (2), (3), (4), (5);
```

```
1 SELECT u.user_id, r.rating
2 FROM ml_users_managed u
3 LEFT OUTER JOIN ml_ratings_ext r
4     ON u.user_id = r.user_id
5 WHERE r.rating IS NULL;
```

#### 5. LEFT SEMI JOIN (Simulación de IN/EXISTS):

```
1 SELECT u.user_id
2 FROM ml_users_managed u
3 LEFT SEMI JOIN ml_ratings_ext r
4     ON u.user_id = r.user_id
5 WHERE u.occupation = 'programmer';
```

u.user\_id

17

29

45

53

55

58

82