



## Ejercicio SQL 3

### Introducción

Decides convertirte en desarrollador de aplicaciones móviles, pero antes de comenzar a crearlas, necesitas conocer a la competencia. Para esto, has buscado y encontrado una base de datos de aplicaciones móviles. Quieres usar esta información para tomar algunas decisiones de diseño. Algunas de estas decisiones son:

- Versión de Android a utilizar
- Hacer una aplicación gratuita o de pago
- Frecuencia recomendada de las actualizaciones
- Qué aplicación recomendarías a algún usuario

Para conocer la respuesta a esta y otras preguntas, decides hacer un programa que consulte a la base de datos utilizando Python y SQL.

### Instrucciones

Este laboratorio se divide en tres partes: procesamiento de la información, modelación y consultas. Además, se incluye una sección inicial con información de la base de datos. Cada una de ellas se describe a continuación.

### Base de datos

#### Archivos CSV

Los archivos CSV (del inglés *comma-separated values*) son un tipo de documento en formato abierto sencillo para representar datos en forma de tabla, en las que las columnas se separan por comas (o punto y coma en donde la coma es el separador decimal: `Chile,Perú,Argentina,España,Brasil...`) y las filas por saltos de línea.<sup>1</sup> En caso que algún elemento de una columna posea coma en el texto, este elemento será encerrado entre comillas, por ejemplo:

---

<sup>1</sup>Cita de wikipedia

PLANK!,GAME,4.7,7196,38M,"500,000+",Free,0,Everyone,Arcade,"July 17, 2018",1.0.2,4.3 and up

En este ejemplo el “500,000+” y “July 17, 2018” poseen comas y por lo tanto, en el archivo se le pone comillas para diferenciar la coma del texto y la coma como separador. Para la lectura de estos archivos se recomienda utilizar alguna librería como **Pandas** o **CSV**.

## Archivos entregados

Para este laboratorio, dispondrás de 2 archivos CSV con la información inicial del sistema. El primero se llama *googleplaystore.csv*, el cual contiene la información de cada aplicación. Sus columnas son las siguientes:

- **App:** nombre de la aplicación
- **Category:** categoría de la aplicación. Entre estas se encuentran *ART\_AND\_DESIGN*, *FAMILY*, *GAME*, etc.
- **RatingOverall:** *rating* de la aplicación
- **Reviews:** número de comentarios que tiene la aplicación. Esta columna debe ser ignorada. En otro archivo estan todos los comentarios y esa cantidad será la utilizada.
- **Size:** tamaño de la aplicación
- **Installs:** indicador que representa cantidad de descarga de la aplicación.
- **Type:** indica si es pagada o gratuita.
- **Price:** precio de la aplicación
- **Content Rating:** grupo de edad a la que va dirigida la aplicación.
- **Genres:** género de la aplicación, se puede considerar como una sub-categoría y cada aplicación puede tener más de un género. En ese caso, cada género está dividido por un punto y coma (“,”).
- **Last Updated:** fecha de la última vez que fue actualizada la aplicación
- **Current Ver:** versión de la aplicación.
- **Android Ver:** versión de *Android* requerida por la aplicación.

El segundo archivo se llama *googleplaystore\_user\_reviews.csv* que contiene diferentes comentarios de los usuarios para cada aplicación. Las columnas son:

- **App:** nombre de la aplicación.
- **Translated\_Review:** Comentario de la aplicación. Todas traducidas al ingles.

- **Sentiment:** Indica si el comentario es positivo, negativo o neutral.
- **Sentiment\_Polarityws:** puntaje asociado a la polaridad del comentario, valor entre -1 y 1. Esta columna debe ser ignorada.
- **Sentiment\_Subjectivity:** puntaje asociado al sentimiento entregado en el comentario. Esta columna debe ser ignorada.

## Procesamiento

Dentro de los archivos, hay varias aplicaciones o comentarios con valores **NaN** los cuales deberán eliminar antes de pasar a la siguiente etapa.

## Modelación

Luego del procesamiento, deberán pasar la información a un base de datos relacional en SQL. Para esto deben crear múltiples entidades y relaciones que **aseguren las restricciones de integridad**. Es requisito que existan al menos **3 tablas que representen entidades y 2 que representen relaciones entre entidades**, y todas las entidades posean un **id numérico único** de manera que las comparaciones para realizar *joins* entre tablas se haga solamente basado en estos ids. El modelo debe incluir la misma información entregada en los .CSV.

## main.py

Para este laboratorio, se les adjunta un archivo de nombre *main.py* y *main.ipynb* que deberán completar con la solución. El archivo contiene el formato esperado de las consultas descritas en la siguiente sección. Es obligatorio seguir el formato allí descrito, pero queda a su criterio el desarrollo de cualquier otra función o clase que sea necesaria. Podrán notar del archivo, que este ya incluye ejecuciones de las consultas y se basa en que estas retornen la solución en el formato solicitado. En caso de no respetar dicho formato, no será evaluada la funcionalidad.

## Ediciones a la base de datos

Antes de las consultas, su desarrollo del laboratorio debe incluir 4 funciones para agregar, editar y eliminar información de su base de datos. Las funciones son:

- **add\_app(app\_data):** esta función recibe un diccionario con la información de la aplicación. El formato del diccionario es:

Llave del diccionario	Valor del diccionario	Condición del valor
name	nombre de la aplicación	Debe ser un <b>string</b> .
category	categoría de la aplicación	Debe ser un <b>string</b> .
rating	rating de la aplicación.	Debe ser un número $> 0$ .
size	tamaño de la aplicación	Debe ser un número $> 0$ .
price	precio de la aplicación.	Debe ser un número $\geq 0$ .
version	versión de la aplicación	Deber ser un número $> 0$ .
android	versión de <i>android</i> de la aplicación	Deber ser un número $> 0$ .
genres	géneros de la aplicación	Debe ser una lista de <b>string</b> .

Cuando se agrega la aplicación, debe suponer que esta tiene 0 comentarios, 0 instalaciones y que la última actualización fue cuando se agregó la aplicación a la base de datos. Para la lista de géneros, si algún género no existe, este debe ser agregado a la base de datos. En caso que un valor del diccionario no sea correcta, debe imprimir en consola “No es posible agregar la aplicación” y retornar **None**.

- **add\_comment(app\_name, text, sentiment)**: agregar un nuevo comentario de una aplicación. Debe recibir el nombre de la aplicación, un texto y un sentimiento (positivo, neutro o negativo). Puede suponer que siempre va a existir el nombre de la aplicación, *text* será un *string*. Esta función debe retornar **None**.
- **download\_app(app\_name)**: descarga una aplicación. Esto incrementa en 1 el número del atributo descargas. Puede suponer que siempre va a existir el nombre de la aplicación. Esta función debe retornar **None**.
- **delete\_app(app\_name)**: elimina toda la información asociada a la aplicación con el nombre ingresado. Puede suponer que siempre va a existir el nombre de la aplicación. Esta función debe retornar **None**.

## Consultas

Las consultas mínimas que debe soportar su sistema son:

- **get\_info(app)**: dado el nombre de una aplicación (app), retorna la información de dicha aplicación en forma del diccionario. En particular, la información retornada debe seguir el siguiente formato en términos de *key*:

```
{
    "name": "nombre de la aplicación",
    "category": "categoría de la aplicación",
    "ratingOverall": "rating de la aplicación.",
    "reviews": "cantidad de comentarios",
    "size": "tamaño de la aplicación",
    "installs": "cantidad de instalaciones",
}
```

En caso que el nombre de la aplicación no exista, debe retornar el siguiente diccionario:

```
{
    "error": "La aplicación {nombre de la aplicación} no existe en la base de datos"
}
```

- **best\_by\_genre(n, genre)**: esta consulta recibe un número entero (n) y un género (genre) de la base de datos. Imprime en consola los nombres de las n mejores aplicaciones de dicho género. En caso de que el genero no exista, debe imprimir en consola “Este género no existe”.
- **count\_by\_version(date\_1, date\_2)**: dadas dos fechas (date\_1, date\_2), deben imprimir en consola cuantas aplicaciones por versión de Android hay y cada versión, indicar la categoría mejor evaluada. Para esta última parte, deben considerar el *rating* promedio de las aplicaciones de dicha categoría para encontrar la mejor evaluada. Las fechas vendrán como string y su formato será “YYYY-MM-DD”.
- **price\_of\_the\_best\_by\_genre(n, genre)**: dado un número entero (n) y un género (genre), debe retornar el precio promedio de las n aplicaciones mejor evaluadas que posean dicho género. En caso de que no exista ese género, debe retornar -1. Esta consulta **debe ser realizada completamente con SQL**. Se espera que ninguna acción (aparte del **return**) sea realizada con comandos propios de Python, solo con comandos de SQLite.
- **recommend(genre, size)**: Dado un género (genre) y un número (size), debes recomendar hasta 5 aplicaciones de dicho género que tengan un tamaño menor al indicado en size. Para recomendar, debes ordenar las aplicaciones de posean dicho género en base a la cantidad de comentarios con sentimiento positivo. Esta consulta debe retornar una lista con los nombres de las aplicaciones recomendadas.
- **need\_update(app)**: dado un **string** con el nombre de una aplicación (app), se debe calcular el fecha promedio de la última actualización de las aplicaciones que pertenecen a dicha categoría y que poseen una versión de Android igual o superior a la aplicación dada. La función retorna True o False de acuerdo a si la fecha de la última actualización de la aplicación es más antigua que el promedio de su categoría. En esta función, **debe ser realizada completamente con SQL**, toda la búsqueda y filtrado debe realizarse mediante consultas SQL, puede utilizar Python para calcular el promedio y retornar lo correspondiente.
- **app\_with\_more\_income()**: se debe retornar el nombre de la aplicación que ha recibido más ingresos. En caso de empate, se deben retornar hasta cinco nombres. Haga todos los supuestos que estime conveniente y especifíquelos en el informe.