



## IIC2115 – Programación como Herramienta para la Ingeniería (I/2023)

### Ejercicio Capítulo 2a

#### Aspectos generales

- **Objetivos:** Aplicar los contenidos de análisis exploratorio de datos para completar y expandir una base de datos incompleta y responder consultas sobre la misma.
- **Lugar de entrega:** lunes 10 de abril a las 16:50 hrs. en sitio del curso.
- **Formato de entrega:** archivo Python Notebook (**C2a.ipynb**) con el avance logrado durante la sesión. Utilice múltiples celdas de texto y código para facilitar el trabajo del cuerpo docente.

#### Introducción

Con el fin de ejercitar los contenidos de análisis exploratorio de datos en Python, en este ejercicio deberá realizar los pasos básicos para expandir y completar una base de datos con información faltante. El cómo hacerlo en cada caso será una decisión de uds., que deberá ser tomada y **JUSTIFICADA** en base a las características de los datos analizados. Además de esto, una vez teniendo la base de datos completa, deberá contestar una serie de consultas con respecto a los datos, que requerirán el uso de técnicas de agregación, agrupación y visualización.

#### Descripción del problema

Netflix es un servicio de *streaming* muy conocido con vasto catálogo de películas y series. La empresa ha hecho un concurso de “Director/a por 1 día!!” y usted ha ganado el ticket dorado. Como persona cinéfila, se ha propuesto crear una película de taquilla y para eso, decide analizar las tendencias en el mundo del

*streaming*. Ha descubierto que en la plataforma de datasets públicos Kaggle hay un dataset disponible con el catálogo de Netflix hasta el 2021. Su objetivo es analizarlo para ayudarle con su nueva película.

## La base de datos

La base de datos se encuentra disponible en el Syllabus, en la carpeta de material de clases. En el archivo `netflix.csv`, encontrará información del catálogo de series y películas de Netflix por medio de las siguientes columnas:

<b>show_id</b>	Identificador único del programa
<b>type</b>	Indica si es una película o serie
<b>title</b>	Título del programa
<b>director</b>	Nombre de quién(es) dirigieron la producción
<b>country</b>	País de producción
<b>date_added</b>	Año en que se añade al catálogo de Netflix
<b>release_year</b>	Año de estreno
<b>rating</b>	Clasificación por edades
<b>duration</b>	Duración en minutos
<b>listed_in</b>	Lista de géneros/clasificaciones a las que pertenece la película

Table 1: Descripción de las columnas del dataset de Netflix

## Misiones

1. **Carga y exploración:** cargue el archivo con los datos y describa su contenido, indicando qué columnas tienen información incompleta. Finalmente, visualice las variables, con el fin de evaluar la existencia de *outliers*.
2. **Imputación y eliminación:** en esta tarea hay varias cosas que revisar. Primero, hay que revisar la presencia de elementos duplicados y manejarlos. Segundo, hay que resolver el problema de la información faltante. Para cada una de las columnas con elementos faltantes, impute los valores en base a algún criterio basado en los datos. Además de esto, analice la posible eliminación de filas en base a los valores faltantes.

3. **Expandir el dataset:** considerando los datos que contiene la columna `listed_in`, genere una nueva tabla que para cada `show_id` tenga columnas binarias indicando si el programa pertenece a una de las siguientes categorías:

- Action & Adventure
- Comedy
- Drama
- Horror
- Mystery
- Sci-Fi & Fantasy
- Thrillers
- Independent
- Children & Family

Se recomienda usar el comando `str.contains()` para revisar presencia de texto en los valores de la columna. Hay que tener cuidado con categorías que sean sinónimos o palabras que se escriban plural. En el archivo `categories.csv` puede ver una lista de las 42 categorías que existen y que aparecen en la columna `listed_in`. Como último paso, **debe unir esta tabla con la original**.

4. **Consultas:** conteste cada una de las siguientes consultas, justificando los análisis y supuestos realizados:

- a) ¿Cuál es el género con más presencia en los últimos 5 años? ¿En qué década hay más películas de ese género? ¿Y qué país ha producido más películas de este tipo?
- b) En promedio, ¿cuánto demora en llegar un programa al catálogo de Netflix? Para películas anteriores a la creación de Netflix, utilice la fecha de la creación del catálogo de *streaming* (2008) para sus cálculos.
- c) Visualice y analice cómo ha cambiado en el tiempo la clasificación por edades del género de **Horror**.