

▼ Actividad - Estadística básica

- **Nombre:**
- **Matrícula:**

Entregar: Archivo PDF de la actividad, así como el archivo .ipynb en tu repositorio. **Nota:** Recuerda habrá una penalización de **50** puntos si la actividad fue entregada fuera de la fecha límite.

Carga el conjunto de datos `bestsellers with categories.csv` (se encuentra en el repositorio de la clase) y realiza un análisis estadístico de las variables.

```
1 # Carga las librerías necesarias.
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns

1 # Carga el conjunto de datos al ambiente de Google Colab y muestra los primeros
2 from google.colab import files
3
4 uploaded = files.upload()
5
6 for fn in uploaded.keys():
7     print('User uploaded file "{name}" with length {length} bytes'.format(
8         name=fn, length=len(uploaded[fn])))
9
10 # 6 renglones.
11
```

Choose Files

bestsellers ...tegories.csv

- **bestsellers with categories.csv**(text/csv) - 51161 bytes, last modified: 3/21/2023 - 100% done

Saving bestsellers with categories.csv to bestsellers with categories.csv

User uploaded file "bestsellers with categories.csv" with length 51161 bytes

```
1 # 6 renglones.
2
3 df = pd.read_csv('bestsellers with categories.csv')
4 df.head(6)
```

	Name	Author	User Rating	Reviews	Price	Year	Genre
0	10-Day Green Smoothie Cleanse	JJ Smith	4.7	17350	8	2016	Non Fiction
1	11/22/63: A Novel	Stephen King	4.6	2052	22	2011	Fiction
2	12 Rules for Life: An Antidote to Chaos	Jordan B. Peterson	4.7	18979	15	2018	Non Fiction
3	1984 (Signet Classics)	George Orwell	4.7	21424	6	2017	Fiction
4	5,000 Awesome Facts (About)	National	4.8	7777	10	2010	Non

El conjunto de datos es una tabla que contiene el top 50 de los libros más vendidos por Amazon por año desde 2009 hasta 2019. Cada libro está clasificado como Ficción o No ficción.

Las variables que contiene son:

- **Name:** Nombre del libro.
- **Author:** Autor.
- **User Rating:** Calificación promedio que los usuarios asignaron al libro (1-5).
- **Reviews:** Número de reseñas.
- **Price:** Precio del libro.
- **Year:** Año de publicación.
- **Genre:** Género literario (ficción/no ficción).

```
1 # Crea una tabla resumen con las estadísticas generales de las variables
```

```

1 # Crea una tabla resumen con los estadísticos generales de las variables
2 # numéricas.
3 df.describe()
4
5

```

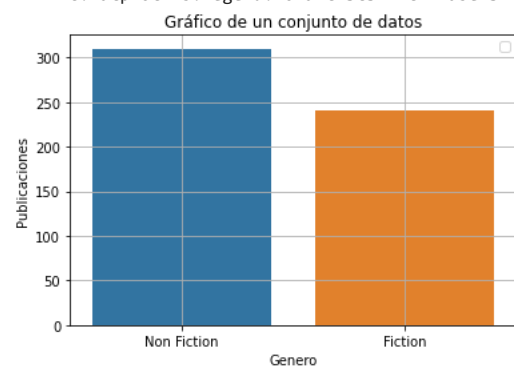
	User Rating	Reviews	Price	Year
count	550.000000	550.000000	550.000000	550.000000
mean	4.618364	11953.281818	13.100000	2014.000000
std	0.226980	11731.132017	10.842262	3.165156
min	3.300000	37.000000	0.000000	2009.000000
25%	4.500000	4058.000000	7.000000	2011.000000
50%	4.700000	8580.000000	11.000000	2014.000000
75%	4.800000	17253.250000	16.000000	2017.000000
max	4.900000	87841.000000	105.000000	2019.000000

```

1 ## ¿Cuál es el género con más publicaciones? Muéstralo en un gráfico.
2 # Configuramos el tamaño de imagen
3 fig = plt.figure(figsize=(6,4))
4
5 # Graficamos
6 sns.countplot(data=df, x = 'Genre')
7
8 # Agregamos títulos a los ejes y al gráfico
9 plt.xlabel('Genero')
10 plt.ylabel('Publicaciones')
11 plt.title('Gráfico de un conjunto de datos')
12
13 # Aquí la leyenda hace mucho más sentido
14 plt.legend(loc='best')
15
16 # Agregamos la cuadrícula para que se vea mejor
17 plt.grid(True)

```

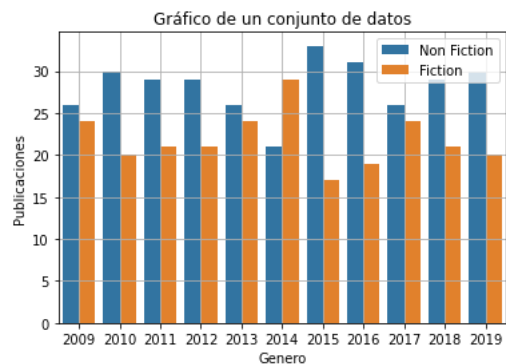
WARNING:matplotlib.legend:No artists with labels found to put in legend. Note that arti



```

1 # ¿Cuántos libros del top 50 se publicaron por género en cada año? ¿Hay algún
2 # año donde hubo más libros de ficción en el top 50?. Muéstralo en un gráfico.
3 # Configuramos el tamaño de imagen
4 fig = plt.figure(figsize=(6,4))
5
6 # Graficamos
7 sns.countplot(data=df, x = 'Year', hue = 'Genre')
8
9 # Agregamos títulos a los ejes y al gráfico
10 plt.xlabel('Genero')
11 plt.ylabel('Publicaciones')
12 plt.title('Gráfico de un conjunto de datos')
13
14 # Aquí la leyenda hace mucho más sentido
15 plt.legend(loc='best')
16
17 # Agregamos la cuadrícula para que se vea mejor
18 plt.grid(True)

```

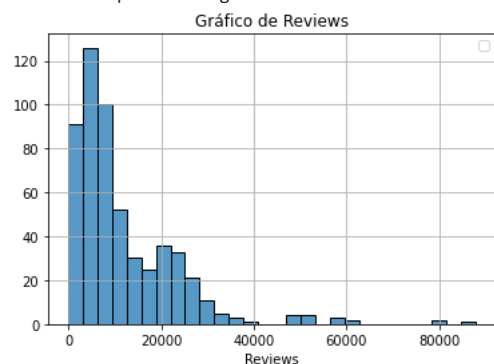


```

1 # ¿Cómo se distribuye la variable Review? Muestra el histograma.
2
3 # Configuramos el tamaño de imagen
4 fig = plt.figure(figsize=(6,4))
5
6 # Graficamos
7 sns.histplot(data=df, x='Reviews')
8
9 # Agregamos títulos a los ejes y al gráfico
10 plt.xlabel('Reviews')
11 plt.ylabel('')
12 plt.title('Gráfico de Reviews')
13
14 # Aquí la leyenda hace mucho más sentido
15 plt.legend(loc='best')
16
17 # Agregamos la cuadrícula para que se vea mejor
18 plt.grid(True)

```

WARNING:matplotlib.legend:No artists with labels found to put in legend. Note that arti



```

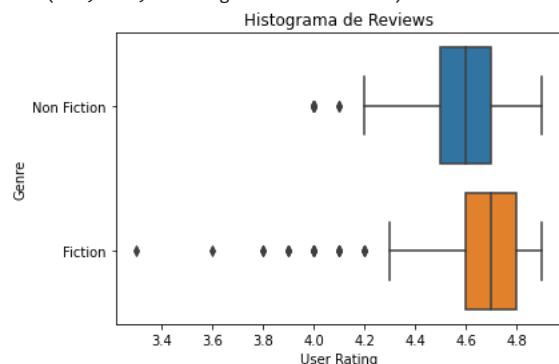
1 # Ahora muéstralo en un gráfico de caja y bigote.
2 # Configuramos el tamaño de imagen
3 fig = plt.figure(figsize=(6,4))
4
5 # Graficamos
6 sns.boxplot(data=df, x='Reviews')
7
8 # Ejes y título
9 plt.title('Histograma de Reviews')

```

```
Text(0.5, 1.0, 'Histograma de Reviews')
Histograma de Reviews
```

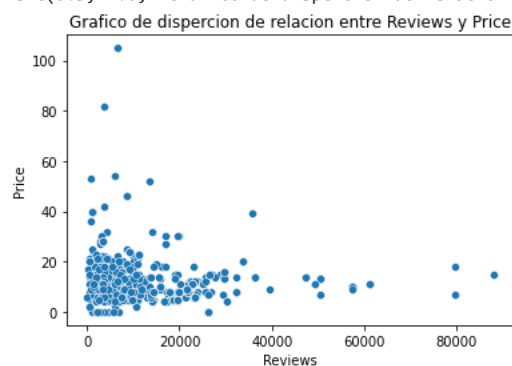
```
1 # ¿Cómo se compara la evaluación del libro por género? ¿Qué genero es mejor
2 # evaluado por los lectores? Muéstralo en un solo gráfico de caja y bigote.
3 fig = plt.figure(figsize=(6,4))
4
5 # Graficamos
6 sns.boxplot(data=df, x='User Rating', y = 'Genre')
7
8 # Ejes y título
9 plt.title('Histograma de Reviews')
```

```
Text(0.5, 1.0, 'Histograma de Reviews')
```

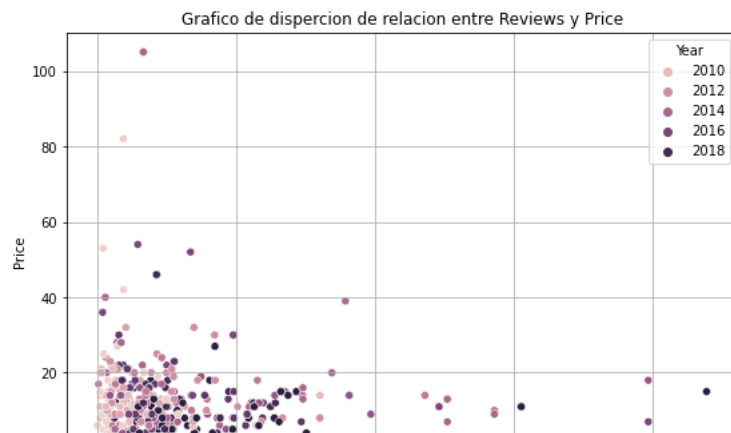


```
1 # ¿Cuál es la relación entre el número de reseñas y precios? Muéstralo en un
2 # gráfico de dispersión.
3
4 fig = plt.figure(figsize=(6,4))
5
6 # Graficamos
7 sns.scatterplot(data=df, x='Reviews', y = 'Price')
8
9 # Ejes y título
10 plt.title('Grafico de dispersion de relacion entre Reviews y Price')
```

```
Text(0.5, 1.0, 'Grafico de dispersion de relacion entre Reviews y Price')
```



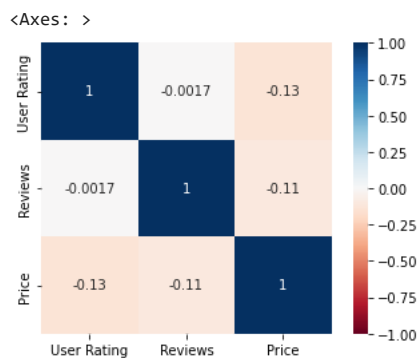
```
1 # De la pregunta anterior, ¿influye algo el año de publicación? ¿Cuál es la
2 # relación entre el número de reseñas, el precio y el año de publicación?
3 # IMPORTANTE: Selecciona una paleta de colores adecuada.
4
5 fig = plt.figure(figsize=(9,6))
6
7 # Graficamos
8 sns.scatterplot(data=df, x='Reviews', y = 'Price', hue = 'Year')
9
10 # Ejes y título
11 plt.title('Grafico de dispersion de relacion entre Reviews y Price')
12
13 # Agregamos la cuadrícula para que se vea mejor
14 plt.grid(True)
```



```

1 # ¿Cuál es la correlación entre las variables numéricas? Muéstralo en un
2 # gráfico. La variable año, a pesar de ser numérica, la vamos a considerar como
3 # cualitativa, así que la eliminaremos del análisis.
4
5 # Vamos a graficar la matriz de correlación del dataset Iris
6
7 iris_corr = df[['User Rating', 'Reviews', 'Price']].corr()
8
9 # Gráfico heatmap. Seleccionamos los valores extremos con vmin y vmax.
10 # El mapa de color que usaremos es de un extremo azul y del otro rojo.
11 # Con annot podemos desplegar el valor de cada celda
12 # Con square hacemos que el gráfico sea simétrico en tamaño de ejes
13 sns.heatmap(data=iris_corr, vmin=-1, vmax=1, cmap = 'RdBu', annot=True, square = True)

```



¿Cuáles variables tiene una fuerte relación positiva entre sí y cuáles tienen una fuerte relación negativa? (Esta pregunta no es de código)
 Responde la pregunta en la siguiente celda de texto.

**** Realmente no hay variables que tengan una relación positiva, solamente entre ellas mismas, sin embargo las que tienen un valor más acercado al 0 son entre user rating y price, reviews y price. Para aquellas negativas más alejadas del 0 se encuentran las relaciones entre user ratings y reviews****

```

1 # Haz una gráfica donde podemos comparar la relación entre las tres variables
2 # numéricas (User Rating, Reviews y Price) y que, además, podamos ver el efecto
3 # del libro. La variable año, a pesar de ser numérica, la vamos a considerar como
4 # cualitativa, así que la eliminaremos del análisis.
5
6 sns.pairplot(df[['User Rating', 'Reviews', 'Price']])
7
8

```



