

Tarea S8.02. Power BI con Python

Esta labor consiste en la elaboración de un informe de Power BI, aprovechando las capacidades analíticas de Python. Se utilizarán los scripts de Python creados previamente en la Tarea 1 para generar visualizaciones personalizadas con las bibliotecas Seaborn y Matplotlib. Estas visualizaciones estarán integradas en el informe de Power BI para ofrecer una comprensión más profunda de la capacidad del lenguaje de programación en la herramienta Power BI.

Entrega:

Almacena en un repositorio de tu GitHub una carpeta que contenga:

- El archivo .pbix de Power BI con el informe solicitado.
- Un PDF del informe exportado.
- Un archivo .ipynb de Jupyter Notebook que contenga cada script correspondiente a cada visualización del informe junto al enunciado correspondiente.

Nivel 1

Los 7 ejercicios del nivel 1 de la tarea 01

Ejercicio 1

Una variable numérica.

Procedimiento seguido:

- Se importan en Power BI las tablas deseadas. Con esto se evita la necesidad de acceder a la base de datos con un conector de Python.
- Se crea un objeto de tipo Python
- Se arrastra en este objeto el campo 'amount' de la tabla 'transactions'. Al hacer esto se genera automáticamente en segundo plano un archivo csv y un dataframe de pandas, que se pueden ver al editar el código en VScode:
- `dataset = pandas.read_csv('input_df_a8a8192f-07b7-4329-a75e-f43713f6afdb.csv')`
- En el editor de scripts se introduce y ejecuta el siguiente script, similar al original, salvo la línea `data = dataset['amount']` que se usa para acceder al campo deseado y la línea `matplotlib.pyplot.figure(figsize=(10,6), dpi=72)` que se usa para ajustar el tamaño del gráfico:

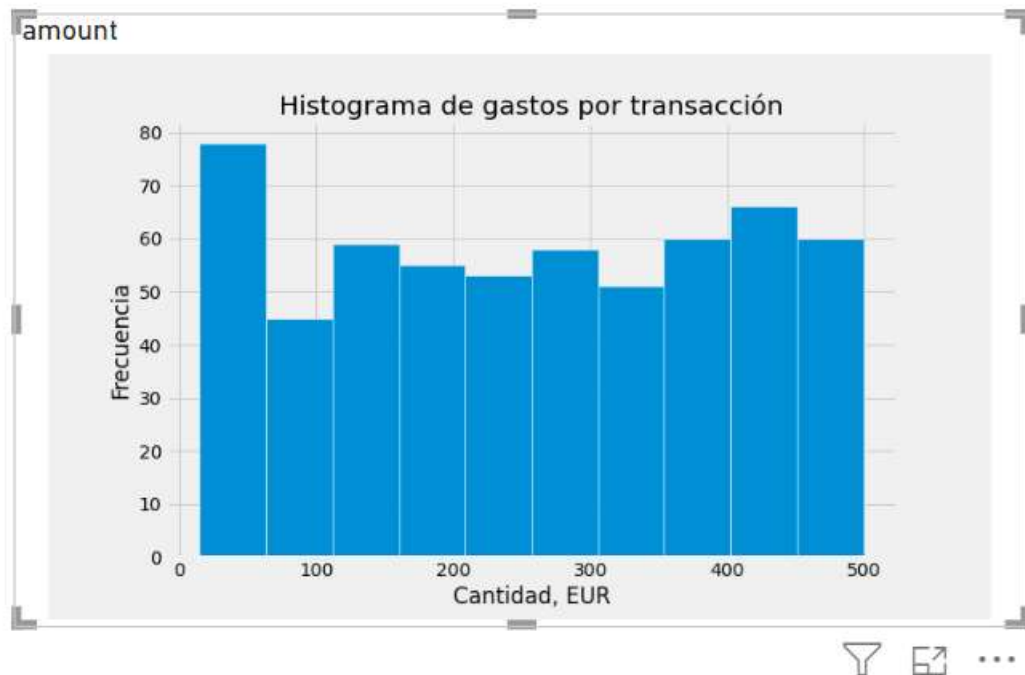
```
import matplotlib.pyplot as plt
data = dataset['amount']
```

```
matplotlib.pyplot.figure(figsize=(10,6), dpi=72)
plt.style.use('fivethirtyeight')
```

```
plt.hist(data, bins=10, edgecolor='white')
plt.xlabel('Cantidad, EUR')
plt.ylabel('Frecuencia')
plt.title('Histograma de gastos por transacción')

plt.show()
```

El resultado es:



Editor de scripts de Python

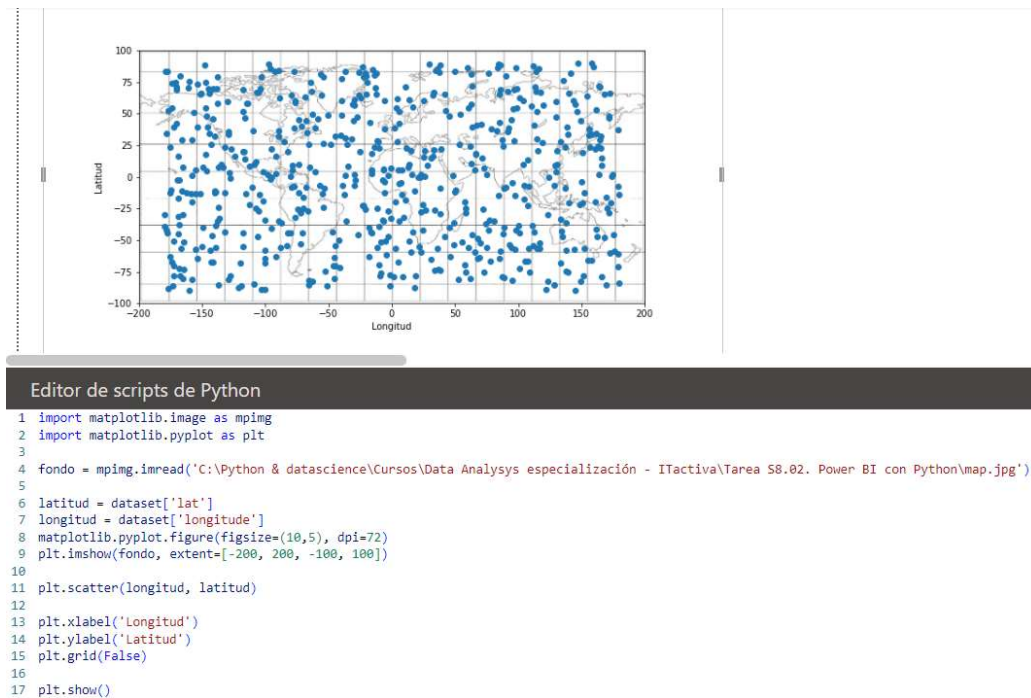
⚠ Las filas duplicadas se quitarán de los datos.

```
1 import matplotlib.pyplot as plt
2 cantidad = dataset['amount']
3
4 matplotlib.pyplot.figure(figsize=(10,6), dpi=72)
5 plt.style.use('fivethirtyeight')
6 plt.hist(cantidad, bins=10, edgecolor='white')
7 plt.xlabel('Cantidad, EUR')
8 plt.ylabel('Frecuencia')
9 plt.title('Histograma de gastos por transacción')
10
11 plt.show()
```

Ejercicio 2

Dos variables numéricas.

Se sigue el mismo procedimiento del Ejercicio 1 con la única diferencia que se arrastran las variables lat y longitude y se modifica el código original para acceder a ambas ellas. El resultado es:



Ej. 3 Una variable categórica.

Power BI elimina por defecto las filas duplicadas. La variable categórica escogida (user_id) tiene duplicados y de hecho el objetivo es contar el número de repeticiones y representarla en un piechart. Para evitar la eliminación de duplicados de este campo, se ha añadido otro campo (id) cuyos valores son todos diferentes, de modo que ya no hay filas duplicadas y se mantienen todos los id-s de usuario. El código introducido ha sido idéntico al original, salvo el cambio del nombre del dataframe de df a dataset:

```
import matplotlib.pyplot as plt
```

```
user_id_counts = dataset['user_id'].value_counts()
top_seven_counts = user_id_counts.head(7)
```

```
# Calcular el porcentaje del total de transacciones para los top 7 usuarios
total_top_seven = top_seven_counts.sum()
percentages_top_seven = top_seven_counts / total_top_seven * 100
```

```
# Crear etiquetas para los IDs de usuario
labels = ['id=' + str(user_id) for user_id in user_id_counts.index]
```

```
# Crear datos para todos los IDs de usuario
sizes = user_id_counts.tolist()
```

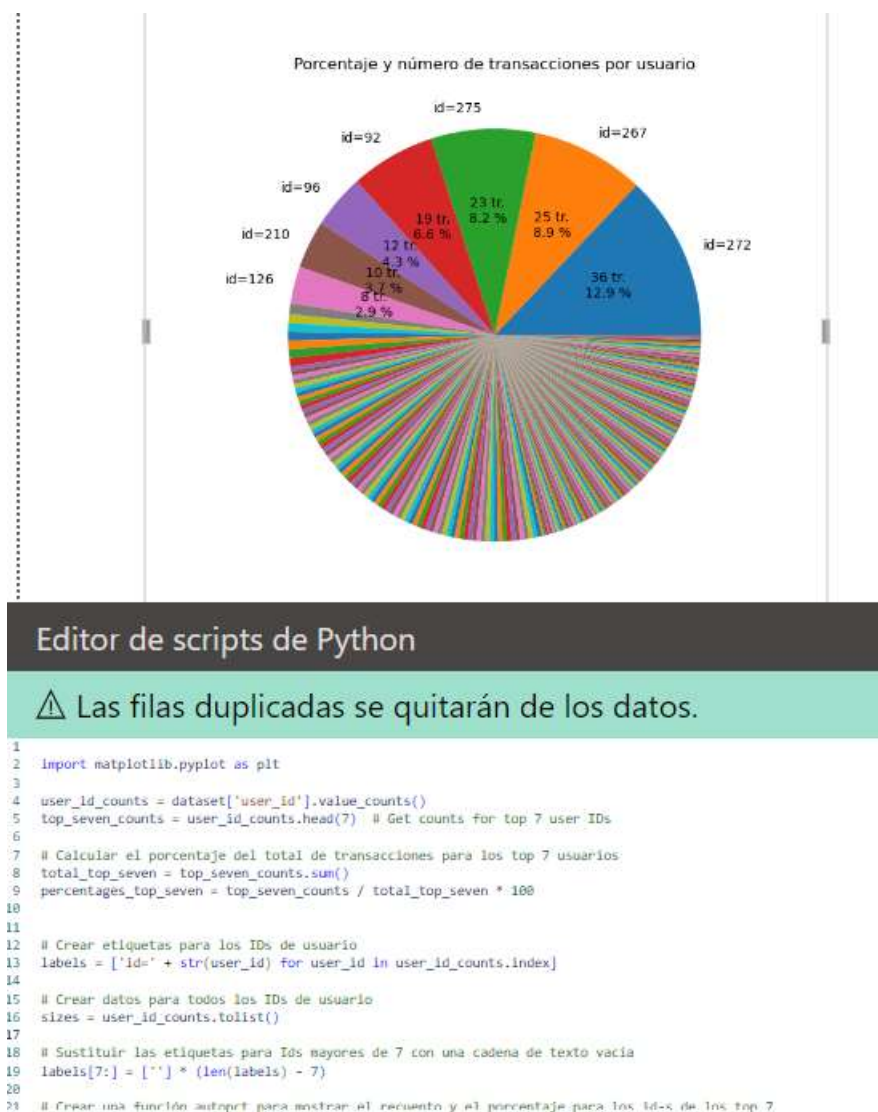
```
# Sustituir las etiquetas para IDs mayores de 7 con una cadena de texto vacía
labels[7:] = [""] * (len(labels) - 7)
```

```
# Crear una función autopct para mostrar el recuento y el porcentaje para los id-s de los top 7
def autopct_func(pct):
    if pct > 2.8:
        count = total_top_seven * pct / 100
        return '{:.0f} tr.\n{:.1f} %'.format(count, pct)
    else:
        return ""

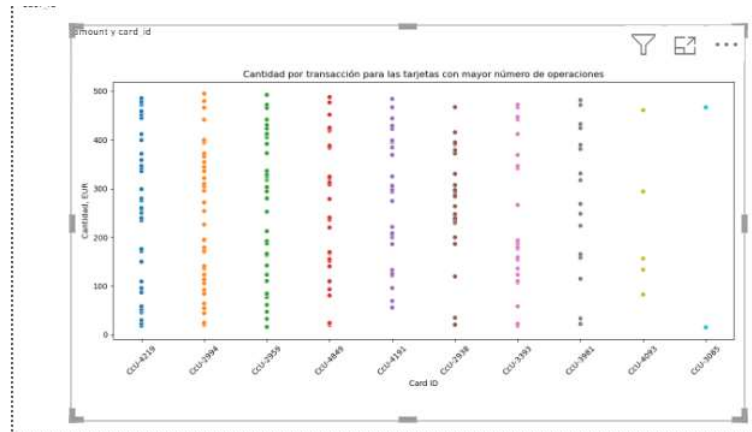
plt.figure(figsize=(7, 7))
plt.pie(sizes, labels=labels, autopct=autopct_func)
plt.title("Porcentaje y número de transacciones por usuario")

plt.show()
```

De este modo el resultado obtenido es idéntico al obtenido en S8.01:



Ejercicio 4 Una variable categórica y una numérica.



Editor de scripts de Python

⚠ Las filas duplicadas se quitarán de los datos.

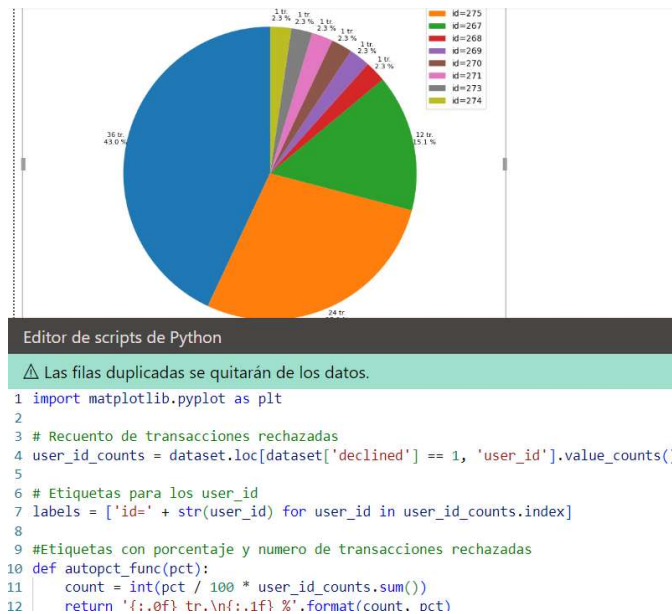
```
1 #Debido a que en Tarea S8.02. no se puede representar el plot interactivo,  
2 # se genera un plot no interactivo, solo con los card_id-s que tienen el mayor número operaciones.  
3  
4 import seaborn as sns  
5 import matplotlib.pyplot as plt  
6 import pandas as pd  
7  
8 #Top N card_id-s por número de transacciones  
9 top_n = 10  
10 top_card_ids = dataset['card_id'].value_counts().nlargest(top_n).index  
11  
12 # Filtrado del DataFrame dejando solo las filas que corresponden a los top N  
13 df_top_n = dataset[dataset['card_id'].isin(top_card_ids)]  
14 df_top_n  
15 # Número de puntos(transacciones) por cada card_id  
16 num_puntos_por_card_id = df_top_n.groupby('card_id').size()
```

Ejercicio 5 Dos variables categóricas.

En este caso también se arrastra un campo extra para evitar la eliminación de los duplicados.

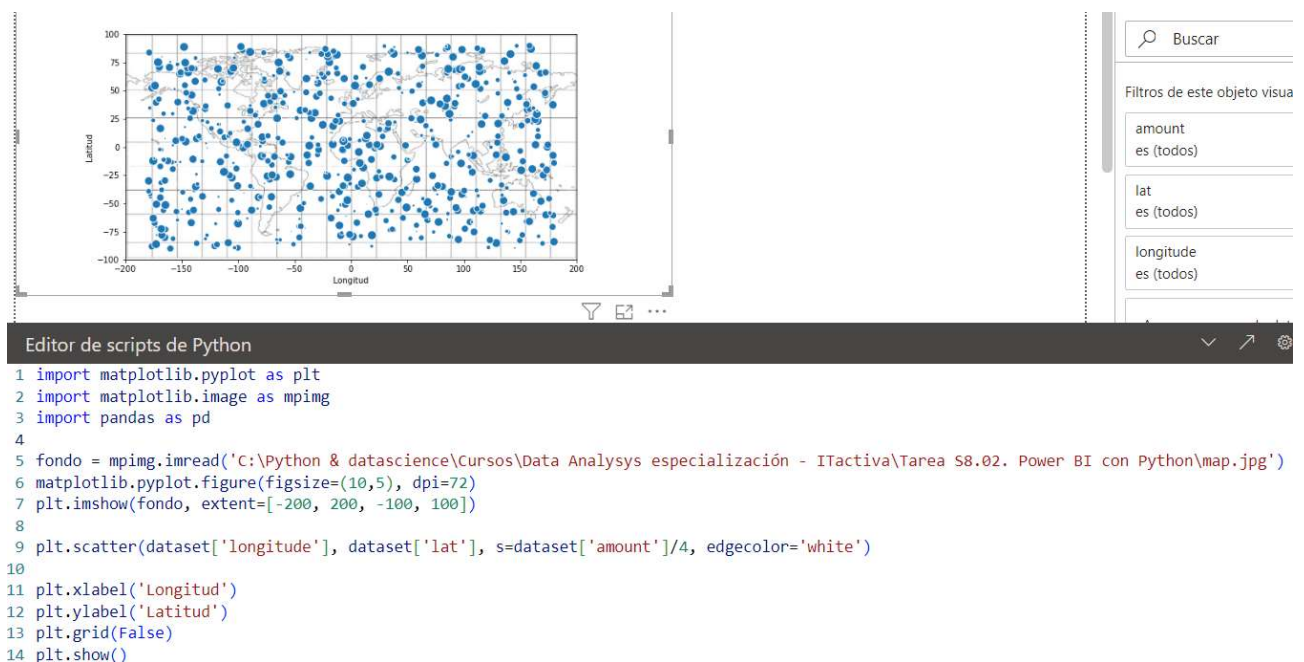
De nuevo, se utiliza el código original, sustituyendo el nombre del dataframe a 'dataset'.

El resultado es:



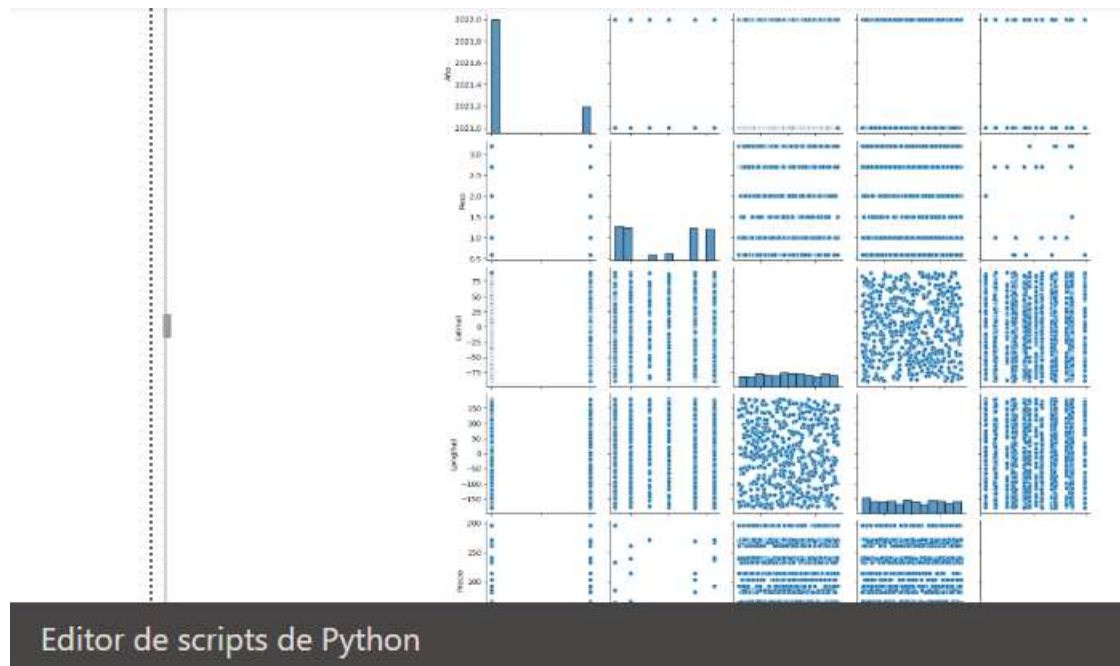
Ejercicio 6 Tres variables.

Se emplea el código original de la tarea 8.01, sin la parte de conversión de object a float, ya que al tomar los datos del dataframe llamado dataset ya están en el formato correcto. El resultado es:



Ejercicio 7 Graficar un Pairplot.

Para que funcione el código original que genera el pairplot, tan solo hay que eliminar el símbolo de \$ del campo de precio y convertir el campo a numérico. También se cambian los nombres de los campos que están en inglés para que todos estén en castellano.



```
Editor de scripts de Python
1 import seaborn as sns
2 import matplotlib.pyplot as plt
3 import pandas as pd
4
5 dataset['Precio'] = dataset['price'].str.replace('$', '')
6 dataset['Precio'] = dataset['Precio'].apply(pd.to_numeric)
7
8
9 sns.pairplot(dataset)
10 plt.show()
```

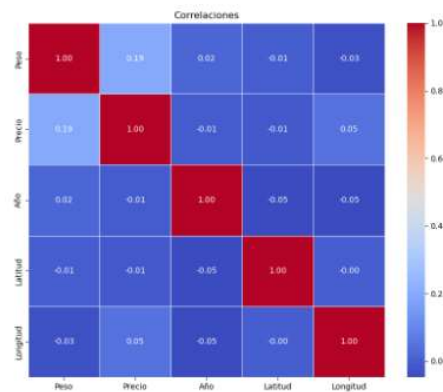
Nivel 2

Los 2 ejercicios del nivel 2 de la tarea 01

Ejercicio 1 Correlación de todas las variables numéricas.

La variable precio que está en formato texto se ha convertido en número mediante Power Query, ya que la conversión con Python fallaba por alguna razón difícil de aclarar.

El resultado es:



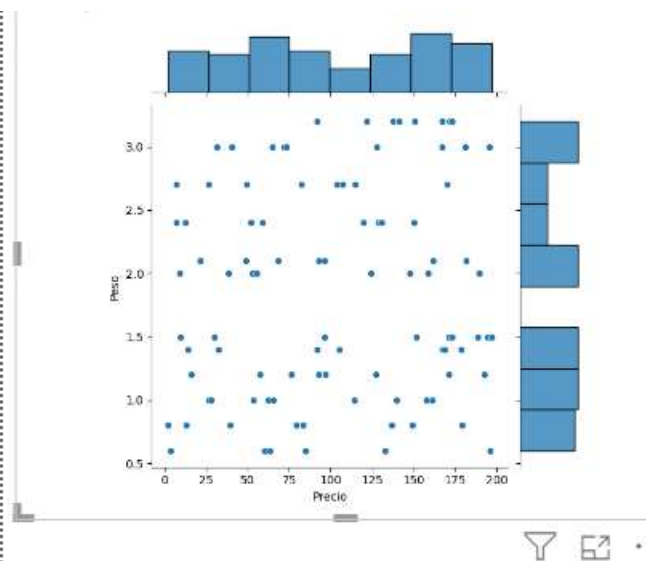
Editor de scripts de Python

```

1
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4
5 corr_matrix=dataset.corr()
6
7 plt.figure(figsize=(10, 8))
8 sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', fmt=".2f", linewidths=0.5)
9 plt.title('Correlaciones')
10 plt.show()

```

Ejercicio 2 Implementa un jointplot.



Editor de scripts de Python

```

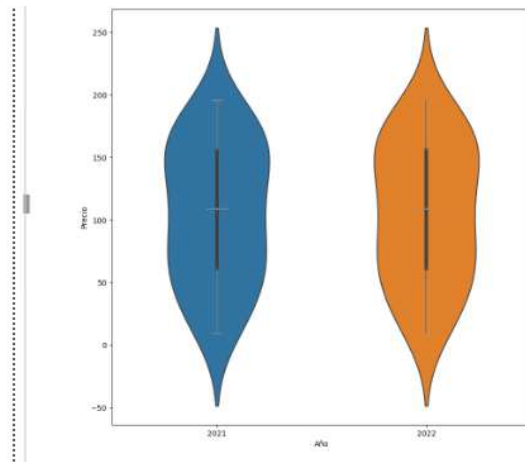
1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4 sns.jointplot(x='Precio', y='Peso', data=dataset)
5
6 plt.show()

```


Nivel 3

Los 2 ejercicios del nivel 3 de la tarea 01

Ejercicio 1 Implementa un violinplot combinado con otro tipo de gráfico.



Editor de scripts de Python

⚠ Las filas duplicadas se quitarán de los datos.

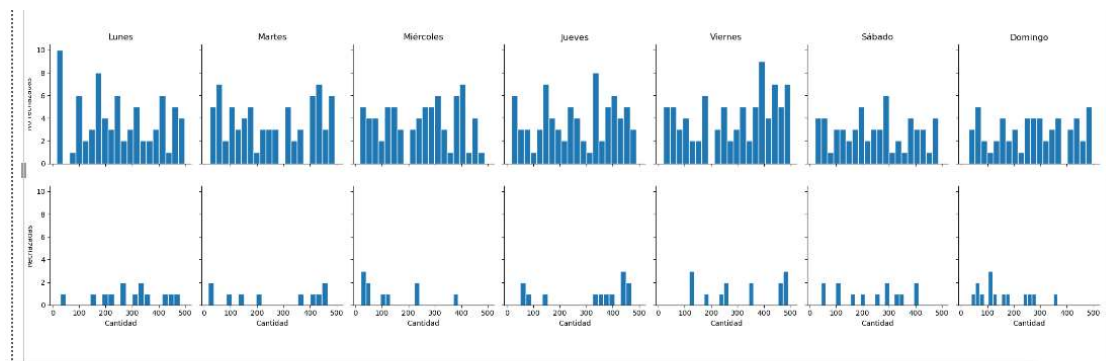
```
1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4 plt.figure(figsize=(10, 10))
5 sns.violinplot(x='Año', y='Precio', data=dataset, width=0.5)
5 sns.boxplot(x='Año', y='Precio', data=dataset, width=0.1, color='white', linewidth=1)
7
8 plt.show()
```

Ejercicio 2 Genera un FacetGrid para visualizar múltiples aspectos de datos simultáneamente.

Primero se crea un campo de día de la semana mediante Power Query, ya que el dataframe llamado dataset contiene los elementos de la jerarquía y no el propio timestamp que se había empleado en el código original. El código M usado para generar el campo nuevo es el siguiente:

```
= Date.DayOfWeekName([Timestamp])
```

A continuación se adapta el código a los nombres de los campos y del dataframe (dataset). El resultado es:



Editor de scripts de Python

⚠ Las filas duplicadas se quitarán de los datos.

```
1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4 dataset['Rechazada'] = dataset['declined'].map({0: 'No', 1: 'Si'})
5
6 g = sns.FacetGrid(dataset, col="Día de la semana", row='Rechazada')
7
8 g.map(plt.hist, "Cantidad", bins=20, edgecolor='white')
9
10 g.set_titles("", "")
11
12 for ax, label in zip(g.axes[:, 0], ['No rechazadas', 'Rechazadas']):
13     ax.set_ylabel(label)
```