

# Proyecto Data Science:

# “Global sales prediction”

Autor: Gambin Mariano  
Profesor: Germán Rodríguez  
Tutor: Mateo Bongiorno  
Comisión: 60895

<b>Video Games Dataset:</b>	<b>3</b>
Descripción columnas del dataset	3
Columnas interesantes para análisis	4
<b>Abstract</b>	<b>5</b>
<b>Exploración y limpieza:</b>	<b>6</b>
Imputación de valores *NaN* con sklearn	7
Imputación valores columnas ['Critic_Count'], ['Critic_Score'], ['User_Count'], ['Year_of_Release']	8
Transformación columna ['User Score']	9
Imputación de las columnas categóricas	10
<b>EDA (Exploratory Data Analysis):</b>	<b>11</b>
Gráficos y observaciones	11
Cantidad de videojuegos por fecha de lanzamiento	11
Distribución de la cantidad de géneros por fecha de lanzamiento	12
Cantidad de videojuegos por género	14
Ventas globales por género	17
Proporción de las ventas por región	18
Ventas por año por región	20
Outliers	21

# Video Games Dataset

Dataset que contiene una lista de videojuegos. Cada uno de los mismos están caracterizados por un género, un año de lanzamiento, publisher(editor), plataforma, así también como las ventas que representaron esos juegos en correspondientes regiones (NA, EU, JP, otros) y también encontraremos las valoraciones promedio de los críticos y de los usuarios.

## Descripción columnas del dataset

[Name]: Nombre del juego.

[Year]: Año en el que fue lanzado.

[Genre]: Género oficial del juego.

[Publisher]: Editora del juego.

[Platform]: Plataforma para la que el juego fue lanzado

[NA\_Sales]: Ventas en Norte América por millón.

[EU\_Sales]: Ventas en la Unión Europea por millón.

[JP\_Sales]: Ventas en Japón por millón.

[Other\_Sales]: Ventas del resto del mundo por millón.

[Global\_Sales]: Ventas totales por millón (Suma [NA\_Sales], [EU\_Sales], [JP\_Sales], [Other\_Sales]).

[Critic\_Score]: Promedio de los puntajes del juego recibido por los críticos.

[Critic\_Count]: Cantidad de críticos que realizaron una valoración.

[User\_Score]: Promedio de los puntajes del juego recibido por los usuarios.

[User\_Count]: Cantidad de usuarios que realizaron una valoración.

[Developer]: Entidad que desarrolló el juego.

[Rating]: Clasificación del juego (E, Everyone, T-Teens, etc).

## Columnas interesantes para análisis

[NA\_Sales], +[EU\_Sales], +[JP\_Sales], +[Other\_Sales]\*= [Global\_Sales]. Las columnas de ventas pueden sernos de mucha utilidad, por ejemplo, al cruzarla con el género o los años, podríamos obtener insights sobre cómo evolucionaron las ventas de las distintas regiones por año, como evolucionó la cantidad de géneros por año.

\*[Platform]\*. Si sumamos esta columna a las ventas podríamos determinar que plataformas fueron las más rentables por región y también cuantos juegos han salido por plataforma por año y la diferencia, si es que la hay, de los géneros más vendidos entre las plataformas.

\*[Developer]\*. Puede ayudarnos a recolectar insights sobre la relación entre los desarrolladores y las ventas globales ¿Cuales son los desarrolladores que más ventas globales lograron?

## Abstract:

En un mercado tan competitivo como el de los videojuegos, comprender las tendencias de ventas y los factores que influyen en ellas es crucial para el éxito empresarial. Un modelo predictivo preciso puede proporcionar a los desarrolladores y editores una ventaja al permitirles anticipar la demanda y ajustar estrategias de lanzamiento, marketing y desarrollo de productos de manera más efectiva.

El análisis de datos en esta industria puede revelar patrones y comportamientos del consumidor que pueden ser difíciles de detectar de otra manera. Al comprender mejor las preferencias de los jugadores y cómo estas se relacionan con las características de los juegos, las empresas pueden tomar decisiones más informadas sobre qué tipos de juegos desarrollar y cómo dirigirse a diferentes segmentos de mercado.

La aplicación de técnicas de aprendizaje automático en este contexto puede proporcionar valiosos conocimientos sobre la relación entre variables, permitiendo a las empresas optimizar sus procesos de toma de decisiones y maximizar su retorno de inversión en el desarrollo y comercialización de videojuegos.

Este proyecto puede ayudar a identificar oportunidades de inversión y crecimiento dentro de la industria de los videojuegos. Al analizar qué características y géneros de juegos tienden a tener un mejor desempeño en términos de ventas globales, los inversores pueden tomar decisiones más fundamentadas sobre dónde asignar recursos y capital.

Como objetivo buscamos predecir las ventas globales por millon utilizando un algoritmo de regresión, teniendo en cuenta los puntos ya establecidos como afirmaciones que guían nuestra investigación.

## Exploración y limpieza

Queres determinar cuál es el % de los valores NaN de cada columna para poder tener un ideal de proporción y así especificar las estrategias de imputación para cada uno de los casos

<Fragmento de código>

```
#Porcentaje de valores nulos por columna (Sumo los valores y los divido
por el largo del dataframe * 100)
percentage_null = round(df_games.isnull().sum() / len(df_games)*100,2)

percentage_null
```

Nombre columna	%Nulos
Name	0.01%
Platform	0.00%
Year_of_Release	1.61%
Genre	0.01%
Publisher	0.32%
NA_Sales	0.00%
EU_Sales	0.00%
JP_Sales	0.00%
Other_Sales	0.00%
Global_Sales	0.00%
Critic_Score	51.33%
Critic_Count	51.33%
User_Score	40.10%

User_Count	54.60%
Developer	39.61%
Rating	40.49%

Podemos observar con esta variable aplicada el porcentaje de valores nulos por columna. Existen columnas en nuestro data frame que cuentan con más del 35% de sus valores nulos. Dependiendo su tipo de dato e importancia utilizaremos distintos métodos para imputar dichos valores descritos más adelante.

## Imputación de valores \*NaN\* con sklearn

Considerando que la idea es predecir las ventas globales existen columnas como "Year\_of\_Release", "Genre" y las distintas de "Sales", para poder hacer una exploración y buscar patrones a través de dichas columnas.

- Las columnas `**[Critic_Score], [Critic_Count], [User_Count], [Rating], [Developer]**` tienen más de un 35% de valores nulos. Estas columnas pueden asistirnos para la predicción que estamos buscando, por ende, optamos por imputar los valores nulos utilizando distintos métodos proporcionados por la librería sklearn
- La columna `**[User_Score]**` tiene una situación un poco más particular. La misma columna se encuentra como tipo objeto ya que cuenta con un valor que es 'tbd' (to be determined, "ha ser determinado") para recuperar esos valores vamos a realizar 3 acciones:
  - 1) Reemplazar los valores 'tbd' por 0 utilizando un bucle for y el método `.replace`
  - 2) Cambiar el tipo de dato de la columna a float
  - 3) Por último utilizar Sklearn para imputar los valores por la mediana ya que la media es más sensible a los valores extremos y el reemplazo por 0 afectaría dicha media

- La columna "Year\_of\_Release" que podría darnos insights sobre cómo se comportaron las ventas y los géneros a lo largo del tiempo, presenta un problema al detectar que, además de los valores \*NaN\*, la misma contiene valores flotantes. Vamos a imputar la poca cantidad de valores \*NaN\* utilizando la media y convertir el tipo de dato en entero (int)
- Las columnas \*\*[Publisher], [Rating], [Developer]\*\*, los valores \*NaN\* serán reemplazados por un valor "Otro", entendiendo que el mismo no podemos asumirlo o aproximarlos con los métodos utilizados en las columnas numéricas anteriores.

## Imputación valores columnas ['Critic\_Count'], ['Critic\_Score'], ['User\_Count'], ['Year\_of\_Release']

<Fragmento de código>

```
#Instanciamos la estrategia de nuestro imputador en una variable
llamada "imputer" (esta misma es para las columnas con variables
numéricas)
imputer_mean = SimpleImputer(missing_values=np.nan, strategy='mean')

#Hacemos reshape a cada una de las columnas a las que vamos a pasar por
esta estrategia
Critic_Count = df_games['Critic_Count'].values.reshape(-1,1)
Critic_Score = df_games['Critic_Score'].values.reshape(-1,1)
User_Count = df_games['User_Count'].values.reshape(-1,1)
Year_of_Release = df_games['Year_of_Release'].values.reshape(-1,1)

#Asignamos en un variable los valores de la columna imputada
Critic_Count_imputed = imputer_mean.fit_transform(Critic_Count)
Critic_Score_imputed = imputer_mean.fit_transform(Critic_Score)
User_Count_imputed = imputer_mean.fit_transform(User_Count)
Year_of_Release_imputed = imputer_mean.fit_transform(Year_of_Release)

#Asignó cada columna imputada con su columna correspondiente
df_games['Critic_Count'] = Critic_Count_imputed
df_games['Critic_Score'] = Critic_Score_imputed
df_games['User_Count'] = User_Count_imputed
df_games['Year_of_Release'] = Year_of_Release_imputed
```



## Transformación columna ['User Score']

<Fragmento de código>

```
#Busco los valores de la columna user score

user_score_valores = df_games['User_Score'].values

user_score_valores_reemplazados = [] #Nueva lista para los valores
reemplazados

for valor in user_score_valores:
    if valor == 'tbd': #Si el valor es tbd
        nuevo_valor = valor.replace('tbd','0') #Reemplazalo por '0'
        user_score_valores_reemplazados.append(nuevo_valor) #Apendealo a la
nueva lista
    else:
        user_score_valores_reemplazados.append(valor) #Si no, apendea el
valor original

df_games['User_Score'] = user_score_valores_reemplazados
```

Al hacer esta transformación y efectivamente cambiar el valor 'tbd' por '0', vemos que hay una gran cantidad de valores en '0', reemplazar por la media los valores nan podría afectar negativamente ya que la misma es sensible a los valores extremos. Por ello optamos por reemplazar los valores 'NaN' por la mediana

<Fragmento de código>

```
#Instanciamos la estrategia de nuestro imputador en una variable
llamada "imputer" (esta misma es para las columnas con variables
numéricas)
imputer_median = SimpleImputer(missing_values=np.nan, strategy='median')

#Hacemos reshape a cada una de las columnas a las que vamos a pasar por
esta estrategia
User_Score = df_games['User_Score'].values.reshape(-1,1)

#Asignamos en un variable los valores de la columna imputada
User_Score_imputed = imputer_median.fit_transform(User_Score)

df_games['User_Score'] = User_Score_imputed
```

## Imputación de las columnas categóricas

```
#Imputación manual de las variables categóricas

df_games["Rating"].fillna("Otro", inplace=True)
df_games["Developer"].fillna("Otro", inplace=True)
df_games["Publisher"].fillna("Otro", inplace=True)
df_games["Genre"].fillna("Misc", inplace=True)
df_games["Name"].fillna("Otro", inplace=True)
```

Volvemos a mostrar nuestro dataframe para poder identificar si, efectivamente, hemos imputado la totalidad de los valores NaN.

Name	0
Platform	0
Year_of_Release	0
Genre	0
Publisher	0
NA_Sales	0
EU_Sales	0
JP_Sales	0
Other_Sales	0
Global_Sales	0
Critic_Score	0
Critic_Count	0
User_Score	0
User_Count	0
Developer	0
Rating	0

# EDA (Exploratory Data Analysis)

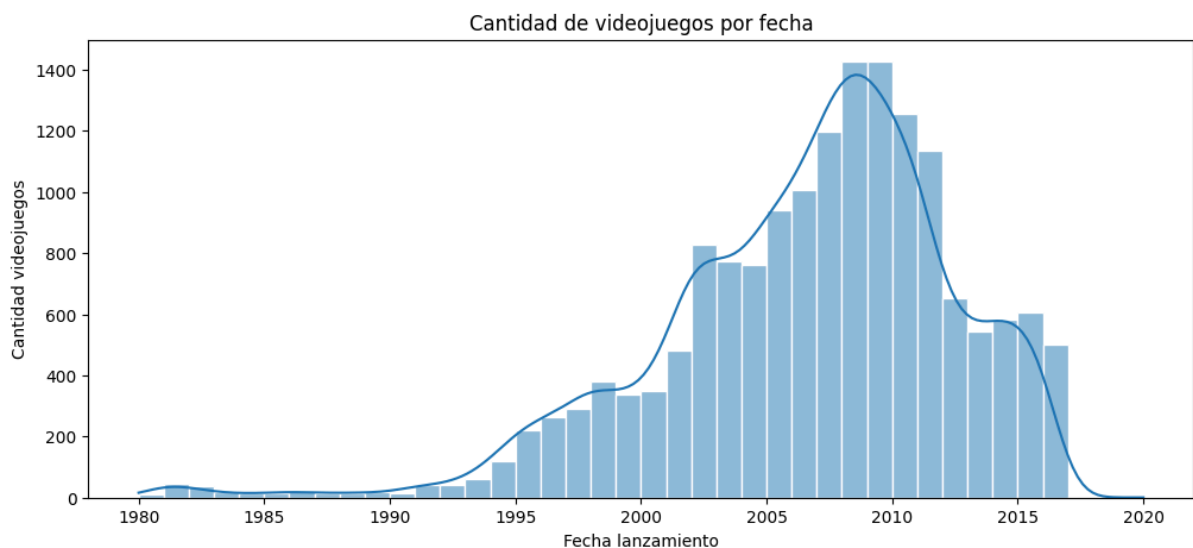
## Gráficos y observaciones

### Cantidad de videojuegos por fecha de lanzamiento

- 1) Como primera observación quisiera ver cómo se distribuyen los videojuegos según su fecha de lanzamiento. Creo una variable que va a ser mi histograma, donde le voy a pasar mi df, en el eje x quiero los años de lanzamiento y muestra el conteo por año

<Fragmento de código>

```
fig, ax = plt.subplots(figsize=(12,5))
frecuencia_años = sns.histplot(data=df_games, x="Year_of_Release",
binwidth = 1, edgecolor = "white", linewidth = 1, kde="True")
frecuencia_años.set(xlabel="Fecha lanzamiento", ylabel="Cantidad
videojuegos" ,title="Cantidad de videojuegos por fecha")
plt.show()
```



Podemos observar una distribución asimétrica positiva y podríamos decir en base a este gráfico que la mayoría de mis videojuegos aparecen entre el 2005 y el 2012 aproximadamente

## Distribución de la cantidad de géneros por fecha de lanzamiento

2) Podríamos evaluar como se comportan los distintos géneros por año.

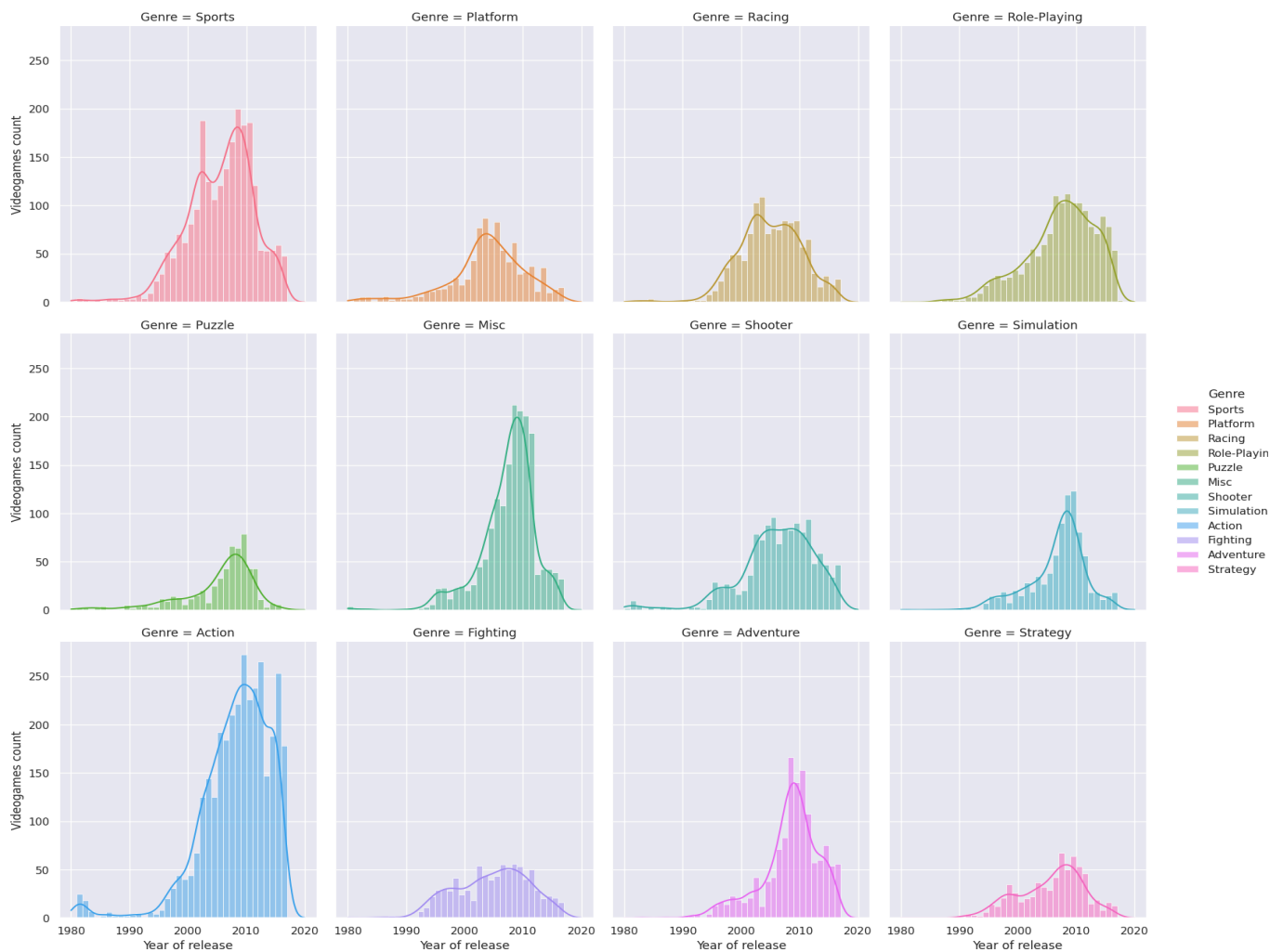
Creamos un plot donde graficamos. Estipulamos el tamaño y creamos un gráfico para determinar la frecuencia

Pasaremos el "año de lanzamiento" al eje x, el conteo será el eje y, y dividiremos los gráficos por la columna "genre".

La misma al tener 12 géneros nos mostrará 12 gráficos distintos. Pasamos col wrap para que solo muestre 4 gráficos por columna y algunos aspectos más de formato.

<Fragmento de código>

```
plt.figure(figsize=(16, 16))
#Tamaño de cada uno de los gráficos
sns.set(rc={'figure.figsize': (25, 15)})
displot_for_genre= sns.displot(data=df_games, x="Year_of_Release",
hue="Genre", col="Genre", col_wrap=4, edgecolor = "white",binwidth = 1,
aspect=.8, kde="True")
# Itera a través de cada gráfico y establece etiquetas para los ejes
for ax in displot_for_genre.axes:
    ax.set_ylabel("Videogames count")
    ax.set_xlabel("Year of release")
plt.show()
```



Cada uno de los géneros parecen estar comportándose bastante parecido donde alcanzan un pico entre los años 2005 y 2012. Y además parecen presentar similitudes en su distribución, marcando una clara asimetría positiva.

## Cantidad de videojuegos por género

- 3) Quisiera ver cómo se distribuyen mis videojuegos por géneros ¿Cuántos videojuegos tengo por género? Creo una variable que va a almacenar en el índice los géneros y en una columna el conteo total de ese género.

<Fragmento de código>

```
#Conteo de la cantidad de veces que aparece cada género
genre_count = df_games["Genre"].value_counts()

genre_count
```

Genre	count
Action	3370
Sports	2348
Misc	1750
Role-Playing	1500
Shooter	1323
Adventure	1303
Racing	1249
Platform	888
Simulation	874
Fighting	849
Strategy	683
Puzzle	580

Con esta información creamos un gráfico de barras para visualizar de mejor manera la cantidad de videojuegos según el género.

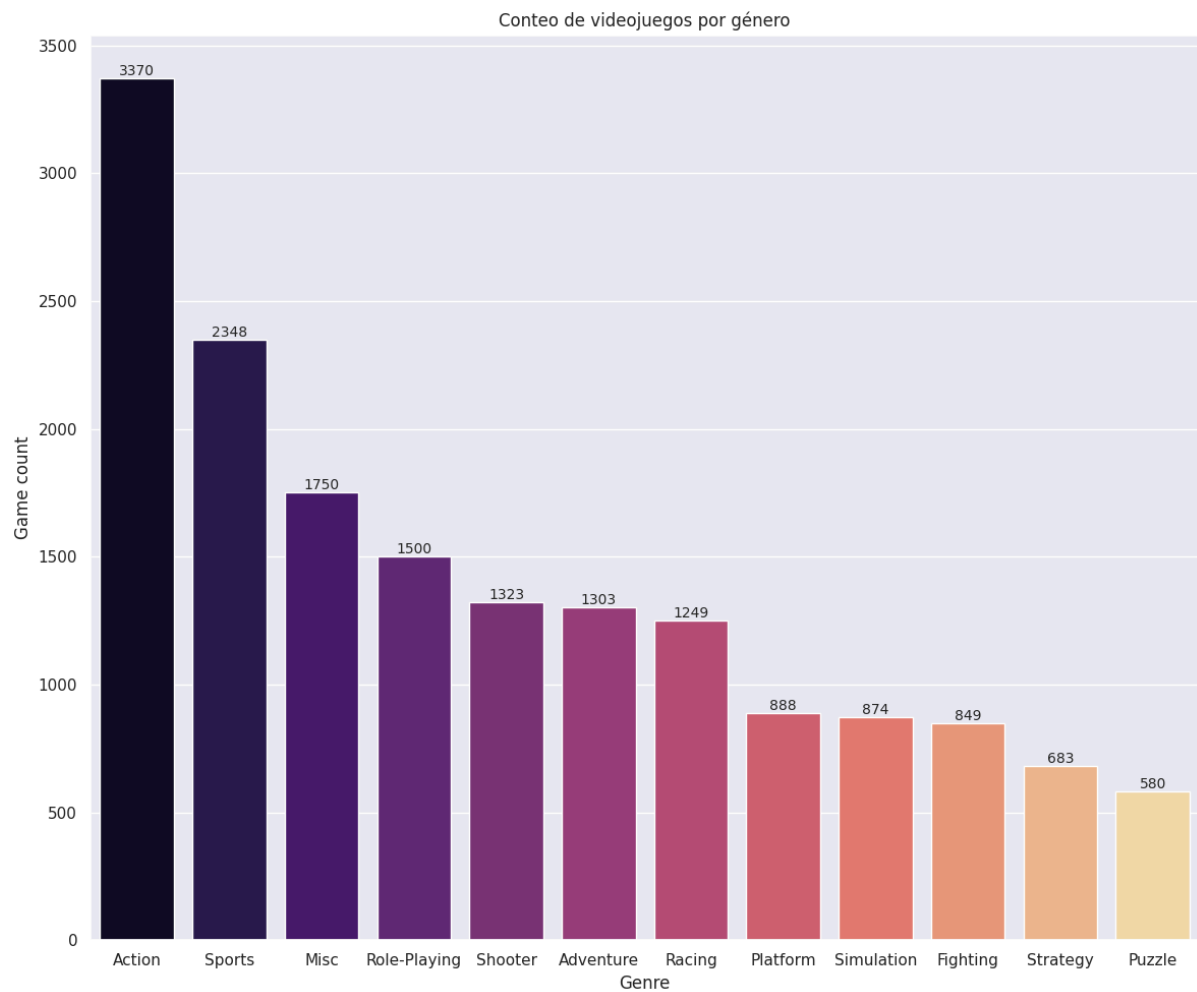
Creamos un plot donde se graficara. Con seaborn creamos el gráfico de barras, pasamos los géneros como eje x y el conteo de

los mismos como eje y. Agregamos una distinción de color por cada género para hacer mas amena la lectura.

#### <Fragmento de código>

```
#Como la columna género y videojuegos coinciden en su cantidad de datos
puedo asumir que según la cantidad de veces que aparece un género es la
cantidad de videojuegos que tengo del mismo
plt.figure(figsize=(12,10))
barplot_genre = sns.barplot(x=genre_count.index, y=genre_count,
hue=genre_count.index, palette="magma")
barplot_genre.set(xlabel="Genre", ylabel="Game count", title="Conteo de
videojuegos por género")

for container in barplot_genre.containers: #Para cada columna, en el
grafico
    barplot_genre.bar_label(container, fontsize=10) #Ingreso el valor y
el tamaño de la fuente
plt.show()
```



Se puede observar una predominancia de videojuegos del género de "acción" dentro de este dataset, seguido de "deportes" y "variados"



## Ventas globales por género

- 4) Quiero entender cuántas fueron las ventas globales por género y ver cuales han generado el mayor número de ventas.

<Fragmento de código>

```
#Utilizo el metodo groupby para agrupar las ventas totales por género
genre_global_sales =
df_games.groupby("Genre")["Global_Sales"].sum().sort_values(ascending=False)

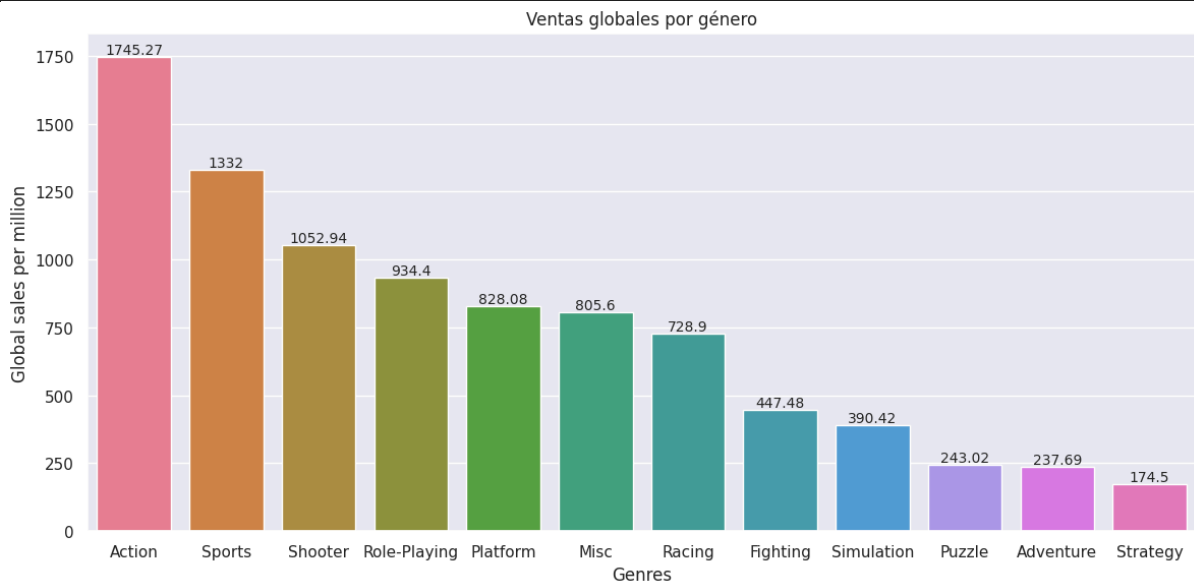
genre_global_sales
```

Genre	Global_Sales
Action	1745.27
Sports	1332.00
Shooter	1052.94
Role-Playing	934.40
Platform	828.08
Misc	805.60
Racing	728.90
Fighting	447.48
Simulation	390.42
Puzzle	243.02
Adventure	237.69

<Fragmento de código>

```
plt.figure(figsize=(12,10))
barplot_global_sales = sns.barplot(x=genre_global_sales.index,
y=genre_global_sales, hue=genre_global_sales.index)
barplot_global_sales.set(xlabel= "Genres", ylabel="Global sales per
million", title="Ventas globales por género")
```

```
for container in barplot_global_sales.containers:
    barplot_global_sales.bar_label(container, fontsize=10)
plt.tight_layout()
plt.show()
```



Rápidamente podemos identificar como los géneros de acción, deportes, disparos y de rol son los que más ventas globales han tenido

## Proporción de las ventas por región

5) ¿Cómo se comportó el porcentaje de las ventas por región?

Agrupamos las ventas de cada región por año utilizando el método `.groupby`

<Fragmento de código>

```
sum_sales_region = df_games.groupby("Year_of_Release") [ ["NA_Sales",
"EU_Sales", "JP_Sales", "Other_Sales"] ].sum()
sum_sales_region
```

Creamos una variable que va a sumar la cantidad de ventas, que se llamara "total\_sales"

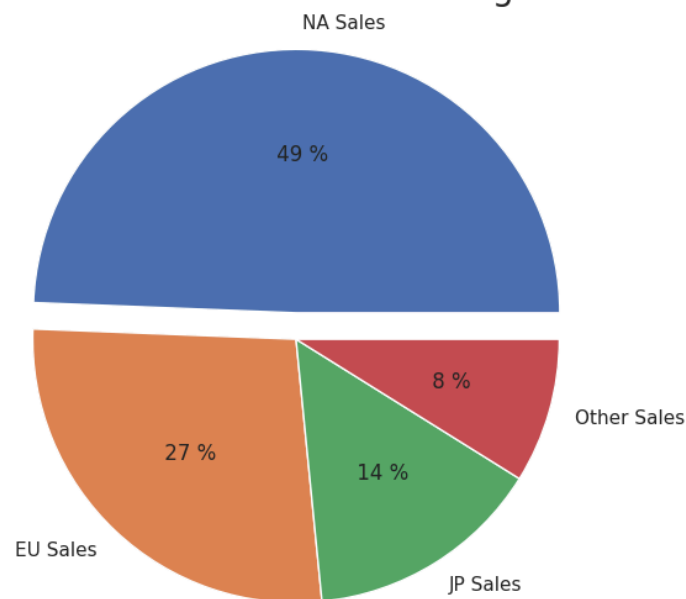
y luego creamos el gráfico de torta. Pasaremos las ventas totales y que los cortes sean las columnas de mi agrupamiento anterior, en este caso las regiones

<Fragmento de código>

```
fig, ax = plt.subplots(figsize=(10,6))
total_sales = sum_sales_region.sum()
```

```
# Cambiar nombres de los labels
labels = ["NA Sales", "EU Sales", "JP Sales", "Other Sales"]
plt.pie(total_sales, labels=labels, explode=(0.1,0,0,0), autopct="% 0i
%%")
ax.set_title('Ventas acumuladas totales de las distintas Regiones de
1980 a 2020', fontsize=20)
ax.set_label(labels)
plt.tight_layout()
plt.show()
```

Ventas acumuladas totales de las distintas Regiones de 1980 a 2020



Podemos observar que el 49% de las ventas corresponden a la región de NorteAmérica seguido de la Unión Europea.

## Ventas por año por región

6)¿Cómo han evolucionado las ventas de las distintas regiones por año?

Antes analizamos el porcentaje total de ventas ¿Que región obtuvo la mayor cantidad de ventas?

Ahora nos puede interesar cómo se han movido las ventas por región a lo largo de los años

Establecemos nuestra figura, y un grid para formato. Creamos 4 líneas donde pasamos los años de lanzamiento en el eje x y cada una de suma de las ventas de las regiones en el eje y, aplicamos una etiqueta para cada línea. Sumamos a esto un título y una etiqueta para cada eje.

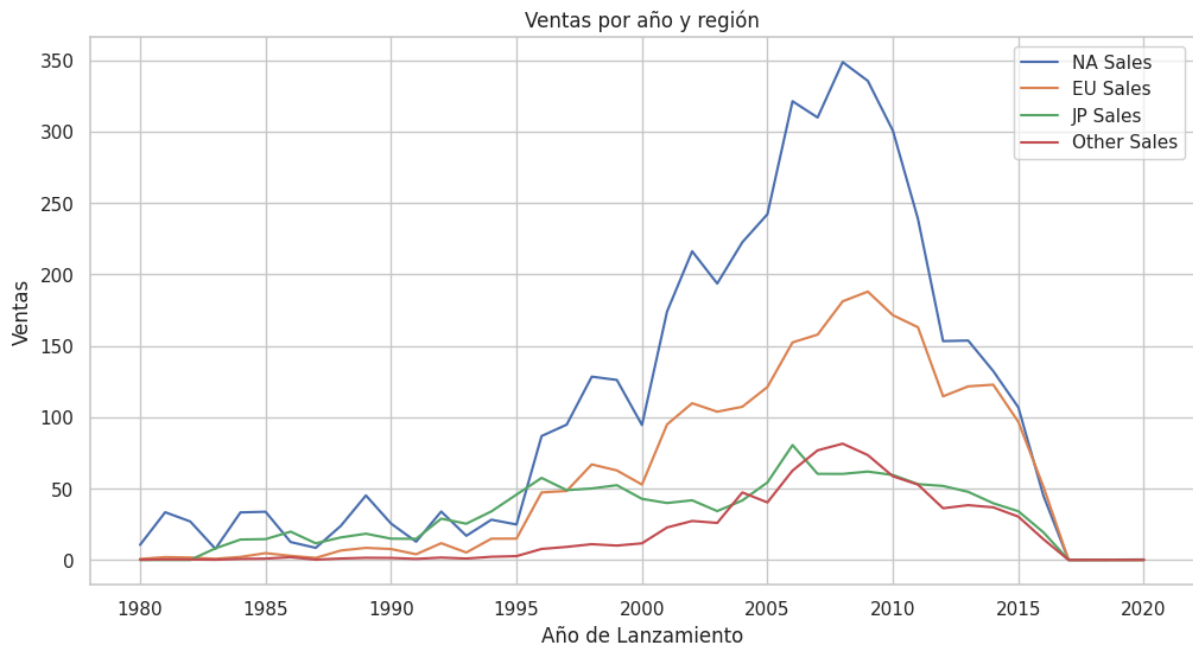
<Fragmento de código>

```
plt.figure(figsize=(12, 6))
sns.set(style="whitegrid")

sns.lineplot(data=sum_sales_region, x="Year_of_Release", y="NA_Sales",
label="NA Sales")
sns.lineplot(data=sum_sales_region, x="Year_of_Release", y="EU_Sales",
label="EU Sales")
sns.lineplot(data=sum_sales_region, x="Year_of_Release", y="JP_Sales",
label="JP Sales")
sns.lineplot(data=sum_sales_region, x="Year_of_Release",
y="Other_Sales", label="Other Sales")

plt.title("Ventas por año y región")
plt.xlabel("Año de Lanzamiento")
plt.ylabel("Ventas")

plt.show()
```



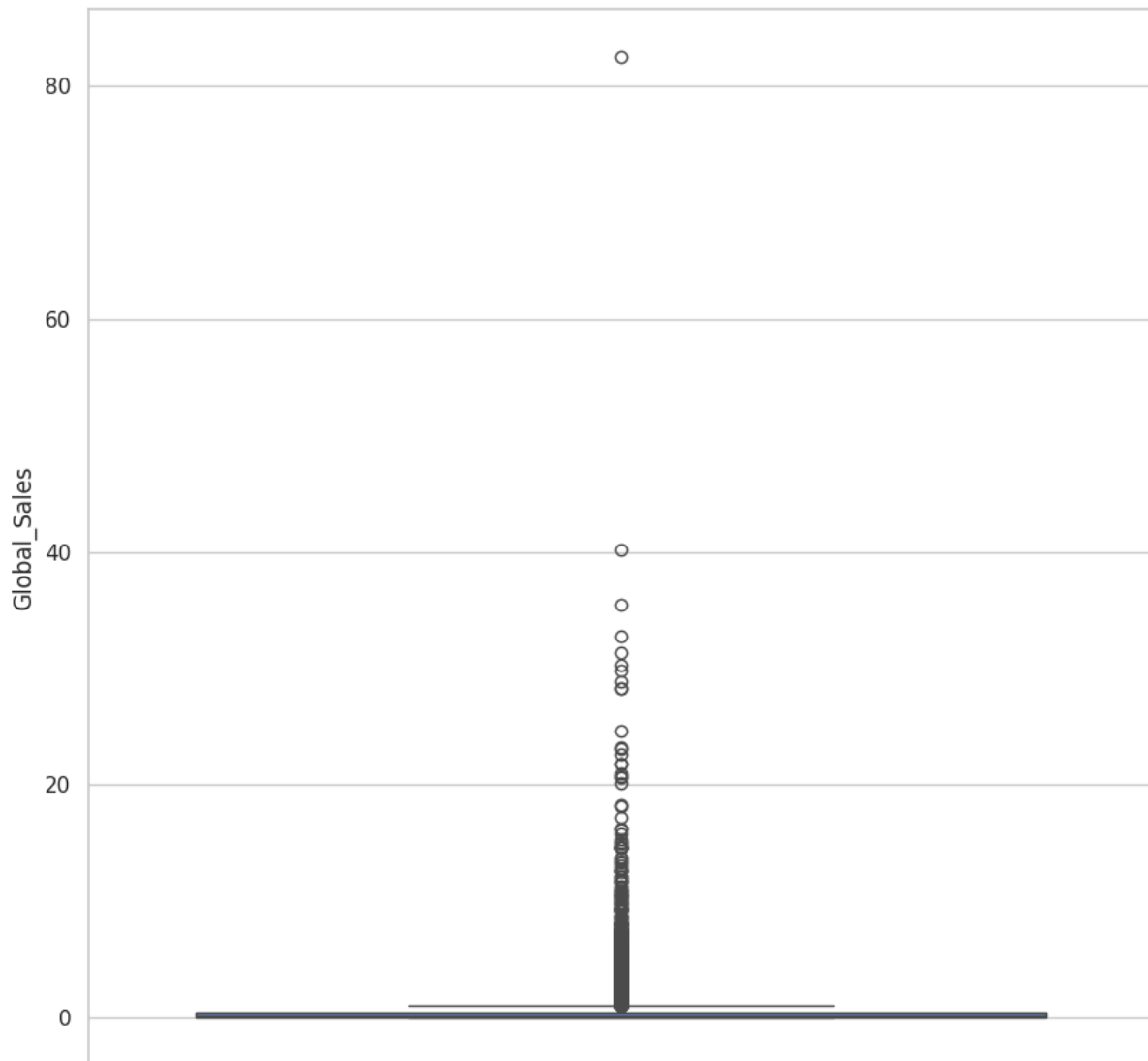
Podemos observar cómo si bien las 4 regiones parecen cumplir un mismo patrón, en donde la mayor cantidad de ventas se da entre el 2000 y el 2012, el volumen de ventas ha sido mayor en la región de Norteamérica y la Unión Europea que en las otras dos en ese mismo periodo

## Outliers

Creamos un gráfico de cajas para identificar la cantidad de valores outliers presentes en nuestra columna [Global\_Sales]

<Fragmento de código>

```
plt.figure(figsize=(10,10))
sns.boxplot(data=df_games['Global_Sales'])
plt.show()
```



Observamos una gran cantidad de valores outliers en nuestra variable a predecir. Crearemos dos funciones, una que calcule los límites superiores e inferiores del boxplot para determinar los valores outliers y otra que traiga los índices de dichos outlier para que los elimine.

<Fragmento de código>

```
# Quiero crear una función que calcule los límites del boxplot y me
# devuelva el límite superior e inferior en una variable con nombre ls
def outliers(data, feature):
    q1 = data[feature].quantile(0.25)
    q3 = data[feature].quantile(0.75)
    iqr = q3 - q1

    lower_limit = q1 - iqr * 1.5
    upper_limit = q3 + iqr * 1.5
```

```

    ls = data.index[(data[feature]<lower_limit) |
(data[feature]>upper_limit)]
    return ls

#Y otra función para que remueva los valores que están por encima o por
debajo de esos límites. Como parámetros va a tomar el data frame y el
ls calculado anteriormente
def remove(df_games,ls):
    ls = sorted(set(ls))
    df_games = df_games.drop(ls)
    return df_games

#obtención de outliers
index_list_global = []
index_list_global.extend(outliers(df_games,'Global_Sales'))

#Borramos dichos outliers y los asigno a una nueva variable que va a
almacenar nuestro dataset limpio
df_games_clean = remove(df_games,index_list_global)

```

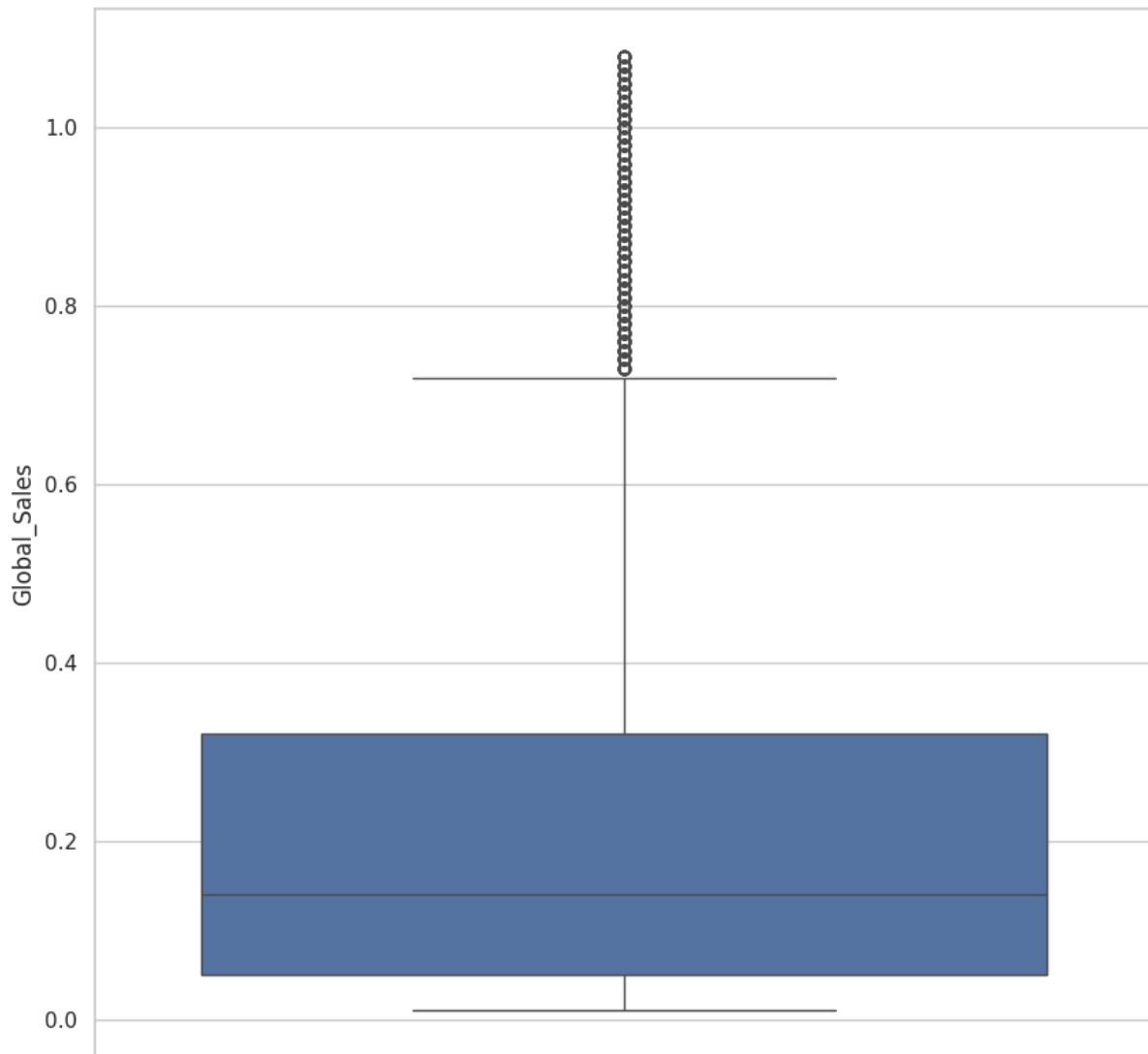
Al volver a graficar la los valores de columna en un diagrama de cajas notamos la diferencia en el gráfico con la eliminación de dichos outliers

<Fragmento de código>

```

fig, ax = plt.subplots(figsize=(10,10))
sns.boxplot(df_games_clean['Global_Sales'])
plt.show()

```



### Próximas fases del proyecto

- Normalización e imputación de valores
- Elección de algoritmos de regresión (Lineal, KNN, RGBRegressor,...)
- Separación de data (entrenamiento, testing)
- Prueba de algoritmos
- Medición de resultados
- Validación de hipótesis / Cambio de variable para subsiguientes pruebas