

Análisis del Dataset Netflix Movies and TV Shows

```
In [32]: 1 print("Dataset obtenido de https://www.kaggle.com/shivamb/netflix-shows")
```

Dataset obtenido de <https://www.kaggle.com/shivamb/netflix-shows> (<https://www.kaggle.com/shivamb/netflix-shows>)

1. Importación de Librerías

```
In [58]: 1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
```

2. Importación de Datos

```
In [175]: 1 #Cargamos los datos.
          2 ruta = "C:/Users/Usuario/Downloads/netflix_titles.csv"
          3 datos = pd.read_csv(ruta)
          4 datos.head()
```

Out[175]:

	show_id	type	title	director	cast	country	date_added	release_year	rating	duration	listed_in	description
0	s1	TV Show	3%	NaN	João Miguel, Bianca Comparato, Michel Gomes, R...	Brazil	August 14, 2020	2020	TV-MA	4 Seasons	International TV Shows, TV Dramas, TV Sci-Fi &...	In a future where the elite inhabit an island ...
1	s2	Movie	7:19	Jorge Michel Grau	Demián Bichir, Héctor Bonilla, Oscar Serrano, ...	Mexico	December 23, 2016	2016	TV-MA	93 min	Dramas, International Movies	After a devastating earthquake hits Mexico Cit...
2	s3	Movie	23:59	Gilbert Chan	Tedd Chan, Stella Chung, Henley Hii, Lawrence ...	Singapore	December 20, 2018	2011	R	78 min	Horror Movies, International Movies	When an army recruit is found dead, his fellow...
3	s4	Movie	9	Shane Acker	Elijah Wood, John C. Reilly, Jennifer Connelly...	United States	November 16, 2017	2009	PG-13	80 min	Action & Adventure, Independent Movies, Sci-Fi...	In a postapocalyptic world, rag-doll robots hi...
4	s5	Movie	21	Robert Luketic	Jim Sturgess, Kevin Spacey, Kate Bosworth, Aar...	United States	January 1, 2020	2008	PG-13	123 min	Dramas	A brilliant group of students become card-coun...

3. Pre-Procesamiento de Datos

```
In [176]: 1 #Obtenemos información del dataset.  
          2 datos.info()
```

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 7787 entries, 0 to 7786  
Data columns (total 12 columns):  
#   Column          Non-Null Count  Dtype  
---  -  
0   show_id         7787 non-null   object  
1   type            7787 non-null   object  
2   title           7787 non-null   object  
3   director        5398 non-null   object  
4   cast            7069 non-null   object  
5   country         7280 non-null   object  
6   date_added      7777 non-null   object  
7   release_year    7787 non-null   int64  
8   rating          7780 non-null   object  
9   duration        7787 non-null   object  
10  listed_in       7787 non-null   object  
11  description     7787 non-null   object  
dtypes: int64(1), object(11)  
memory usage: 730.2+ KB
```

```
In [177]: 1 #Verificamos la cantidad de valores faltantes por columna.
2 print("Cantidad de Valores Faltantes por Columna")
3 print(datos.isnull().sum())
4
5 print("")
6
7 print("Cantidad de Valores Faltantes por Columna (%)")
8 print((datos.isnull().sum() / len(datos)) * 100)
9
10 print("")
11
12 print('''Debido a la cantidad de valores faltantes en ciertas columnas, se tendrá precaución uso de la columna director,
13 ya que, al no tener el 30% de los valores de esta columna, podría llegarse a resultados erróneos, y no es conveniente
14 eliminar una cantidad tan grande de observaciones.
15 Por otro lado, las demás columnas con valores faltantes presentan porcentajes de valores faltantes menores al 10%.
16 Aunque no es lo ideal, puede haber más flexibilidad en el uso de estas columnas.''')
```

Cantidad de Valores Faltantes por Columna

show_id	0
type	0
title	0
director	2389
cast	718
country	507
date_added	10
release_year	0
rating	7
duration	0
listed_in	0
description	0

dtype: int64

Cantidad de Valores Faltantes por Columna (%)

show_id	0.000000
type	0.000000
title	0.000000
director	30.679337
cast	9.220496
country	6.510851
date_added	0.128419
release_year	0.000000
rating	0.089893

```

duration      0.000000
listed_in     0.000000
description    0.000000
dtype: float64

```

Debido a la cantidad de valores faltantes en ciertas columnas, se tendrá precaución uso de la columna director, ya que, al no tener el 30% de los valores de esta columna, podría llegarse a resultados erróneos, y no es conveniente eliminar una cantidad tan grande de observaciones. Por otro lado, las demás columnas con valores faltantes presentan porcentajes de valores faltantes menores al 10%. Aunque no es lo ideal, puede haber más flexibilidad en el uso de estas columnas.

```

In [178]: 1 #Omitimos las columnas show_id y description.
          2 datos = datos.drop(columns = ["show_id", "description"], axis = 1)
          3 datos.head()

```

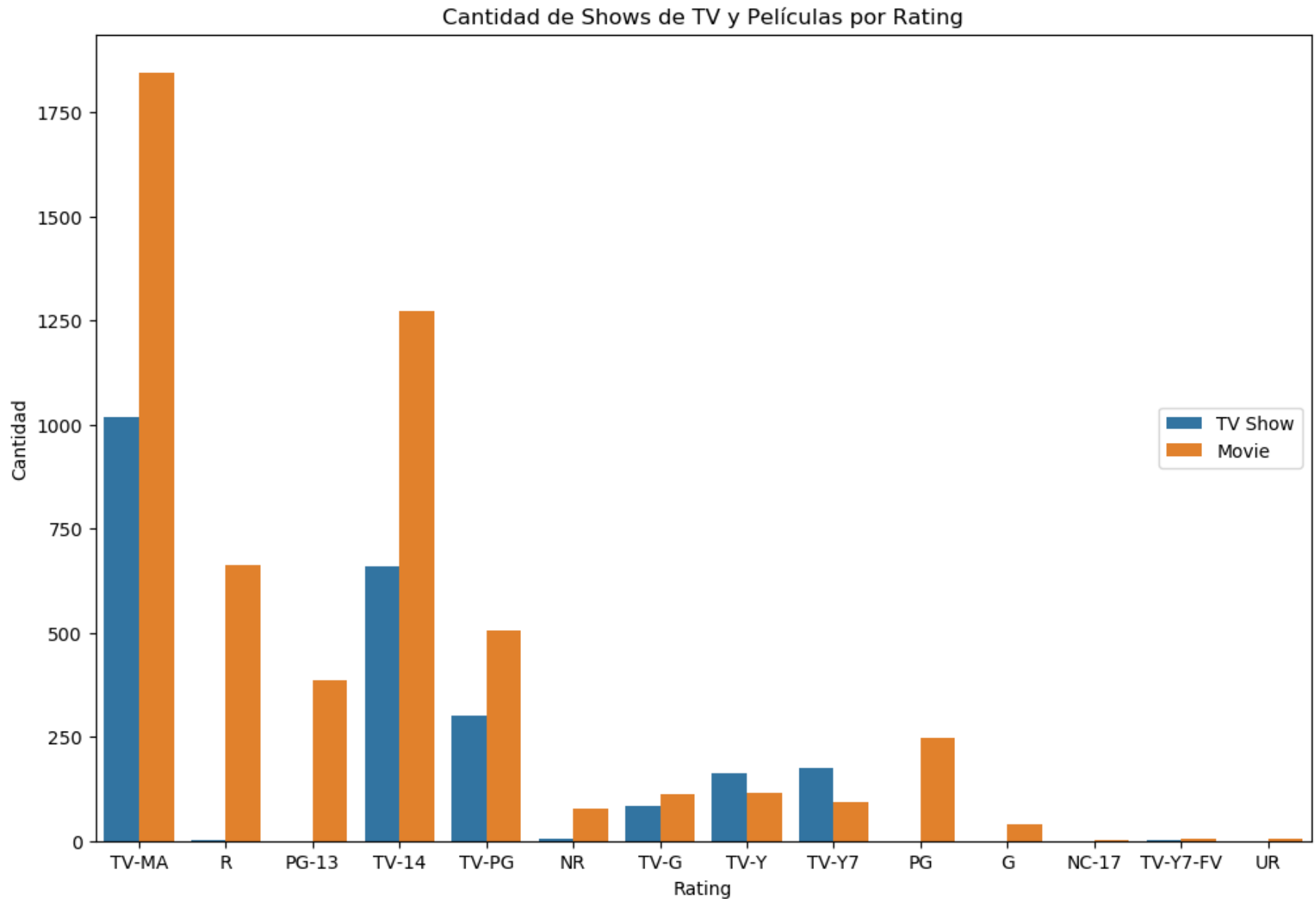
Out[178]:

	type	title	director	cast	country	date_added	release_year	rating	duration	listed_in
0	TV Show	3%	NaN	João Miguel, Bianca Comparato, Michel Gomes, R...	Brazil	August 14, 2020	2020	TV-MA	4 Seasons	International TV Shows, TV Dramas, TV Sci-Fi &...
1	Movie	7:19	Jorge Michel Grau	Demián Bichir, Héctor Bonilla, Oscar Serrano, ...	Mexico	December 23, 2016	2016	TV-MA	93 min	Dramas, International Movies
2	Movie	23:59	Gilbert Chan	Tedd Chan, Stella Chung, Henley Hii, Lawrence ...	Singapore	December 20, 2018	2011	R	78 min	Horror Movies, International Movies
3	Movie	9	Shane Acker	Elijah Wood, John C. Reilly, Jennifer Connelly...	United States	November 16, 2017	2009	PG-13	80 min	Action & Adventure, Independent Movies, Sci-Fi...
4	Movie	21	Robert Luketic	Jim Sturgess, Kevin Spacey, Kate Bosworth, Aar...	United States	January 1, 2020	2008	PG-13	123 min	Dramas

4. Análisis de Datos

```
In [227]: 1 print("Análisis 1: Cantidad de Shows de TV y Películas por Rating")
2
3 fig, ax = plt.subplots()
4 fig.set_size_inches(12, 8)
5
6 sns.countplot(x = "rating", data = datos, hue = "type")
7 ax.set_xlabel("Rating")
8 ax.set_ylabel("Cantidad")
9 ax.set_title("Cantidad de Shows de TV y Películas por Rating")
10 ax.legend(loc = "center right")
11
12 plt.show()
13
14 print('''Observación del Análisis 1: Por lo general, hay más películas que shows de TV en casi todas las categorías.
15 La mayoría de series y películas se encuentran dentro de la categoría TV-MA, la cual apunta a público mayor de 17 años,
16 por lo que podría considerarse que este es el mercado al que más atención se le presta en la actualidad.''' )
17 print("")
```

Análisis 1: Cantidad de Shows de TV y Películas por Rating

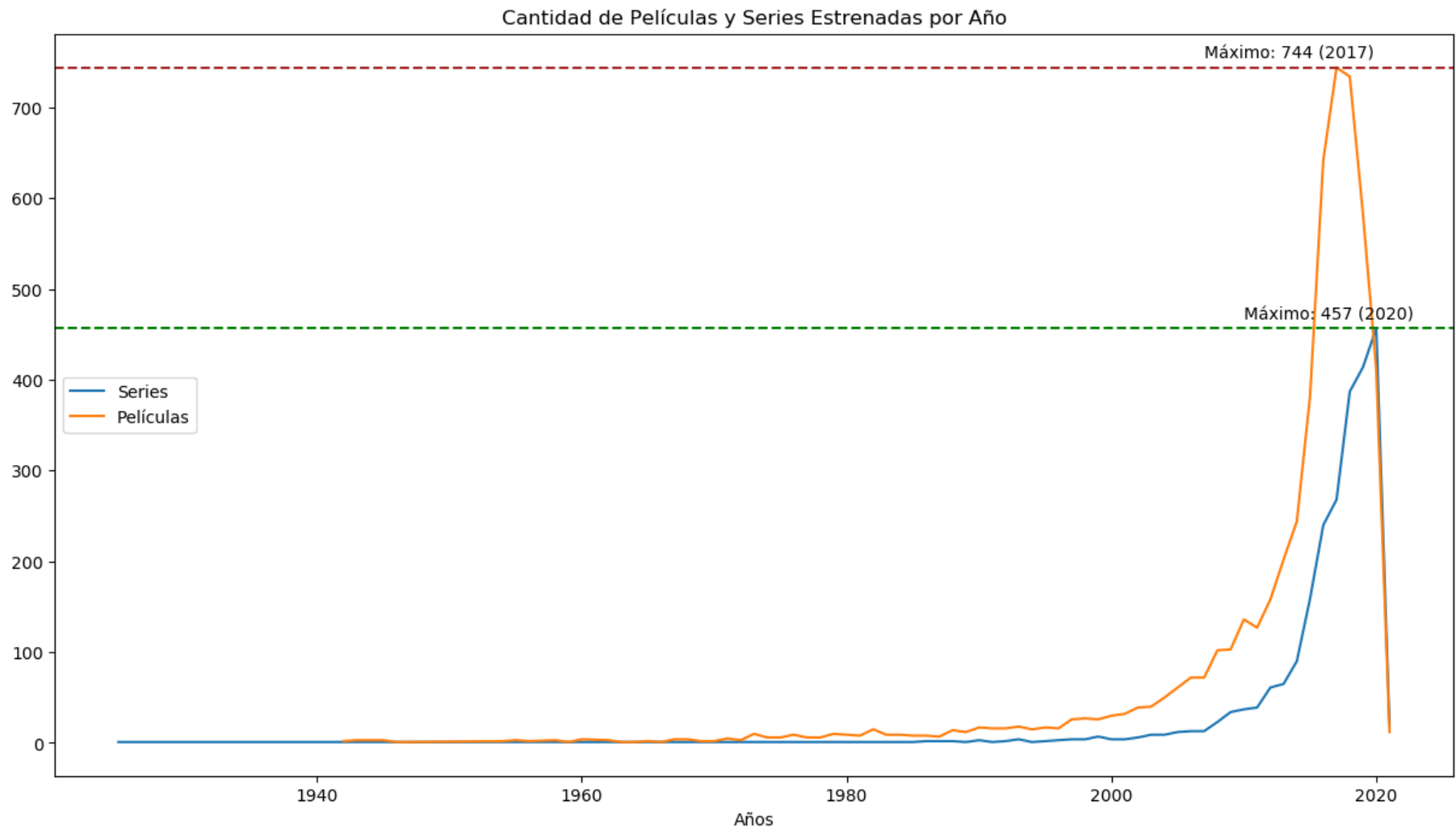


Observación del Análisis 1: Por lo general, hay más películas que shows de TV en casi todas las categorías. La mayoría de series y películas se encuentran dentro de la categoría TV-MA, la cual apunta a público mayor de 17 años, por lo que podría considerarse que este es el mercado al que más atención se le presta en la actualidad.

In [202]:

```
1 print("Análisis 2: Cantidad de Películas y Series Estrenadas por Año")
2
3 plt.style.use("default")
4
5 contador1 = pd.value_counts(datos.loc[datos["type"] == "TV Show", "release_year"]).sort_index(ascending = True)
6 contador2 = pd.value_counts(datos.loc[datos["type"] == "Movie", "release_year"]).sort_index(ascending = True)
7
8 fig, ax = plt.subplots()
9 fig.set_size_inches(15, 8)
10
11 ax.plot(contador1, label = "Series")
12 ax.plot(contador2, label = "Películas")
13 ax.set_xlabel("Años")
14 ax.set_title("Cantidad de Películas y Series Estrenadas por Año")
15 ax.axhline(y = max(contador1), linestyle = "--", color = "green")
16 ax.annotate("Máximo: " + str(max(contador1)) + " (" + str(contador1[contador1 == max(contador1)].index[0]) + ")",
17            xy = [contador1[contador1 == max(contador1)].index[0] - 10, max(contador1) + 10])
18 ax.axhline(y = max(contador2), linestyle = "--", color = "brown")
19 ax.annotate("Máximo: " + str(max(contador2)) + " (" + str(contador2[contador2 == max(contador2)].index[0]) + ")",
20            xy = [contador2[contador2 == max(contador2)].index[0] - 10, max(contador2) + 10])
21 plt.legend()
22
23 plt.show()
24
25 print('''Observación del Análisis 2: La producción de series y películas ha incrementado drásticamente desde mediados
26 del siglo XX hasta la actualidad, teniendo las películas el mayor volumen de producción y, por lo tanto, de estrenos.
27 Sin embargo, el impacto de la pandemia del COVID-19 se hace evidente al final del gráfico, en donde puede apreciarse
28 el desplome en estrenos tanto de películas como de series, debido a los paros en muchas producciones.''' )
29 print("")
```

Análisis 2: Cantidad de Películas y Series Estrenadas por Año

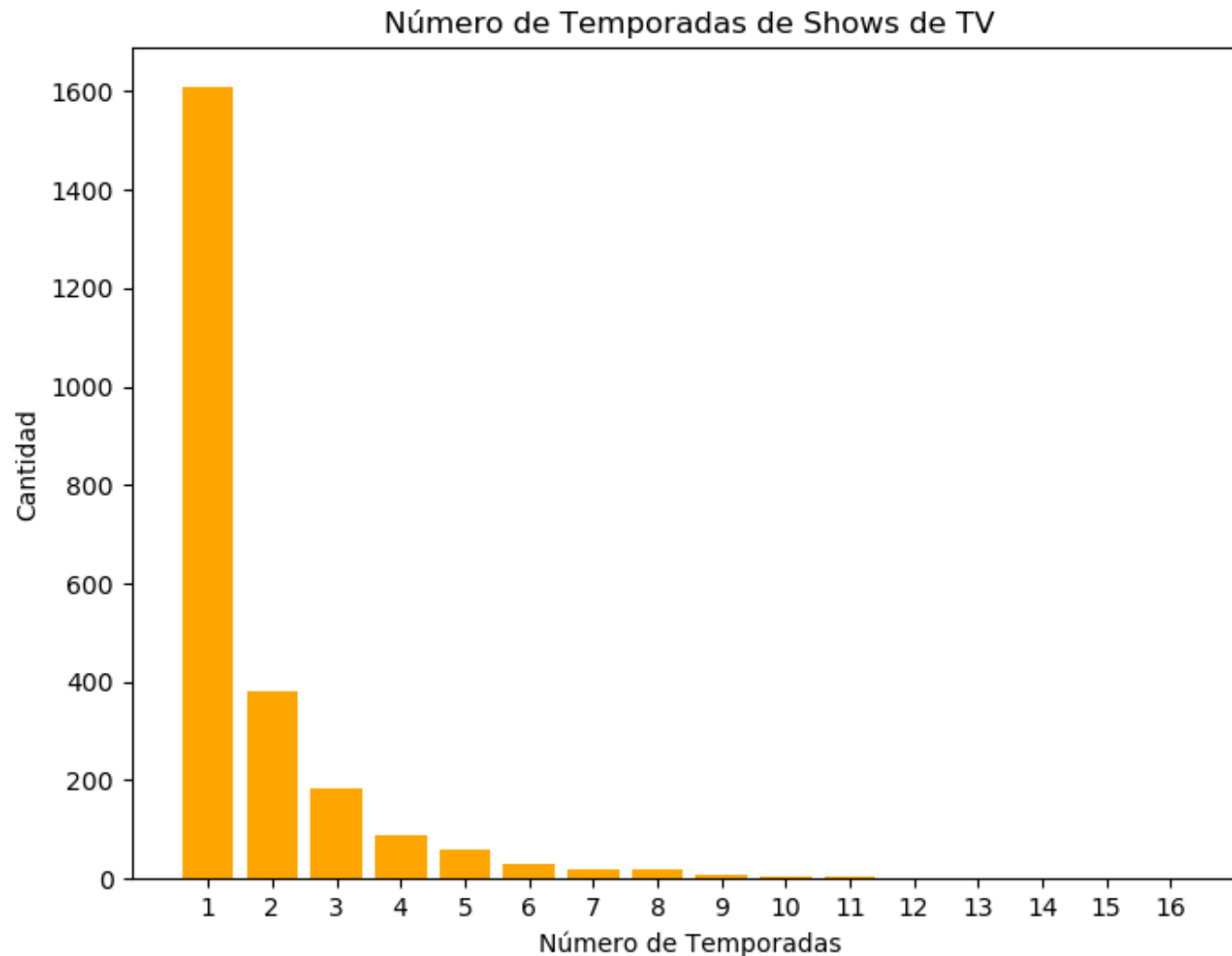


Observación del Análisis 2: La producción de series y películas ha incrementado drásticamente desde mediados del siglo XX hasta la actualidad, teniendo las películas el mayor volumen de producción y, por lo tanto, de estrenos. Sin embargo, el impacto de la pandemia del COVID-19 se hace evidente al final del gráfico, en donde puede apreciarse el desplome en estrenos tanto de películas como de series, debido a los paros en muchas producciones.

In [225]:

```
1 print("Análisis 3: Distribución del Número de Temporadas en Shows de TV")
2
3 fig, ax = plt.subplots()
4 fig.set_size_inches(8, 6)
5
6 subset1 = datos.loc[datos["type"] == "TV Show", ["duration"]]
7 subset1["duration"] = subset1["duration"].apply(lambda x: x.replace(" Seasons", ""))
8 subset1["duration"] = subset1["duration"].apply(lambda x: int(x.replace(" Season", "")))
9 subset1 = subset1.sort_values("duration", ascending = True)
10
11 contador = pd.value_counts(subset1["duration"]).sort_index(ascending = True)
12 contador = contador.append(pd.Series([0], index = [14]))
13
14 ax.bar(contador.index, contador.values, tick_label = contador.index, color = "orange")
15 ax.set_xlabel("Número de Temporadas")
16 ax.set_ylabel("Cantidad")
17 ax.set_title("Número de Temporadas de Shows de TV")
18
19 plt.show()
20
21 print('''Observación del Análisis 3: La enorme mayoría de las series suelen tener una sola temporada actualmente.
22 Esto puede deberse a tres posibles causas:
23 1) Algunas de estas series fueron recientemente producidas y estrenadas (como se evidenció en el Análisis 2) y por lo
24 tanto no ha transcurrido suficiente tiempo para que una "siguiente temporada" ocurra aún.
25 2) Algunas series no tuvieron la recepción esperada y se quedaron en una sola temporada, sin renovación.
26 3) Otras series podrían haber sido renovadas para una siguiente temporada, la cual está todavía en producción, lo cual
27 llevaría a esperar lanzamientos aplazados debido al efecto de la pandemia.''' )
28 print("")
```

Análisis 3: Distribución del Número de Temporadas en Shows de TV



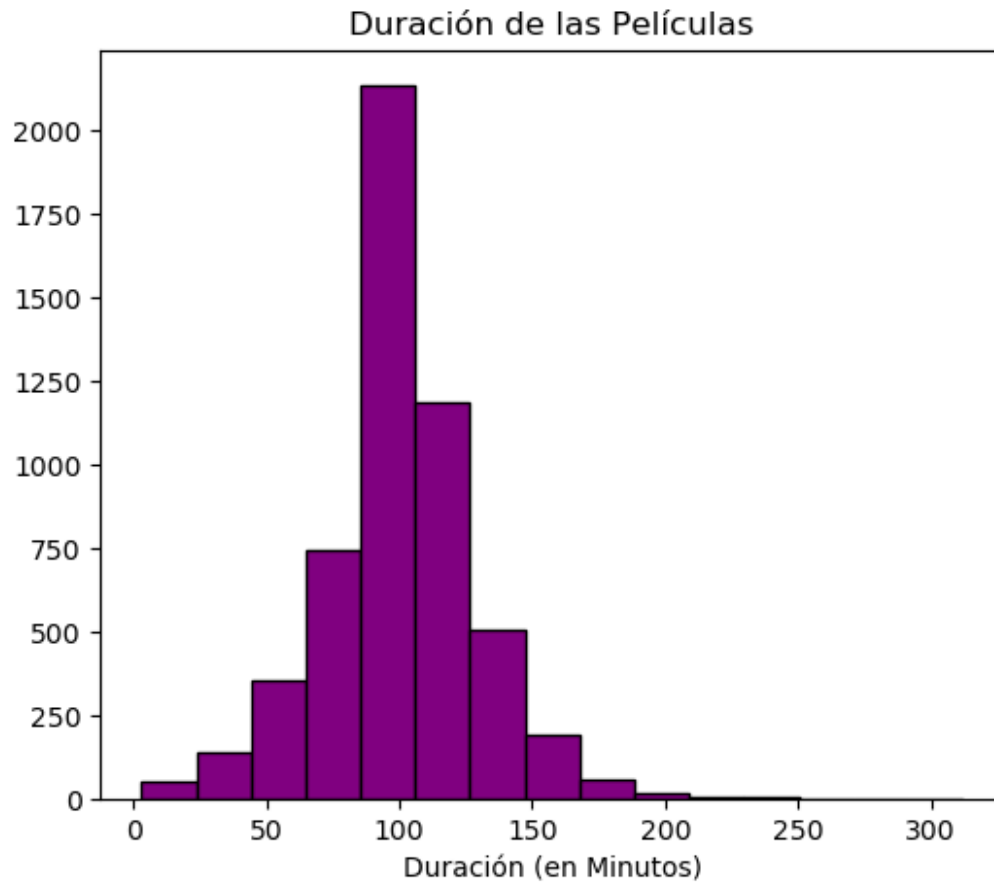
Observación del Análisis 3: La enorme mayoría de las series suelen tener una sola temporada actualmente.

Esto puede deberse a tres posibles causas:

- 1) Algunas de estas series fueron recientemente producidas y estrenadas (como se evidenció en el Análisis 2) y por lo tanto no ha transcurrido suficiente tiempo para que una "siguiente temporada" ocurra aún.
- 2) Algunas series no tuvieron la recepción esperada y se quedaron en una sola temporada, sin renovación.
- 3) Otras series podrían haber sido renovadas para una siguiente temporada, la cual está todavía en producción, lo cual llevaría a esperar lanzamientos aplazados debido al efecto de la pandemia.

```
In [228]: 1 print("Análisis 4: Distribución de la Duración de las Películas")
2
3 subset2 = datos.loc[datos["type"] == "Movie", ["duration"]]
4 subset2["duration"] = subset2["duration"].apply(lambda x: int(x.replace(" min", "")))
5
6 fig, ax = plt.subplots()
7 fig.set_size_inches(6, 5)
8
9 contador = pd.value_counts(subset1["duration"]).sort_index(ascending = True)
10 contador.head()
11
12 ax.hist(subset2["duration"], bins = 15, color = "purple", ec = "black")
13 ax.set_xlabel("Duración (en Minutos)")
14 ax.set_title("Duración de las Películas")
15
16 plt.show()
17
18 print('''Observación del Análisis 4: La mayoría de películas parece tener una duración de alrededor de 100 minutos, es
19 decir, cerca de dos horas. Curiosamente, hay una cierta cantidad de películas con duraciones por debajo de los 50
20 minutos; esto podría tratarse de cortos. Por otra parte, hay muy pocas películas que duren más de 200 minutos.
21 Podría decirse que el rango habitual de duración de películas a nivel histórico se sitúa entre una hora y dos horas y
22 media.''' )
23 print("")
```

Análisis 4: Distribución de la Duración de las Películas



Observación del Análisis 4: La mayoría de películas parece tener una duración de alrededor de 100 minutos, es decir, cerca de dos horas. Curiosamente, hay una cierta cantidad de películas con duraciones por debajo de los 50 minutos; esto podría tratarse de cortos. Por otra parte, hay muy pocas películas que duren más de 200 minutos. Podría decirse que el rango habitual de duración de películas a nivel histórico se sitúa entre una hora y dos horas y media.

5. Conclusiones

```
In [230]: 1 print('''El análisis realizado permitió llegar a las siguientes conclusiones:
2 - El mercado de películas y series con contenido para público maduro es el más explorado históricamente.
3 - La pandemia del Coronavirus afectó fuertemente los estrenos de películas y series en 2020.
4 - La mayoría de series tiene una sola temporada hasta el momento. Podría hacerse un análisis más profundo para validar
5 la existencia de las tres causas hipotéticas de esto (series no renovadas, series en producción de nuevas temporadas y
6 series en espera de la crítica).
7 - Las películas por lo general duran entre una y dos horas y media.
8 ''')
```

El análisis realizado permitió llegar a las siguientes conclusiones:

- El mercado de películas y series con contenido para público maduro es el más explorado históricamente.
- La pandemia del Coronavirus afectó fuertemente los estrenos de películas y series en 2020.
- La mayoría de series tiene una sola temporada hasta el momento. Podría hacerse un análisis más profundo para validar la existencia de las tres causas hipotéticas de esto (series no renovadas, series en producción de nuevas temporadas y series en espera de la crítica).
- Las películas por lo general duran entre una y dos horas y media.