

## Análisis del Perfil de Datos

Tras analizar el perfil de datos generado para los datos del 1 al 31 de enero de la API de TVMAZE, se identificaron varios problemas en el conjunto de datos, tales como variables con altas correlaciones, valores perdidos, datos no soportados y desequilibrios en variables categóricas. A continuación, se detallan las conclusiones.

1. Se han detectado múltiples variables con altas correlaciones entre sí, columnas como:

- `_embedded.show.averageRuntime` con `_embedded.show.network.country.code` y otras 6 variables.
- `_embedded.show.externals.thetvdb` con `_embedded.show.externals.tvrage` y otras 9 variables.

Estas columnas pueden duplicar la información y ser redundantes.

<code>_embedded.show.averageRuntime</code> is highly overall correlated with <code>_embedded.show.network.country.code</code> and 6 other fields	High correlation
<code>_embedded.show.externals.thetvdb</code> is highly overall correlated with <code>_embedded.show.externals.tvrage</code> and 9 other fields	High correlation
<code>_embedded.show.externals.tvrage</code> is highly overall correlated with <code>_embedded.show.externals.thetvdb</code> and 13 other fields	High correlation
<code>_embedded.show.id</code> is highly overall correlated with <code>_embedded.show.externals.thetvdb</code> and 6 other fields	High correlation
<code>_embedded.show.language</code> is highly overall correlated with <code>_embedded.show.externals.tvrage</code> and 10 other fields	High correlation
<code>_embedded.show.network.country.code</code> is highly overall correlated with <code>_embedded.show.averageRuntime</code> and 23 other fields	High correlation
<code>_embedded.show.network.country.name</code> is highly overall correlated with <code>_embedded.show.averageRuntime</code> and 23 other fields	High correlation
<code>_embedded.show.network.country.timezone</code> is highly overall correlated with <code>_embedded.show.averageRuntime</code> and 23 other fields	High correlation
<code>_embedded.show.network.id</code> is highly overall correlated with <code>_embedded.show.externals.tvrage</code> and 13 other fields	High correlation
<code>_embedded.show.network.name</code> is highly overall correlated with <code>_embedded.show.averageRuntime</code> and 25 other fields	High correlation
<code>_embedded.show.network.officialSite</code> is highly overall correlated with <code>_embedded.show.averageRuntime</code> and 25 other fields	High correlation
<code>_embedded.show.rating.average</code> is highly overall correlated with <code>_embedded.show.network.country.code</code> and 5 other fields	High correlation
<code>_embedded.show.runtime</code> is highly overall correlated with <code>_embedded.show.averageRuntime</code> and 8 other fields	High correlation

2. Varias variables presentan un alto porcentaje de valores perdidos:

- `airtime`: 51.4% valores perdidos.
- `runtime`: 9.5% valores perdidos.
- `summary`: 69.2% valores perdidos.
- `rating.average`: 92.9% valores perdidos.

Hay más variables con más del 30% de valores perdidos, lo que puede distorsionar los resultados y las conclusiones del análisis.

<code>airtime</code> has 2460 (51.4%) missing values	Missing
<code>runtime</code> has 453 (9.5%) missing values	Missing
<code>summary</code> has 3312 (69.2%) missing values	Missing
<code>rating.average</code> has 4448 (92.9%) missing values	Missing
<code>image.medium</code> has 3561 (74.4%) missing values	Missing
<code>image.original</code> has 3561 (74.4%) missing values	Missing
<code>_embedded.show.language</code> has 330 (6.9%) missing values	Missing
<code>_embedded.show.runtime</code> has 3582 (74.8%) missing values	Missing
<code>_embedded.show.averageRuntime</code> has 309 (6.5%) missing values	Missing
<code>_embedded.show.ended</code> has 3066 (64.0%) missing values	Missing
<code>_embedded.show.officialSite</code> has 473 (9.9%) missing values	Missing
<code>_embedded.show.rating.average</code> has 4048 (84.6%) missing values	Missing
<code>_embedded.show.network</code> has 4787 (100.0%) missing values	Missing
<code>_embedded.show.webChannel.id</code> has 112 (2.3%) missing values	Missing
<code>_embedded.show.webChannel.name</code> has 112 (2.3%) missing values	Missing
<code>_embedded.show.webChannel.country.name</code> has 1602 (33.5%) missing values	Missing

### 3. Desequilibrio en variables categóricas:

La variable `type` y `schedule.time` se encuentra desequilibradas con una categoría predominando sobre las demás. Lo que puede provocar que los modelos se inclinen a la categoría mayoritaria.

<code>type</code> is highly imbalanced (96.3%)	Imbalance
<code>_embedded.show.schedule.time</code> is highly imbalanced (50.6%)	Imbalance

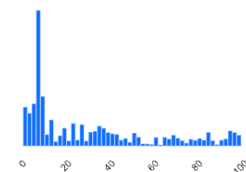
### 4. Variables con valores en cero como: `_embedded.show.weight`, donde el 2.8% de los valores son cero

`_embedded.show.weight`

Real number (R)

High correlation Zeros

Distinct	101	Minimum	0
Distinct (%)	2.1%	Maximum	100
Missing	0	Zeros	134
Missing (%)	0.0%	Zeros (%)	2.8%
Infinite	0	Negative	0
Infinite (%)	0.0%	Negative (%)	0.0%
Mean	30.856695	Memory size	37.5 KiB



More details

5. Varias variables están marcadas como no soportadas o requieren limpieza. Entre ellas se encuentran las siguientes, las cuales son fundamentales para la información de los shows y episodios:

- `_embedded.show.genres`
- `_embedded.show.schedule.days`
- `_embedded.show.network`

<code>_embedded.show.genres</code>	<code>_embedded.show.schedule.days</code>	<code>_embedded.show.network</code>
[Comedy, Crime]	[Monday, Tuesday, Wednesday, Thursday]	NaN
[Thriller, Mystery]	[Thursday]	NaN
[]	[Thursday]	NaN
[]	[Monday, Tuesday, Wednesday, Thursday, Friday, Saturday, Sunday]	NaN
[Drama]	[]	NaN
[Drama]	[]	NaN
[Drama]	[]	NaN
[Drama]	[]	NaN
[Drama]	[]	NaN
[Drama]	[]	NaN

- Como se puede observar las variables: `embedded.show.genres` y `_embedded.show.schedule.days` almacenan datos en forma de listas o vectores, por lo que deben ser transformadas.
- El campo NaN puede significar valores faltantes.

6. `id`, `url`, `_links.self.href`: Tienen valores únicos para cada registro. Siendo útiles para rastrear los episodios y todos los datos. No aportan información analítica, pero son esenciales para realizar las relaciones de los elementos.

<code>id</code> has <b>unique</b> values	Unique
<code>url</code> has unique values	Unique
<code>_links.self.href</code> has unique values	Unique

7. Columnas como season presentan valores atípicos, se puede observar que los números van del 1 al 25 y una distorsión en 2024.

Statistics			Histogram			Common values			Extreme values		
Value			Count		Frequency (%)						
1	2		2519		52.6%						
2024			694		14.5%						
2			559		11.7%						
3			259		5.4%						
5			120		2.5%						
4			114		2.4%						
6			73		1.5%						
8			66		1.4%						
25			36		0.8%						
11			33		0.7%						
Other values (24)			314		6.6%						

### Conclusión

El conjunto de datos presenta inconsistencias que deben solucionarse antes de incorporarlos en un modelo estructurado de base de datos o utilizarlos en análisis posteriores.

- Redundancia de Datos: Se debe manejar la redundancia de datos para evitar problemas de multicolinealidad. Es recomendable eliminar o combinar variables altamente correlacionadas.
- Valores Faltantes: Es crucial manejar adecuadamente los valores faltantes para mantener la integridad del análisis.
- Desequilibrio en Variables Categóricas: Las variables categóricas desequilibradas requieren ajustes para garantizar su confiabilidad.
- Preprocesamiento de Variables No Estructuradas: Es necesario preprocesar y transformar variables no estructuradas y no soportadas para su uso efectivo. En particular, las variables que contienen arrays, como géneros y días de emisión, deben ser transformadas en un formato adecuado para la carga y el análisis.

## Acciones realizadas

1. faltantes o contienen datos no soportados e irrelevantes para el análisis. Se eliminaron estas columnas identificadas para reducir la redundancia y mejorar la calidad del conjunto de datos.
2. Se eliminaron las filas donde el valor de la columna season es 2024, considerado un valor atípico no esperado, ya que las temporadas suelen numerarse de manera secuencial, adicional se puede evidenciar que van del 1 al 25, ese valor es incorrecto.
3. Se seleccionaron las columnas `_embedded.show.genres`, `_embedded.show.schedule.days` y `_embedded.show.network` que contienen arrays o estructuras complejas. Estas columnas se convirtieron en cadenas de texto separadas por comas para facilitar su manipulación y análisis. Aunque los datos de la columna `_embedded.show.network` son escasos, se decidió incluirlos en el modelo para no perder información valiosa.
4. Se definió un diccionario `days_mapping` para convertir los nombres de los días de la semana a números, con el fin de estandarizar los datos y facilitar los análisis numéricos y el modelado.
5. Se rellenaron los valores faltantes de columnas numéricas como `runtime` y `_embedded.show.averageRuntime` con la mediana de cada columna, manteniendo así la consistencia en la distribución de los datos.
6. Se eliminaron las filas duplicadas en el DataFrame para asegurar que cada registro sea único y evitar sesgos en el análisis.
7. Se ajustaron las variables categóricas, reemplazando las categorías con pocas apariciones por 'Other', con el objetivo de equilibrar la distribución de categorías y mejorar la confiabilidad de los modelos.
8. Se convierten variables categóricas en variables dummy para análisis y modelado.
9. Se eliminaron las filas del conjunto de datos que tienen más del 20% de valores faltantes (NaN), mejorando así la calidad y confiabilidad de los datos restantes para el análisis.