

Análisis de Sentimiento en Reseñas de Disneyland

Diana Cordero

Febrero 2025

Abstract

Este reporte presenta un análisis de sentimiento aplicado a reseñas de visitantes de los parques Disneyland. Se utilizó la librería **TextBlob** para clasificar los comentarios en positivos, negativos y neutros. Los resultados permiten identificar tendencias en la percepción de los usuarios sobre su experiencia en los parques.

1 Introducción

Los parques de Disneyland son destinos turísticos populares en todo el mundo. Las opiniones de los visitantes reflejan su satisfacción y pueden proporcionar información valiosa para mejorar la experiencia del cliente. En este estudio, se analizan reseñas extraídas de un conjunto de datos disponible en Kaggle. El conjunto de datos contiene 42,656 registros y las siguientes columnas:

- **Review_ID**: Identificador único de la reseña.
- **Rating**: Calificación otorgada por el usuario.
- **Year_Month**: Fecha de la reseña en formato año-mes.
- **Reviewer_Location**: Ubicación del usuario que escribió la reseña.
- **Review_Text**: Texto completo de la reseña.
- **Branch**: Sucursal de Disneyland a la que se refiere la reseña.

2 Descripción del Conjunto de Datos

El conjunto de datos utilizado en este análisis proviene de Kaggle y contiene un total de 42,656 registros con 18 columnas. A continuación, se describen las principales características del conjunto:

- **Review_ID**: Identificador único de cada reseña (`int64`).
- **Rating**: Calificación otorgada por el usuario en una escala de 1 a 5 (`int64`).

- **Year_Month:** Fecha de la reseña en formato YYYY-MM (`object`).
- **Reviewer_Location:** Ubicación del usuario que escribió la reseña (`object`).
- **Review_Text:** Texto completo de la reseña (`object`).
- **Branch:** Sucursal de Disneyland a la que se refiere la reseña (`object`).
- **Sentiment:** Valor numérico de sentimiento asignado a la reseña (`float64`).
- **Sentiment_Label:** Clasificación del sentimiento en positivo, negativo o neutro (`object`).
- **Cleaned_Review:** Versión preprocesada del texto de la reseña (`object`).
- **Sentiment_Score:** Puntuación de sentimiento obtenida (`float64`).
- **TextBlob_Polarity:** Polaridad del sentimiento calculada con `TextBlob` (`float64`).
- **TextBlob_Subjectivity:** Nivel de subjetividad de la reseña (`float64`).
- **VADER_Label:** Clasificación de sentimiento utilizando el modelo VADER (`object`).
- **TextBlob_Label:** Clasificación de sentimiento utilizando `TextBlob` (`object`).
- **Review_Segments:** Segmentación del texto de la reseña (`object`).
- **Truncated_Review:** Versión truncada de la reseña (`object`).
- **Tokenized_Review:** Representación tokenizada del texto de la reseña (`object`).
- **Review_Length:** Longitud del texto de la reseña en número de caracteres (`int64`).

El dataset tiene un tamaño aproximado de 5.9 MB y no presenta valores nulos en ninguna de sus columnas, lo que permite un análisis completo de las reseñas sin necesidad de imputación de datos.

3 Metodología

Se realizó un análisis exploratorio de los datos, que incluyó:

- Conteo de la distribución de reseñas de acuerdo con la calificación.
- Conteo de reseñas por mes y año.
- Identificación del top 10 de ubicaciones de los reseñadores.
- Análisis de la distribución de la longitud de las reseñas.

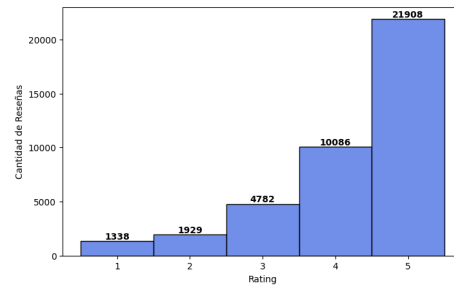


Figure 1: Distribución de las reseñas

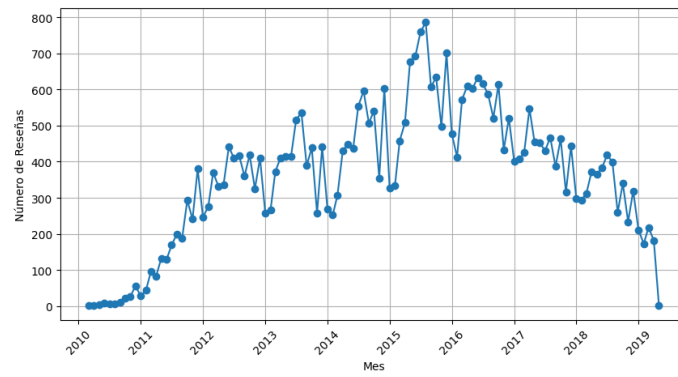


Figure 2: Conteo de las reseñas a través del tiempo

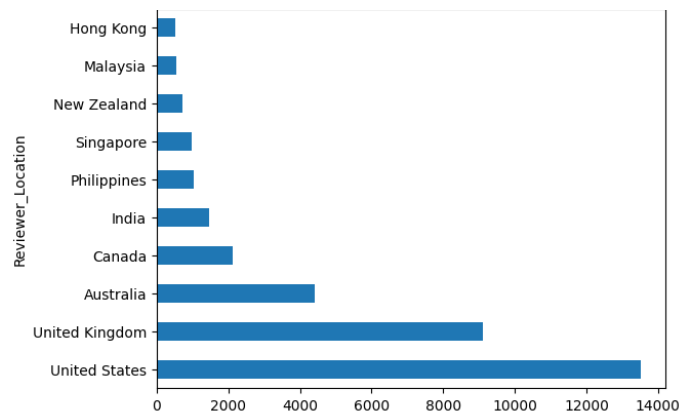


Figure 3: Ubicaciones de los reseñadores

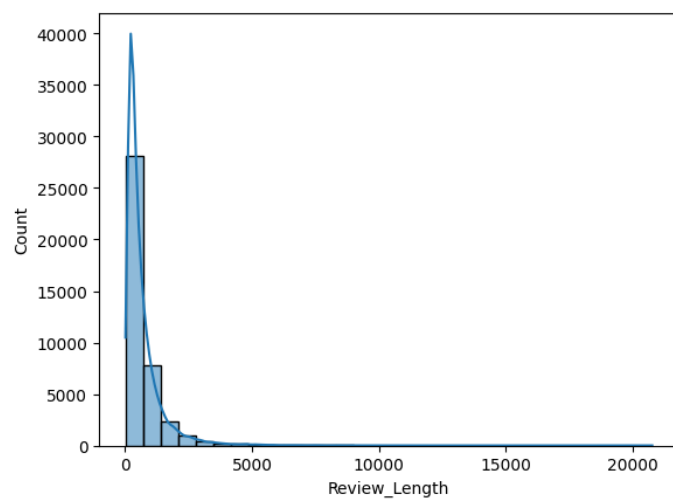


Figure 4: Longitud de las reseñas

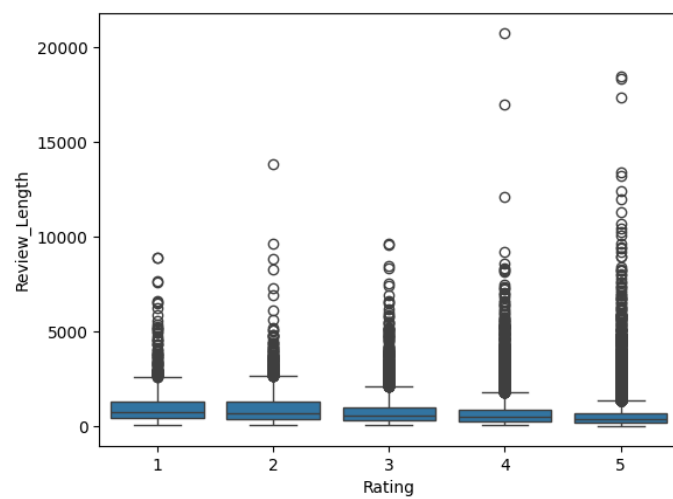


Figure 5: Relación entre la longitud de la reseña y la calificación dada

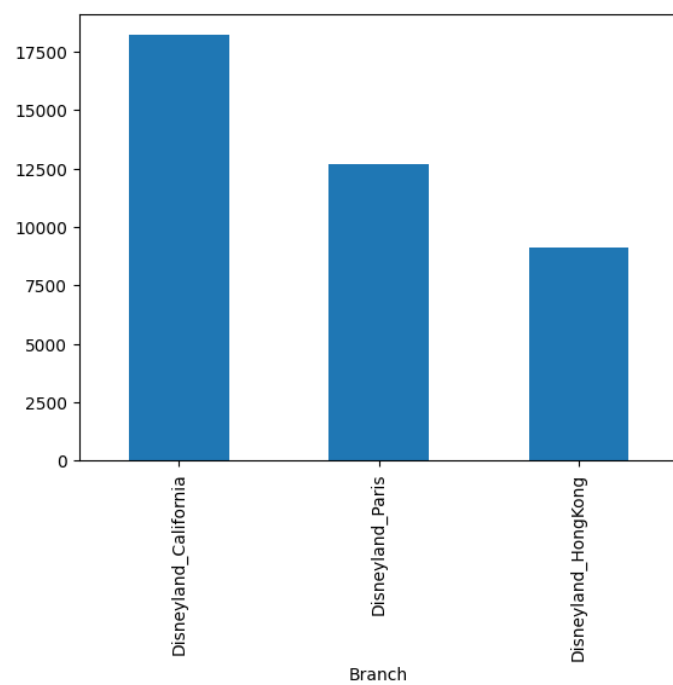


Figure 6: Reseñas por sucursal de Disneyland

- Comparación de la longitud de la reseña vs la calificación dada.
- Conteo de reseñas por sucursal de Disneyland.

Posteriormente, se llevó a cabo la limpieza de los textos mediante:

- Eliminación de stopwords.
- Eliminación de signos de puntuación.
- Conversión del texto a minúsculas.
- Lematización de las palabras.
- Conteo de palabras en las reseñas.

Finalmente, se realizó el análisis de sentimiento utilizando las herramientas **VADER** y **TextBlob** para clasificar los comentarios en positivos, negativos y neutros.

4 Resultados

4.1 Distribución de Sentimientos

A continuación, se muestra la distribución de los sentimientos en las reseñas analizadas:

Table 1: Distribución de sentimientos en las reseñas de Disneyland

Categoría	Porcentaje
Positivo	88.31%
Neutro	9.99%
Negativo	1.69%

La Figura 7 muestra la distribución de calificaciones en función del sentimiento de las reseñas. Se observa que las reseñas etiquetadas como positivas tienen una mediana de calificación de 5, con la mayoría de los valores concentrados entre 4 y 5, lo que indica que las experiencias bien valoradas suelen estar asociadas con sentimientos positivos. Por otro lado, las reseñas clasificadas como negativas presentan una mayor dispersión en las calificaciones, con una mediana cercana a 3 y valores distribuidos entre 1 y 5, lo que sugiere que algunas experiencias con sentimientos negativos no necesariamente recibieron las calificaciones más bajas. Finalmente, las reseñas etiquetadas como neutrales tienen una mediana de 4, con valores distribuidos en todo el rango de calificaciones, lo que indica que el análisis de sentimiento puede haber clasificado algunas reseñas como neutrales a pesar de tener calificaciones altas o bajas.

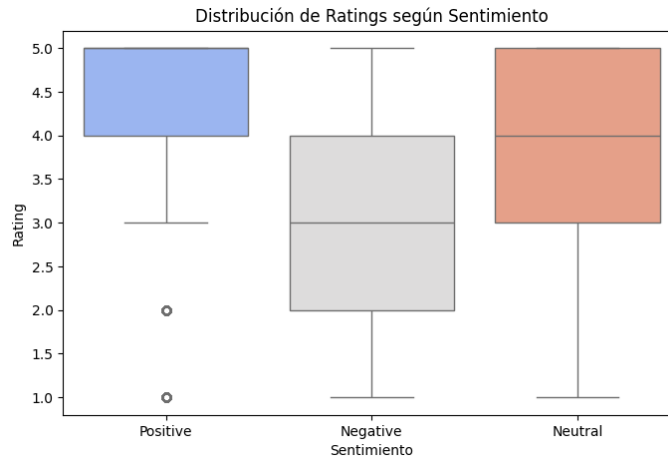


Figure 7: Distribución de calificaciones según el sentimiento de la reseña

La Figura 8 muestra la distribución de las puntuaciones de sentimiento obtenidas con VADER. Se observa que la mayoría de las reseñas tienen puntuaciones cercanas a 1.0, lo que indica una fuerte tendencia hacia sentimientos positivos. En contraste, hay pocas reseñas con puntuaciones negativas, lo que sugiere que las opiniones desfavorables son menos frecuentes en el conjunto de datos. Además, se puede notar una menor concentración de valores alrededor de 0, lo que indica que pocas reseñas fueron clasificadas como neutrales. Esta distribución sesgada hacia valores positivos sugiere que la mayoría de las experiencias reflejadas en las reseñas son favorables.

La Figura 9 presenta la cantidad de reseñas clasificadas según su sentimiento. Se observa que la mayoría de las reseñas corresponden a la categoría "Positive", con una cantidad significativamente mayor en comparación con las categorías "Negative" y "Neutral". Esto sugiere que la percepción general de los visitantes es predominantemente positiva. La cantidad de reseñas negativas es relativamente baja, y las reseñas neutrales son las menos frecuentes. Esta distribución refuerza la tendencia observada en la Figura 8, donde las puntuaciones de sentimiento estaban sesgadas hacia valores positivos.

Para evaluar la coherencia entre las calificaciones otorgadas por los usuarios y el análisis de sentimiento, se definieron umbrales para considerar una calificación como alta (≥ 4) o baja (≤ 2). Posteriormente, se identificaron casos en los que la etiqueta de sentimiento no coincidía con la calificación del usuario. Se consideraron inconsistentes aquellas reseñas en las que un usuario otorgó una calificación alta, pero el análisis de sentimiento fue negativo, o cuando la calificación fue baja y el análisis de sentimiento resultó positivo. Estos casos pueden indicar errores en la clasificación de sentimiento o usuarios que asignaron una puntuación que no refleja el contenido de su reseña. Finalmente, se contabilizaron y examinaron algunas de estas reseñas inconsistentes con el objetivo de

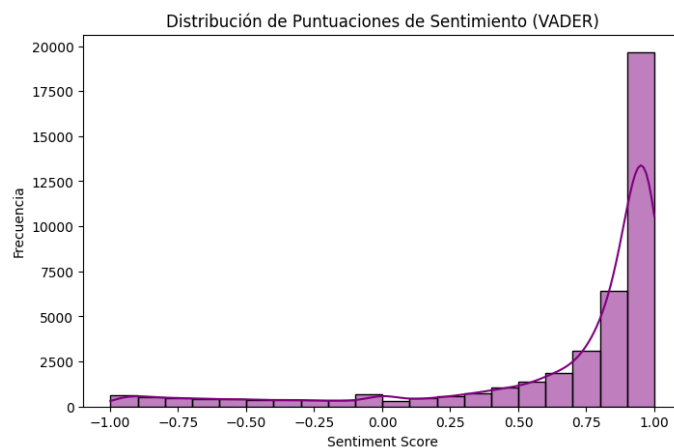


Figure 8: Distribución de puntuaciones de sentimiento según VADER

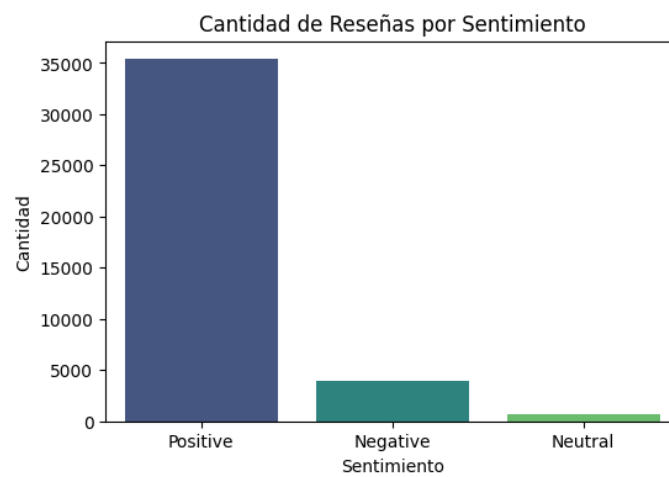


Figure 9: Cantidad de reseñas según el sentimiento clasificado

mejorar la precisión del análisis y detectar posibles anomalías en la clasificación de los sentimientos.

4.2 Análisis de Nube de Palabras para Reseñas Negativas con Calificación Alta

Como parte del análisis de sentimiento, se generó una **nube de palabras** para identificar los términos más frecuentes en reseñas que presentan una discrepancia entre la calificación numérica y la evaluación de sentimiento. Se consideraron aquellas reseñas con una calificación alta (≥ 4) pero clasificadas con un sentimiento negativo.

4.2.1 Filtrado de Datos

Se estableció un umbral de calificación alta (`high_rating_threshold = 4`) y se filtraron las reseñas donde:

- La calificación otorgada por el usuario es mayor o igual a 4.
- El modelo de análisis de sentimiento clasificó el texto como **negativo**.

4.2.2 Preprocesamiento del Texto

Para analizar el contenido de estas reseñas, se aplicaron técnicas de procesamiento de lenguaje natural (NLP), incluyendo:

- **Tokenización:** división de cada reseña en palabras individuales.
- **Eliminación de stopwords:** exclusión de palabras comunes sin carga informativa (ej. "el", "de", "y").
- **Filtrado de caracteres no alfabéticos:** eliminación de números y símbolos.

4.2.3 Cálculo de Frecuencia de Palabras

Se utilizó la función `Counter` de la librería `collections` para contar la frecuencia de cada palabra en el conjunto de datos filtrado. Posteriormente, se extrajeron las 20 palabras más repetidas.

4.2.4 Generación de la Nube de Palabras

Con los datos obtenidos, se generó una **nube de palabras** mediante la librería `WordCloud`, donde:

- Las palabras más frecuentes aparecen con un mayor tamaño.
- Se configuró un fondo blanco para mejorar la visibilidad.
- Se utilizaron los valores de frecuencia previamente calculados.

VADER	TextBlob	Cantidad
Positive	Positive	34,332
Positive	Neutral	2,545
Negative	Positive	1,545
Negative	Negative	1,376
Negative	Neutral	1,350
Positive	Negative	789
Neutral	Positive	369
Neutral	Neutral	211
Neutral	Negative	139

Table 2: Comparación de Clasificación de Sentimiento entre VADER y TextBlob

ocurre en la clasificación de sentimientos positivos (34,332 casos), mientras que en otras categorías existen discrepancias, especialmente en la clasificación de sentimientos negativos y neutrales.

Este análisis permite identificar diferencias entre los enfoques de cada modelo y evaluar su efectividad en la detección de sentimientos en los textos analizados.

5 Referencias

Disneyland Reviews. (2021, 19 enero). Kaggle. <https://www.kaggle.com/datasets/arushchillar/disneyland-reviews>

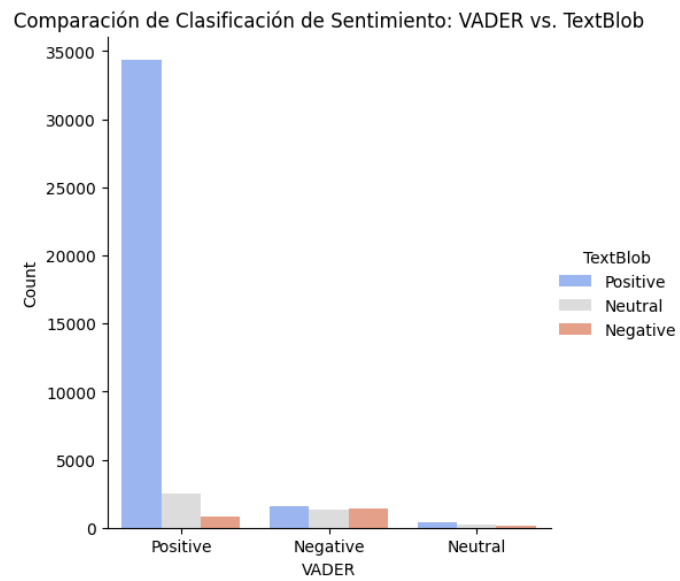


Figure 11: Gráfico de comparación de etiquetas de sentimiento entre VADER y TextBlob