

CODING BOOTCAMPS ESPOL: Data-Driven Decisions Specialist

GenAI for Data Analytics

Grupo 7

Proyecto: Análisis de Reseñas (Gestor de Contraseñas)

Entregable Final

Integrantes:

Forero Villota Katherine Sheila
Vallejo Vásquez Aquiles Anders
Gonzalez Gavilanes Luis Alberto
Peralta Baidal Darwin Vicente
Rubira Espinoza Ivonne Bethsabe

Introducción.....	3
Metodología.....	4
Estado Inicial del Dataset.....	4
Estado Actual del Dataset con Mejoras Aplicadas.....	8
Análisis Exploratorio y Estadísticas Descriptivas del Dataset.....	11
Visualizaciones Generadas.....	13
Insights Obtenidos.....	18
Resultados.....	20
Conclusiones.....	23

Introducción

El presente proyecto analiza un conjunto de datos compuesto por reseñas de usuarios de una aplicación de gestión de contraseñas publicadas en Google Play Store. La base incluye información relevante como fechas de publicación, calificaciones otorgadas, contenido textual de los comentarios y posibles respuestas del desarrollador, lo que permite estudiar tanto la percepción de los usuarios como la interacción del equipo responsable de la aplicación.

El objetivo principal es evaluar el nivel de satisfacción de los usuarios, identificar patrones temporales en las calificaciones, y generar indicadores clave que aporten información accionable para la mejora continua del producto y la experiencia del usuario. Básicamente, entender qué está funcionando, qué no, y dónde se están quejando más. Porque siempre se están quejando en algún lado.

Metodología

Estado Inicial del Dataset

Al cargar el archivo `GooglePlay_App_Data.xlsx`, el dataset presenta 206 registros y 12 columnas, lo que lo clasifica como un conjunto de datos pequeño. Debido a su tamaño, cada fila tiene un peso importante, por lo que cualquier eliminación debe analizarse con cuidado y justificarse técnicamente.

Al revisar la estructura general del dataset, se identifican problemas de calidad, tipado y completitud, los cuales afectan directamente la confiabilidad del análisis.

En cuanto a las columnas, el estado inicial es el siguiente:

- `review_id`
 - Tipo: object
 - 206 valores no nulos
 - Funciona como identificador único
 - No presenta duplicados
- `user_name`
 - Tipo: object
 - 206 valores completos
 - No requiere tratamiento
- `review_title`
 - Tipo: float64
 - 206 valores nulos (100%)
 - Una columna textual no puede ser float

- Evidencia error en el origen del archivo
 - No contiene información recuperable
- review_description
 - Tipo: object
 - 206 valores completos
 - Columna clave para análisis cualitativo
- rating
 - Tipo: int64
 - Valores entre 1 y 5
 - No existen valores fuera de rango
- thumbs_up
 - Tipo: float64
 - 161 valores no nulos y 45 nulos
 - Representa un conteo, pero aparece como decimal
 - El tipo incorrecto se debe a la presencia de nulos
- review_date
 - Tipo: object
 - 205 valores no nulos
 - Mezcla interna de strings, fechas y NaN
 - Presenta caracteres invisibles en algunos registros

- developer_response
 - Tipo: object
 - 204 valores no nulos
 - Los valores faltantes representan ausencia de respuesta
- developer_response_date
 - Tipo: object
 - 187 valores no nulos
 - Solo tiene sentido cuando existe respuesta del desarrollador
- appVersion
 - Tipo: object
 - 146 valores no nulos
 - Alto porcentaje de valores faltantes
- language_code
 - Tipo: object
 - Error ortográfico en el nombre de la columna
- country_code
 - Tipo: object
 - 206 valores completos

El análisis de valores nulos revela que el dataset no se encuentra limpio y que no todos los nulos tienen el mismo significado:

- Columnas totalmente inutilizables (review_title)
- Columnas con nulos estructurales (developer_response)
- Columnas con nulos críticos (review_date)
- Columnas con nulos imputables (thumbs_up, appVersion)

No se detectan registros duplicados ni duplicados lógicos por review_id, lo que confirma que el problema del dataset no es duplicación, sino calidad estructural.

En este estado inicial, el dataset no es apto para análisis temporal ni estadístico confiable, debido a problemas de tipado, nulos y coherencia semántica.

Estado Actual del Dataset con Mejoras Aplicadas

Tras aplicar el proceso completo de limpieza y validación, el dataset queda con 205 registros y 11 columnas, manteniendo prácticamente todo el volumen original, pero con una estructura mucho más sólida y consistente.

Los principales cambios estructurales aplicados son los siguientes:

- Eliminación definitiva de review_title, ya que:
 - Estaba 100% vacía
 - No existía posibilidad de imputación válida
 - Solo introducía ruido
- Corrección del nombre de columna:
 - language_code → language_code
- Conversión correcta de tipos de datos:
 - review_date → datetime64[ns]
 - developer_response_date → datetime64[ns]
 - Eliminación de caracteres invisibles antes de la conversión
- Tratamiento diferenciado de valores nulos:
 - Se elimina solo un registro sin review_date (variable crítica)
 - thumbs_up:
 - Nulos imputados con 0
 - Conversión de float64 a int64
 - developer_response:

- Nulos imputados como "No Response"
- appVersion:
 - Nulos imputados como "Unknown"
- developer_response_date:
 - Se mantiene como NaT únicamente cuando no hubo respuesta
- Creación de nueva variable estructural:
 - has_response:
 - 1 si el desarrollador respondió
 - 0 si no respondió
 - Permite análisis sin depender de valores nulos
- Validación de coherencia temporal:
 - Se detectaron fechas de respuesta anteriores a la reseña
 - Las fechas incoherentes fueron invalidadas (NaT)
 - No se eliminaron registros innecesariamente

En el estado final, el dataset presenta:

- 0 valores nulos en variables críticas (review_date, rating)
- Tipos de datos coherentes con el significado de cada columna
- Variables listas para análisis temporal, estadístico y visual
- Conservación del mayor volumen posible de información

Desde el punto de vista analítico, el dataset pasa de ser estructuralmente inconsistente a completamente utilizable, permitiendo:

- Análisis de tendencias en el tiempo
- Evaluación de percepción del usuario mediante ratings
- Medición del engagement con thumbs_up
- Análisis del comportamiento del desarrollador con has_response

Análisis Exploratorio y Estadísticas Descriptivas del Dataset

Una vez finalizado el proceso de limpieza y validación del dataset, se realizó un análisis exploratorio con el objetivo de comprender el comportamiento general de las variables numéricas y temporales, así como identificar patrones relevantes que permitan contextualizar los datos antes de etapas posteriores del proyecto.

A partir del resumen estadístico generado con `df.describe()`, se analizó la variable `rating`, la cual presenta una media de 3.98 y una mediana de 5, con valores comprendidos entre 1 y 5. Estos resultados indican una percepción mayoritariamente positiva de la aplicación, ya que la mayor concentración de reseñas se encuentra en calificaciones altas.

En concreto, el comportamiento de `rating` muestra que:

- La mayoría de los usuarios otorgan 5 estrellas.
- Las calificaciones bajas (1 y 2) representan una proporción menor.
- La desviación estándar (1.35) indica cierta variabilidad, pero sin valores extremos fuera del rango esperado.

En cuanto a la variable `thumbs_up`, que representa el nivel de interacción de otros usuarios con las reseñas, se observa una media de 1.52 y una mediana de 0. Esto evidencia una distribución altamente sesgada, donde la mayoría de reseñas no reciben votos adicionales.

Este comportamiento se puede resumir de la siguiente manera:

- Predominan valores cercanos a cero.
- Existen algunos casos aislados con alta cantidad de votos (hasta 31).
- Solo ciertas reseñas logran generar mayor visibilidad o relevancia.

Desde el punto de vista temporal, las fechas de reseña abarcan un período aproximado entre 2016 y 2024, lo que permite realizar análisis de evolución en el tiempo. El análisis por año muestra un pico significativo de reseñas en 2016, seguido de una disminución progresiva y una posterior estabilización en años posteriores.

Respecto a la gestión del desarrollador, la variable `has_response` presenta una media cercana a 0.99, lo que indica que el desarrollador respondió a más del 98% de las reseñas. Este resultado refleja una política activa de atención al usuario.

Finalmente, las visualizaciones descriptivas refuerzan los hallazgos estadísticos:

- El gráfico de distribución de rating confirma la concentración en calificaciones altas.
- El histograma de `thumbs_up` evidencia la asimetría de la variable.
- El gráfico de reseñas por año permite observar la evolución del volumen de opiniones.

En conjunto, este análisis exploratorio permite validar que el dataset final es coherente, consistente y adecuado para continuar con análisis más avanzados, modelado o visualización de resultados.

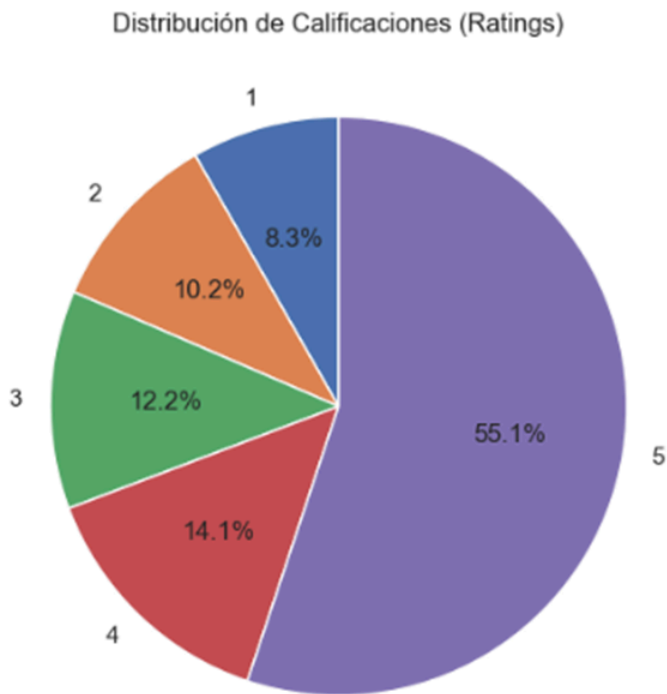
Visualizaciones Generadas

Se crearon visualizaciones en ChatGPT que nos ayuden a explicar los patrones encontrados en los datos. Algunas visualizaciones sugeridas para nuestro caso pueden ser:

1. *Gráfico Circular*

El gráfico circular muestra la distribución porcentual de las calificaciones otorgadas por los usuarios, permitiendo identificar la proporción de reseñas positivas, neutrales y negativas.

Visualización



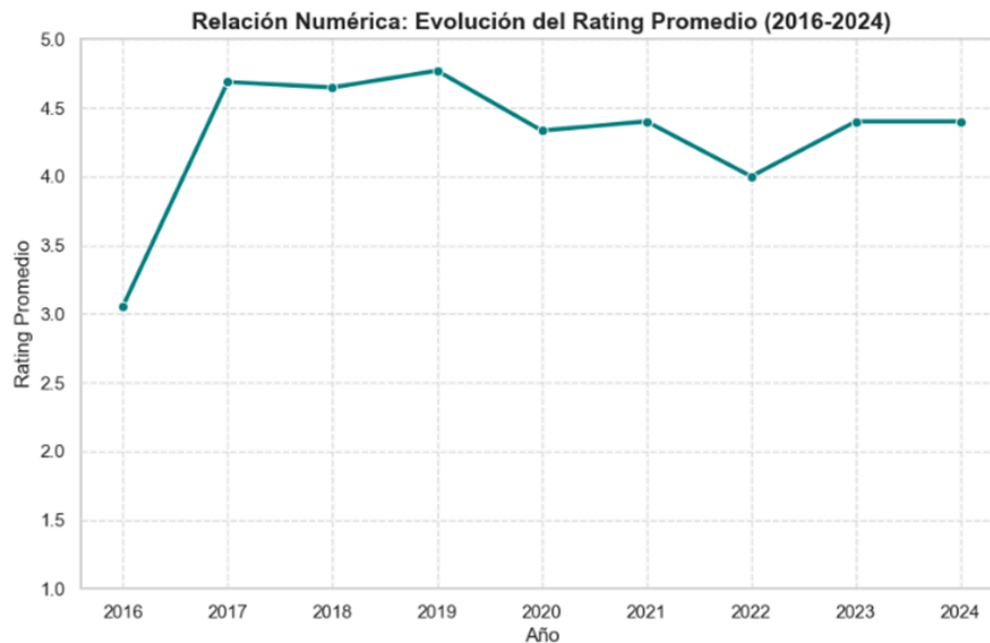
2. Gráfico de Líneas: Relación Numérica Temporal

Este gráfico muestra la relación entre dos variables numéricas: el Año y el Promedio de Calificación obtenido en ese periodo.

Detalle: Explica la evolución de la calidad a través del tiempo.

Patrón: La línea muestra una tendencia ascendente, lo que sugiere que la app ha mejorado desde su lanzamiento en 2016.

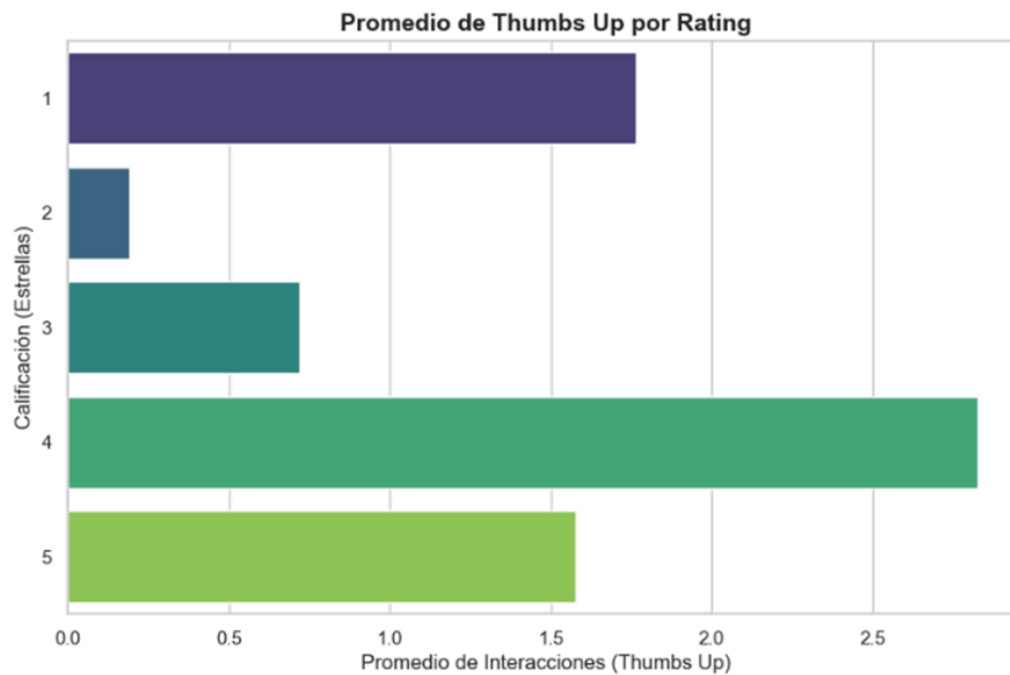
Visualización



3. Gráfico Engagement de Usuarios (Barras Horizontales)

Este gráfico muestra la relación entre la calificación (Rating) y el nivel de interacción, medido por el promedio de "Thumbs Up" (Me gusta). Al usar barras horizontales, es más fácil comparar cuál categoría de calificación despierta más interés o debate entre los usuarios.

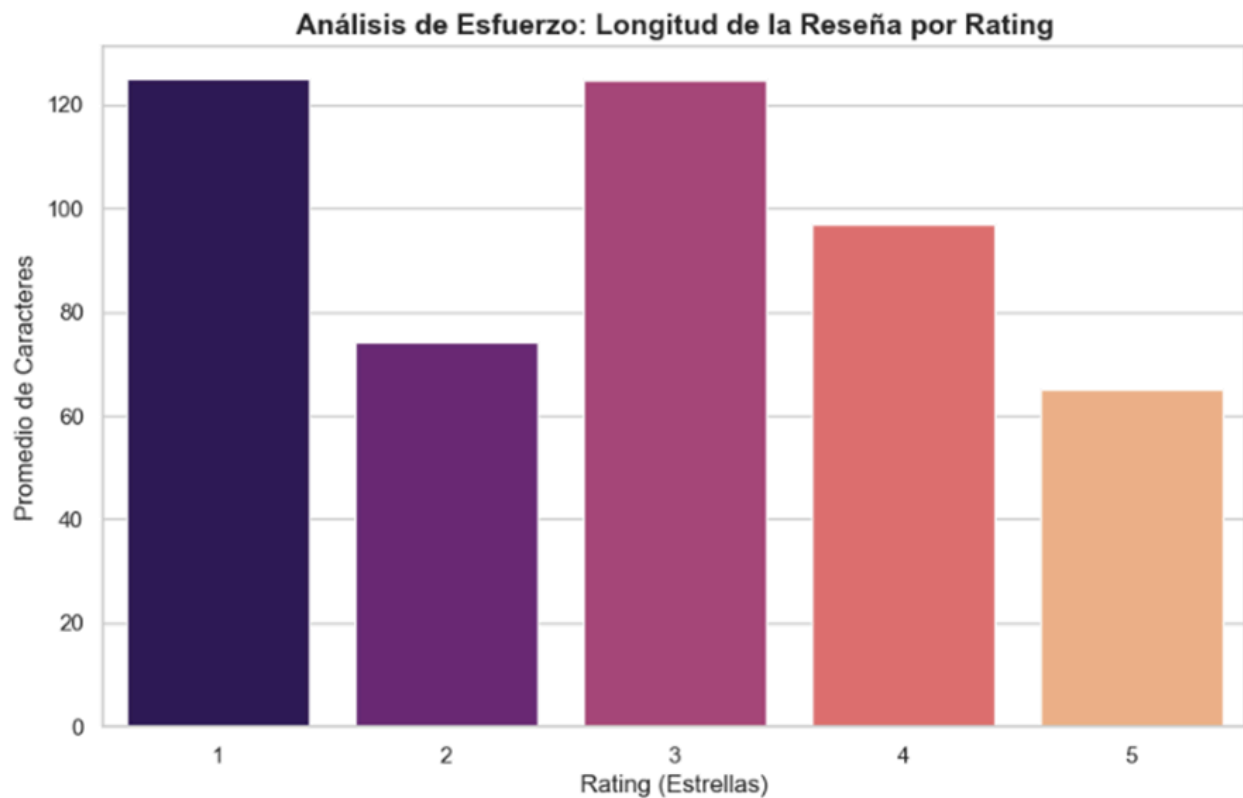
Visualización



4. Gráfico de Barras: Longitud de la Reseña vs. Rating.

Este nos ayuda a entender si los usuarios se esfuerzan más en escribir cuando están molestos o cuando están felices (análisis de sentimiento y esfuerzo).

Visualización



Validación

Para garantizar la confiabilidad del análisis, se realizaron procesos de validación de los gráficos generados por herramientas de inteligencia artificial. Se utilizó Python en Jupyter Notebook, empleando las librerías Pandas, Matplotlib.

Los resultados obtenidos fueron consistentes, lo que valida la precisión de los gráficos generados y confirma la confiabilidad del análisis exploratorio de datos realizado.


```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

rating_counts = df['rating'].value_counts().sort_index()

plt.figure(figsize=(6,6))
plt.pie(rating_counts, labels=rating_counts.index, autopct='%1.1f%%', startangle=90)
plt.title("Distribución de Calificaciones (Ratings)")
plt.show()
```

```
# Preparar datos: promedio de rating por año
df_evolucion = df.groupby('year')['rating'].mean().reset_index()

plt.figure(figsize=(10, 6))

# Creación del gráfico de líneas
sns.lineplot(data=df_evolucion, x='year', y='rating', marker='o', color='teal', linewidth=2.5)

plt.title('Relación Numérica: Evolución del Rating Promedio (2016-2024)', fontsize=14, fontweight='bold')
plt.xlabel('Año', fontsize=12)
plt.ylabel('Rating Promedio', fontsize=12)
plt.ylim(1, 5) # Escala real de calificación
plt.grid(True, linestyle='--', alpha=0.6)
plt.show()
```

```
# 2. Preparar los datos: Agrupamos por rating y sacamos el promedio de thumbs_up
thumbs_rating = df.groupby('rating')['thumbs_up'].mean().reset_index()

# 3. Configurar el estilo y tamaño
plt.figure(figsize=(10, 6))
sns.set_theme(style="whitegrid")

# 4. Crear el gráfico de barras HORIZONTALES
# fíjate que ahora 'rating' está en el eje Y
sns.barplot(data=thumbs_rating, x='thumbs_up', y='rating', palette='viridis', orient='h')

# 5. Títulos y etiquetas
plt.title('Promedio de Thumbs Up por Rating', fontsize=14, fontweight='bold')
plt.xlabel('Promedio de Interacciones (Thumbs Up)', fontsize=12)
plt.ylabel('Calificación (Estrellas)', fontsize=12)

plt.show()
```

```
# Calcular la longitud de cada reseña
df['review_length'] = df['review_description'].str.len()

# Configurar gráfico
plt.figure(figsize=(10, 6))
sns.barplot(data=df, x='rating', y='review_length', palette='magma', ci=None)

plt.title('Análisis de Esfuerzo: Longitud de la Reseña por Rating', fontsize=14, fontweight='bold')
plt.xlabel('Rating (Estrellas)', fontsize=12)
plt.ylabel('Promedio de Caracteres', fontsize=12)
plt.show()
```

Insights Obtenidos

Insight 1. Alta satisfacción general, con un segmento relevante de clientes insatisfechos

- ***Cómo se llegó a este insight:*** Se analizó la distribución de las calificaciones, clasificándolas en positivas (69.27%), neutrales (12.2%) y negativas (18.54%). Esto permitió evaluar la percepción general de los usuarios sobre el producto.
- ***Lo que dicen los datos:*** La mayoría de los usuarios tiene una percepción positiva del producto. Sin embargo, el 18.54% presenta experiencias negativas, lo que evidencia la existencia de problemas que afectan a una parte importante de los clientes.
- ***Preocupación por el cliente:*** Los clientes insatisfechos pueden perder confianza en el producto y reducir su uso, lo que afecta su experiencia general.
- ***Valor para la empresa:*** Este resultado permite identificar oportunidades de mejora que pueden aumentar la retención de clientes y fortalecer la reputación del producto.

Insight 2. Las reseñas negativas tienen mayor impacto en la percepción de otros usuarios

- ***Cómo se llegó a este insight:*** Se analizó el nivel de interacción (likes) de las reseñas según su calificación, identificando que las reseñas negativas reciben bastante atención.
- ***Lo que dicen los datos:*** Las experiencias negativas generan mayor interacción en promedio, lo que aumenta su visibilidad e influencia en otros usuarios.
- ***Preocupación por el cliente:*** Los usuarios pueden verse influenciados por estas reseñas, afectando su confianza y percepción del producto.
- ***Valor para la empresa:*** Este insight permite priorizar la solución de problemas más visibles, reduciendo el impacto negativo en la reputación del producto.

Insight 3. La experiencia del cliente mejoró significativamente a partir de 2017

- ***Cómo se llegó a este insight:*** Se analizó el rating promedio por año, observando una mejora notable en la satisfacción del cliente desde 2017.
- ***Lo que dicen los datos:*** El nivel de satisfacción aumentó después de 2017 y se ha mantenido estable, lo que indica que las mejoras implementadas fueron efectivas.
- ***Preocupación por el cliente:*** Es importante mantener la calidad del producto para asegurar una experiencia positiva y evitar una disminución en la satisfacción.
- ***Valor para la empresa:*** Este resultado confirma que las mejoras realizadas han sido efectivas y deben mantenerse para fortalecer la satisfacción y fidelización.

Insight 4. Existe una base sólida de clientes altamente satisfechos

- ***Cómo se llegó a este insight:*** Se analizó la distribución de las calificaciones, identificando que la mayoría de las reseñas corresponde a la calificación máxima.
- ***Lo que dicen los datos:*** La alta frecuencia de calificaciones de 5 estrellas indica un alto nivel de satisfacción entre los usuarios.
- ***Preocupación por el cliente:*** Esto refleja que el producto cumple con las expectativas y genera una experiencia positiva para la mayoría de los usuarios.
- ***Valor para la empresa:*** Contar con clientes satisfechos fortalece la fidelización y contribuye al crecimiento y posicionamiento del producto.

KPIs principales (Cards)

El dashboard inicia con cuatro indicadores clave que permiten comprender rápidamente el estado general del producto y el tamaño del conjunto de datos analizado. En primer lugar, se presenta el número total de reseñas, lo que proporciona contexto sobre la cantidad de información utilizada para el análisis y la representatividad de los resultados. En segundo lugar, se muestra la calificación promedio, indicador fundamental del nivel general de satisfacción de los usuarios.

Adicionalmente, se incluyen dos métricas de segmentación: el porcentaje de reseñas positivas y el porcentaje de reseñas negativas, clasificadas según el número de estrellas otorgadas. Estas métricas permiten identificar de manera inmediata la percepción general del producto y facilitan la comparación entre experiencias favorables y desfavorables.

Gráficos de análisis

El dashboard incorpora visualizaciones que permiten profundizar en la interpretación de los datos desde distintas perspectivas, estas son las antes ya propuestas.

- ★ La distribución por calificación muestra la proporción de reseñas según el número de estrellas, evidenciando una concentración importante en las calificaciones altas, especialmente en 5 estrellas. Esta visualización confirma la percepción general positiva detectada en los indicadores principales.
- ★ El gráfico de evolución del rating promedio permite analizar la tendencia temporal de la satisfacción del usuario desde 2016 hasta 2024. Se observa una mejora significativa después del año 2016 y una posterior estabilidad en niveles altos, lo que sugiere que las mejoras implementadas en la aplicación han tenido un impacto positivo sostenido en la experiencia del usuario.
- ★ El gráfico de “Thumbs Up por calificación” analiza el nivel de interacción de los usuarios con las reseñas según su rating. Se identifica que las reseñas con calificaciones bajas generan una interacción considerable, lo que indica que las experiencias negativas tienden a captar mayor atención de la comunidad, incrementando su impacto en la percepción pública del producto.
- ★ Por otro lado, el gráfico de longitud promedio de reseña permite evaluar el esfuerzo de los usuarios al escribir comentarios según su nivel de satisfacción. Se observa que las reseñas con calificaciones extremas, tanto positivas como negativas, tienden a ser más extensas, lo que sugiere una mayor motivación para expresar experiencias cuando estas son particularmente buenas o malas.

Tabla de detalle de reseñas

El dashboard incluye una tabla de detalle que presenta información individual de cada reseña, incluyendo usuario, fecha, calificación, contenido textual, número de interacciones y versión de la aplicación. Esta sección permite pasar del análisis agregado a la exploración específica de casos individuales, facilitando la validación de los patrones observados en los gráficos y permitiendo identificar ejemplos concretos de experiencias de usuario.

La tabla también aporta transparencia al análisis, ya que el usuario puede verificar directamente los datos fuente que respaldan los indicadores y visualizaciones presentadas.

Filtros interactivos

El panel de filtros permite segmentar la información según fecha de reseña, calificación y versión de la aplicación. Estos controles brindan flexibilidad, permitiendo explorar comportamientos específicos en distintos periodos de tiempo o versiones del producto. La capacidad de filtrar facilita la identificación de cambios asociados a actualizaciones de la aplicación y contribuye a un análisis más profundo orientado a la toma de decisiones.

Insights clave

A partir del análisis del dashboard y las visualizaciones generadas, se ratificaron los siguientes hallazgos principales sobre la percepción de los usuarios:

- **Existe un alto nivel de satisfacción general:** La mayoría de las reseñas corresponden a calificaciones positivas, con un rating promedio cercano a 4 estrellas, lo que indica que el producto cumple las expectativas de la mayoría de los usuarios.
- **Las reseñas negativas tienen mayor impacto en la percepción pública:** Las calificaciones bajas generan más interacción (likes) y visibilidad, lo que puede influir de manera desproporcionada en la reputación del producto frente a otros usuarios.
- **La experiencia del usuario mejoró después de los primeros años de lanzamiento:** El análisis temporal muestra una mejora significativa en las calificaciones a partir de 2017, seguida de una estabilidad en niveles altos, lo que sugiere que las mejoras implementadas han sido efectivas.
- **Existe una base sólida de clientes altamente satisfechos:** La alta proporción de reseñas de 5 estrellas refleja una experiencia positiva para un segmento importante de usuarios, lo que representa una oportunidad para fortalecer la fidelización y el posicionamiento del producto.

Conclusiones

El análisis de las reseñas permitió comprender de manera completa la percepción de los usuarios sobre la aplicación de gestión de contraseñas, identificando tanto fortalezas como oportunidades de mejora. Los resultados evidencian que la mayoría de los usuarios presenta un alto nivel de satisfacción, reflejado en la predominancia de calificaciones positivas, lo que confirma que el producto cumple con las expectativas principales del mercado.

Sin embargo, también se identificó un segmento relevante de usuarios insatisfechos cuyas reseñas generan mayor interacción y visibilidad, lo que implica que las experiencias negativas pueden influir de manera desproporcionada en la percepción pública del producto. Esto sugiere que la gestión proactiva de problemas y la atención prioritaria a incidencias críticas pueden tener un impacto significativo en la reputación de la aplicación.

Desde el punto de vista temporal, la mejora observada en las calificaciones a partir de 2017 indica que las actualizaciones y mejoras implementadas por el equipo de desarrollo han sido efectivas, lo que demuestra la importancia de la evolución continua del producto basada en la retroalimentación de usuarios. Mantener este enfoque orientado al usuario resulta clave para sostener la satisfacción a lo largo del tiempo.

Adicionalmente, la alta tasa de respuesta del desarrollador refleja una estrategia positiva de atención al cliente que contribuye a fortalecer la confianza y fidelización de los usuarios. Este comportamiento representa una ventaja competitiva muy importante, especialmente en aplicaciones relacionadas con seguridad y gestión de información sensible (al ser nosotros una app de gestión de contraseñas).

En resumen, los insights obtenidos pueden aplicarse para priorizar mejoras funcionales, optimizar la experiencia del usuario y diseñar estrategias de comunicación más efectivas. Asimismo, el uso de análisis de datos de reseñas demuestra ser una herramienta valiosa para la toma de decisiones basada en evidencia, permitiendo a las organizaciones identificar necesidades reales del mercado y responder de manera más eficiente a las expectativas de sus clientes.