

Análisis de sentimiento del conflicto Israel-Palestina en Reddit mediante técnicas de procesamiento de lenguaje natural

Paula Andrea Estupiñan Diaz

Deiver Enrique Alfonso Cortes

Asesor

Sebastián Vélez Jaramillo

Universidad Nacional Abierta y a Distancia UNAD

Escuela de Ciencias Básicas, Tecnología e Ingeniería ECBTI

Especialización en Ciencia de Datos y Analítica

2025

Resumen

Las plataformas de redes sociales han transformado la manera en que las personas expresan sus opiniones y sentimientos sobre acontecimientos actuales. Reddit, en particular, se ha consolidado como un espacio activo de discusión, donde los usuarios comparten perspectivas sobre temas de relevancia global, incluido el conflicto entre Israel y Palestina. Las conversaciones generadas en torno a este tema reflejan una amplia variedad de sentimientos que fluctúan según el contexto. Sin embargo, la magnitud, velocidad y complejidad del contenido dificultan su análisis manual, lo que hace necesario recurrir a técnicas automatizadas y escalables.

Este estudio aplica técnicas de Procesamiento de Lenguaje Natural (PLN) y aprendizaje automático para analizar más de 2.4 millones de comentarios en inglés publicados en Reddit sobre dicho conflicto. El proceso metodológico incluyó un análisis exploratorio inicial (EDA) y un preprocesamiento exhaustivo de los textos, mediante una clase personalizada que permitió limpiar y normalizar los datos para su posterior análisis. A partir de este, se implementaron y compararon tres modelos preentrenados de clasificación de sentimientos: DistilBERT, T5 y RoBERTa.

Los resultados permitieron identificar sentimientos predominantes como enojo, tristeza, alegría, miedo y neutralidad, así como analizar su evolución temporal y detectar picos significativos asociados a eventos clave del conflicto. La comparación entre modelos evidenció diferencias en la forma de interpretar el discurso digital, aportando perspectivas complementarias sobre la conversación colectiva en línea.

Palabras claves: Reddit, conflicto Israel-Palestina, análisis de sentimientos, procesamiento de lenguaje natural, modelos preentrenados.

Abstract

Social media platforms have transformed the way people express their opinions and emotions regarding current events. Reddit, in particular, has become an active space for discussion, where users share perspectives on globally relevant issues, including the conflict between Israel and Palestine. The conversations surrounding this topic reflect a wide range of emotions that fluctuate depending on the context. However, the scale, speed, and complexity of the content make manual analysis challenging, requiring the use of automated and scalable techniques.

This study applies Natural Language Processing (NLP) and machine learning techniques to analyze over 2.4 million English-language comments posted on Reddit about the conflict. The methodological process included an initial exploratory data analysis (EDA) and comprehensive text preprocessing using a custom class designed to clean and normalize the data for subsequent analysis. Based on this, three pre-trained sentiment classification models—DistilBERT, T5, and RoBERTa—were implemented and compared.

The results allowed for the identification of predominant emotions such as anger, sadness, joy, fear, and neutrality, as well as the analysis of their temporal evolution and the detection of significant peaks associated with key conflict events. The comparison between models revealed differences in the interpretation of digital discourse, offering complementary perspectives on the collective online conversation.

Keywords: Reddit, Israel-Palestine conflict, sentiment analysis, natural language processing, pre-trained models.

Tabla de Contenido

Introducción	10
Planteamiento del Problema	12
Justificación	13
Objetivos	14
Objetivo General	14
Objetivos Específicos	14
Marco Conceptual	15
Marco Teórico	17
Metodología	21
Recursos Computacionales	22
Preprocesamiento Textual	23
Análisis y Preprocesamiento	23
Limpieza de Datos	26
Clasificación de Sentimientos Mediante Modelos Preentrenados	29
Modelos Implementados	29
Clasificación por Lotes con Aceleración GPU	29
Procedimiento Implementado	30
Agrupación y Análisis Temporal de Sentimientos	31
Agrupación Temporal	32
Visualización de Sentimientos por Período	32

Detección de Picos de Sentimientos	33
Resultados	35
Análisis Exploratorio Inicial (EDA)	35
Preprocesamiento Textual	36
Clasificación de Sentimientos	36
Análisis Temporal de la Evolución de Sentimientos	40
Evolución de sentimientos – Modelo DistilBERT	40
Evolución de sentimientos – modelo T5	42
Evolución de sentimientos – Modelo RoBERTa	43
Comparación de la Evolución de Sentimiento Entre Modelos	45
Proporciones de Sentimiento a lo Largo del Tiempo	45
Proporciones de Sentimiento – Modelo DistilBERT	46
Proporciones de Sentimiento – Modelo T5	47
Proporciones de Sentimiento – Modelo RoBERTa	49
Sentimiento Dominante por Periodo	52
Modelo DistilBERT	52
Modelo T5	54
Modelo RoBERTa	55
Comparación de Sentimiento Dominante entre Modelos	57
Eventos Asociados a los Picos de Sentimiento	58
Conclusiones	61

Recomendaciones 63

Referencias 65

Lista de Tablas

Tabla 1 <i>Tecnologías Utilizadas en el Proyecto</i>	22
Tabla 2 <i>Recursos Computacionales Utilizados para el Análisis</i>	23
Tabla 3 <i>Variables del Conjunto de Datos Analizado</i>	24
Tabla 4 <i>Ejemplos de Entradas de Texto</i>	27
Tabla 5 <i>Etapas del Proceso de Limpieza Textual</i>	27
Tabla 6 <i>Variables Seleccionadas para el Análisis de Sentimientos</i>	35
Tabla 7 <i>Comparación entre el Texto Original y el Texto Limpio</i>	36
Tabla 8 <i>Conteo de Sentimientos por Modelo de Clasificación</i>	39
Tabla 9 <i>Comparación de Proporciones Promedio Mensuales de Sentimientos por Modelo</i> . .	50
Tabla 10 <i>Comparación de Sentimientos Dominantes por Modelo y Periodo</i>	57

Lista de Figuras

Figura 1 <i>Representación Estructural de Componentes</i>	20
Figura 2 <i>Flujo Metodológico del Análisis de Sentimientos</i>	21
Figura 3 <i>Diagrama de Flujo del Proceso de Clasificación de Sentimientos</i>	31
Figura 4 <i>Conteo Total de Sentimientos Clasificados por el Modelo DistilBERT</i>	37
Figura 5 <i>Conteo Total de Sentimientos Clasificados por el Modelo T5</i>	38
Figura 6 <i>Conteo Total de Sentimientos Clasificados por el Modelo RoBERTa</i>	38
Figura 7 <i>Evolución Diaria de Sentimientos Clasificados por DistilBERT</i>	40
Figura 8 <i>Evolución Semanal de Sentimientos Clasificados por DistilBERT</i>	41
Figura 9 <i>Evolución Mensual de Sentimientos Clasificados por DistilBERT</i>	41
Figura 10 <i>Evolución Diaria de Sentimientos Clasificados por T5</i>	42
Figura 11 <i>Evolución Semanal de Sentimientos Clasificados por T5</i>	42
Figura 12 <i>Evolución Mensual de Sentimientos Clasificados por T5</i>	43
Figura 13 <i>Evolución Diaria de Sentimientos Clasificados por RoBERTa</i>	43
Figura 14 <i>Evolución Semanal de Sentimientos Clasificados por RoBERTa</i>	44
Figura 15 <i>Evolución Mensual de Sentimientos Clasificados por RoBERTa</i>	44
Figura 16 <i>Proporción Diaria de Sentimiento Según DistilBERT</i>	46
Figura 17 <i>Proporción Semanal de Sentimiento Según DistilBERT</i>	46
Figura 18 <i>Proporción Mensual de Sentimiento Según DistilBERT</i>	47
Figura 19 <i>Proporción Diaria de Sentimiento Según T5</i>	47
Figura 20 <i>Proporción Semanal de Sentimiento Según T5</i>	48

Figura 21 <i>Proporción Mensual de Sentimiento Según T5</i>	48
Figura 22 <i>Proporción Diaria de Sentimiento Según RoBERTa</i>	49
Figura 23 <i>Proporción Semanal de Sentimiento Según RoBERTa</i>	49
Figura 24 <i>Proporción Mensual de Sentimiento Según RoBERTa</i>	50
Figura 25 <i>Distribución Diaria de Sentimientos Dominantes Clasificados por DistilBERT</i> . .	52
Figura 26 <i>Distribución Semanal de Sentimientos Dominantes Clasificados por DistilBERT</i> .	53
Figura 27 <i>Distribución Mensual de Sentimientos Dominantes Clasificados por DistilBERT</i> .	53
Figura 28 <i>Distribución Diaria de Sentimientos Dominantes Clasificados por T5</i>	54
Figura 29 <i>Distribución Semanal de Sentimientos Dominantes Clasificados por T5</i>	54
Figura 30 <i>Distribución Mensual de Sentimientos Dominantes Clasificados por T5</i>	55
Figura 31 <i>Distribución Diaria de Sentimientos Dominantes Clasificados por RoBERTa</i> . . .	55
Figura 32 <i>Distribución Semanal de Sentimientos Dominantes Clasificados por RoBERTa</i> . .	56
Figura 33 <i>Distribución Mensual de Sentimientos Dominantes Clasificados por RoBERTa</i> . .	56
Figura 34 <i>Evolución Mensual de Sentimientos y Picos Detectados por el Modelo T5</i>	58
Figura 35 <i>Evolución Mensual de Sentimientos y Picos Detectados por el Modelo DistilBERT</i>	59
Figura 36 <i>Evolución Mensual de Sentimientos y Picos Detectados por el Modelo RoBERTa</i> .	59

Introducción

El análisis automatizado de sentimientos en plataformas digitales ha adquirido gran relevancia debido a su capacidad para reflejar la percepción pública sobre acontecimientos sociales, políticos y culturales con alto impacto. Reddit, plataforma social estructurada en comunidades temáticas (subreddits), es un entorno ideal para explorar cómo se manifiestan y evolucionan los sentimientos colectivos en torno a temas específicos como el conflicto entre Israel y Palestina. La identificación sistemática y objetiva de estos sentimientos facilita una mejor comprensión de la opinión pública global y permite monitorear en tiempo real el impacto emocional de eventos específicos, brindando así una interpretación más profunda de los discursos digitales (Costola et al., 2023; Reyes & Moreno, 2024).

Diversos estudios han demostrado la utilidad y efectividad del análisis automatizado de sentimientos en contextos similares. Por ejemplo, Maldonado Ramones (2022) evidenció cómo la aplicación de técnicas avanzadas de Procesamiento de Lenguaje Natural (PLN) logra categorizar con precisión sentimientos expresados en plataformas como Twitter. Asimismo, investigaciones como las realizadas por Regal et al. (2019) han utilizado modelos basados en aprendizaje profundo, específicamente redes neuronales recurrentes LSTM, para analizar sentimientos y predecir comportamientos relacionados con eventos sociales relevantes, confirmando la robustez de estos métodos para entender fenómenos sociales complejos desde una perspectiva emocional.

Ante este contexto, el presente estudio propone aplicar técnicas de PLN y aprendizaje automático para analizar los sentimientos expresados en comentarios sobre el conflicto Israel-Palestina en Reddit, con el fin de identificar cuáles son los sentimientos predominantes en

estas discusiones y estudiar su evolución a través del tiempo.

El desarrollo de este proyecto, por tanto, no solo contribuirá al ámbito académico fortaleciendo competencias técnicas en ciencia de datos y análisis de textos, sino que también aportará al ámbito social, proporcionando perspectivas críticas y valiosas sobre la dinámica de los sentimientos en debates digitales frente a conflictos internacionales.

Planteamiento del Problema

El conflicto Israel-Palestina es uno de los temas más debatidos en redes sociales debido al impacto que ha tenido a nivel mundial, las noticias publicadas y las diferentes posturas que existen sobre este tema. Reddit, una plataforma digital ampliamente utilizada para la discusión abierta de eventos globales, alberga diariamente miles de comentarios que reflejan sentimientos como ira, tristeza, miedo, frustración y empatía en reacción directa a eventos del conflicto.

Debido a la gran cantidad y a la rápida generación de estos comentarios, resulta muy complicado llevar a cabo un análisis manual que permita identificar, clasificar y comprender de forma objetiva y sistemática la dinámica de los sentimientos de estas interacciones digitales. Esta dificultad no solo limita la capacidad de obtener una visión clara y oportuna del estado emocional colectivo en respuesta a eventos específicos, sino que también restringe el potencial uso académico y social de esta información para investigaciones más profundas sobre el impacto emocional de conflictos internacionales.

En la actualidad, Reddit carece de mecanismos internos especializados capaces de capturar automáticamente estos sentimientos y analizar su evolución de manera eficiente. Esta brecha metodológica impide a investigadores, analistas y actores interesados aprovechar adecuadamente la información generada, lo que representa un desafío significativo en términos de análisis social digital y comprensión pública de conflictos geopolíticos complejos.

Justificación

Comprender los sentimientos expresados por los usuarios en redes sociales frente a eventos de gran discusión es clave para interpretar la percepción pública, la polarización social y la forma en que se construyen narrativas digitales en distintos contextos. En el caso del conflicto entre Israel y Palestina, Reddit concentra un alto volumen de opiniones cargadas de sentimientos diversos, cuyo análisis manual resulta impracticable debido a su complejidad y escala.

Esta situación plantea la necesidad de soluciones metodológicas automatizadas que permitan interpretar de forma objetiva y a gran escala los sentimientos presentes en la conversación digital. En este contexto, el uso de técnicas de Procesamiento de Lenguaje Natural (PLN) y modelos de aprendizaje automático ofrece una alternativa eficaz. Estas herramientas permiten clasificar los sentimientos expresados en comentarios, identificar su evolución a lo largo del tiempo y reconocer patrones emocionales vinculados a eventos clave, superando así las limitaciones de otros tipos de análisis.

Adicional al aporte técnico, este proyecto contribuye a una comprensión más profunda del comportamiento social en entornos virtuales, ofreciendo insumos relevantes para otras áreas de estudio. Asimismo, abre camino a la posible creación de sistemas de monitoreo emocional capaces de detectar reacciones colectivas frente a acontecimientos sensibles, lo cual puede ser útil en escenarios de prevención de conflictos, respuesta comunicacional y gestión del discurso público.

Al combinar ciencia de datos, análisis textual y comprensión sociopolítica, esta investigación no solo aborda una problemática relevante, sino que también aporta a la consolidación de una ciencia de datos aplicada a fenómenos humanos.

Objetivos

Objetivo General

Analizar los sentimientos expresados en comentarios de Reddit sobre el conflicto Israel-Palestina mediante técnicas de procesamiento de lenguaje natural, para identificar sentimientos predominantes y su evolución temporal en la opinión pública digital.

Objetivos Específicos

Preparar los datos recopilados de comentarios de Reddit, mediante técnicas de limpieza y preprocesamiento textual para asegurar la calidad del análisis de sentimiento.

Clasificar los comentarios previamente procesados en categorías de sentimiento definidas, utilizando técnicas de Procesamiento de Lenguaje Natural (PLN) con modelos ya construidos, con el propósito de cuantificar los sentimientos predominantes en la conversación.

Analizar la evolución temporal de los sentimientos predominantes identificados en los comentarios, utilizando series de tiempo, para comprender cambios significativos relacionados con eventos específicos del conflicto.

Comparar los resultados obtenidos por los modelos en la clasificación de sentimientos, con el fin de identificar similitudes, diferencias y complementariedades en su detección.

Marco Conceptual

El análisis de sentimientos en redes sociales se ha consolidado como una herramienta clave para comprender la relación entre la opinión pública y diversos fenómenos sociales, económicos o políticos. En este proyecto se abordan los siguientes conceptos fundamentales:

Sentimientos en Redes Sociales

Hace referencia a las opiniones, emociones y percepciones expresadas por los usuarios en plataformas como Twitter, Reddit o Facebook. Reyes y Moreno (2024) destacan que el análisis de sentimientos, apoyado en técnicas de procesamiento de lenguaje natural (PLN) y aprendizaje automático (machine learning), permite identificar patrones dentro de grandes volúmenes de datos no estructurados, así como clasificar con precisión las opiniones emitidas por los usuarios.

Procesamiento de Lenguaje Natural (PLN)

Es un subcampo de la inteligencia artificial orientado al análisis y generación de lenguaje humano. Herramientas como TextBlob y VADER han demostrado ser útiles para la clasificación básica de sentimientos en textos breves, como los publicados en redes sociales. No obstante, técnicas más avanzadas, como los modelos basados en BERT (Bidirectional Encoder Representations from Transformers), han mejorado considerablemente la interpretación contextual del lenguaje, logrando mayor precisión en tareas de análisis emocional (Costola et al., 2023).

Minería de Datos en Redes Sociales

La minería de datos permite extraer patrones y tendencias relevantes a partir de grandes volúmenes de información. En el contexto de las redes sociales, esta técnica es esencial para

convertir datos no estructurados en hallazgos útiles, como la identificación de palabras clave, temas recurrentes y sentimientos predominantes. Montiel Flores y Valenzuela Valenzuela (2023) señalan que el uso de herramientas como RapidMiner facilita tanto la interpretación como la clasificación de datos masivos, mejorando así la capacidad de análisis en tiempo real.

Marco Teórico

El análisis de sentimientos se ha convertido en una herramienta clave para comprender cómo las personas expresan sus sentimientos en entornos digitales, como las redes sociales, especialmente ante situaciones de alto impacto social o político. En este contexto, Reddit destaca por su estructura basada en comunidades temáticas (*subreddits*), que permite organizar las conversaciones según temas específicos. Esta organización favorece la participación de los usuarios y facilita discusiones más argumentadas y contextualizadas que las que suelen encontrarse en otras redes (Baumgartner et al., 2020; González-González et al., 2022).

Diversos estudios han demostrado que Reddit representa una fuente valiosa para investigaciones sociales y políticas, dado que los comentarios en esta red contienen mayor contexto y riqueza expresiva que los mensajes breves de otras plataformas. Esta característica permite identificar con mayor precisión tanto la emocionalidad como la intención de los usuarios al interactuar en conversaciones públicas (Medvedev et al., 2019). Asimismo, el anonimato que ofrece Reddit promueve una expresión más sincera de opiniones y sentimientos, aspecto especialmente útil al abordar temas sensibles como los conflictos internacionales. Un ejemplo de ello se encuentra en investigaciones recientes que han analizado sentimientos como miedo y esperanza durante el conflicto entre Rusia y Ucrania utilizando datos de Reddit (Guerra & Karakuş, 2023).

Para identificar estas emociones de manera automatizada se recurre al procesamiento de lenguaje natural (PLN), un campo interdisciplinar que combina informática, lingüística e inteligencia artificial. Herramientas tradicionales como VADER o TextBlob permiten clasificar

textos en categorías generales como positivo, negativo o neutral. Sin embargo, presentan limitaciones ante textos complejos con sarcasmo, ambigüedad o dobles sentidos (Viteri, 2021).

Con el desarrollo del aprendizaje profundo se han creado modelos más sofisticados capaces de interpretar con mayor precisión el lenguaje humano. Uno de los más influyentes es BERT (Bidirectional Encoder Representations from Transformers), desarrollado por Google, que introdujo un enfoque bidireccional para entender el contexto completo de una oración. BERT ha logrado avances significativos en tareas como la clasificación de texto y el análisis de sentimientos (Devlin et al., 2019). Posteriormente, se crearon versiones mejoradas como RoBERTa, que optimizó el proceso de entrenamiento y logró resultados aún más precisos (Liu et al., 2019).

En el presente estudio se emplean tres modelos preentrenados disponibles en la plataforma Hugging Face, cada uno ajustado para tareas de clasificación emocional con enfoques arquitectónicos diversos:

- Distilbert-base-uncased-emotion (DistilBERT): desarrollado por Bhadrish Savani, es una versión compacta y eficiente de BERT entrenada con el dataset GoEmotions de Google. Este modelo reconoce seis sentimientos principales: alegría, tristeza, enojo, miedo, amor y sorpresa. Destaca por su equilibrio entre precisión y bajo requerimiento computacional, lo que lo hace adecuado para procesar grandes volúmenes de datos textuales (Savani, 2024).
- T5-base-finetuned-emotion (T5): basado en la arquitectura T5 (Text-to-Text Transfer Transformer), este modelo fue afinado por Manuel Romero para tareas de clasificación de sentimientos tratados como problemas de generación de texto. Su enfoque permite una mayor flexibilidad y capacidad de generalización, particularmente útil en contextos ambiguos o no

estructurados como Reddit (Mendoza, 2024).

- Emotion-english-distilroberta-base (RoBERTa): creado por Jan Hartmann, se fundamenta en una versión distilada del modelo RoBERTa y ha sido entrenado con el conjunto de datos de sentimientos de SemEval2018. Su objetivo es identificar una gama más amplia de matices emocionales con alta precisión en inglés, siendo especialmente efectivo para analizar conversaciones naturales y complejas (Hartmann, 2022).

El uso complementario de estos tres modelos permite realizar un análisis comparativo y robusto de los sentimientos expresados en los comentarios de Reddit sobre el conflicto Israel-Palestina, mejorando así la confiabilidad y profundidad del estudio.

Además del desarrollo de modelos preentrenados para la clasificación de sentimientos, el diseño de estructuras computacionales modulares permite escalar y adaptar el análisis de sentimientos a distintos contextos y volúmenes de datos. En este sentido, es posible implementar componentes orientados a objetos que encapsulan las funciones principales de los modelos emocionales dentro de una interfaz común, favoreciendo su reutilización y extensibilidad.

Estos componentes suelen compartir atributos clave, lo que permite una configuración flexible en distintos entornos. A través de métodos estandarizados, como la predicción de emociones a partir de un texto o la aplicación sobre columnas completas de datos estructurados, se puede automatizar y escalar la tarea de detección emocional en textos de forma eficiente.

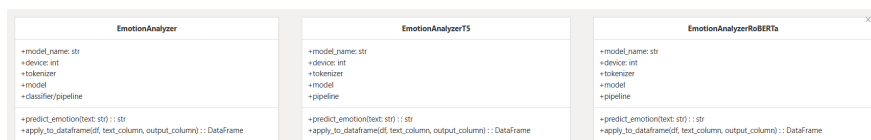
La existencia de estas estructuras modulares favorece la interoperabilidad entre modelos, facilita la comparación de resultados y promueve buenas prácticas en el diseño de sistemas aplicados al procesamiento de lenguaje natural. Este enfoque modular ha sido ampliamente

promovido por plataformas como Hugging Face, que integran modelos mediante pipelines reutilizables y fácilmente integrables en distintas aplicaciones (Wolf et al., 2020).

Este enfoque puede representarse mediante estructuras orientadas a objetos, donde cada componente actúa como una envoltura del modelo subyacente. La Figura 1 ilustra un ejemplo general de cómo podrían organizarse estos elementos para encapsular modelos como DistilBERT, T5 y RoBERTa, facilitando así su aplicación escalable y flexible en tareas de análisis de sentimientos.

Figura 1

Representación Estructural de Componentes



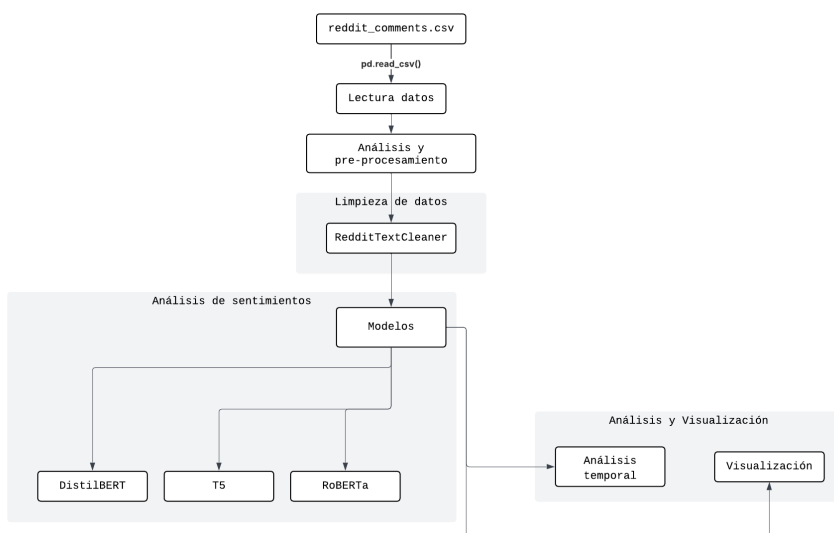
Nota. Cada componente encapsula un modelo orientado a la predicción de sentimientos basado en modelos preentrenados específicos, e implementa métodos comunes que facilitan su aplicación de manera estandarizada en tareas de procesamiento de lenguaje natural.

Metodología

La metodología implementada en este proyecto para el análisis de sentimientos relacionados con el conflicto Israel-Palestina en Reddit se estructuró en una serie de etapas secuenciales, abarcando desde la recolección inicial del conjunto de datos hasta el análisis detallado de los resultados. Cada fase fue realizada para garantizar la calidad, validez y reproducibilidad del proceso analítico. La Figura 2 presenta una visión general del flujo metodológico seguido en esta investigación.

Figura 2

Flujo Metodológico del Análisis de Sentimientos



Nota. Se presenta una visualización esquemática de las etapas implementadas en el proyecto, desde la recolección de datos hasta el análisis de sentimientos de los comentarios.

Recursos Computacionales

Para el desarrollo del análisis de sentimientos, se implementó un conjunto de herramientas tecnológicas especializadas en procesamiento de lenguaje natural, visualización de datos y despliegue de resultados. La Tabla 1 resume las tecnologías empleadas y su propósito específico dentro del proyecto.

Asimismo, el procesamiento de datos a gran escala requirió el uso de recursos computacionales adecuados para garantizar eficiencia y tiempos de ejecución razonables. La Tabla 2 detalla la configuración del equipo utilizado para el entrenamiento y ejecución de los modelos.

Tabla 1

Tecnologías Utilizadas en el Proyecto

Tecnología	Propósito
Python	Lenguaje principal de programación
Pandas	Manipulación y procesamiento de datos
Transformers (Hugging Face)	Modelos de PLN preentrenados
Matplotlib / Seaborn	Visualizaciones estáticas
Streamlit	Panel interactivo de visualización

Nota. Herramientas empleadas para el desarrollo, procesamiento, visualización e interacción en el análisis de sentimientos.

Tabla 2*Recursos Computacionales Utilizados para el Análisis*

Componente	Especificación
Procesador (CPU)	AMD Ryzen 5 3600 (6 núcleos, 12 hilos, 3.6 GHz)
Memoria RAM	32 GB DDR4
Tarjeta gráfica (GPU)	Nvidia GeForce GTX 1650 (4 GB VRAM)
Almacenamiento	Disco SSD
Sistema Operativo	Windows 10 Pro

Nota. Configuración del equipo utilizada para el preprocesamiento, clasificación de sentimientos y análisis de datos masivos.

Preprocesamiento Textual*Análisis y Preprocesamiento*

El proceso de adquisición de datos constituye el punto de partida del flujo metodológico desarrollado para el análisis de sentimientos de comentarios en Reddit. Esta etapa se ejecutó antes de cualquier procedimiento de limpieza textual o aplicación de técnicas automatizadas de clasificación de sentimientos.

En este caso se utilizó el conjunto de datos *Daily Public Opinion on Israel-Palestine War*, disponible en la plataforma Kaggle.¹ Con el fin de automatizar la descarga y facilitar su

¹ <https://www.kaggle.com/datasets/asaniczka/reddit-on-israel-palestine-daily-updated/data>

integración al flujo de procesamiento, se desarrolló una clase en Python denominada `KaggleDatasetHandler`. Esta clase permite autenticar y descargar datasets desde Kaggle mediante su API oficial, asegurando una adquisición estructurada, reproducible y eficiente del conjunto de datos.

El conjunto de datos contiene aproximadamente 2.459.013 comentarios publicados en Reddit entre octubre de 2023 y marzo de 2025, distribuidos en 24 columnas con información detallada sobre los comentarios, sus autores y las publicaciones asociadas. La Tabla 3 presenta la descripción de las variables contenidas en el conjunto de datos original.

Tabla 3

Variables del Conjunto de Datos Analizado

Variable	Descripción
<code>comment_id</code>	ID único del comentario
<code>score</code>	Puntuación total del comentario
<code>self_text</code>	Texto del comentario
<code>subreddit</code>	Comunidad donde fue publicado
<code>created_time</code>	Fecha y hora del comentario
<code>post_id</code>	ID único de la publicación
<code>author_name</code>	Nombre del autor
<code>controversiality</code>	Indicador de controversia
<code>ups</code>	Votos positivos

downs	Votos negativos
user_is_verified	Usuario verificado
user_account_created_time	Fecha de cuenta
user_awardee_karma	Karma por otorgar premios
user_awarder_karma	Karma por recibir premios
user_link_karma	Karma por compartir enlaces
user_comment_karma	Karma en comentarios
user_total_karma	Karma total
post_score	Puntuación de la publicación
post_self_text	Texto de la publicación
post_title	Título de la publicación
post_upvote_ratio	Proporción de votos positivos
post_thumbs_ups	Votos positivos obtenidos
post_total_awards_received	Premios recibidos
post_created_time	Fecha de publicación

Nota. Variables extraídas directamente del dataset original de Kaggle, utilizadas para el análisis exploratorio inicial.

Con el análisis exploratorio del conjunto de datos realizado y seleccionadas las variables relevantes, se procedió a la estructuración formal del proyecto. Con este fin, se creó un repositorio

en GitHub², para centralizar el desarrollo, facilitar el seguimiento del progreso, habilitar el trabajo colaborativo y gestionar de manera eficiente las versiones y actualizaciones del código fuente.

Limpieza de Datos

Se avanzó hacia la preparación del conjunto de datos. Para ello, se generó una muestra aleatoria de 1000 registros a partir del dataset original, con el propósito de probar y ajustar el proceso de limpieza textual antes de aplicarlo a la totalidad de los comentarios. Esta muestra permitió validar la efectividad de las transformaciones implementadas y realizar ajustes en la lógica de preprocesamiento sin comprometer recursos computacionales excesivos.

Para llevar a cabo esta tarea, se diseñó una clase en Python denominada `RedditTextCleaner`, orientada al tratamiento automatizado de textos provenientes de Reddit. Esta clase implementa un conjunto de funciones de preprocesamiento que permiten transformar los comentarios originales, frecuentemente cargados de elementos irrelevantes o ruido textual, en entradas limpias, coherentes y listas para el análisis de sentimientos.

Los comentarios en Reddit suelen contener una amplia variedad de elementos no estructurados, como menciones a usuarios, enlaces, emojis, hashtags, puntuación innecesaria y caracteres no imprimibles. Estos componentes dificultan el análisis automático, ya que introducen ruido y distorsionan los resultados obtenidos por los modelos de clasificación. Por esta razón, la etapa de limpieza textual es crítica dentro del flujo metodológico.

La Tabla 4 presenta ejemplos representativos de comentarios antes de aplicar el proceso de limpieza textual. En ellos se identifican elementos como símbolos, URLs, emojis, menciones y

² <https://github.com/dalfonso75/reddit-emotion-analysis-israel-palestine.git>

cadena vacías, los cuales no aportan valor semántico al análisis y deben ser eliminados durante el preprocesamiento.

Tabla 4

Ejemplos de Entradas de Texto

Nº	Input Text	Descripción
1	@User Hello World! 😊 #hashtag	Texto con mención, URL, emoji y hashtag
2	12345 \$\$\$ % % hello!!!	Texto con números, símbolos y palabras
3	(empty string)	Cadena vacía
4	None	Valor nulo

Nota. Ejemplos representativos de entradas textuales con ruido que fueron procesadas durante la limpieza.

La clase `RedditTextCleaner` se diseñó con un enfoque modular y reutilizable. Su método principal, `clean()`, toma un texto como entrada y retorna una versión preprocesada. Para mejorar la eficiencia, se implementaron expresiones regulares precompiladas y tablas de traducción que permiten filtrar rápidamente los caracteres no deseados.

La Tabla 5 resume las transformaciones aplicadas durante este proceso.

Tabla 5

Etapas del Proceso de Limpieza Textual

Transformación	Descripción
----------------	-------------

Eliminación de URLs	Se eliminan enlaces web para evitar sesgos semánticos.
Eliminación de menciones	Se retiran referencias a otros usuarios como @usuario.
Eliminación de hashtags	Se eliminan etiquetas, aunque son poco comunes en Reddit.
Eliminación de emojis	Se remueven usando la librería emoji.
Filtrado de puntuación	Se eliminan todos los signos, excepto ?, ! y ' por su valor expresivo.
Normalización de espacios	Se unifican y limpian los espacios en blanco.
Conversión a minúsculas	Se transforma todo el texto para evitar duplicidades.
Eliminación de caracteres no imprimibles	Se garantiza compatibilidad con bibliotecas de modelado.

Nota. Transformaciones implementadas por la clase `RedditTextCleaner` durante el preprocesamiento textual.

Una vez validado el funcionamiento de esta clase sobre la muestra de prueba, se procedió a aplicar el proceso de limpieza al conjunto completo de datos. Esta operación se ejecutó en bloque sobre los más de 2.4 millones de comentarios disponibles, generando una nueva columna denominada `clean_text` que almacena la versión normalizada de cada comentario.

Clasificación de Sentimientos Mediante Modelos Preentrenados

Luego del preprocesamiento y limpieza del conjunto de datos, se procedió a la etapa de clasificación de los sentimientos. Esta fase tuvo como objetivo identificar los sentimientos predominantes expresados por los usuarios en comentarios relacionados con el conflicto Israel-Palestina. Para ello, se utilizaron modelos preentrenados basados en arquitecturas tipo transformer, disponibles en la plataforma Hugging Face.

Modelos Implementados

El sistema implementa tres modelos distintos, cada uno encapsulado en una clase Python especializada, permitiendo comparar arquitecturas y enfoques para la predicción de sentimientos:

- **EmotionAnalyzerDistilBERT**: modelo principal basado en bhadresh-savani/distilbert-base-uncased-emotion. clasifica los sentimientos: enojo, alegría, tristeza, miedo, sorpresa, amor, desconocido.
- **EmotionAnalyzerT5**: utiliza el modelo mrm8488/t5-base-finetuned-emotion. Predice: neutral, enojo, sorpresa, miedo, alegría, tristeza, repulsión, desconocido.
- **EmotionAnalyzerRoBERTa**: basado en el modelo j-hartmann/emotion-english-distilroberta-base. clasifica: alegría, tristeza, enojo, miedo, sorpresa, disgusto y neutralidad.

Clasificación por Lotes con Aceleración GPU

Para dar inicio al procesamiento del conjunto de datos y dado el tamaño de este (más de 2.4 millones de comentarios), se implementó un sistema de procesamiento por lotes con aceleración por GPU, lo que permitió realizar la inferencia de sentimientos de manera eficiente y

en tiempos razonables. El sistema divide el conjunto de datos en fragmentos manejables (chunks), aplica el modelo sobre cada uno utilizando la GPU disponible y luego recompone los resultados en un solo archivo consolidado.

Esta estrategia resolvió tres desafíos clave:

- Limitaciones de memoria RAM y VRAM, ya que procesar todos los comentarios de forma simultánea superaba la capacidad disponible.
- Tiempo de ejecución, el cual se redujo significativamente al aprovechar la paralelización en GPU frente a una ejecución secuencial en CPU.
- Flexibilidad de ejecución, ya que esta técnica permitió ejecutar el proceso en momentos específicos, sin necesidad de mantener la máquina encendida hasta completar la tarea. Era posible detener el procesamiento en cualquier punto y retomarlo desde el último bloque procesado al reanudar la ejecución.

Procedimiento Implementado

La validación inicial del pipeline de clasificación se realizó sobre una muestra de 100.000 comentarios. Una vez verificada su efectividad, se aplicó el modelo principal de forma masiva sobre el conjunto completo. Para cada comentario, el sistema generó una nueva columna denominada `predicted_emotion`, que contiene el sentimientos identificado como predominante por el modelo.

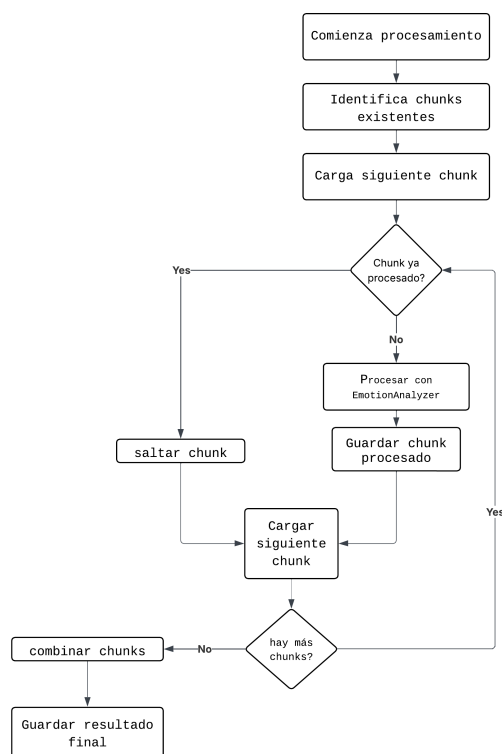
Todo el proceso fue automatizado mediante la clase `EmotionAnalyzer`, que encapsula la carga del modelo, la tokenización del texto, la inferencia y el etiquetado final.

La Figura 3 muestra el flujo general del proceso de clasificación implementado en este

estudio. En él se visualiza cómo el texto preprocesado es transformado mediante modelos preentrenados, aplicando procesamiento por lotes con aceleración GPU para asignar etiquetas de manera eficiente.

Figura 3

Diagrama de Flujo del Proceso de Clasificación de Sentimientos



Nota. Se muestra el flujo general del proceso del proyecto.

Agrupación y Análisis Temporal de Sentimientos

Para estudiar la evolución de los sentimientos a lo largo del tiempo, se diseñó un procedimiento específico que permite agrupar y analizar los registros según distintas escalas temporales. Este enfoque busca identificar tendencias emocionales globales, así como variaciones

significativas en momentos específicos del conflicto Israel-Palestina.

Agrupación Temporal

Se implementó una función denominada `process_emotion_counts_by_period`, cuyo objetivo es transformar un conjunto de comentarios clasificados por sentimiento en una estructura agregada por periodos temporales. Esta función toma como entrada un DataFrame con dos columnas clave: `created_time` (fecha del comentario) y `predicted_emotion` (sentimiento asignado).

Su procedimiento interno consiste en:

- Convertir la columna `created_time` al tipo `datetime` para habilitar operaciones temporales.
- Generar una nueva columna llamada `period`, que agrupa las fechas en intervalos estándar definidos como diario D, semanal W o mensual M).
- Aplicar un conteo agrupado por sentimiento y periodo mediante `groupby + size()`.
- Despivotar los resultados con `unstack()` para obtener una tabla donde las filas representan los periodos y las columnas los sentimientos.

Este proceso permitió generar una matriz tiempo-sentimiento, base fundamental para las visualizaciones y análisis posteriores. La transformación es crucial para detectar tendencias, cambios estructurales y correlaciones temporales entre sentimientos y eventos reales.

Visualización de Sentimientos por Período

A partir de los datos agregados, se generaron distintos tipos de visualizaciones para cada modelo y escala temporal. Estas permitieron explorar tanto la intensidad como la composición relativa de los sentimientos a lo largo del tiempo:

- **Conteo absoluto de sentimientos:** Se graficaron las curvas temporales de cada sentimiento utilizando líneas de color diferenciadas.
- **Proporción relativa por periodo:** Se normalizaron los conteos por fila, mostrando qué porcentaje de comentarios expresaba cada sentimiento en cada periodo.
- **Sentimiento dominante:** Se determinó el sentimiento más frecuente por periodo utilizando `idxmax()` y se visualizó mediante puntos de color sobre una línea temporal.

Estos gráficos permitieron analizar la evolución de los sentimientos en la conversación en Reddit, observando tanto aumentos sostenidos como variaciones abruptas en respuesta a eventos del contexto internacional.

Detección de Picos de Sentimientos

Con el fin de identificar eventos puntuales que generaron respuestas de sentimientos inusuales, se diseñó una función denominada `detectar_picos_multiples`, basada en el algoritmo `find_peaks` de la biblioteca `scipy`. Esta función se aplicó sobre los conteos mensuales de cada sentimiento y realiza las siguientes operaciones:

- Itera por cada columna del DataFrame, correspondiente a un sentimiento.
- Aplica `find_peaks` sobre la serie de tiempo, estableciendo un valor mínimo de prominencia (por defecto, 10.000 menciones) para considerar solo picos significativos.
- Registra la fecha, valor del pico y sentimiento correspondiente.
- Devuelve una lista de tuplas que luego se convierte en un nuevo DataFrame para su análisis.

Este procedimiento permitió identificar momentos clave en los que un sentimiento tuvo un

comportamiento anómalo, proporcionando indicios de posibles eventos del mundo real que provocaron esas respuestas colectivas en la plataforma.

Resultados

Análisis Exploratorio Inicial (EDA)

Como parte del proceso de preprocesamiento, se tomó la decisión de generar un subconjunto del conjunto de datos original, con el fin de optimizar los tiempos de procesamiento y análisis. Para ello, se seleccionaron únicamente las columnas consideradas esenciales para el objetivo de clasificación de sentimiento y análisis temporal. La Tabla 6 presenta las variables que fueron conservadas en este dataset reducido.

Tabla 6

Variables Seleccionadas para el Análisis de Sentimientos

Variable	Justificación de selección
comment_id	Evita duplicados
self_text	Entrada principal para la clasificación de sentimientos
created_time	Permite análisis temporal
subreddit	Permite comparación entre comunidades
score	Representa interacción del usuario
post_title	Contextualiza la publicación

Nota. Subconjunto de variables utilizadas específicamente para clasificación de sentimientos y análisis de series de tiempo.

Preprocesamiento Textual

Como resultado del proceso de limpieza, se obtuvo una nueva versión de cada comentario en la columna `clean_text`, la cual contiene el texto sin ruido estructural ni elementos no lingüísticos. Esta transformación permitió garantizar la uniformidad, legibilidad y compatibilidad del contenido con los modelos de clasificación de sentimientos utilizados en las etapas posteriores.

La Tabla 7 presenta algunos ejemplos representativos de comentarios procesados, mostrando cómo se eliminaron correctamente menciones, enlaces, emojis, hashtags, caracteres especiales y otros elementos que no aportan valor semántico al análisis.

Tabla 7

Comparación entre el Texto Original y el Texto Limpio

Nº	Texto original (<code>input_text</code>)	Texto limpio (<code>clean_text</code>)
1	@User Hello World! 😊 #hashtag	hello world visit
2	12345 \$\$\$ % % hello!!!	12345 hello
3	(empty string)	(empty string)
4	None	(empty string)

Nota. Comparación entre los comentarios originales y sus versiones limpias generadas mediante preprocesamiento textual.

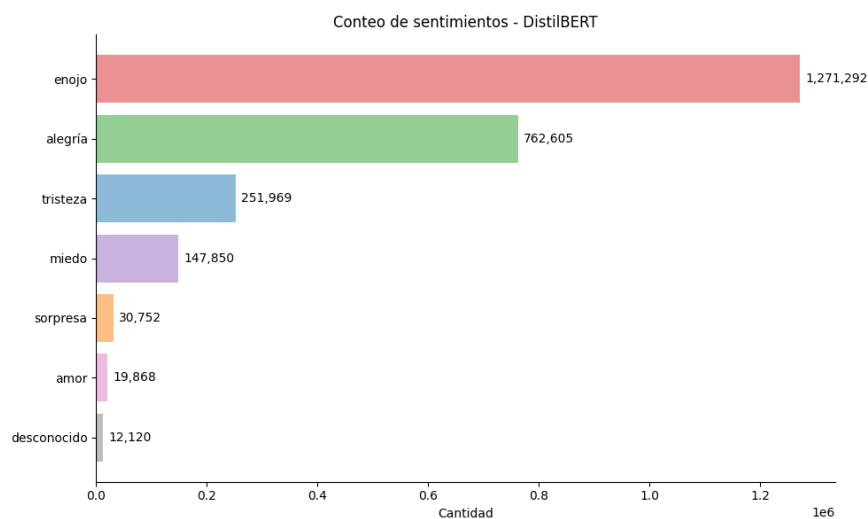
Clasificación de Sentimientos

Una vez aplicada la inferencia de sentimientos sobre los comentarios de Reddit, se obtuvieron las distribuciones globales de sentimientos para cada modelo. En esta etapa se

presentan los conteos absolutos de sentimientos detectados por los modelos DistilBERT, T5 y RoBERTa, con el fin de identificar diferencias en la sensibilidad y comportamiento clasificatorio de cada arquitectura.

Figura 4

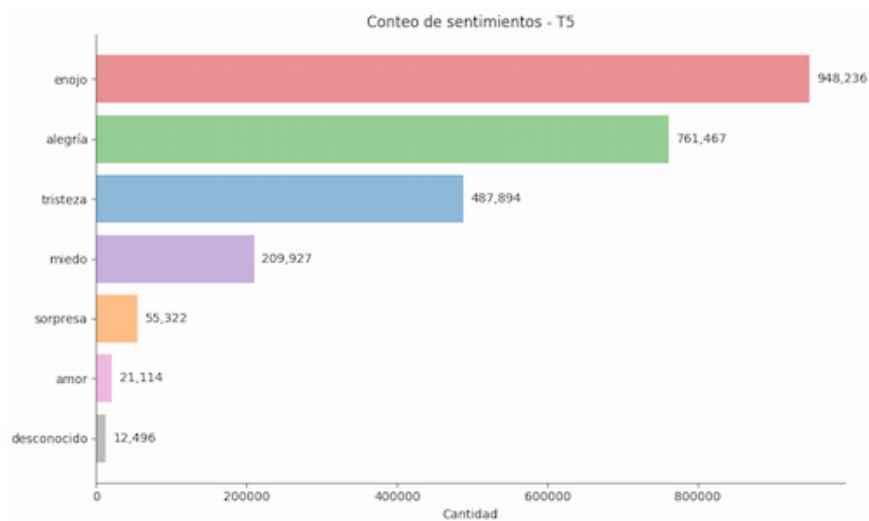
Conteo Total de Sentimientos Clasificados por el Modelo DistilBERT



Nota. El modelo DistilBERT clasifica mayoritariamente sentimientos negativos, siendo enojo y tristeza las más frecuentes, con una notable presencia también de alegría.

Figura 5

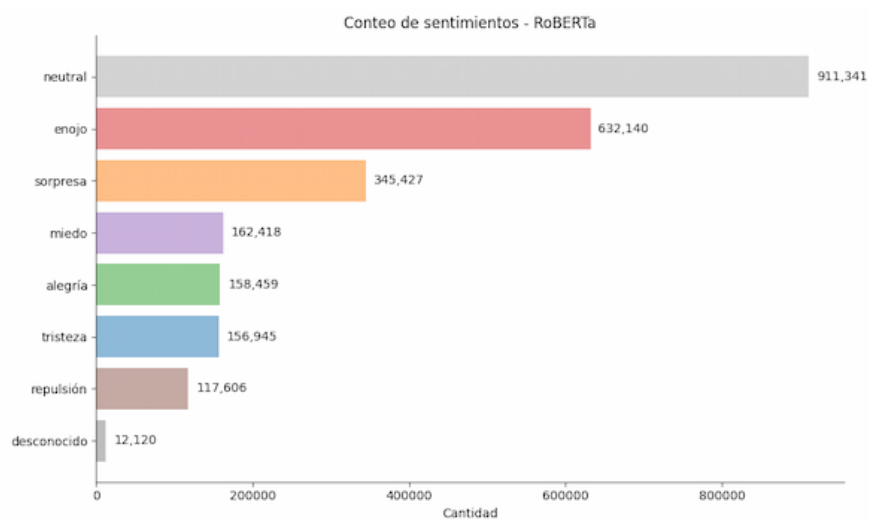
Conteo Total de Sentimientos Clasificados por el Modelo T5



Nota. El modelo T5 presenta una distribución más balanceada.

Figura 6

Conteo Total de Sentimientos Clasificados por el Modelo RoBERTa



Nota. RoBERTa introduce el sentimiento neutral como el más frecuente.

Tabla 8*Conteo de Sentimientos por Modelo de Clasificación*

Sentimiento	DistilBERT	T5	RoBERTa
Alegría	762,605	761,467	158,459
Amor	19,868	21,114	–
Desconocido	12,120	12,496	12,120
Enojo	1,271,292	948,236	632,140
Miedo	147,850	209,927	162,418
Neutral	–	–	911,341
Repulsión	–	–	117,606
Sorpresa	30,752	55,322	345,427
Tristeza	251,969	487,894	156,945

Nota. Todos los modelos coinciden en identificar al enojo como el sentimiento predominante.

DistilBERT muestra una fuerte inclinación hacia sentimientos intensos como enojo y tristeza, mientras que T5 distribuye mejor las emociones, manteniendo niveles altos de alegría. RoBERTa, por su parte, introduce el sentimiento neutral como dominante, lo cual refleja una diferencia en el esquema de etiquetas.

Estas diferencias permiten una comprensión más amplia del espectro de sentimientos presente en los comentarios y sientan la base para el análisis temporal detallado presentado en la siguiente sección.

Análisis Temporal de la Evolución de Sentimientos

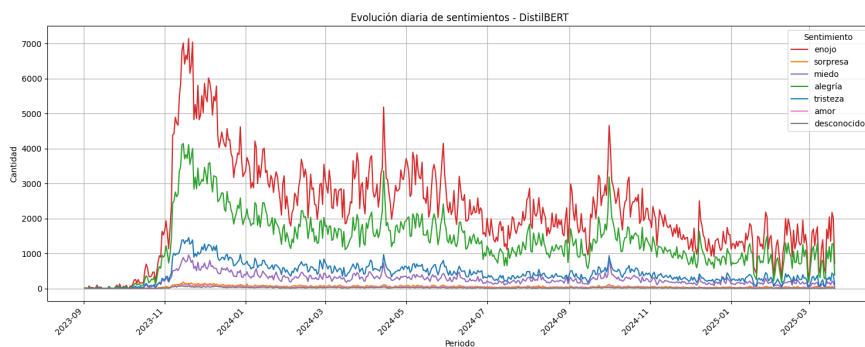
El análisis temporal permite observar cómo varían los sentimientos a lo largo del tiempo, revelando patrones, picos de intensidad emocional y posibles relaciones con eventos específicos del conflicto Israel-Palestina. Se evaluaron tres escalas temporales: diaria, semanal y mensual, para cada uno de los modelos utilizados.

Evolución de sentimientos – Modelo DistilBERT

Evolución Diaria.

Figura 7

Evolución Diaria de Sentimientos Clasificados por DistilBERT

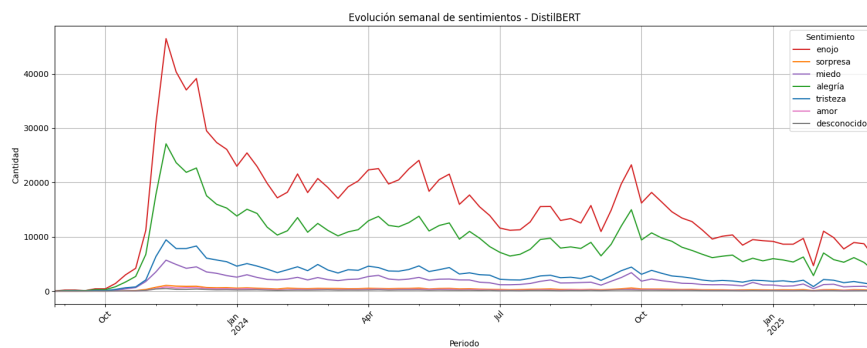


Nota. Se observa alta variabilidad en los sentimientos, con predominancia de enojo y alegría.

Evolución Semanal.

Figura 8

Evolución Semanal de Sentimientos Clasificados por DistilBERT

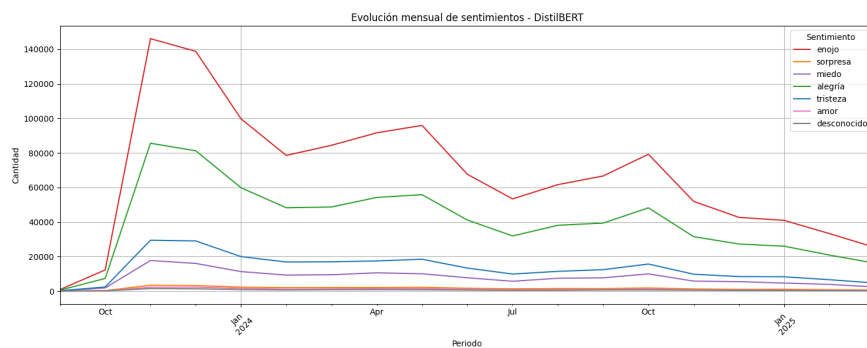


Nota. La agregación semanal permite suavizar los picos de sentimientos.

Evolución Mensual.

Figura 9

Evolución Mensual de Sentimientos Clasificados por DistilBERT



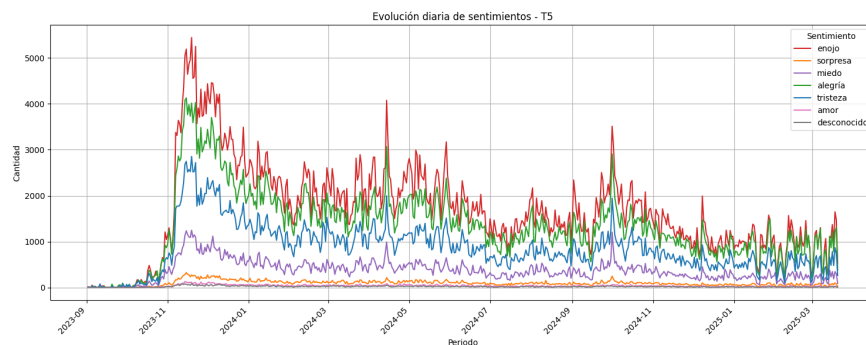
Nota. Disminución general de todos los sentimientos, con una persistente prevalencia del enojo.

Evolución de sentimientos – modelo T5

Evolución diaria.

Figura 10

Evolución Diaria de Sentimientos Clasificados por T5

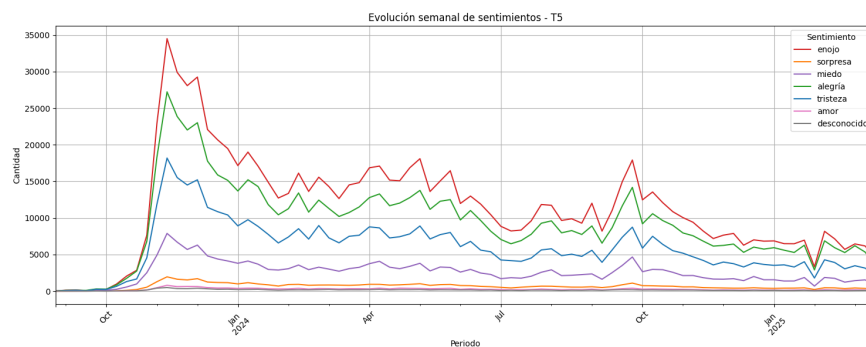


Nota. El modelo T5 muestra una distribución más balanceada entre enojo, alegría y tristeza.

Evolución Semanal.

Figura 11

Evolución Semanal de Sentimientos Clasificados por T5

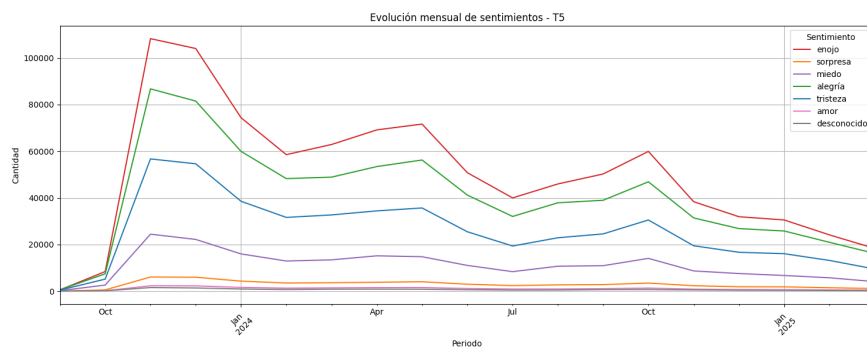


Nota. Alternancia entre enojo y alegría como sentimientos dominantes durante gran parte del año.

Evolución Mensual.

Figura 12

Evolución Mensual de Sentimientos Clasificados por T5



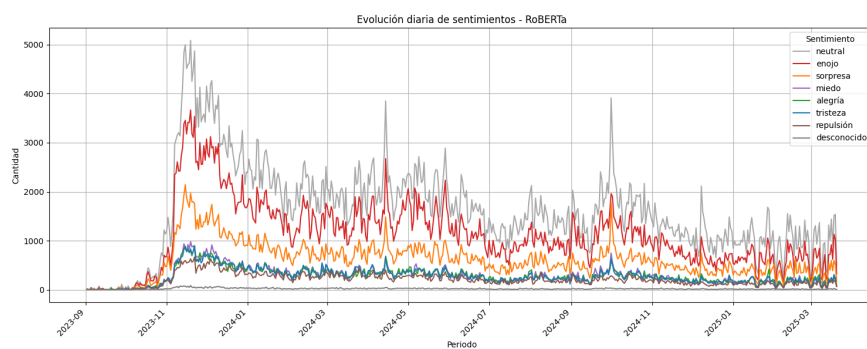
Nota. A nivel mensual, el sentimiento enojo fue predominante en la mayoría de los periodos.

Evolución de sentimientos – Modelo RoBERTa

Evolución Diaria.

Figura 13

Evolución Diaria de Sentimientos Clasificados por RoBERTa

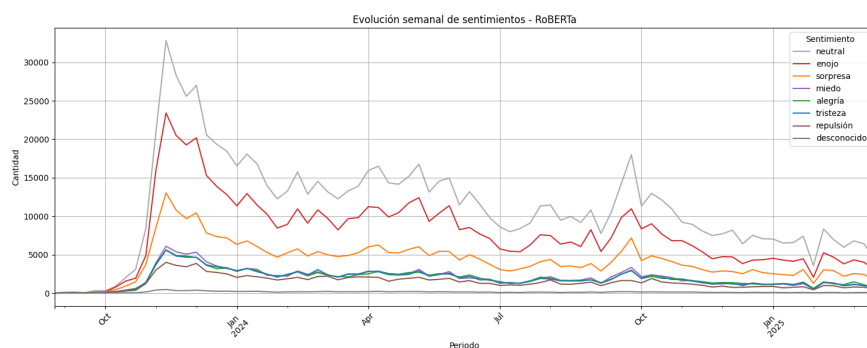


Nota. El sentimiento neutral domina ampliamente en este modelo, seguida por enojo y sorpresa.

Evolución Semanal.

Figura 14

Evolución Semanal de Sentimientos Clasificados por RoBERTa

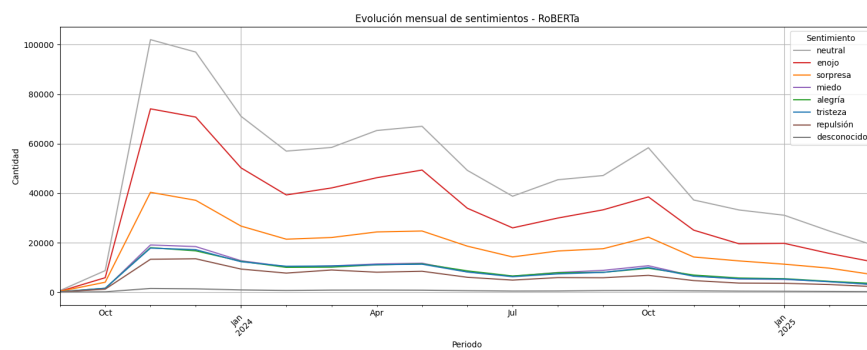


Nota. La evolución semanal revela un descenso paulatino del sentimiento neutral, mientras que enojo y sorpresa mantienen una presencia destacada.

Evolución Mensual.

Figura 15

Evolución Mensual de Sentimientos Clasificados por RoBERTa



Nota. En la escala mensual, se observa una persistente predominancia del sentimiento neutral, lo que diferencia a RoBERTa de los otros modelos.

Comparación de la Evolución de Sentimiento Entre Modelos

Resumen Comparativo. La evolución de sentimientos identificada mediante los modelos DistilBERT, T5 y RoBERTa presenta tanto patrones comunes como diferencias específicas en la distribución temporal de los sentimientos.

- En los tres modelos se observa un aumento abrupto de los sentimientos al inicio del periodo analizado, seguido por una disminución progresiva en los meses posteriores.
- DistilBERT evidencia una frecuencia elevada de sentimientos como enojo y alegría, con notoria variabilidad en el análisis diario y estabilidad relativa en los agregados semanales y mensuales.
- T5 muestra un comportamiento más equilibrado entre enojo, tristeza y alegría, sin una dominancia extrema de un solo sentimiento.
- RoBERTa se diferencia al presentar el sentimiento neutral como la más frecuente a lo largo del tiempo, así como una mayor proporción de sorpresa y repulsión, sentimientos que no están presentes en los otros modelos.

Las diferencias observadas pueden atribuirse a las etiquetas disponibles en cada modelo, así como a sus respectivos métodos de entrenamiento. En las siguientes secciones se continúa con el análisis de proporciones sentimentales y sentimientos dominantes, con el fin de profundizar en la evolución del sentimiento expresado en la conversación digital.

Proporciones de Sentimiento a lo Largo del Tiempo

A diferencia del conteo absoluto, la proporción de sentimiento permite analizar el peso relativo de cada sentimiento respecto al total de comentarios emitidos en un periodo determinado.

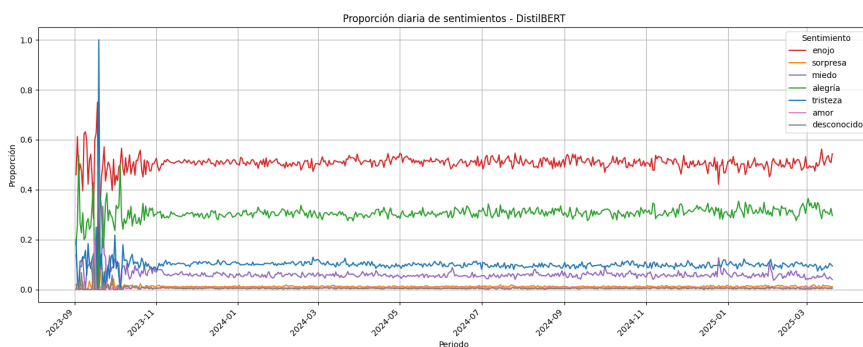
Esto resulta útil para interpretar cambios independientemente del volumen de comentarios.

Proporciones de Sentimiento – Modelo DistilBERT

Proporción Diaria.

Figura 16

Proporción Diaria de Sentimiento Según DistilBERT

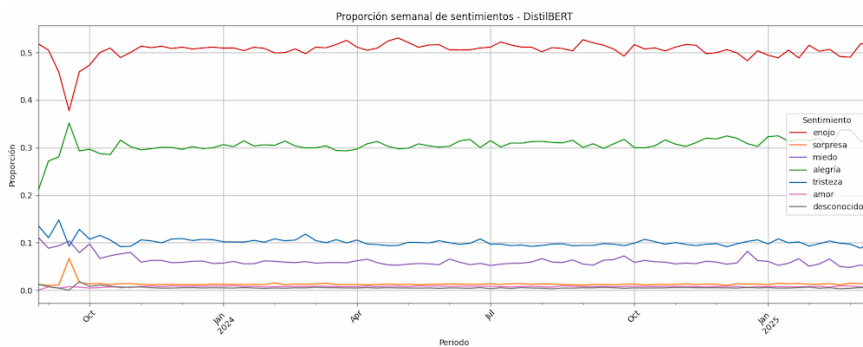


Nota. El sentimiento enojo mantiene una proporción elevada durante todo el periodo.

Proporción Semanal.

Figura 17

Proporción Semanal de Sentimiento Según DistilBERT

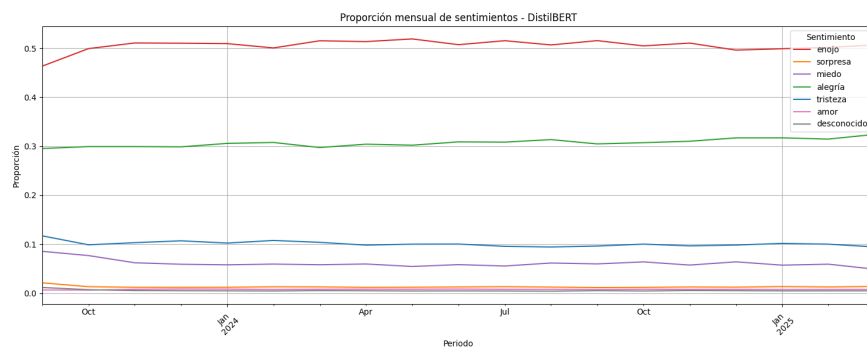


Nota. La proporción de enojo supera de forma consistente a los demás sentimientos.

Proporción Mensual.

Figura 18

Proporción Mensual de Sentimiento Según DistilBERT



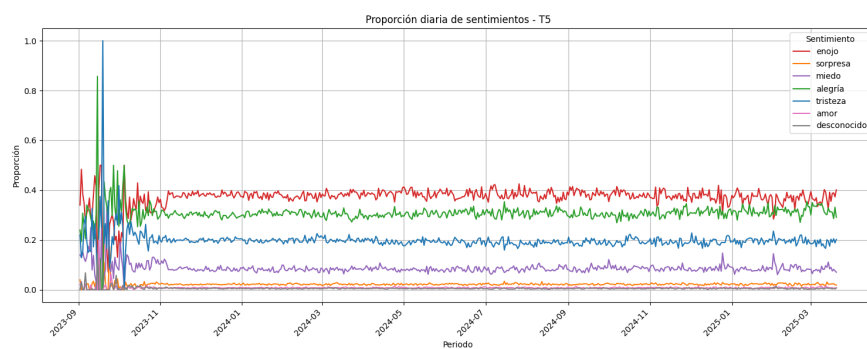
Nota. El sentimiento enojo representa cerca del 50 % de los sentimientos clasificados.

Proporciones de Sentimiento – Modelo T5

Proporción Diaria.

Figura 19

Proporción Diaria de Sentimiento Según T5

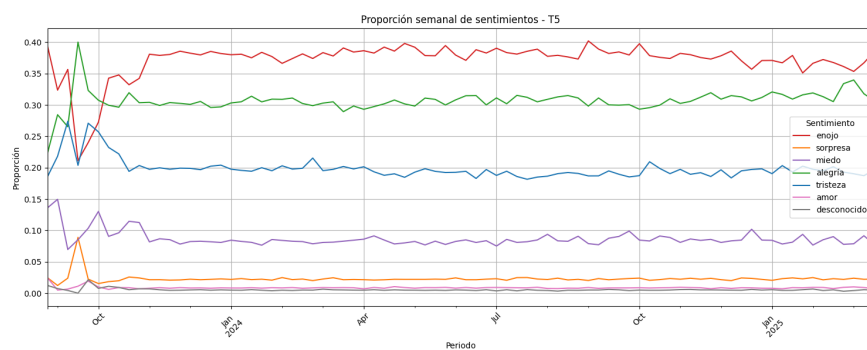


Nota. Se destacan tres sentimientos principales: enojo, alegría y tristeza.

Proporción Semanal.

Figura 20

Proporción Semanal de Sentimiento Según T5

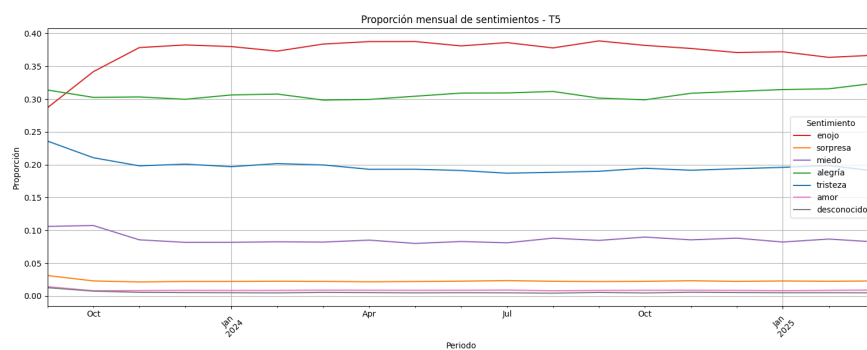


Nota. El sentimiento enojo alcanza proporciones superiores al 35 % en varios tramos del periodo.

Proporción Mensual.

Figura 21

Proporción Mensual de Sentimiento Según T5



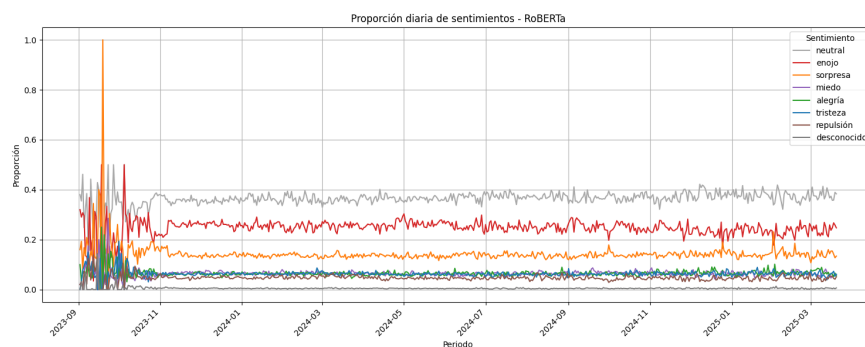
Nota. Se mantiene una proporción elevada de enojo.

Proporciones de Sentimiento – Modelo RoBERTa

Proporción Diaria.

Figura 22

Proporción Diaria de Sentimiento Según RoBERTa

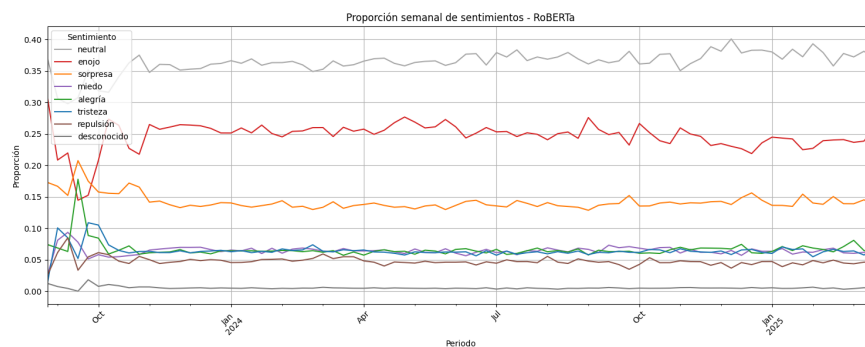


Nota. El sentimiento neutral domina en proporción diaria durante todo el periodo.

Proporción Semanal.

Figura 23

Proporción Semanal de Sentimiento Según RoBERTa

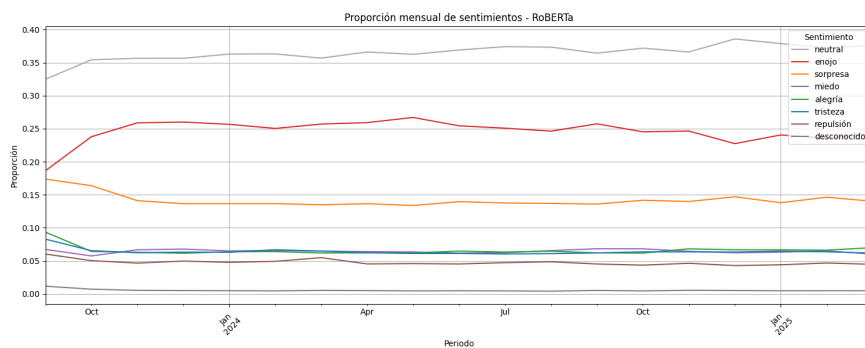


Nota. Se observa una estabilidad en la proporción del sentimiento neutral.

Proporción Mensual.

Figura 24

Proporción Mensual de Sentimiento Según RoBERTa



Nota. RoBERTa mantiene una proporción mensual superior al 35 % para el sentimiento neutral.

Tabla 9

Comparación de Proporciones Promedio Mensuales de Sentimientos por Modelo

Sentimiento	DistilBERT	T5	RoBERTa
Enojo	0.5054	0.3719	0.2463
Alegría	0.3070	0.3073	0.0657
Tristeza	0.1007	0.1974	0.0640
Miedo	0.0609	0.0864	0.0643
Sorpresa	0.0129	0.0228	0.1418
Amor	0.0077	0.0087	—
Repulsión	—	—	0.0475
Neutral	—	—	0.3650

Desconocido	0.0053	0.0055	0.0053
-------------	--------	--------	--------

Nota. DistilBERT presenta un fuerte predominio del enojo, mientras que T5 ofrece una distribución más balanceada, destacando alegría y tristeza. RoBERTa se diferencia por asignar una proporción significativa al sentimiento neutral, y una mayor sorpresa en comparación con los demás.

Comparación de Proporciones Entre Modelos. La Tabla 9 presenta las proporciones promedio mensuales de sentimiento clasificados por cada modelo. Los resultados permiten observar diferencias relevantes en la distribución relativa de los sentimientos:

- DistilBERT asigna una alta proporción al sentimiento enojo (50.5 %), seguida por alegría (30.7 %) y tristeza (10.1 %). Sentimientos como sorpresa, amor y desconocido tienen menor representación.
- T5 muestra una distribución más equilibrada entre enojo (37.2 %), alegría (30.7 %) y tristeza (19.7 %), manteniendo proporciones más altas para miedo y sorpresa en comparación con DistilBERT.
- RoBERTa se distingue por la inclusión del sentimiento neutral como dominante (36.5 %), seguida de enojo (24.6 %) y sorpresa (14.2 %). Además, incluye repulsión, que no está presente en los otros modelos.

Estas diferencias pueden atribuirse a las etiquetas disponibles en el entrenamiento de cada modelo y a sus respectivas arquitecturas, lo que refuerza la complementariedad de enfoques en el análisis de sentimiento.

Sentimiento Dominante por Periodo

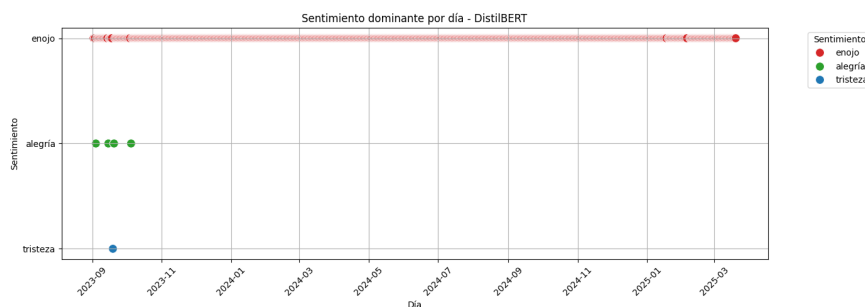
El análisis del sentimiento dominante permite identificar, para cada periodo temporal (día, semana y mes), el sentimiento con mayor frecuencia relativa. Esta visualización resulta útil para detectar cambios abruptos o patrones sostenidos en la opinión pública digital a lo largo del tiempo. En esta sección se presentan los resultados de cada modelo por separado.

Modelo DistilBERT

Sentimiento Dominante Diario.

Figura 25

Distribución Diaria de Sentimientos Dominantes Clasificados por DistilBERT

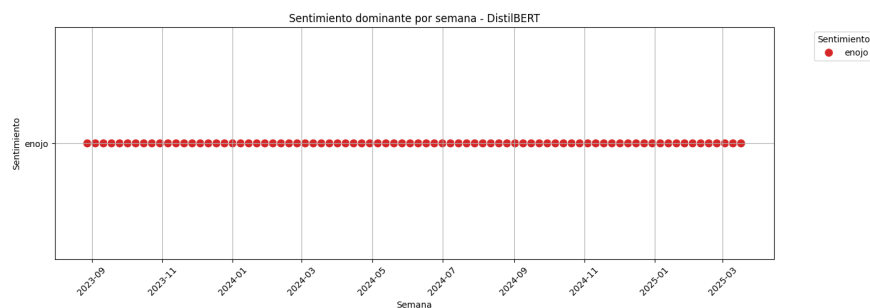


Nota. Se observa una fuerte presencia de enojo como sentimiento dominante durante la mayor parte del periodo. Algunos días aislados muestran predominancia de alegría y tristeza.

Sentimiento Dominante Semanal.

Figura 26

Distribución Semanal de Sentimientos Dominantes Clasificados por DistilBERT

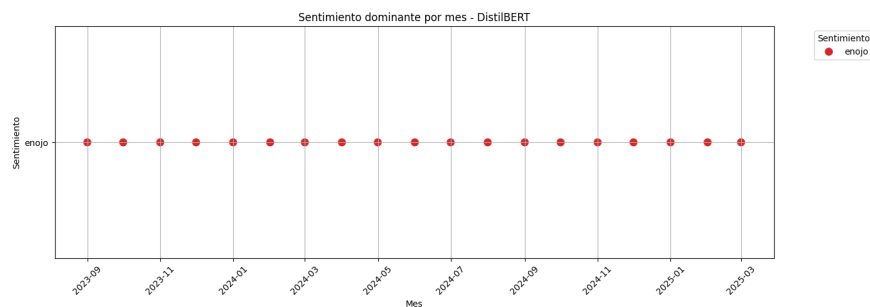


Nota. El modelo mantiene una consistencia en el sentimiento enojo como predominante.

Sentimiento Dominante Mensual.

Figura 27

Distribución Mensual de Sentimientos Dominantes Clasificados por DistilBERT



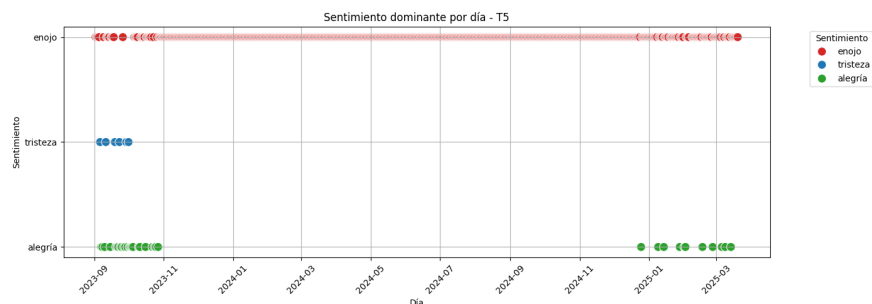
Nota. En todos los meses analizados, el sentimiento enojo se mantiene como el dominante según el modelo DistilBERT.

Modelo T5

Sentimiento Dominante Diario.

Figura 28

Distribución Diaria de Sentimientos Dominantes Clasificados por T5

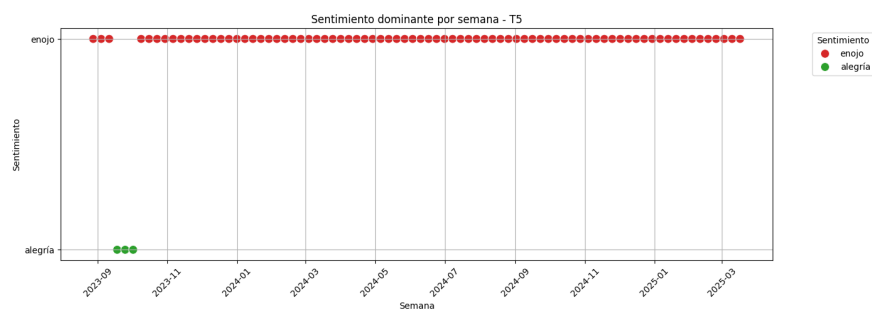


Nota. El modelo T5 identifica una variabilidad más marcada en el sentimiento dominante diario.

Sentimiento Dominante Semanal.

Figura 29

Distribución Semanal de Sentimientos Dominantes Clasificados por T5

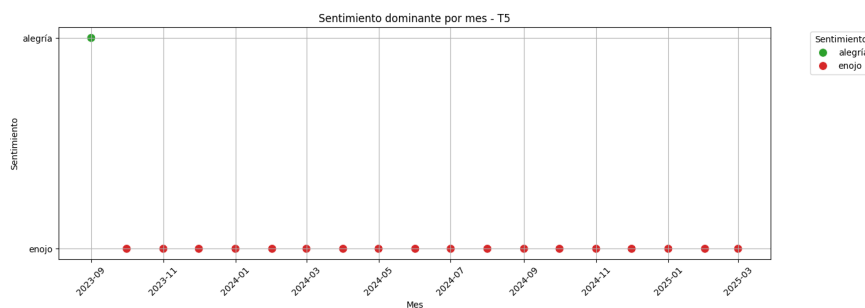


Nota. El sentimiento enojo predomina en la mayoría de las semanas, aunque se presentan algunas semanas con alegría como dominante.

Sentimiento Dominante Mensual.

Figura 30

Distribución Mensual de Sentimientos Dominantes Clasificados por T5



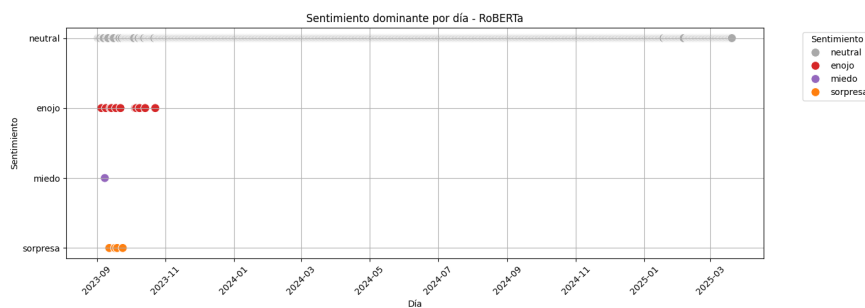
Nota. Se observa predominancia mensual del sentimiento enojo, excepto en un mes donde el dominante es alegría.

Modelo RoBERTa

Sentimiento Dominante Diario.

Figura 31

Distribución Diaria de Sentimientos Dominantes Clasificados por RoBERTa

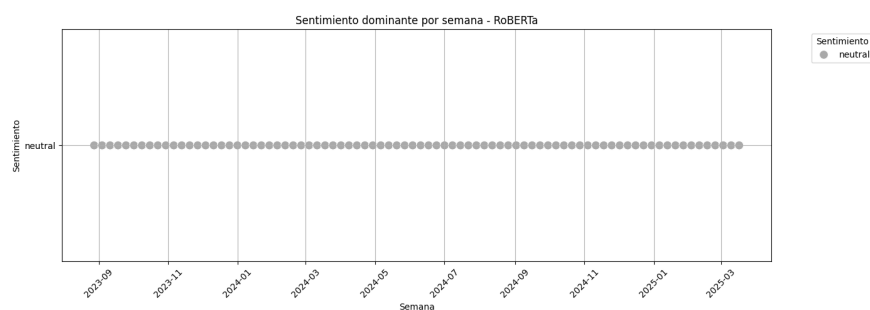


Nota. Aunque neutral es el sentimiento más frecuente, se evidencian días en los que enojo, sorpresa y miedo alcanzan valores dominantes.

Sentimiento Dominante Semanal.

Figura 32

Distribución Semanal de Sentimientos Dominantes Clasificados por RoBERTa

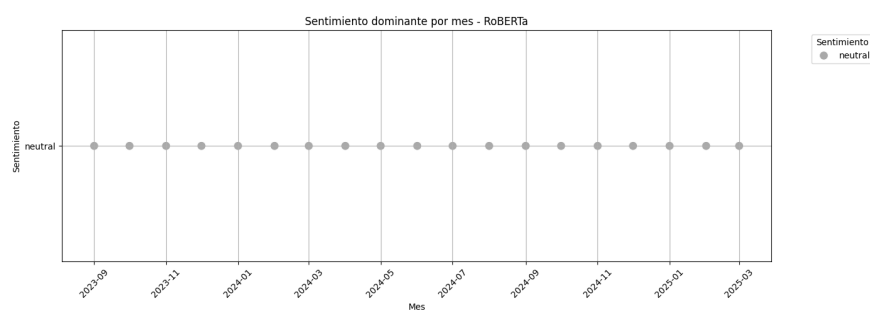


Nota. El sentimiento neutral domina de forma sostenida a lo largo de todas las semanas.

Sentimiento Dominante Mensual.

Figura 33

Distribución Mensual de Sentimientos Dominantes Clasificados por RoBERTa



Nota. RoBERTa clasifica sistemáticamente neutral como sentimiento dominante para todos los meses del periodo estudiado.

Comparación de Sentimiento Dominante entre Modelos

Con el fin de sintetizar los resultados del análisis de sentimientos, se presenta una comparación entre los tres modelos utilizados (DistilBERT, T5 y RoBERTa), identificando cuál fue el sentimiento dominante en cada escala temporal. Esta comparación permite observar consistencias, divergencias y características particulares de cada modelo frente al mismo conjunto de datos.

Tabla 10

Comparación de Sentimientos Dominantes por Modelo y Periodo

Periodo	DistilBERT	T5	RoBERTa
Día	Enojo (predominante, con algunos días de alegría y tristeza)	Enojo (con días de alegría y tristeza)	Neutral (con días de enojo, miedo y sorpresa)
Semana	Enojo (consistente)	Enojo (con semanas de alegría)	Neutral (uniforme)
Mes	Enojo (todos los meses)	Enojo (excepto un mes con alegría)	Neutral (en todo el periodo)

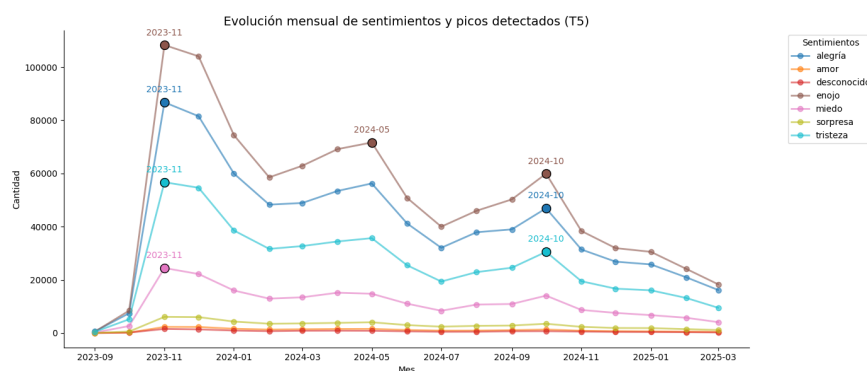
Nota. La tabla resume el sentimiento predominante en distintos niveles temporales. Mientras que DistilBERT y T5 muestran al enojo como dominante en casi todos los periodos, RoBERTa destaca por identificar la neutralidad como sentimiento recurrente.

Eventos Asociados a los Picos de Sentimiento

En las siguientes figuras se presentan los picos mensuales de sentimiento detectados por los modelos T5, DistilBERT y RoBERTa. Las visualizaciones reflejan una mayor intensidad en tres periodos clave: noviembre de 2023, mayo de 2024 y octubre de 2024.

Figura 34

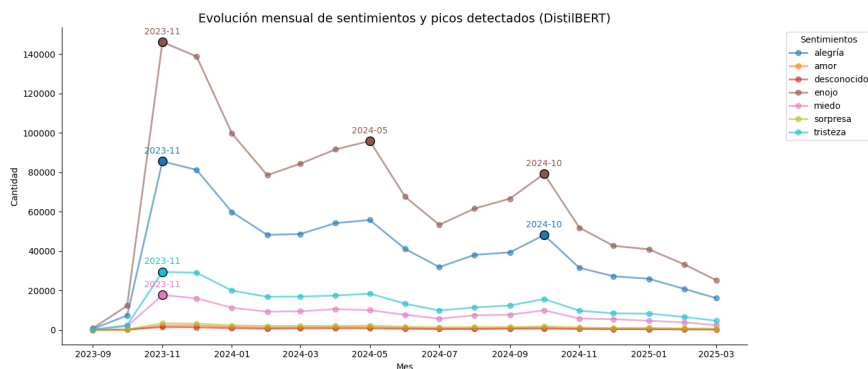
Evolución Mensual de Sentimientos y Picos Detectados por el Modelo T5



Nota. El modelo T5 presenta una clara concentración de sentimientos en noviembre de 2023, con repuntes en mayo y octubre de 2024.

Figura 35

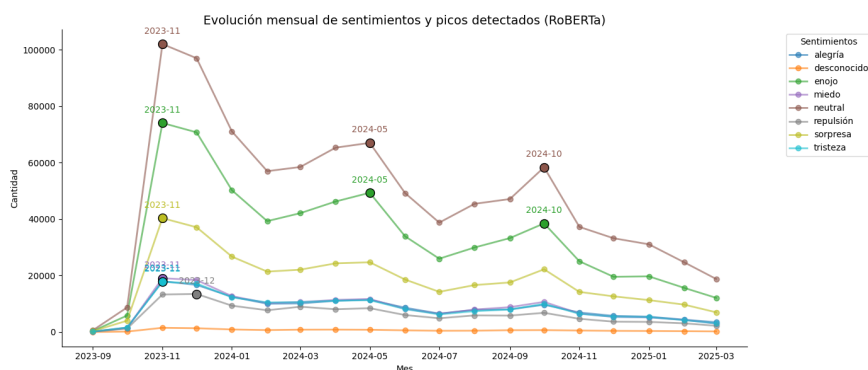
Evolución Mensual de Sentimientos y Picos Detectados por el Modelo DistilBERT



Nota. DistilBERT muestra patrones similares a T5, con énfasis en enojo y tristeza en los mismos meses críticos del conflicto.

Figura 36

Evolución Mensual de Sentimientos y Picos Detectados por el Modelo RoBERTa



Nota. RoBERTa coincide en los picos, con predominancia de neutralidad, enojo y sorpresa como sentimientos destacados.

A continuación se resumen los posibles eventos asociados a los picos de sentimiento

observados en los modelos:

- Noviembre de 2023: Se observa el pico sentimental más alto en los tres modelos. Este incremento coincide con la ofensiva militar intensificada de Israel en Gaza tras el ataque de Hamás del 7 de octubre, así como con los bombardeos a hospitales, campos de refugiados y escuelas (France24, 2023a, 2023b). Predominaron los sentimientos de enojo, tristeza y miedo.

- Mayo de 2024: El segundo pico puede relacionarse con la ofensiva sobre Rafah y los ataques a zonas que habían sido designadas como seguras para la población civil (El País, 2024; France24, 2025). La continuidad del bloqueo y el deterioro de las condiciones humanitarias aumentaron la atención internacional. También se registraron reacciones diplomáticas significativas en la región, como la ruptura de relaciones entre Colombia e Israel (France24, 2024). Se reflejan aumentos en enojo, tristeza y neutralidad.

- Octubre de 2024: Este tercer pico ocurre al cumplirse un año del inicio del conflicto. Coincide con la reanudación de enfrentamientos, el frente abierto con Hezbolá y el anuncio del reconocimiento del Estado palestino por parte de países como España (RTVE, 2024).

Este análisis contextual contribuye a entender mejor la relación entre la evolución de los sentimientos y el desarrollo de eventos claves durante el conflicto, sirviendo como insumo para la discusión posterior.

Conclusiones

El análisis de sentimientos en torno al conflicto Israel-Palestina permitió comprender cómo se manifiestan los sentimientos colectivos en redes sociales ante eventos de alto impacto geopolítico. A través de una estrategia metodológica integral, se abordaron desde el procesamiento inicial de los datos hasta el análisis comparativo entre arquitecturas de modelos de lenguaje, permitiendo una lectura profunda del fenómeno sentimental digital.

La etapa de preparación del conjunto de datos, centrada en la limpieza y normalización de más de 2.4 millones de comentarios de Reddit, fue clave para garantizar que la información estuviera en condiciones óptimas para su análisis automatizado. Esta labor permitió reducir significativamente el ruido textual característico de este tipo de plataformas y transformar una estructura compleja e informal en una base coherente, procesable por modelos de lenguaje.

La clasificación de sentimientos mediante tres modelos preentrenados DistilBERT, T5 y RoBERTa los sentimientos más frecuentes en la conversación digital, y también diferencias relevantes en cómo cada arquitectura interpreta el lenguaje. Se identificaron sentimientos predominantes como enojo, tristeza y alegría, junto con otros menos visibles pero igualmente significativos como repulsión, miedo o neutralidad, lo que aportó una visión más amplia del espectro sentimental presente.

La comparación entre modelos mostró que cada uno ofrece una aproximación particular al fenómeno, lo que permitió descubrir patrones comunes, divergencias y matices que enriquecen la interpretación global. Esta diversidad de salidas evidencia cómo distintas estructuras algorítmicas pueden capturar aspectos complementarios de una misma realidad discursiva, lo que amplía las

posibilidades analíticas en estudios similares.

Finalmente, el análisis temporal de los sentimientos permitió asociar variaciones con momentos específicos del conflicto, revelando reacciones colectivas inmediatas ante eventos como bombardeos, treguas, ofensivas o declaraciones diplomáticas. Estas fluctuaciones mostraron cómo la narrativa sentimental se construye en paralelo con el desarrollo de los acontecimientos, consolidando al análisis de sentimientos como una herramienta útil para explorar las dinámicas sociopolíticas contemporáneas.

En conjunto, este proyecto demuestra que es posible articular herramientas de ciencia de datos, aprendizaje profundo y análisis contextual para interpretar fenómenos humanos complejos a gran escala, abriendo nuevas rutas para el estudio de sentimientos colectivos en entornos digitales.

Recomendaciones

A partir de los resultados obtenidos, se proponen las siguientes recomendaciones para investigaciones futuras y desarrollos más robustos en este campo:

Reentrenamiento de modelos con datos contextualizados. Se recomienda construir un conjunto de datos etiquetado manualmente con comentarios reales sobre el conflicto, a fin de afinar los modelos utilizados. Esto permitiría reducir errores en la clasificación, como la asignación incorrecta del sentimiento alegría en contextos de violencia, lo cual podría obedecer a ambigüedad semántica o ruido en los datos. Un modelo afinado sobre un corpus específico proporcionaría mayor precisión y sensibilidad al contexto.

Integración de modelos complementarios y ensambles. Los tres modelos utilizados muestran diferencias significativas en la clasificación de sentimientos. Se recomienda, por tanto, explorar métodos de ensembling o votación ponderada entre modelos, para combinar sus fortalezas y obtener resultados más equilibrados y confiables. Esta estrategia puede mitigar sesgos individuales y mejorar la capacidad de generalización del análisis.

Ampliación del espectro de sentimientos y su tipificación. Incorporar modelos que manejen un espectro más amplio de sentimientos (por ejemplo, culpa, esperanza, indignación) podría enriquecer la comprensión emocional de la conversación digital.

Uso de servicios en la nube y procesamiento distribuido. Dado el volumen de datos analizado (más de 2.4 millones de comentarios), se recomienda migrar el procesamiento y análisis a entornos escalables en la nube (como Google Cloud, AWS o Azure). Esto facilitaría entrenar modelos propios, manejar datasets aún más grandes y reducir significativamente los tiempos de

ejecución.

Desarrollo de herramientas interactivas de monitoreo. Finalmente, se recomienda implementar dashboards interactivos que permitan monitorear en tiempo real la evolución de sentimientos en redes sociales. Esto podría ser útil para periodistas, investigadores, organizaciones humanitarias o entidades gubernamentales interesadas en identificar cambios en la percepción pública ante eventos clave.

Referencias

- Baumgartner, J., Zannettou, S., Keegan, B., Squire, M., & Blackburn, J. (2020). The Pushshift Reddit Dataset [arXiv preprint arXiv:2001.08435]. <https://arxiv.org/abs/2001.08435>
- Costola, M., Hinz, O., Nofer, M., & Pelizzon, L. (2023). Machine learning sentiment analysis, COVID-19 news and stock market reactions. *Research in International Business and Finance*, 64, 101881. <https://doi.org/10.1016/j.ribaf.2023.101881>
- Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *arXiv preprint arXiv:1810.04805*. <https://arxiv.org/abs/1810.04805>
- El País. (2024). *Israel bombardea Rafah pese a las advertencias de la comunidad internacional*. Consultado el 10 de mayo de 2025, desde <https://elpais.com/internacional/2024-05-26/israel-bombardea-rafah-pese-a-las-advertencias-de-la-comunidad-internacional.html>
- France24. (2023a). *El Consejo de Seguridad de la ONU aprueba pausas humanitarias en Gaza*. Consultado el 10 de mayo de 2025, desde <https://www.france24.com/es/medio-oriente/20231115-%F0%9F%94%B4-en-directo-israel-toma-el-hospital-al-shifa-la-oms-pierde-comunicaci%C3%B3n-con-su-personal-en-el-centro-m%C3%A9dico>
- France24. (2023b). *Hamás reporta al menos 195 muertos en Jabalia; la ONU dice que ataques israelíes podrían ser crímenes de guerra*. Consultado el 10 de mayo de 2025, desde <https://www.france24.com/es/medio-oriente/20231101-%F0%9F%94%B4-en-directo-abre-el-paso-de-rafah-para-la-evacuaci%C3%B3n-limitada-de-palestinos-heridos-y-extranjeros>

- France24. (2024). *Gustavo Petro anuncia que Colombia romperá relaciones diplomáticas con Israel por ofensiva en Gaza*. Consultado el 10 de mayo de 2025, desde <https://www.france24.com/es/am%C3%A9rica-latina/20240501-gustavo-petro-anuncia-que-colombia-romper%C3%A1-relaciones-diplom%C3%A1ticas-con-israel-por-ofensiva-en-gaza>
- France24. (2025). *23 muertos en Gaza por ataques aéreos israelíes y crece preocupación por distribución de ayuda*. Consultado el 10 de mayo de 2025, desde <https://www.france24.com/es/medio-orient/20250510-23-muertos-en-gaza-por-ataques-a%C3%A9reos-israel%C3%ADes-y-crece-preocupaci%C3%B3n-por-distribuci%C3%B3n-de-ayuda>
- González-González, R., Pérez, L., & Salinas, D. (2022). El análisis del discurso emocional en Reddit: una aproximación desde la lingüística computacional. *Revista Latinoamericana de Ciencias Sociales*. <https://revistalatinacs.org/>
- Guerra, J., & Karakuş, M. (2023). Tracking emotions during the Russia-Ukraine conflict on Reddit. *International Journal of Social Media Research*. <https://example.org/guerra2023emotions>
- Hartmann, J. (2022). emotion-english-distilroberta-base [Modelo de PLN] [Hugging Face]. <https://huggingface.co/j-hartmann/emotion-english-distilroberta-base>
- Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., & Stoyanov, V. (2019). RoBERTa: A Robustly Optimized BERT Pretraining Approach. *arXiv preprint arXiv:1907.11692*. <https://arxiv.org/abs/1907.11692>

- Maldonado Ramones, E. S. (2022). *Análisis de sentimientos en la red social Twitter mediante el procesamiento de lenguaje natural* [bachelorThesis]. Universidad Nacional de Chimborazo. <http://dspace.unach.edu.ec/handle/51000/10103>
- Medvedev, A., Lambiotte, R., & Delvenne, J.-C. (2019). Anatomy of Reddit: Understanding the evolution of online communities. *arXiv preprint arXiv:1810.07772*.
<https://arxiv.org/abs/1810.07772>
- Mendoza, M. R. (2024). t5-base-finetuned-emotion [Modelo de PLN] [Hugging Face].
<https://huggingface.co/mrm8488/t5-base-finetuned-emotion>
- Montiel Flores, J. C., & Valenzuela Valenzuela, J. P. (2023). Minería de datos para un análisis del sentimiento: Estudio de caso. *New Trends in Qualitative Research*, 19, 1-13.
<https://doi.org/10.36367/ntqr.19.2023.e927>
- Regal, A., Morzán, J., Fabbri, C., Herrera, G., Yaulli, G., Palomino, A., & Gil, C. (2019). Proyección del precio de criptomonedas basado en Tweets empleando LSTM. *Ingeniare. Revista chilena de ingeniería*, 27(4), 696-706.
<https://doi.org/10.4067/S0718-33052019000400696>
- Reyes, N. S., & Moreno, G. E. T. (2024). Análisis de sentimientos en datos de redes sociales: Aplicación de técnicas de procesamiento de lenguaje natural y machine learning. *Dominio de las Ciencias*, 10(1). <https://doi.org/10.23857/dc.v10i1.3714>
- RTVE. (2024). *España reconoce oficialmente al Estado de Palestina junto con Irlanda y Noruega*. Consultado el 10 de mayo de 2025, desde <https://www.rtve.es/noticias/20240528/espana-reconoce-oficialmente-estado-palestina-junto-irlanda-noruega/2451234.shtml>

Savani, B. (2024). distilbert-base-uncased-emotion [Modelo de PLN] [Hugging Face].

<https://huggingface.co/bhadresh-savani/distilbert-base-uncased-emotion>

Viteri, A. (2021). Limitaciones de los modelos tradicionales de análisis de sentimientos en redes sociales. *Revista Iberoamericana de Tecnología y Sociedad*.

<https://revistaiberoamericana.org/>

Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Rault, T., Louf, R., Funtowicz, M., Davison, J., Shleifer, S., von Platen, P., Ma, C., Jernite, Y., Plu, J., Xu, C., Le Scao, T., Gugger, S., . . . Rush, A. M. (2020). Transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*. <https://arxiv.org/abs/1910.03771>