Equipo Deutsche Welle -Presentación del Proyecto: Análisis de Sentimiento

Sarah Armstrong Carlos Gómez Guijarro Martín Pérez Bernal Sophia Lunetto



Resumen de DW

- La emisora internacional de Alemania, fundada en 1953
- Financiada con dinero público, con oficinas principales en Bonn y Berlín

Algunos datos

- Millones de contactos semanales de usuarios a través de TV, radio y en línea (2023)
- DW transmite en 32 idiomas

Misión

Brindar información imparcial para mentes libres





Contenido

Estructura de la presentación

- 1. 🖈 Introducción
- 2. A Motor de búsqueda
- 3. Análisis de sentimiento
- 4. 4 Análisis de sesgo
- 5. Resultados
- 6. Conclusiones
- 7. **Perspectivas**

- **K** Entorno técnico del proyecto
- Lenguaje: Python
- Librerías principales:
 - nltk, pandas, NumPy
 - openai, transformers, PyABSA
 - MLflow



Objetivos



Automatizar el análisis de sentimiento para identificar la tonalidad (positiva, neural o critica) de una entidad especifica mencionada en artículos periodísticos.

Input: Un query o topico

Encontrar los artículos relevantes

Análisis de sentimiento hacia una entidad

Resultados



- Datos de evaluación etiquetados
- Motor de búsqueda
- Pipeline de análisis de sentimiento
- Evaluación de modelos
- Seguimiento de rendimiento con la plataforma MLflow



Análisis Exploratorio de Datos y Limpieza



Data



EDA



Data Processing

Datos sin etiquetar de DW:

ID del artículo, idioma, título, resumen, texto completo

Conjuntos de datos externos (etiquetados por expertos):
•PerSenT

NewsMTSC

Datos de evaluación (tarea de etiquetado)

EDA basico:

- Filtrar articulos cortos
- Artículos unicos
- Frecuencia mensual de publicación

EDA Avanzado:

- · modelado de temas
- NER (Reconocimiento de Entidades Nombradas)



Preprocesamiento General y Específico del corpus



Data



EDA



Data Processing

Datos sin etiquetar de DW:

ID del artículo, idioma, título, resumen, texto completo

Conjuntos de datos externos (etiquetados por expertos): •PerSenT

NewsMTSC

Datos de evaluación (tarea de etiquetado)

EDA basico:

- Filtrar articulos cortos
- Artículos unicos
- Frecuencia mensual de publicacion

EDA Avanzado:

- modelado de temas
- NER (Reconocimiento de Entidades Nombradas)

Preprocesamiento general:

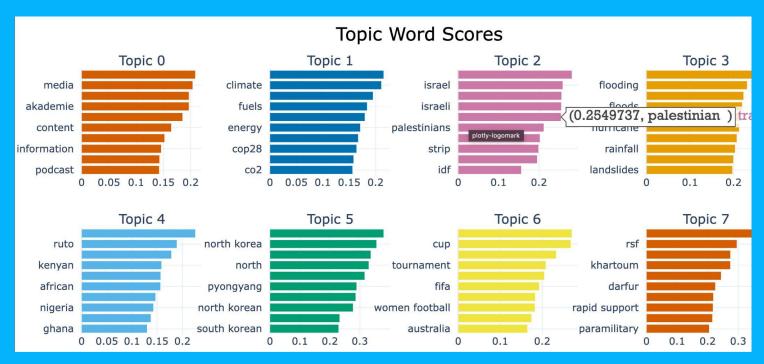
- Manejo de valores faltantes
- Manejo de textos muy cortos o muy largos

Preprocesamiento específico del modelo:

- NLTK library
- Hugging Face



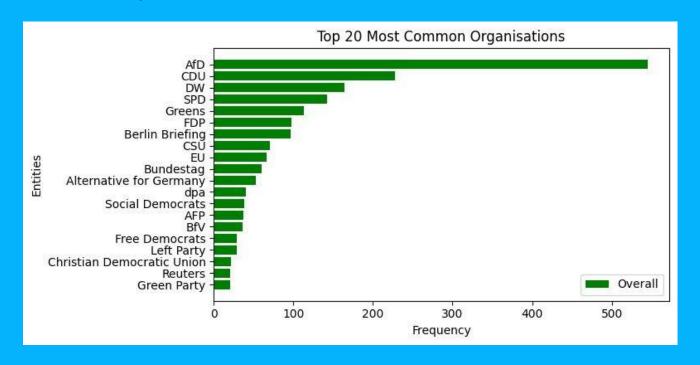
EDA Avanzado: Modelado de Temas



El modelado de temas hace clusters semanticos



Análisis Exploratorio Avanzado: Reconocimiento de Entidades Nombradas



NER (Reconocimiento de Entidades Nombradas) resalta las entidades principales y está integrado en el motor de búsqueda para obtener perspectivas más profundas



Datos de Evaluación

Crear un conjunto de datos de evaluación (entregable clave)

- Ground Truth
- Mejoramiento del Modelo
- Confianza en el Modelo

Desafios:

- Complejidad del sentimiento
- Menciones indirectas
- Lenguaje dependiente del contexto

"Merkel es realmente genial, mientras que Steinmeier es un incompetente..."





Proceso Colaborativo de Etiquetado

Selección de un sub set para evaluación

- → Tema: "Política alemana" (138 artículos)
- → Entidades clave: Olaf Scholz (Canciller alemán), AfD (partido de extrema derecha) y el Papa Francisco

Desarrollo de Guidelines

- → Se crearon instrucciones claras para el etiquetado
- → Se realizaron reuniones de equipo para asegurar la consistencia y obtener retroalimentación

Primera ronda

- → Cada persona etiquetó aproximadamente 46 artículos
- → Cada artículo fue etiquetado dos veces
- → Se presentaron desacuerdos en el 37% de los artículos

Segunda ronda

- → Ronda de desempate: Resolución de discrepancias
- → Artículos sin acuerdo fueron eliminados del conjunto de datos

Final de Evaluación

→ Se obtuvo un conjunto de datos de alta calidad (109 artículos) para la evaluación del modelo



Motor de búsqueda

Ingresar query

Encontrar artículos relevantes

relevantes

Búsqueda por palabra clave

Búsqueda por palabra clave mejorada

Búsqueda Semántica



Entidad nombrada Coincidencias exactas

Conjunto ampliado semánticamente

Coincidencia basada en similitud semántica



Coincidencia con entidades identificadas por el modelo



Motor de búsqueda Búsqueda por palabra clave

Ingresar query

Ingresar query

I. Encontrar artículos relevantes

query = "Germany"

Búsqueda por palabra clave

resultado = articulos que contienen:
"Germany"

Pros: Simple, lightweight.

Cons: Rigido, el tiempo requerido escala linealmente con el tamaño del corpus.

resultado =

"Greenhouse gases"



Ingresar query

Motor de búsqueda Búsqueda por palabra clave mejorada

I. Encontrar artículos Ingresar query relevantes Búsqueda por query "Global palabra Articulos que warming" clave contienen: "Global mejorada warming" O "Climate change" 0

Pros: Captura sinonimos, lightweight.

Cons: Resultados con baja relevancia, el tiempo requerido escala linealmente con el tamaño del corpus.



Motor de búsqueda Búsqueda por entidad nombrada

Ingresar query

Ingresar query

I. Encontrar artículos relevantes

query = "Olaf Scholz"

Selección de la entidad mas relevante desde una lista: "Olaf Scholz" resultados =
Articulos que contienen:
"Olaf Scholz" como entidad

Pros: Precision alta, extrae entitiades especificas.

Cons: Reconocimiento de entidades requiere bastante computacion. Mejora la precisión del análisis de sentimiento.



Motor de búsqueda Búsqueda semántica

Ingresar query

I. Encontrar artículos relevantes

query =
"Pope Francis
tour of Sudan"

Ingresar query

Búsqueda semántica

resultado =
Articulos ordenados por
relevancia (conteniendo
frases "suficientemente"
similares)

Pros: Comparacion al nivel del significado, alta flexibilidad/relevancia

Cons: Requiere Embeddings, los cuales son intensivos en términos de recursos.



Motor de búsqueda

Ingresar query

similaridad semántica seguido de búsqueda por entiedad nombrada

Ingresar query

la. Encontrar artículos relevantes

Ingresar query

Ib. Encontrar artículos relevantes

query =
"German
elections""

Búsqueda Semántica

query = "Olaf Scholz"

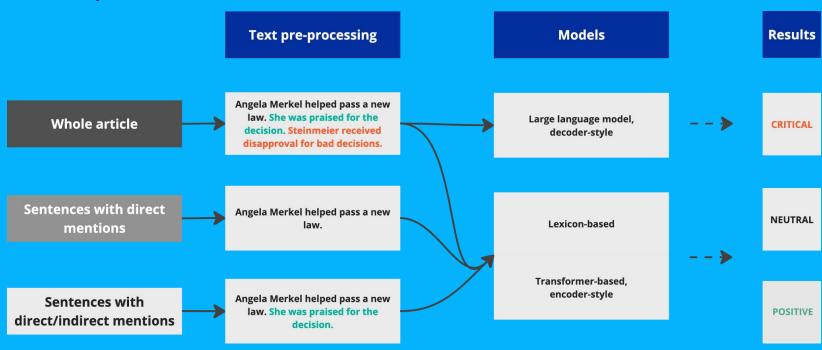
Busqueda por entidad nombrada

resultado =
Articulos sobre
"Olaf Scholz"
en el contexto
de
"German elections"



Modelos Preprocesamiento de texto

Ingresar query





Modelos "Lexicon"

Valence Aware Dictionary for Sentiment Reasoning (Vader)



Ingresar query

Léxico de puntuación desde -4 hasta 4

"horrible": -2.5

"great": 3.1

"okay": 0.9



Heurísticas

intensidad (como"very") punctuación (como"great!!!") capitalización (como"GREAT") Puntaje agregado entre -1 y 1 entre +0.05 y +1 es POSITIVO entre -0.05 y 0.05 es NEUTRAL entre -1 y -0.05 es CRITICO



Modelos "transformer"

Ingresar query

Los transformer son NNs especializados con 100M-300M de parametros.



Pros: Bajos prerrequisitos computacionales.



Cons: Tienen mejor precisión cuando se hace un preprocesamiento adecuado (depende del modelo).

Output: Devuelve un label en forma de Critical, Neutral, o Positive.



Modelos "transformer"

Ingresar query

Bidirectional Encoder Representations from Transformers (BERT).

	NLPtown	NewsSentiment	CardiffNLP	PyABSA
Tipo de Modelo	BERT-based transformer	BERT-based transformer	RoBERTa-based transformer	BERT-based transformer
Uso principal	Analisis de sentimiento general (diferentes idiomas)	Analisis dirigido a entidades en artículos periodisticos	Analisis de sentimiento en social media y tweets	Analisis dirigido a entidades
Training Data	Amazon product reviews, movie reviews, etc.	News datasets	Twitter data	ABSA datasets como SemEval, restaurant reviews
Entity-based	*	V	×	V



Ingresar query

Large Language Models (LLMs)

NNs masivas como GPT y Claude con 1B-14B de parametros.





Pros: Generalizable. Entiende sentimientos de artículos enteros. No hay prerrequisitos locales.

Cons: Inferencia más costosa (\$), Requieren conocimientos técnicos adicionales (API).



Output: Puntuación de polaridad y, en algunos casos, explicación sobre cómo se llegó a la predicción (útil para prompt engineering)



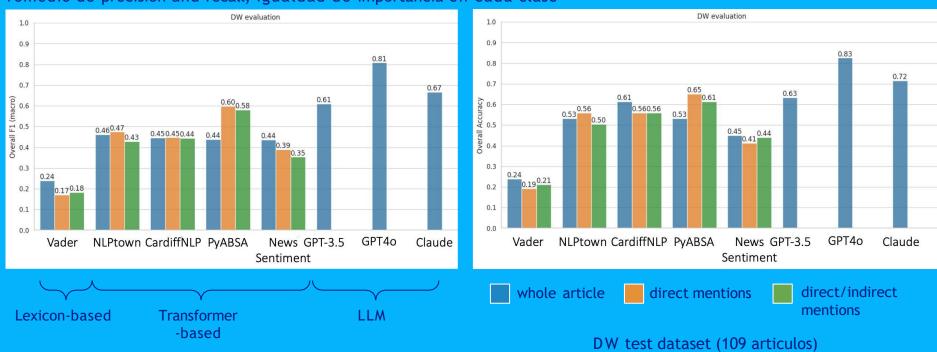


Evaluación de los Modelos

Ingresar Query

F1 score Promedio de precision and recall, igualdad de importancia en cada clase

Accuracy Proporción de predicciones correctas

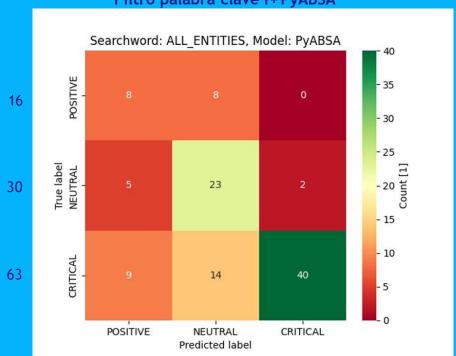




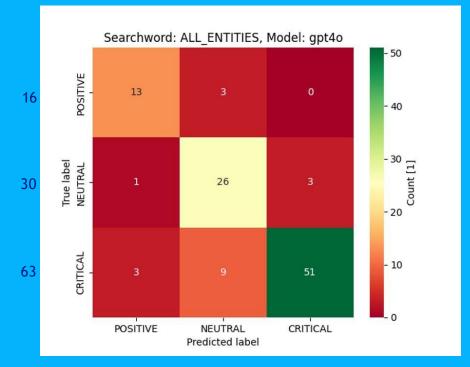
Evaluación de los Modelos

Ingresar Query

Filtro palabra clave f+PyABSA



GPT-40



II. Análisis de

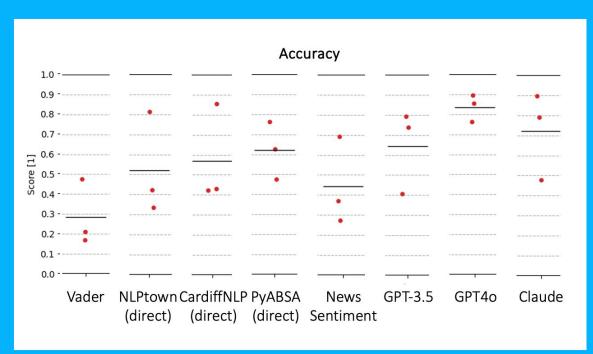
sentimiento



Evaluación de los Modelos

Consistencia entre Entidades

Ingresar query



Entidades:

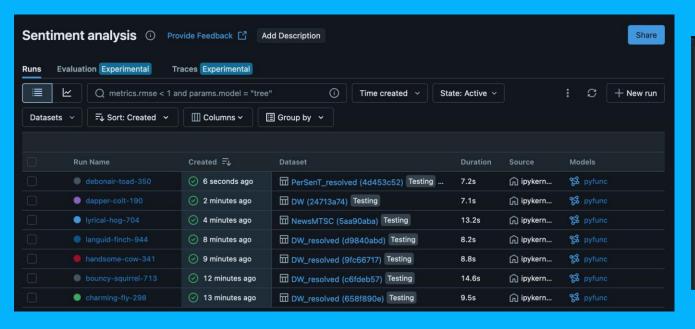
AfD (89.6% CRITICAL) Olaf Scholz (50.0% NEUTRAL) Pope Francis (68.4% POSITIVE)

Hallazgo:

La mayoría de los modelos tienen dificultades para mantener un rendimiento constante entre distintas entidades. Mejor desempeño: GPT-40



Seguimiento de Experimentos en MLflow





Análisis de sesgo

Tema: German politics

Entidades: Olaf Scholz - CDU - AfD

Olaf Scholz (413 articulos)

Búsqueda: Combinación de búsqueda semántica y por entidad nombrada

Modelos: PyABSA - GPT-40

Sentiment Distribution (%)

CRITICAL

19.4%

POSITIVE

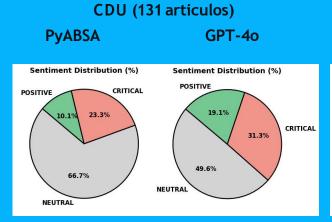
POSITIVE

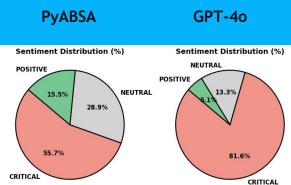
11.7%

29.7%

NEUTRAL

NEUTRAL





AfD (98 articulos)



Análisis

Ingresar query

Tema: German politics

Conclusiones clave:

 Olaf Scholz y CDU están cubiertos predominantemente con un tono neutral

La distribución de tonalidades positiva o crítica varía según el modelo

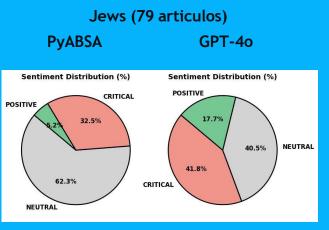
 AfD recibe una cobertura más crítica, lo que destaca la naturaleza polémica de sus políticas y la percepción pública

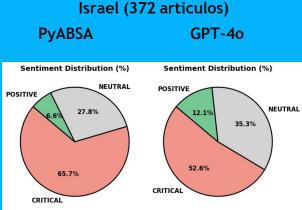


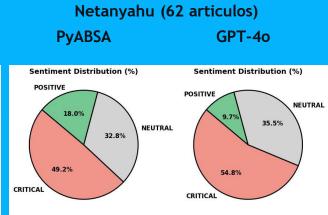
Tema: Israel-Palestine conflict

Entidades: Jews - Arabs - Israel - Netanyahu - Palestinian Authority - Hamas Búsqueda: Combinación de búsqueda semántica y por entidad nombrada

Modelos: PyABSA - GPT-40







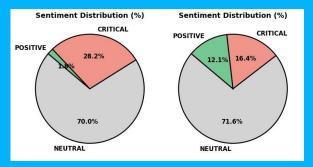


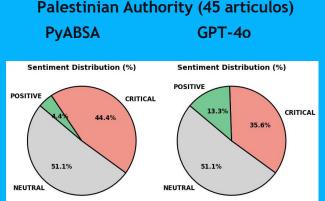
Tema: Israel-Palestine conflict

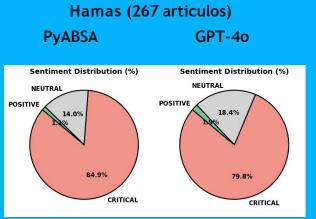
Entidades: Jews - Arabs - Israel - Netanyahu - Palestinian Authority - Hamas Búsqueda: Combinación de búsqueda semántica y por entidad nombrada

Modelos: PyABSA - GPT-40

Arabs (116 articulos)
PyABSA GPT-40









Análisis

Ingresar query

Tema: Israel-Palestine conflict

Conclusiones clave:

• El sentimiento hacia las comunidades judía y árabe depende del modelo utilizado.

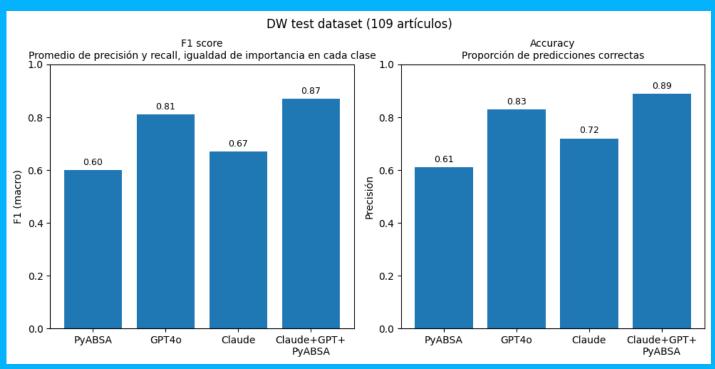
Se necesita más pruebas exhaustivas

•Las entidades con roles de gobierno o conflicto activo (como Netanyahu y Hamas) reciben un sentimiento más **crítico**



Una nueva idea

Ingresar Query





Análisis

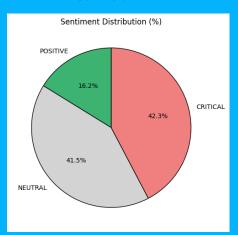
Ingresar query

Tema: Israel-Palestine conflict

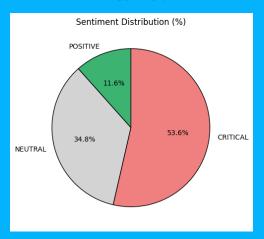
Entidades: Jews - Arabs - Israel - Netanyahu - Palestinian Authority - Hamas **Búsqueda:** Combinación de búsqueda semántica y por entidad nombrada

Modelos: claude+gpt+pyabsa

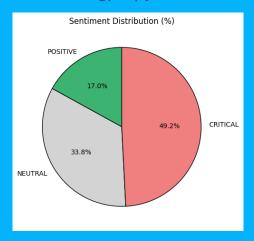
Jews (79 articulos) claude+gpt+pyabsa



Israel (372 articulos)
claude+gpt+pyabsa



Netanyahu (62 articulos) claude+gpt+pyabsa



sentimiento



Análisis

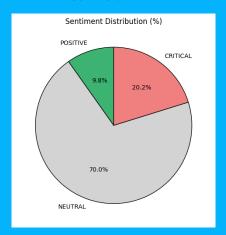
Ingresar query

Tema: Israel-Palestine conflict

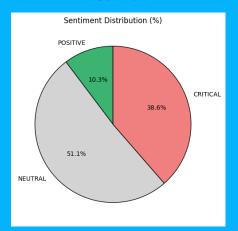
Entidades: Jews - Arabs - Israel - Netanyahu - Palestinian Authority - Hamas Búsqueda: Combinación de búsqueda semántica y por entidad nombrada

Modelos: claude+gpt+pyabsa

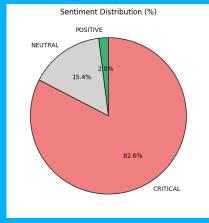
Arabs (116 articulos) claude+gpt+pyabsa



Palestinian Authority (45 articulos) claude+gpt+pyabsa



Hamas (267 articulos) claude+gpt+pyabsa





Resumen

- Datos etiquetados de alta calidad son esenciales para evaluar el rendimiento de los modelos
- Mejor modelo en rendimiento: GPT-40
 Modelo más eficiente y liviano: PyABSA
 Top performer: GPT+Claude+PyABSA
- El sistema de votación de los tres mejores modelos tiene una precisión increíble y esta determino que hay sesgo en algunas entidades



Perspectivas

- Continuar el análisis en temas relevantes y con metadatos
- Mejoras de rendimiento:
 - Ajuste fino de modelos basados en transformadores con corpus de DW
 - Ingeniería de prompts para modelos LLM
 - Estudiar porque el sistema de votación tiene un incremento en precisión
- Facilidad de uso para periodistas:
 - Crear una interfaz de usuario
 - Desarrollar un chat interactivo
- Nuevas funcionalidades:
 - Soporte para múltiples idiomas
 - Mostrar evolución del sentimiento a lo largo del tiempo





Gracias por su atención

