

Proyecto de Ingeniería de Datos

IMEDIA Project



ITESO, Universidad Jesuita de Guadalajara

Esteban Javier Berumen Nieto
Isabel Valladolid Dillanes

Mayo 2025

Introducción

Actualmente, las redes sociales se han convertido en uno de los principales espacios donde las personas expresan sus opiniones, emociones y reacciones respecto a eventos sociales, políticos y culturales. Gracias al gran volumen y velocidad de publicación, es complicado extraer información significativa de manera manual. La necesidad de interpretar el sentimiento colectivo y detectar patrones emocionales es cada vez más importante para empresas, medios de comunicación, etc.

“IMEDIA”, busca contribuir a la demanda de herramientas prácticas que permitan analizar la opinión pública en plataformas o redes sociales. IMEDIA proporciona una solución práctica para extraer, clasificar y presentar insights emocionales a partir de textos publicados por los usuarios.

Objetivo

El objetivo de IMEDIA es ser una plataforma de análisis de sentimientos en redes sociales, que integre extracción de datos, procesamiento NLP y visualización interactiva de insights emocionales, que además, permita explorar patrones de opinión en tiempo real a través de un dashboard dinámico y una API de clasificación.

Metodos

El proyecto realiza la extracción de datos desde subreddits populares, incluyendo publicaciones, comentarios y metadatos, enriquecidos con información de los autores. Los datos en formato JSON se transforman en archivos CSV organizados. Luego, se aplica limpieza textual y tokenización con *bert-base-uncased*. Los textos se clasifican por sentimiento o categoría usando un modelo ajustado previamente, y los resultados se exportan. Finalmente, se visualizan mediante dashboards en Power BI que muestran análisis por subreddit, tipo de publicación, nube de palabras y tendencias temporales.

Resultados obtenidos

El flujo de datos sigue una estructura organizada: los archivos JSON se almacenan en *data/raw/*, se transforman en CSVs en *data/interim/*, luego se limpian y procesan en *data/processed/*, y finalmente se clasifican en *data/analice/*. Los resultados se visualizan mediante dashboards interactivos en Power BI, facilitando el análisis y la interpretación de la información extraída.

```

xteb@teb10d15 MINGW64 ~/apps/IMEDIA_Project/media_project (main)
$ python pipeline.py

Tenemos: 91 subreddits que son:
('OldSchoolCool', 'AskReddit', 'politics', 'LeopardsAteMyFace', 'ChoosingBeggars', 'DunderMifflin', 'BlackPeopleTwitter', 'Dammthatsinteresting', 'dataisbeautiful', 'GetMot
ivated', 'Unexpected', 'Instagramreality', 'movies', 'PublicFreakout', 'WatchPeopleDieInside', 'EarthPorn', 'MadeMeSmile', 'Wellthatsucks', 'Freefolk', 'gaming', 'funny', '
withtheonion', 'comics', 'Inanyparents', 'MallStreetBets', 'trashy', 'PoliticalHumor', 'couplepics', 'Bekeezed', 'maxfuckinglevel', 'gifswhatkeepgoing', 'todaylearned',
'Hongkong', 'meirl', 'Jokes', 'worldnews', 'technology', 'PregueMemes', 'Murderedbywords', 'pokemon', 'dankmemes', 'announcements', 'whatcouldgowrong', 'guityourbullshi
t', 'WhitePeopleTwitter', 'wildlyinteresting', 'Futurology', 'IAmA', 'ListenToThis', 'StarWarsBattlefront', 'BikinBottomTwitter', 'Music', 'facepalm', 'upliftingnews', 'ti
fu', 'rickandmorty', 'IdiotInCars', 'humansdelingros', 'TheAveller', 'PewdiepieSubmissions', 'interestingasfuck', 'Showerthoughts', 'pcasterrace', 'twochromosomes', 'te
chnicallythetruth', 'books', 'videos', 'Memeconomy', 'assholedesign', 'Coronavirus', 'comedyheaven', 'wholesomememes', 'awfuleverything', 'space', 'lotmemes', 'aww', 'Ayy
WD', 'sports', 'teenagers', 'AskMen', 'memes', 'JusticeServed', 'oddlysatisfying', 'science', 'news', 'instant_regret', 'insanepeoplefacebook', 'pics', 'LifeProTips', 'Kids
AreFuckingStupid', 'gifs')

Recopilando datos de Reddit...

-----
Fetching data for subreddit: OldSchoolCool subreddit 1 of 91
-----

-----
Fetching data for subreddit: AskReddit subreddit 2 of 91
-----

```

```

-----
Fetching data for subreddit: KidsAreFuckingStupid subreddit 90 of 91
-----

-----
Fetching data for subreddit: gifs subreddit 91 of 91
-----

Data saved to C:\Users\teb\apps\IMEDIA_Project\data\raw\raw_reddit_data.json
Datos procesados y guardados con éxito.
Token indices sequence length is longer than the specified maximum sequence length for this model (515 > 512). Running this sequence through the model will re
sult in indexing errors.
Comentarios procesados guardados en: C:\Users\teb\apps\IMEDIA_Project\data\processed\comments_processed.csv
Posts procesados (posts_processed.csv) guardados en: C:\Users\teb\apps\IMEDIA_Project\data\processed\posts_processed.csv
Posts procesados (subreddit_info_processed.csv) guardados en: C:\Users\teb\apps\IMEDIA_Project\data\processed\subreddit_info_processed.csv
Posts procesados (top_posts_processed.csv) guardados en: C:\Users\teb\apps\IMEDIA_Project\data\processed\top_posts_processed.csv
Posts procesados (ultimos_posts_processed.csv) guardados en: C:\Users\teb\apps\IMEDIA_Project\data\processed\ultimos_posts_processed.csv
Cargando modelo...
Modelo cargado
Procesando archivo: ../data/processed/comments_processed.csv
Tiene 8800 filas y 8 columnas
Analizando columna: body_tokenized
Analizando columna: subreddit_tokenized
Analizando columna: user_subscriptions_tokenized
Archivo guardado con sentimientos en: ../data/analice\comments_processed_labeled.csv
Procesando archivo: ../data/processed/posts_processed.csv
Tiene 9100 filas y 6 columnas
Analizando columna: title_tokenized
Analizando columna: subreddit_tokenized
Archivo guardado con sentimientos en: ../data/analice\posts_processed_labeled.csv
Procesando archivo: ../data/processed/subreddit_info_processed.csv
Tiene 91 filas y 5 columnas
Analizando columna: subreddit_tokenized
Archivo guardado con sentimientos en: ../data/analice\subreddit_info_processed_labeled.csv
Procesando archivo: ../data/processed/top_posts_processed.csv
Tiene 9100 filas y 6 columnas
Analizando columna: title_tokenized
Analizando columna: subreddit_tokenized
Archivo guardado con sentimientos en: ../data/analice\top_posts_processed_labeled.csv
Procesando archivo: ../data/processed/ultimos_posts_processed.csv
Tiene 2015 filas y 6 columnas
Analizando columna: title_tokenized
Analizando columna: subreddit_tokenized
Archivo guardado con sentimientos en: ../data/analice\ultimos_posts_processed_labeled.csv

```

Conclusiones

Este proyecto desarrolló un sistema automatizado para extraer y procesar datos de Reddit usando técnicas de NLP. Mediante la API de Reddit, librerías de Python y modelos pre entrenados, se recolectaron y limpiaron publicaciones, clasificándolas por sentimiento y almacenados localmente. La estructura modular del código y el uso de archivos CSV facilitan su reutilización y escalabilidad. Este trabajo demuestra cómo el NLP puede transformar datos textuales en información útil para analizar opiniones y comportamientos en redes sociales.

Bibliografía

PRAW 7.7.1 documentation. (s. f.). <https://praw.readthedocs.io/en/stable/>

Tokenizer. (s. f.). https://huggingface.co/docs/transformers/en/main_classes/tokenizer

Repositorio de github: https://github.com/XTEP63/IMEDIA_Project.git