

Machine Learning - Propuesta de Proyecto: ”Sentiment Stream”

Alondra Soto, Angel Santiago, Emilia Couret

7 de marzo del 2025

Este proyecto se estará enfocando en la aplicación de un modelo de procesamiento de texto. Dado una cadena de caracteres, tal modelo se enfocará en identificar la intención de la entrada dada. Es decir, el modelo será capaz de clasificar un texto dado basado en las emociones percibidas. Nuestra intención es proveer una clase correspondiente a la intención de un mensaje (positivo, neutral y negativo), y entrenar nuestro modelo en consecuencia. Queremos identificar palabras claves a distinguir el tono de voz de un texto, con la intención de obtener un modelo predictivo razonablemente preciso.

Sugerimos utilizar este modelo para optimizar el proceso de monitoreo psicológico. Muchas veces, la forma en la que una persona se expresa de manera escrita releva mucho sobre su estado emocional, pero procesar gran cantidades de texto puede ser un proceso limitado y poco eficiente. De tal manera, sugerimos utilizar un modelo tal como el que hemos descrito para proveer una manera de identificar patrones emocionales, permitiendo un análisis más profundo y continuo de un paciente. Este modelo podría mostrar el cambio de estado de ánimo a través del tiempo e incluso detallar las palabras que más contribuyan a texto clasificaciones negativas. Esta aplicación no solo optimizaría el tiempo de un profesional, sino que permite prevención en momentos oportunos y personalizados, lo cual mejora la calidad del tratamiento de un paciente.

Para construir nuestro modelo, utilizaremos una base de datos (<https://www.kaggle.com/datasets/durgeshrao9993/twitter-analysis-dataset-2022>) con aproximadamente 29 mil observaciones. La data contiene publicaciones aleatorias hechas en X (previamente Twitter) últimamente actualizada en el año 2023. En si, la data no tiene características, lo cual es sujeto a cambio con el progreso de nuestro proyecto. Adicionalmente, estaremos conduciendo los pasos de preprocesamiento estándares, y particionando la data de manera esperada (60% para entrenamiento, 20% de validación y 20% de prueba).

Nosotros proponemos un método similar al de la “bolsa de palabras”, donde tenemos una lista de palabras obtenidas de nuestra base de datos. Después de su debido pre-proceso, utilizaremos métodos de clasificación como regresión logística, con intención de interpretar las predicciones dadas entre las 3 clases mencionadas previamente.

Al tener un modelo, mostraremos sus predicciones comparadas a las etiquetas reales para determinar su rendimiento. Finalmente, proponemos también experimentar con “pesos” basados en la intención de cada palabra, en el cual cada entrada se le asigna un valor de peso. Suponemos que tal característica adicional, la cual se puede generar a través de otro modelo predictivo, podría mejorar la precisión del modelo principal.