

Don't Patronize Me!

- Flores Tiburcio Luis Fernando
- Vázquez Rojas José David
- Yáñez Espíndola José Marcos

Tabla de contenido



01

Introducción

02

Descripción del corpus

03

Metodología

04

Resultados

05

Conclusiones

Introducción

Lenguaje Condescendiente

- Denota una actitud superior hacia los demás.
- Describe una situación personal de una manera caritativa.
- Lo que genera un sentimiento de lástima y compasión.
- Es a menudo involuntario e inconsciente.

Detección de PCL

- Es difícil de detectar para los sistemas de PNL.
- No existe un formalismo en los algoritmos para su detección.
- Se usan algoritmos conocidos para aproximar un modelo aceptable
- Es inconsistente pues está sujeto a la naturaleza de las lenguas

Descripción del Corpus

01 10469 total de textos

02 En total, 9476 enunciados sin PCL (90% de los textos)

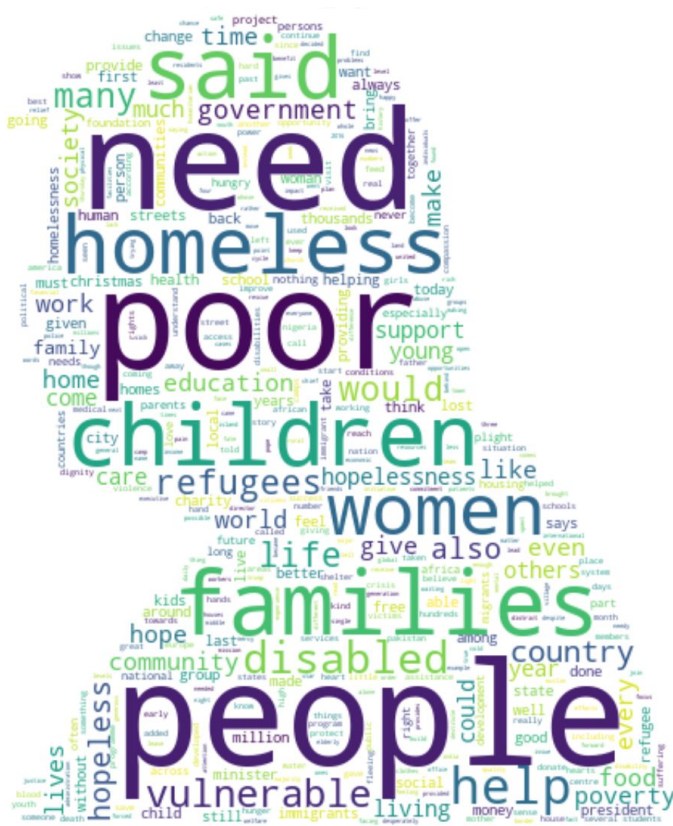
03 En total, 933 enunciados con PCL (10% de los textos)



10 categorías de PCL **04**

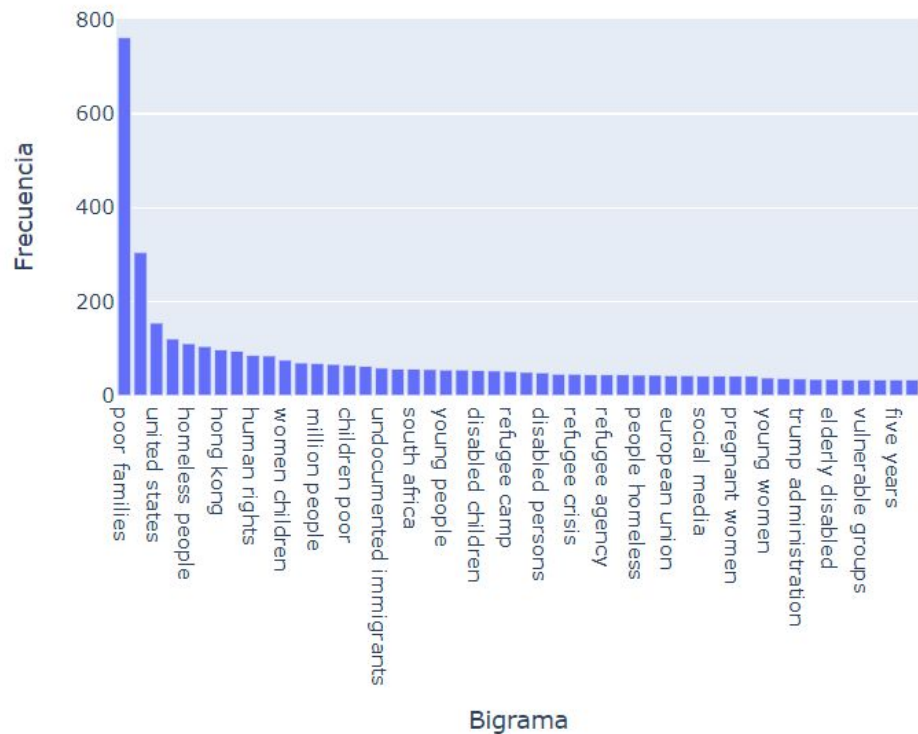
Textos provenientes de **05**
20 países diferentes

Textos obtenidos de **06**
984 artículos
diferentes en total

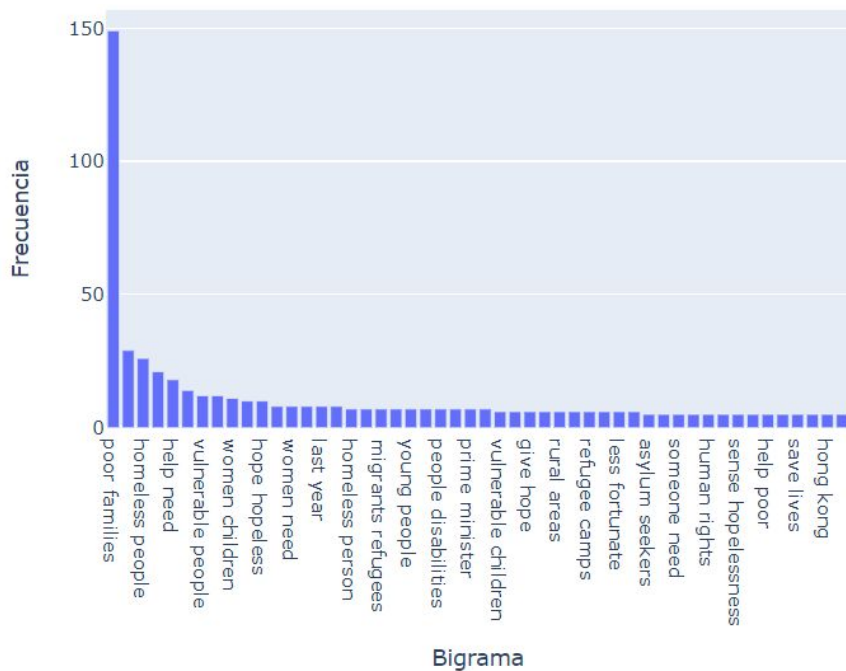


Gráfica de Barras de Bigramas

Lenguaje sin PCL

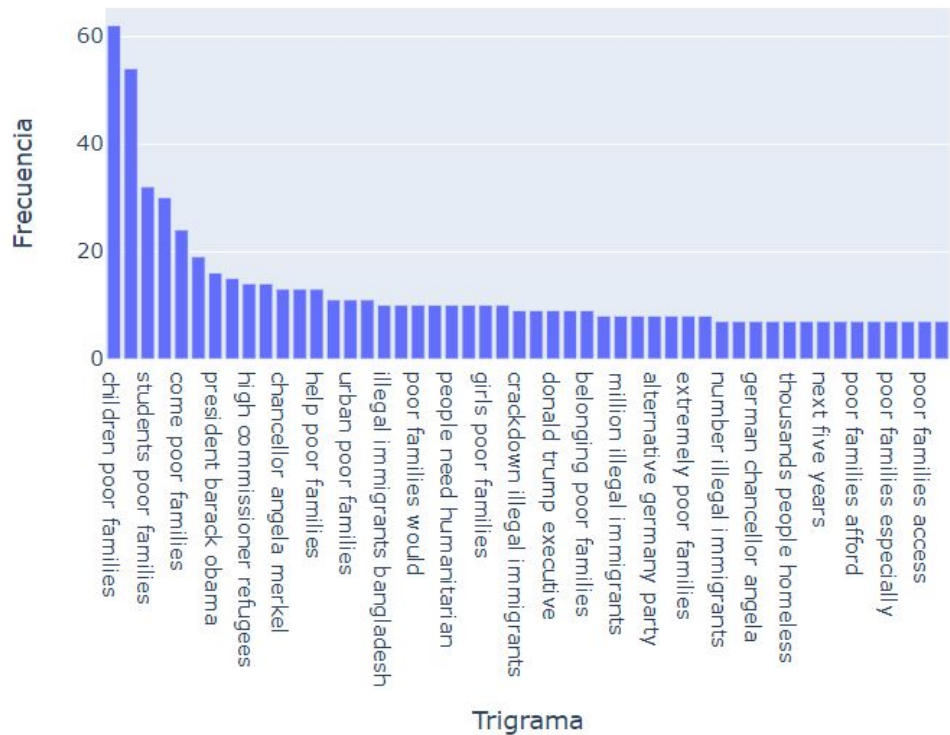


Lenguaje con PCL

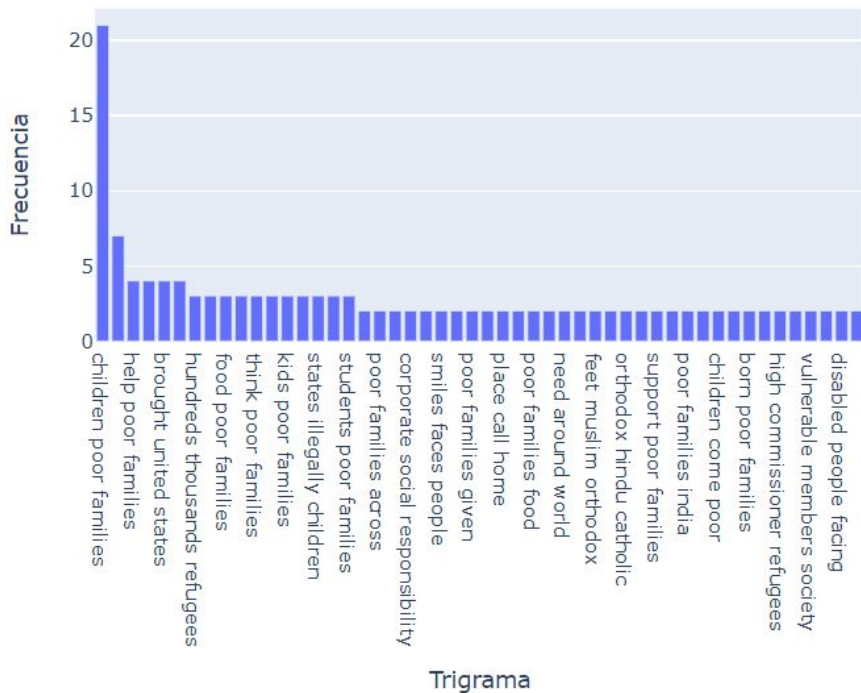


Gráfica de Barras de Trigramas

Lenguaje sin PCL

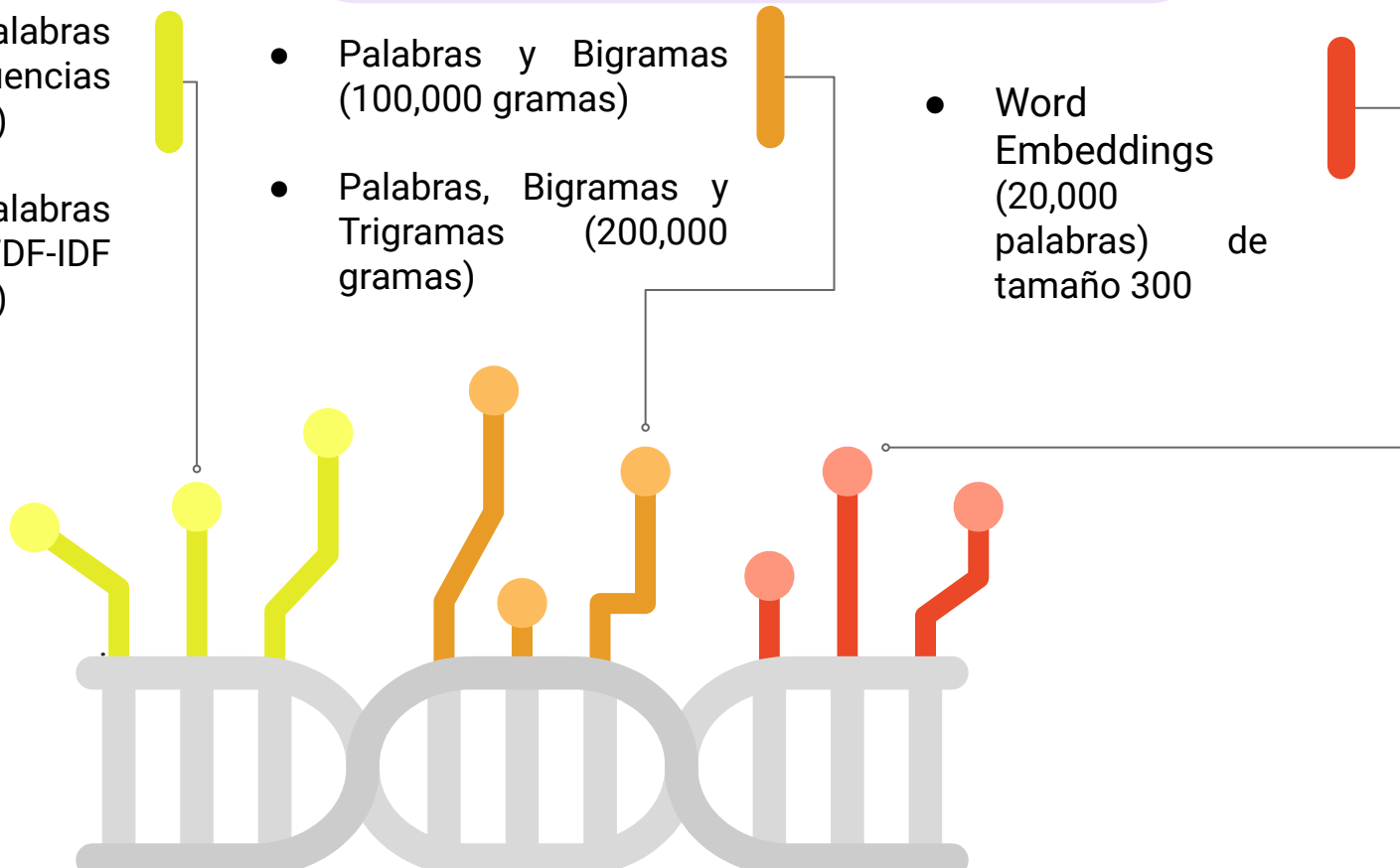


Lenguaje con PCL



Metodología

- Bolsa de Palabras usando las frecuencias (20,000 palabras)
- Bolsa de Palabras usando TDF-IDF (20,000 palabras)
- Palabras y Bigramas (100,000 gramas)
- Palabras, Bigramas y Trigramas (200,000 gramas)
- Word Embeddings (20,000 palabras) de tamaño 300



Conjuntos de Training, Testing y Validacion

67.5%

Training

Conjunto usado para entrenar los modelos

7.5%

Testing

Conjunto usado para probar cual es el mejor modelo

25%

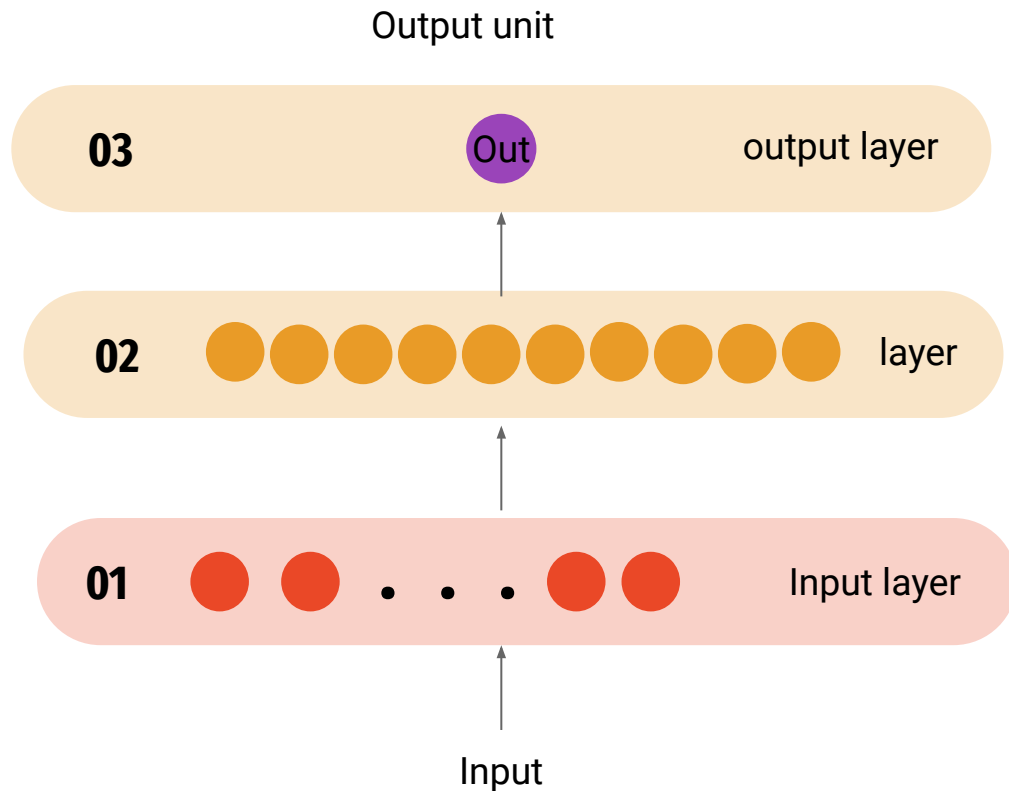
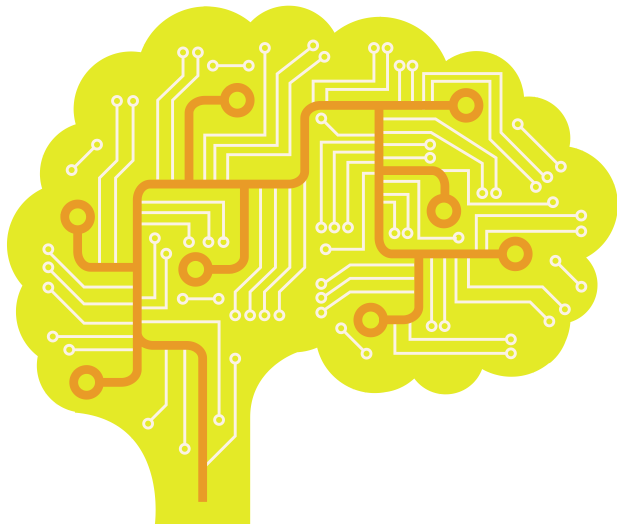
Validacion

Conjunto usado para probar con el mejor modelo



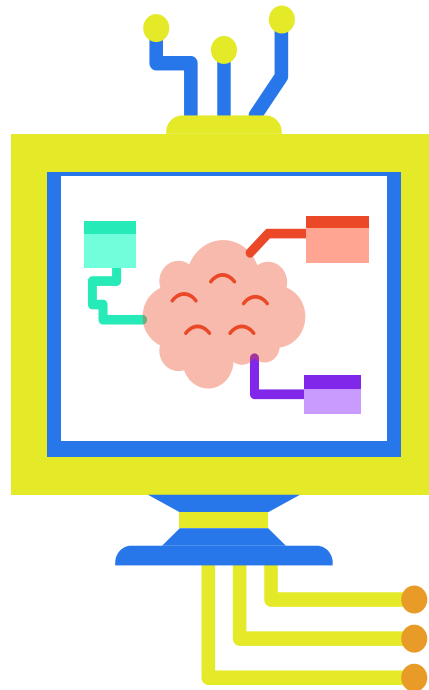
Modelo de la red neuronal

Input: Numero de
características
Hidden Layer: 10
Out: 1



Entrenamiento de la red

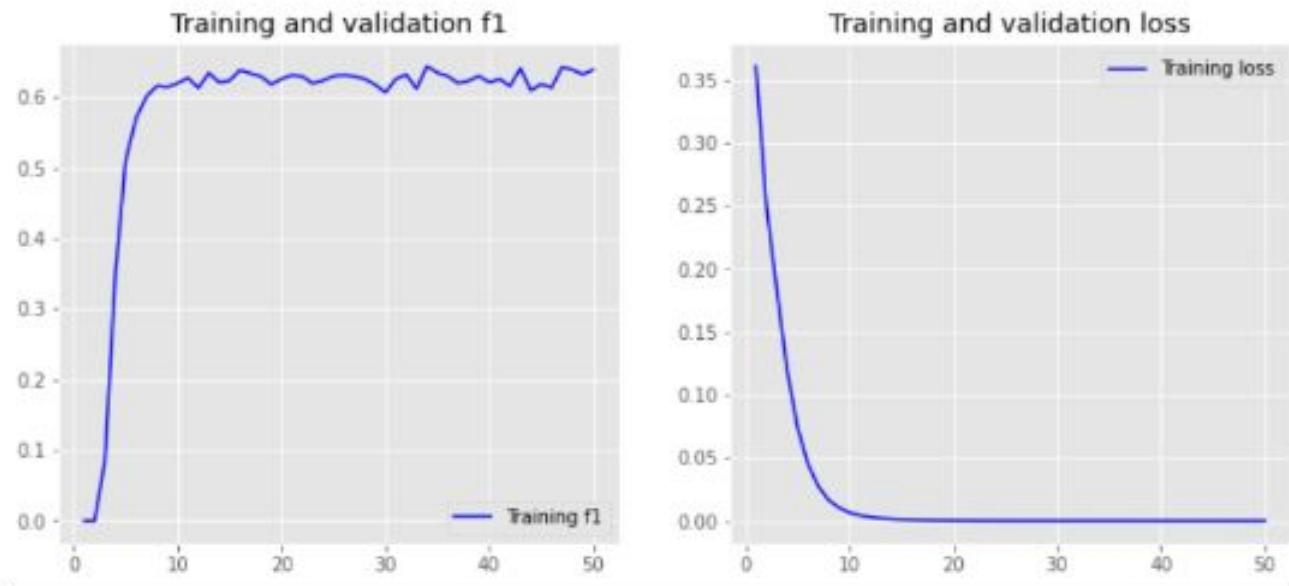
Para el entrenamiento, usamos el optimizador Adam con learning rate de 0.001 y 50 épocas con un batch size de 10.



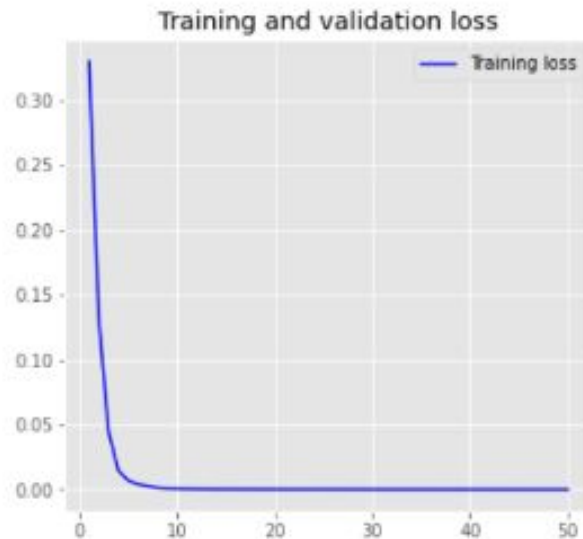
Gráfica de Barras de Trigramas



Bolsa de palabras TF-IDF



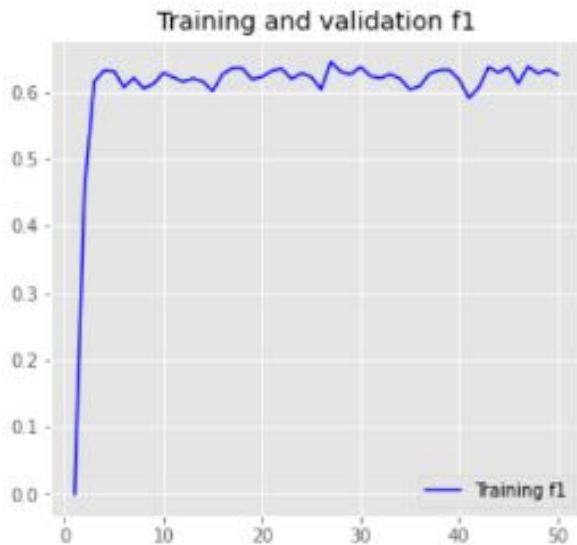
Palabras y Bigramas



Palabras, Bigramas y Trigramas



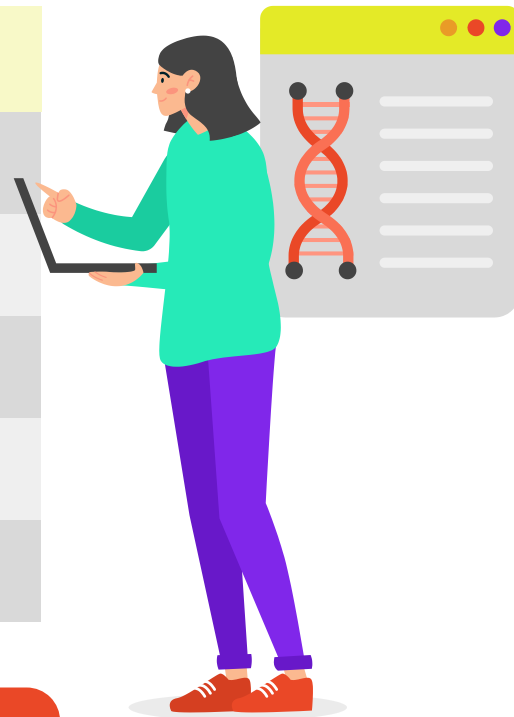
Word Embeddings



Resultados

Models	F1 Score en el test set
Bolsa de palabras simple	0.2730
Bolsa de palabras TDF-IDF	0.2750
Bigramas	0.2494
Trigramas	0.2221
Embeddings	0.2434

F1 Score con TDF-IDF en el test validación: 0.3378



Conclusiones

01 La clasificación de textos no es una tarea sencilla

Sobre todo cuando la diferencia entre las clases es demasiado sutil como en este caso.

02 Resulta difícil identificar PCL (inclusive para nosotros los humanos)

por ello muchas veces lo usamos inconscientemente, sin embargo, considerando que usamos un modelo muy sencillo de aprendizaje de máquina y todos los problemas antes mencionados, el resultado obtenido fue bastante bueno, aunque muy mejorable

03 El uso de n-gramas y word embeddings no mejora el F1 score

Para este caso, se esperaba por naturaleza de la lengua esperado una mejora con bigramas y trigramas, al igual con word embeddings

04 Trabajo a futuro

Como siguiente se probaría con algún modelo de red neuronal más complejo y haciendo un procesamiento de datos más complejo como usar stemming o lematización

Referencias

- Basant Agarwal and Namita Mittal. 2014. Text classification using machine learning methods-a survey. *In Proceedings of the Second International Conference on Soft Computing for Problem Solving(SocProS 2012), December 28-30, 2012, pages 701–709. Springer.*
- Carla Perez Almendros, Luis Espinosa Anke, and Steven Schockaert. 2020. Don't patronize me! An annotated dataset with patronizing and condescending language towards vulnerable communities. *In Proceedings of the 28th International Conference on Computational Linguistics, pages 5891–5902, Barcelona, Spain (Online). International Committee on Computational Linguistics.*
- Real Python. 2021. Practical text classification with python and keras. Keras Team. Keras documentation: Keras api reference.