

Analisis Exploratorio

# ANALISIS DE DATOS DEL LITCOIN NPL CHALLENGE

Analisis de datos del LitCoin NPL Challenge por  
Jeyner Arango, Estefania Elvira, Sofia Escobar, Jose Monzon

Grupo 8



## Analysis Exploratorio



# SITUACION PROBLEMATICA

Es necesario acelerar la investigación científica en el área de la medicina mediante el uso de la ciencia de los datos.



# Analisis Exploratorio

# ASPECTOS CLAVE



## Deteccion de Conceptos cientificos

Se buscan sistemas que puedan reconocer con precision conceptos científicos clave de textos científicos biomédicos.



## Identificacion de Relaciones Biomedicas

Identificar todas las relaciones entre el título y el resumen basadas en un modelo predictivo BioLink.





# PROBLEMA CIENTIFICO

El problema científico consta en desarrollar sistemas de procesamiento de lenguaje natural avanzados que puedan comprender y analizar textos científicos de manera eficiente para extraer información, identificar relaciones y determinar la novedad de los hallazgos.





# OBJETIVOS

Nuestros principales objetivos para la elaboracion del proyecto

- • • • •
- • • • •
- • • • •
- • • • •
- • • • •

## **Desarrollar modelos NLP**

Desarrollar modelos de procesamiento del lenguaje natural, que sean capaces de procesar textos científicos.

## **Determinar novedad científica**

Elaborar sistemas capaces de evaluar si las relaciones identificadas presentan información relevante

## **Reconocimiento de Entidades**

Se quiere lograr que los sistemas puedan reconocer correctamente las entidades biomedicas presentes



- • • • •
- • • • •
- • • • •
- • • • •

# Analisis Exploratorio



## Entidades (13,600 datos)

• `id`, `abstract\_id`, `offset\_start` y `offset\_finish`, `type`, `mention`:  
`entity\_ids`



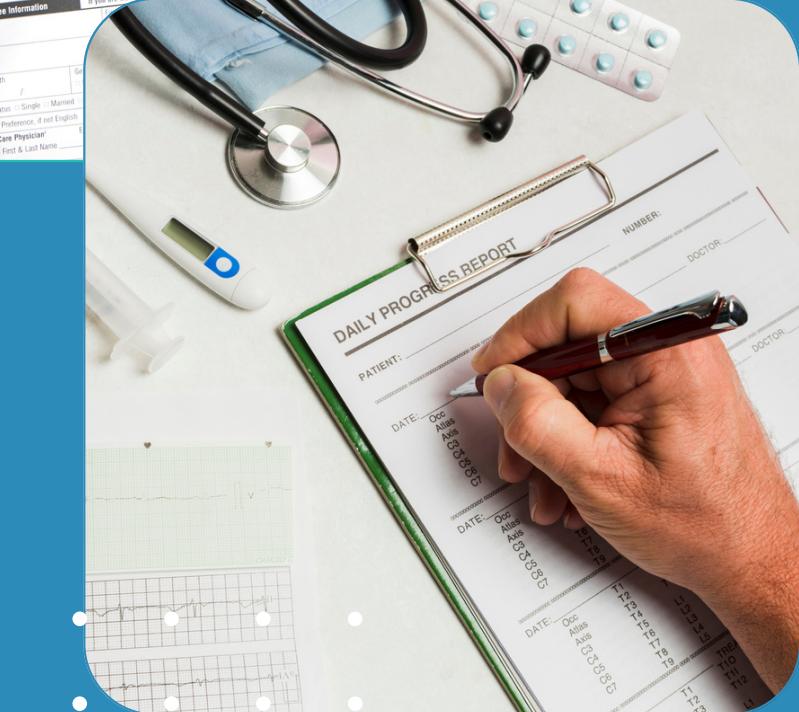
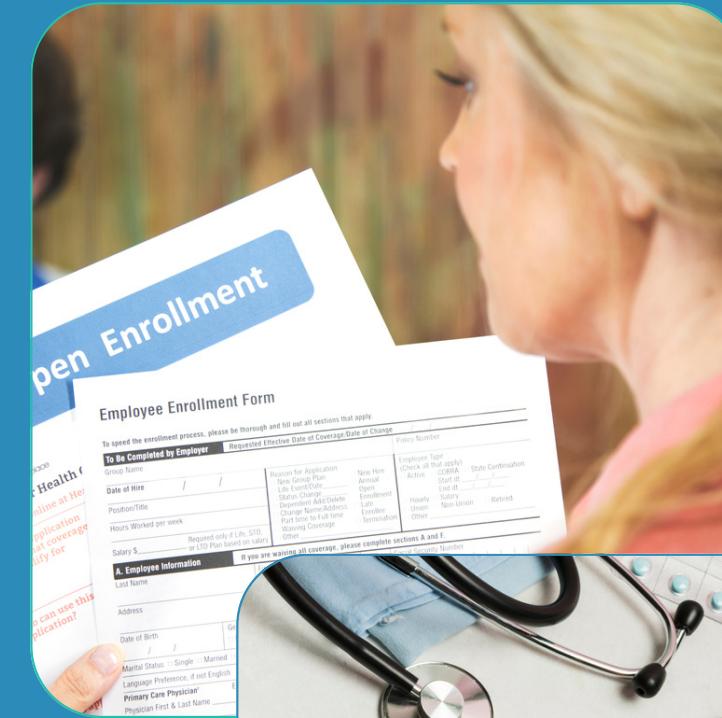
## Relaciones (2,300 datos)

- `id`, `abstract\_id`, `type`, `entity\_1\_id` y `entity\_2\_id`, `novel`



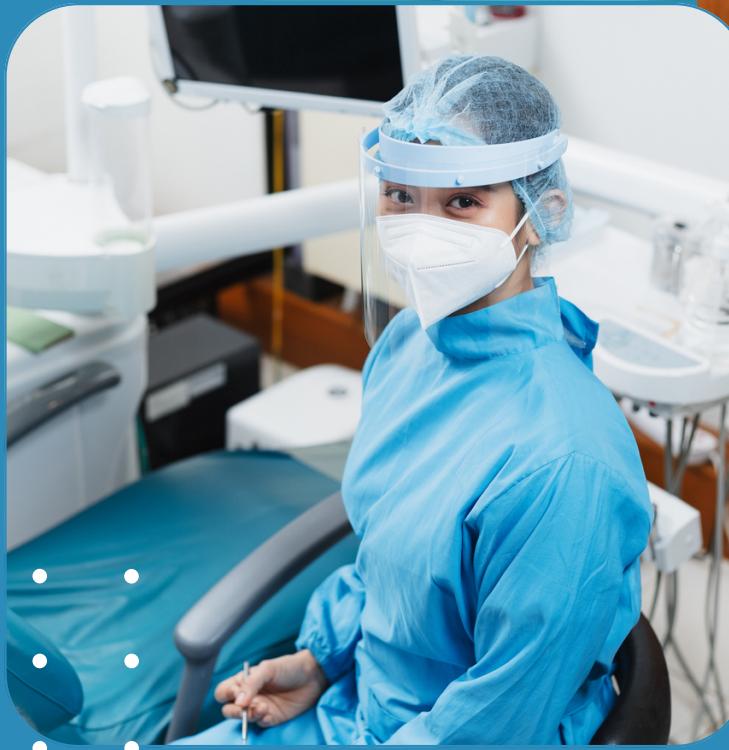
## Abstract (400 datos)

- `abstract\_id`, `title`, `abstract`



# LIMPIEZA

- Se eliminaron los valores Nulos
- Se convirtieron los textos a minusculas para mas uniformidad
- Se eliminaron caracteres especiales y puntuacion
- Se eliminaron las stopwords utilizando la base de datos de ntl.corpus

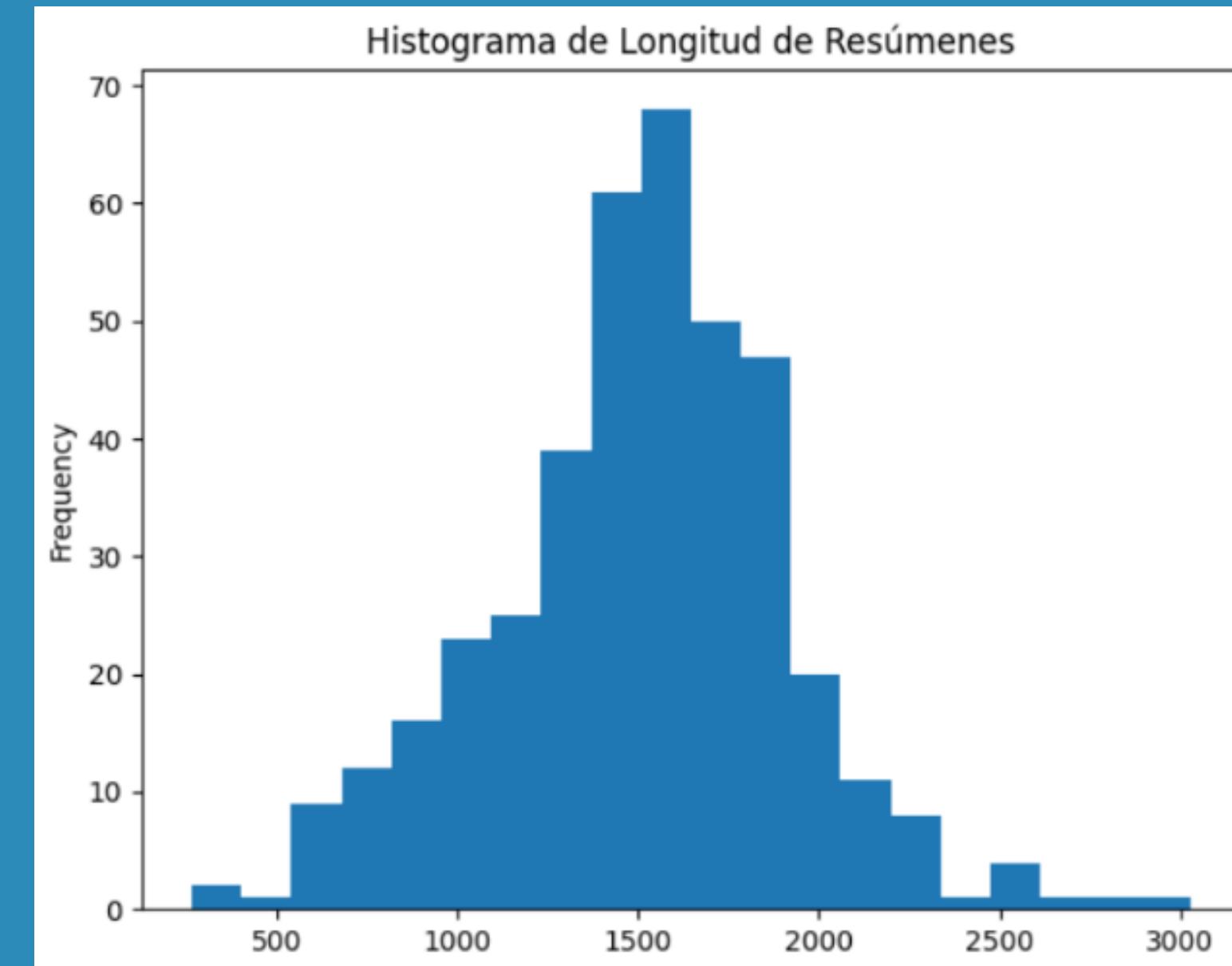
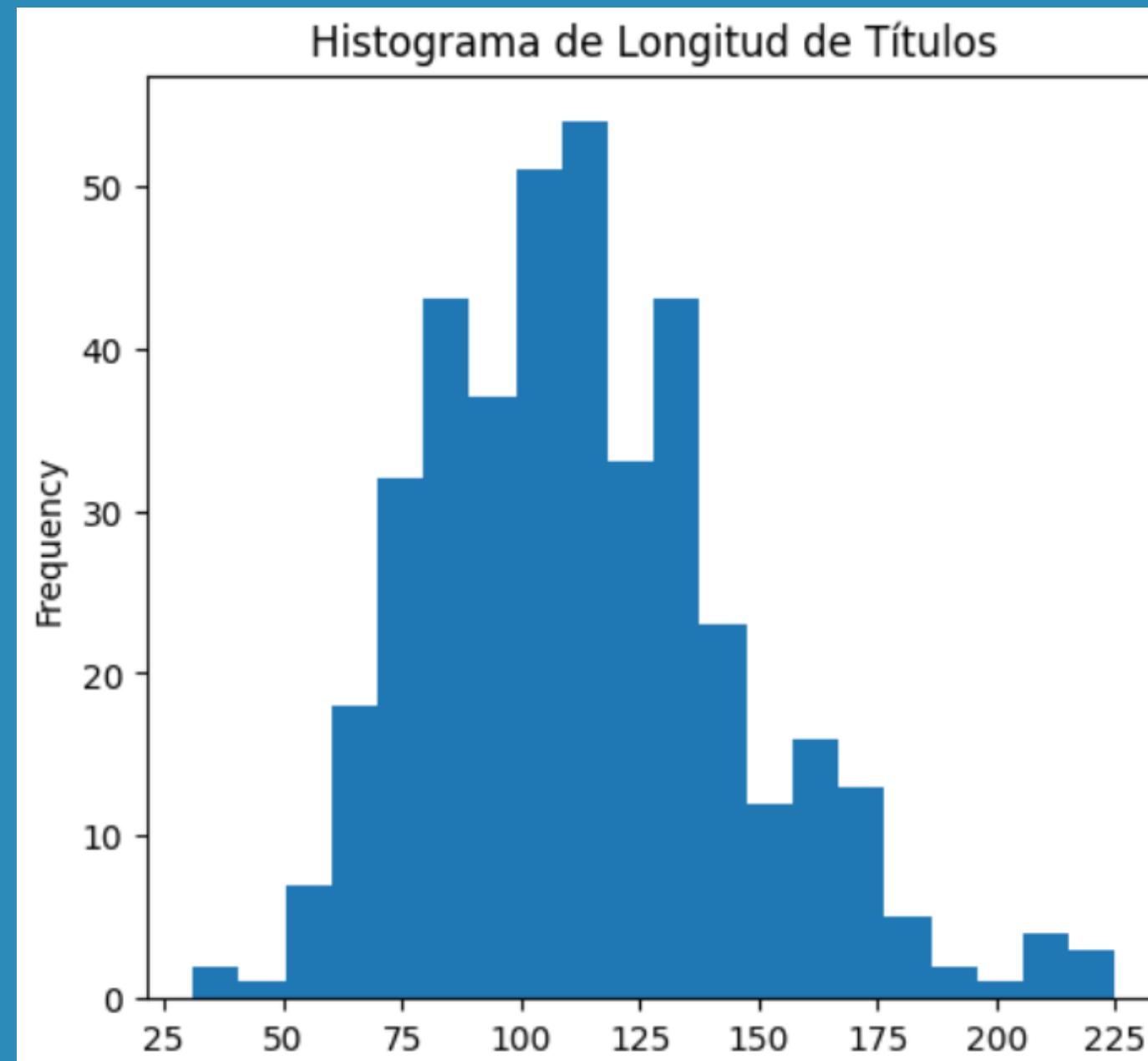




Analisis Exploratorio

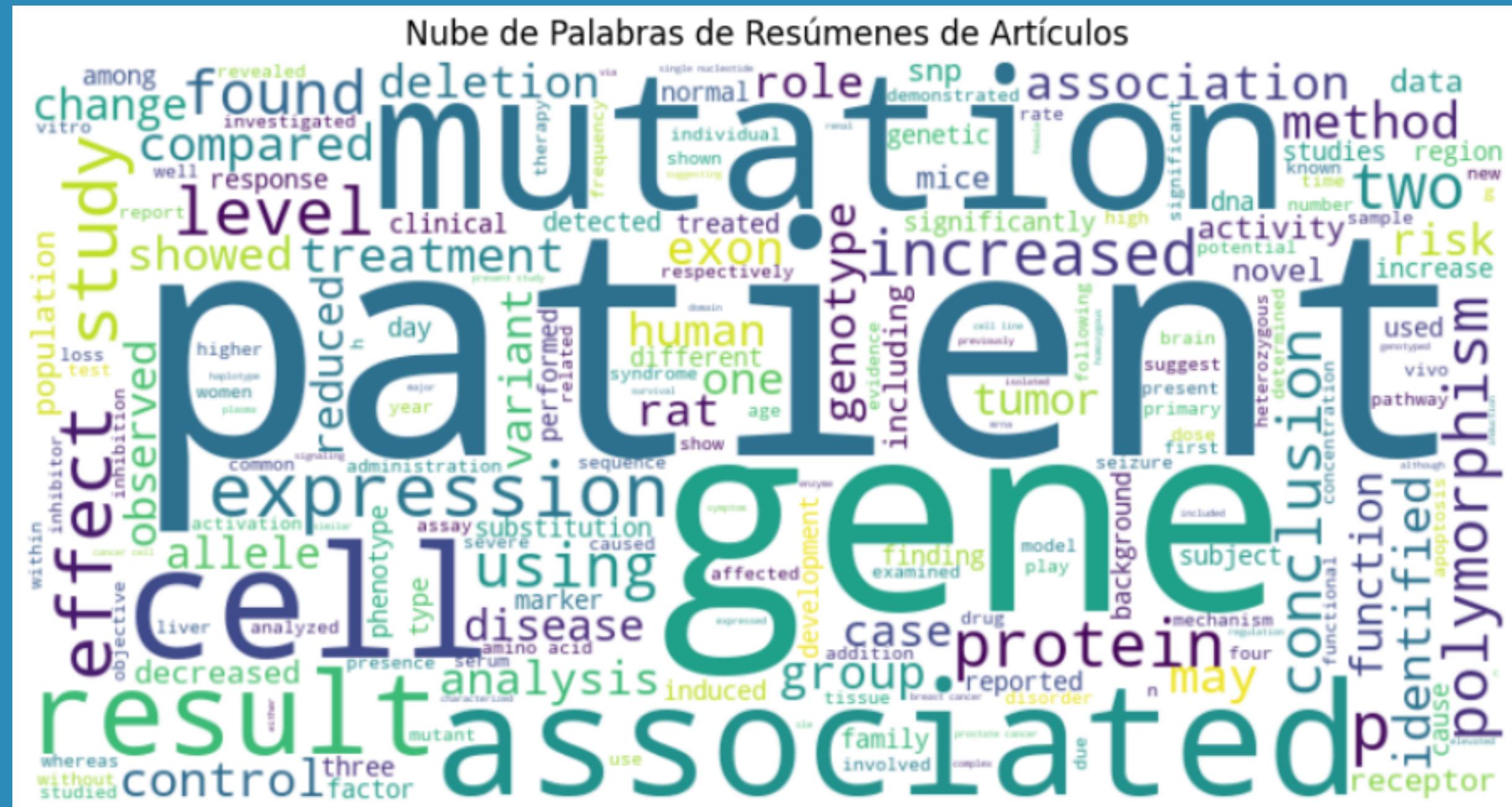


# ABSTRACT





# ABSTRACT





## Analisis Exploratorio

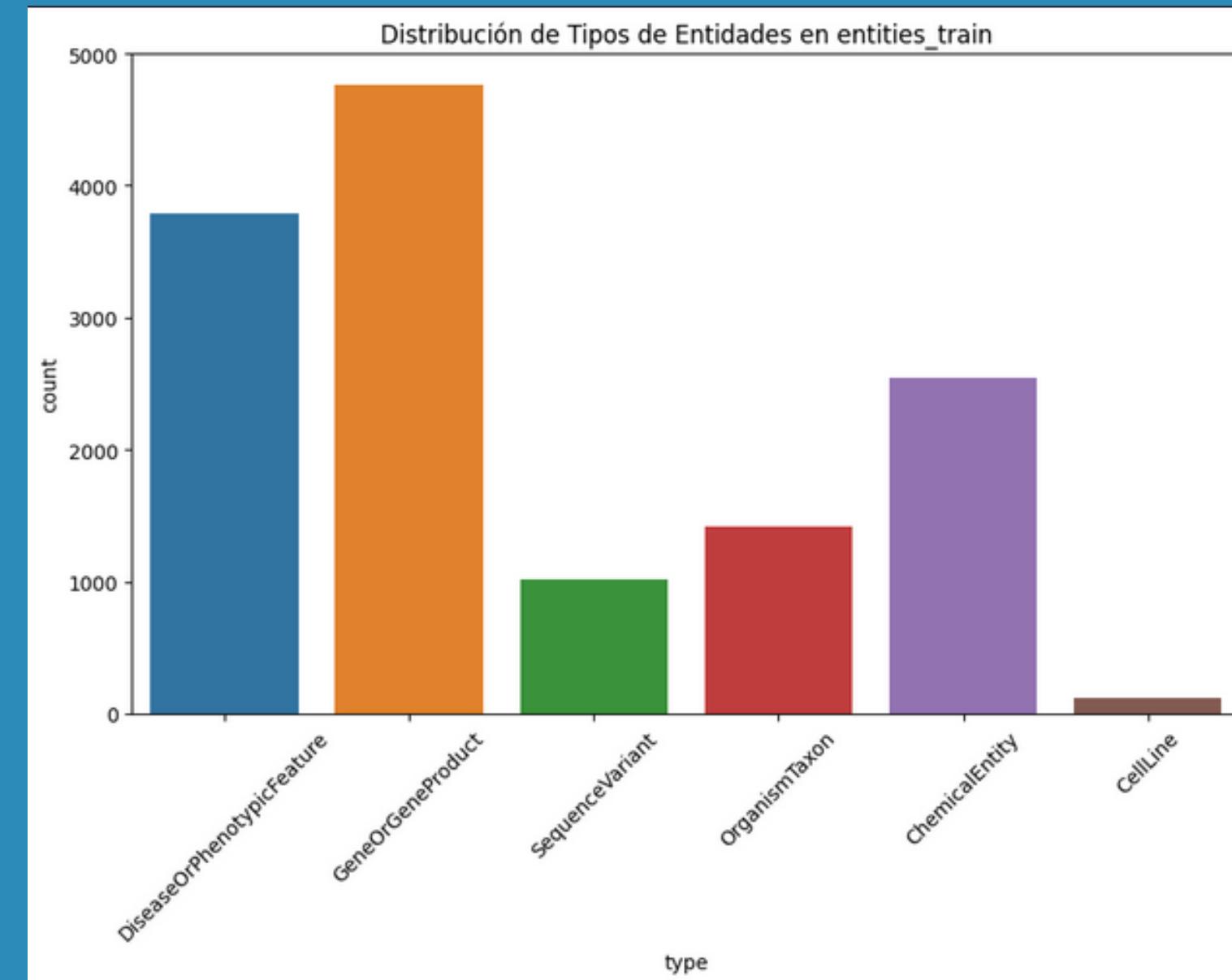


# ENTIDADES

Tabla de frecuencia de tipos de entidades:

GeneOrGeneProduct	4764
DiseaseOrPhenotypicFeature	3784
ChemicalEntity	2540
OrganismTaxon	1420
SequenceVariant	1011
CellLine	117

Name: type, dtype: int64



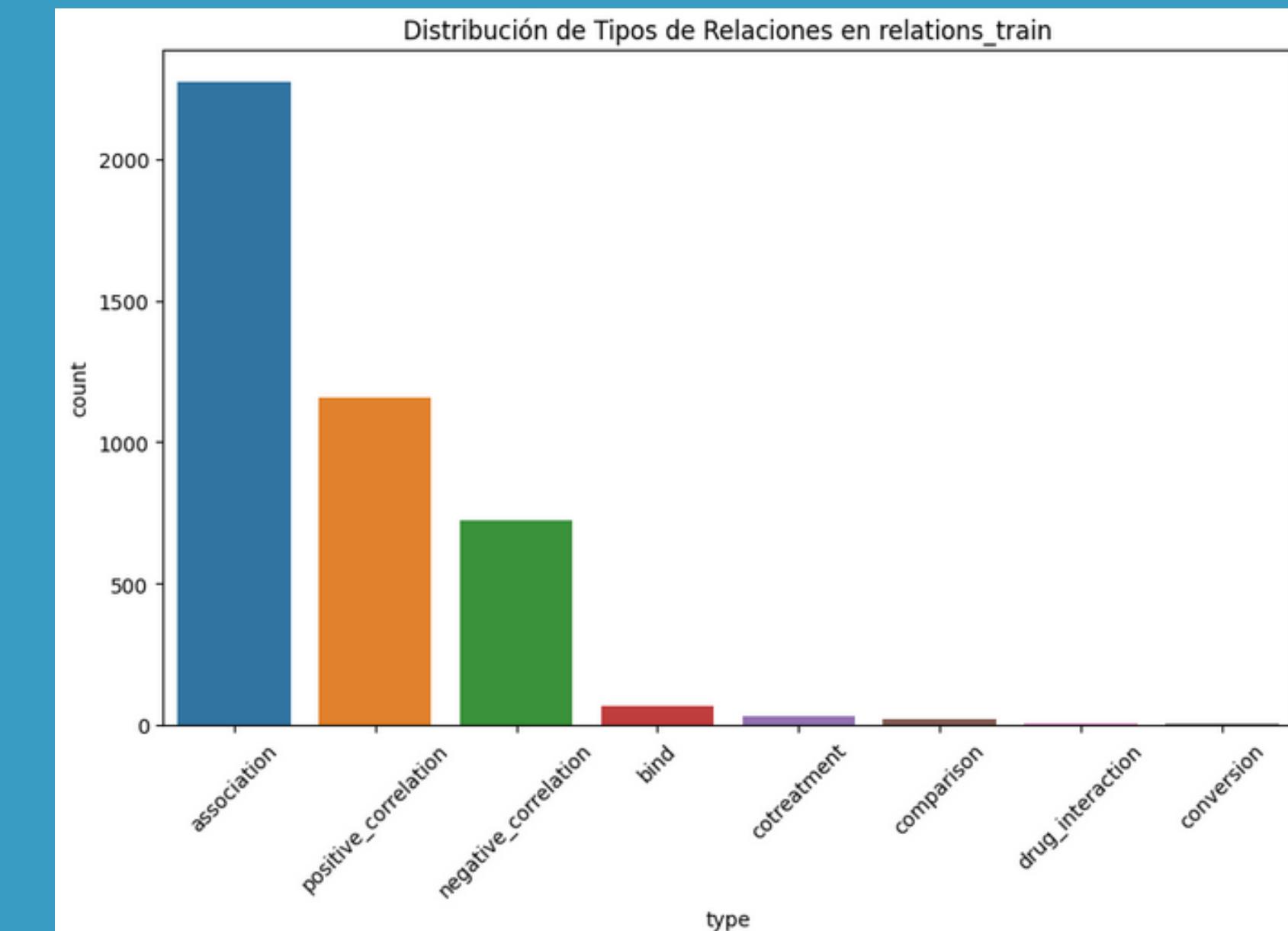


## Analisis Exploratorio



# RELACIONES

```
Tabla de frecuencia de tipos de relaciones:  
association           2274  
positive_correlation 1159  
negative_correlation 721  
bind                 69  
cotreatment          29  
comparison           22  
drug_interaction     3  
conversion            3  
Name: type, dtype: int64
```

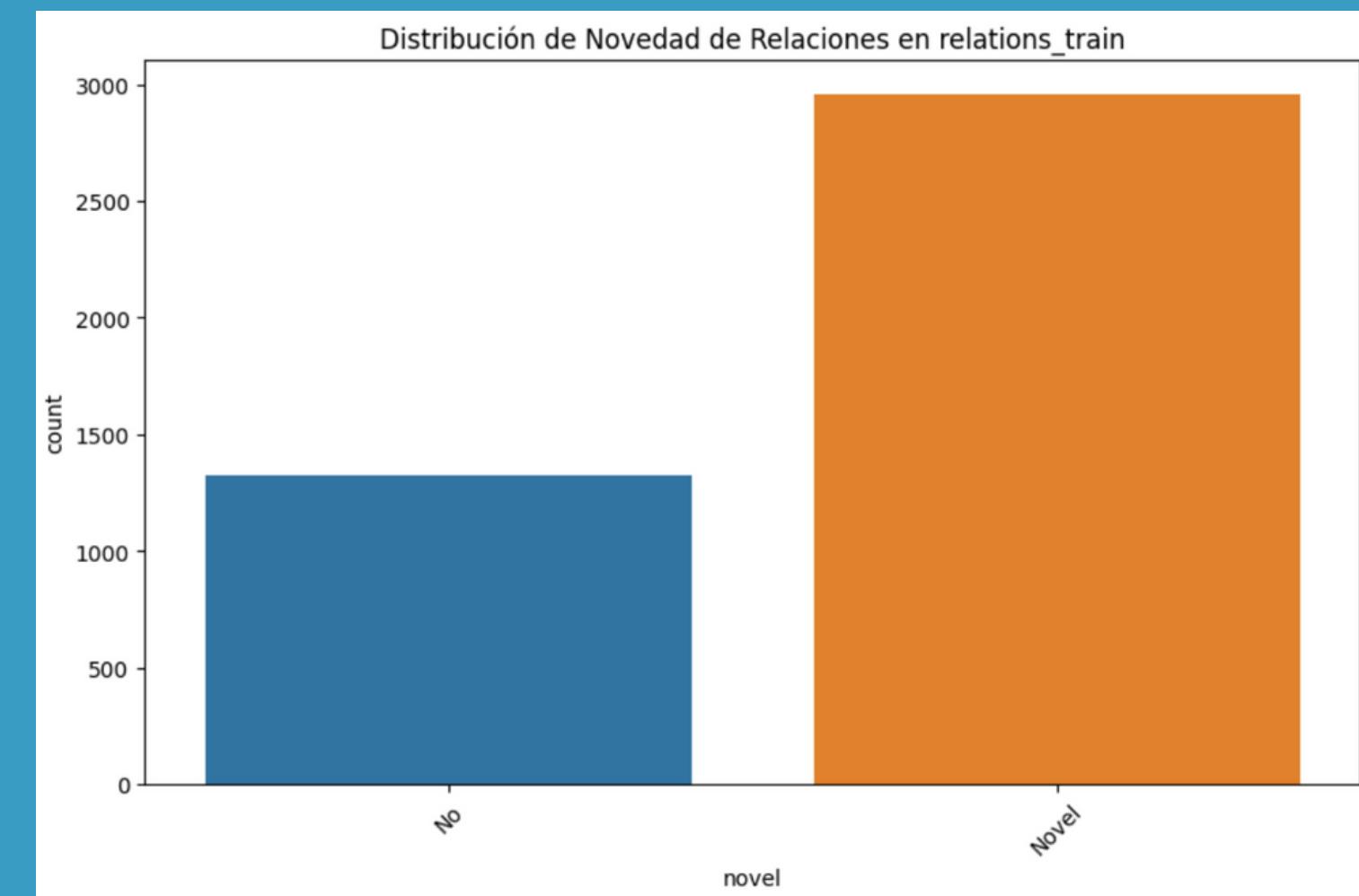




Analisis Exploratorio



# RELACIONES





# CONCLUSIONES

- Los datos de entidades y relaciones proporcionan información importante, lo que es esencial para el problema planteado.
- Los datos de texto contienen una variedad de temas y requieren técnicas avanzadas de NLP.
- Se podría considerar la implementación de la clasificación de texto o la extracción de información, para abordar el desafío de identificar relaciones en los textos de los artículos.





Analisis Exploratorio

# GRACIAS



Grupo 8

