



datos

Procesamiento y

Clasificación de

– 🗆 ×

CLASIFICACIÓN DE TWEETS MISÓGINOS Y NO MISÓGINOS



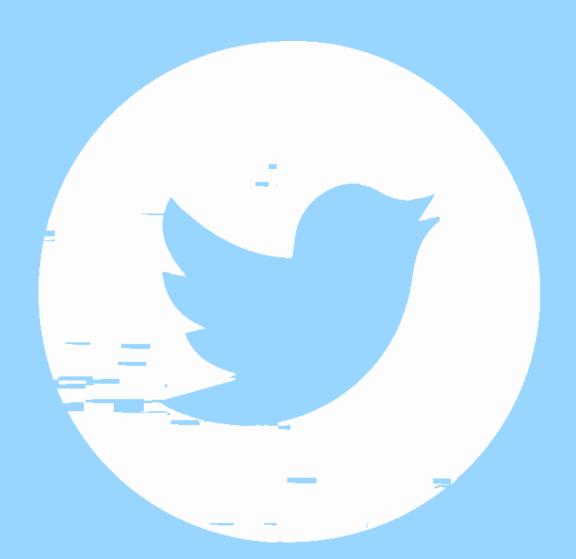
FERNANDO BARAJAS RAMÍREZ

LEOBARDO GARCÍA REYES

DAVID EDUARDO GALLARDO







Introducción

- Según la RAE misoginia se define como la 'Aversión a las mujeres'. Procede del griego misogynía, de miso- 'odio' + gyné 'mujer'.
- Por definición, la misoginia se refiere a "...conductas de odio hacia la mujer y se manifiesta en actos violentos y crueles contra ella por el hecho de ser mujer".
- Estas conductas prevalecen en varias culturas y sociedades, y se refleja en la práctica de subordinación, violencia e incluso crímenes contra la mujer.
- Es importante continuar hablando de este tema para la prevención de la violencia de género. Una de las prácticas relacionadas a este comportamiento es el uso de vocablos, expresiones y sentencias machistas, misóginas y/o que atacan directamente una mujer por su género.





Plantamiento del problema

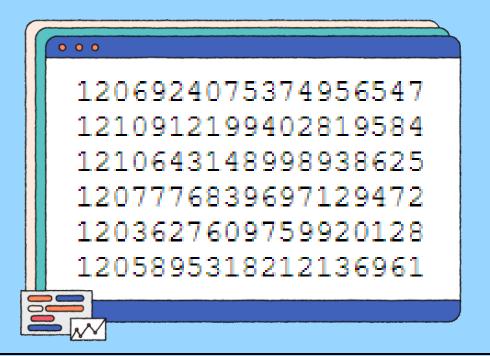
A partir de los tweets extraídos se busca determinar mediante diversos modelos predictivos, si un tweet corresponde a una sentencia que puede ser clasificada como misógina o no y de esta manera poder medir la violencia contra las mujeres en las redes sociales las cuales podrían servir como punto de referencia para medir el "nivel" de misoginia que existe en los países actualmente.





Extracción de los datos

Los IDs de cada tweet, junto con la etiqueta o clasificación inicial, fueron obtenidos del articulo "Detecting misogyny in Spanish tweets. An approach based on linguistics features and word embeddings".



A partir de estos IDs, se utilizó la API de Twitter para extraer los tweets con información extra del usuario que lo publicó. Se extraen 10,244 tweets.

ID	Tweet	Creado_tweet	Locacion_usuario	Seguidores_usuario	Amigos_usuario	Favoritos_usuario	Descripcion_usuario	Verificado_usuario	Idioma	Tipo
1206248589992685568	- Ya todos saben	2019-12-15 16:2		1488	1504	11098	Denme su energía para	FALSO	es	no_misogino
1231950470761078784				82	220	281		FALSO	es	misogino
1204807821637931008	Un condenado p	2019-12-11 16:5	Región de Murcia	148104	1946	5390	El periódico de actualid	VERDADERO	es	no_misogino
1205162379451011072				352281	423	2831	El sitio de los inconforn	VERDADERO	es	no_misogino
1211066868179570688	@GretaThunber	2019-12-28 23:3	España	234	488	3444	El mundo de hoy deja t	FALSO	es	misogino

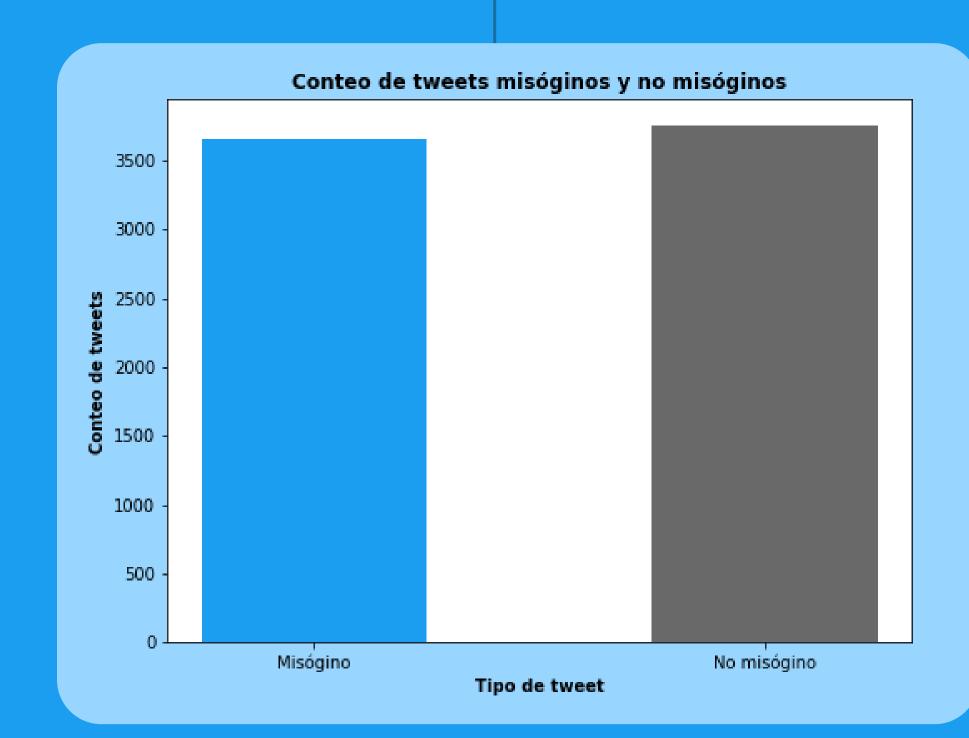




Distribución de tweets misóginos/no misóginos Q



Los 10,244 tweets que se extrajeron, había IDs duplicados y tweets en otro idioma. Se realizó un filtro quedando un total de 7,418 tweets de los cuales 3,661 están etiquetados como misóginos y 3,757 están etiquetados como no misóginos.

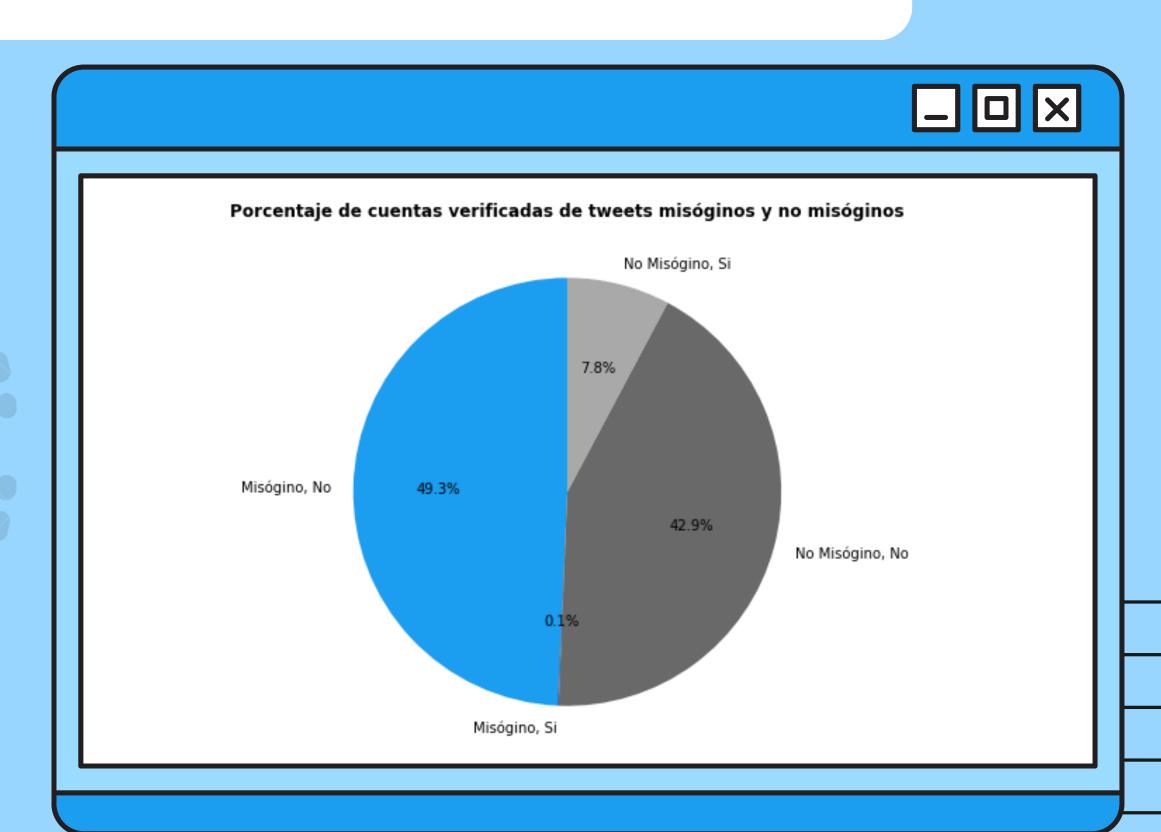






Porcentaje de cuentas verificadas Q

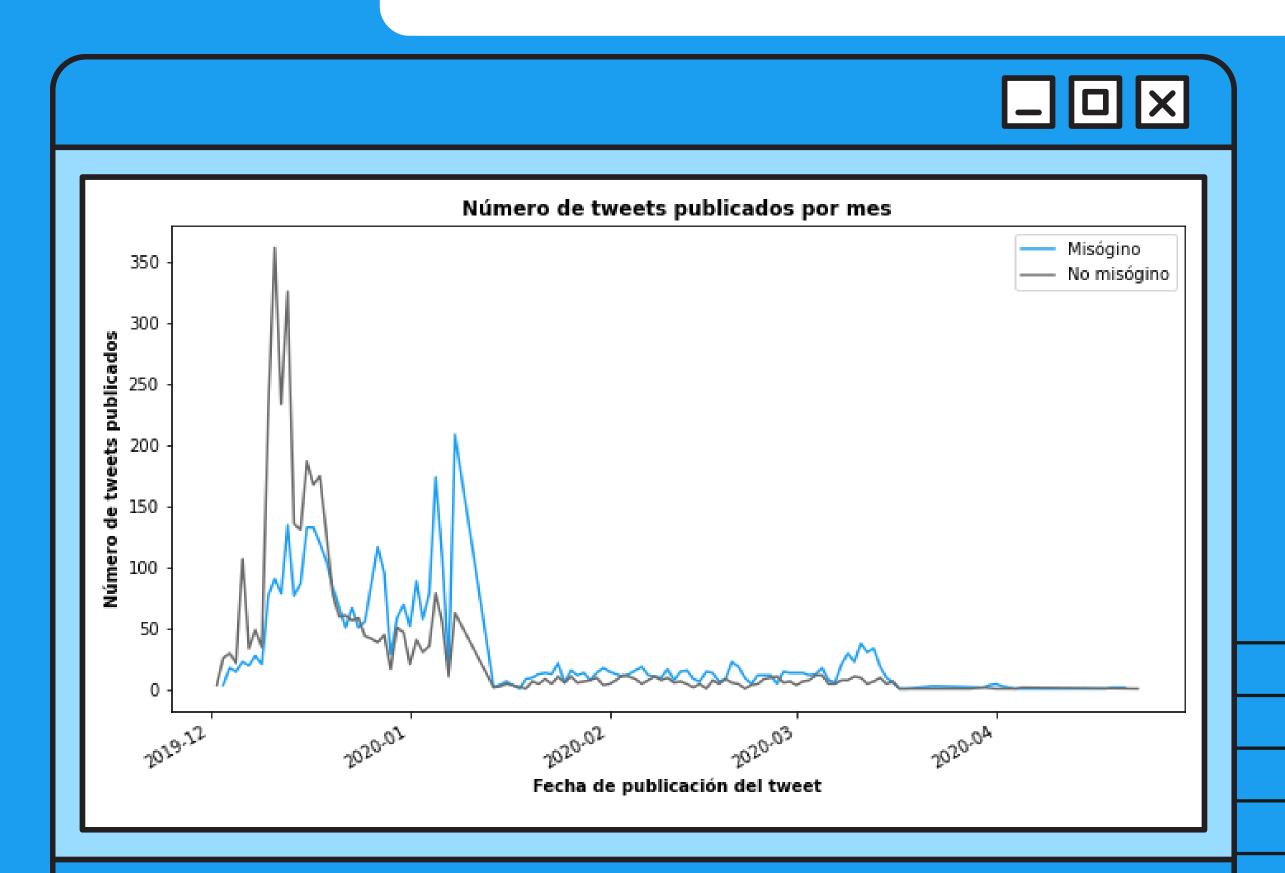
El porcentaje de tweets misóginos es de 49.35% y para los tweets no misóginos es de 50.64%







Publicación de los tweets



El tweet más antiguo es del 02/12/2019 y el más reciente fue del 23/04/2020





Limpieza de texto







- Convertir todo el texto a minúsculas
 Eliminación de usuarios
- 3. Eliminación de páginas web
- 4. Eliminación de texto en paréntesis
- 5. Eliminación de signos de puntuación
- 6. Eliminación de números
- 7. Eliminación de acentos
- 8. Eliminación de emojis
- 9. Eliminación de espacios en blanco múltiples











@Famelica_legion @navedelmisterio@GretaThunberg Greta es una niña histérica ymanipulada (hola) ☎ 10. https://t.co/BBE5XPHAKS



@famelica_legion @navedelmisterio @gretathunberg greta es una niña histérica y manipulada (hola) # 10. https://t.co/bbe5xphaks



greta es una niña histérica y manipulada (hola) 🕁 10.

greta es una niña histerica y manipulada



greta es una niña histerica y manipulada 🕸



greta es una niña histérica y manipulada 🕸 10







Stopwords y Lematización 9



def Toke_StopW_POS(texto limpio):

Tokenizar texto

Eliminar stopwords

Asignación de tipo de palabra

return(texto tokenizados)

Lista de algunas stopwords

de, la, que, el, en, y, a, los, del, se, las, por, un, para, con, no, una, su, al, lo, como, más, pero, sus, le, ya, o, este, sí, porque, esta, entre, ...

También se elimino palabras menores a 1 caractér





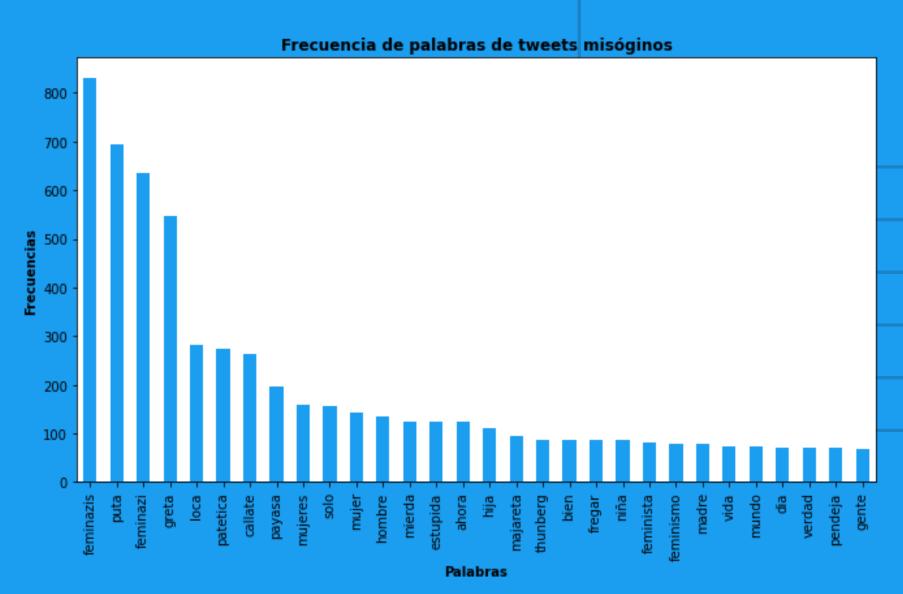
Tweets misóginos











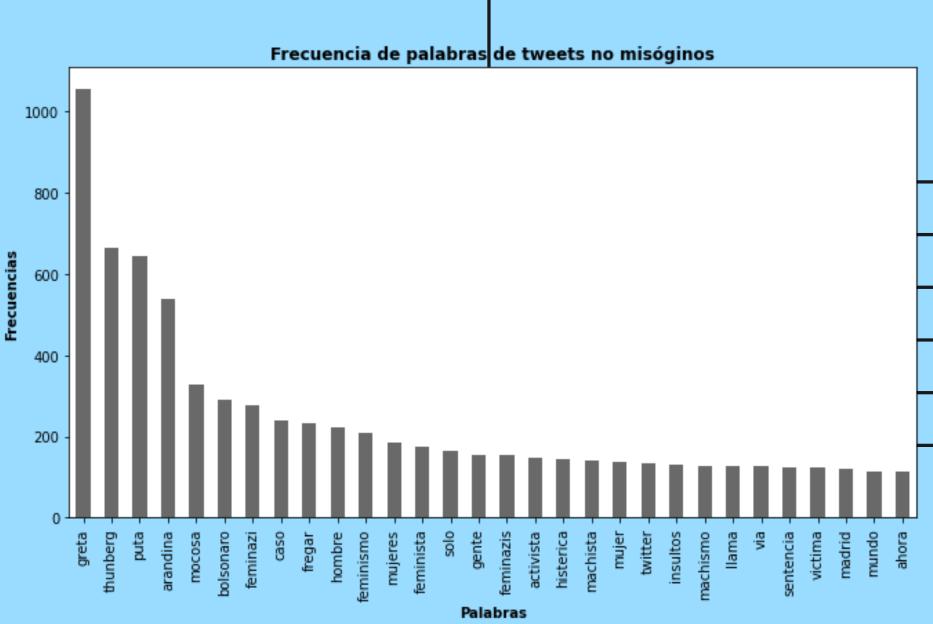




Tweets no misóginos









n-gramas



Q

Se género un unigrama y bigrama y es útil en muchas aplicaciones de análisis textual en que la secuencia de palabras es relevante.

Misógino Unigrama	Misógino Bigrama	No Misógino Unigrama	No Misógino Bigrama
patetica	greta majareta	voto	feminismo radical
callate	greta thunberg	religion	libertad expresion
feminazis	callate puta	silencio	violencia mujer



Modelos



2

Se hizo el entrenamiento con 4 modelos diferentes los cuales dieron resultados bastante buenos que se mencionan mas adelante. Los modelos usados fueron:

- Random Forest.
- Linear SVC.
- Multinomial NB.
- Logistic Regression.

De estos 4 modelos se eligirá el mejor y se utilizará para predecir nuevos valores.

Para predecir se escogió un conjunto de prueba del 35% del conjunto completo.

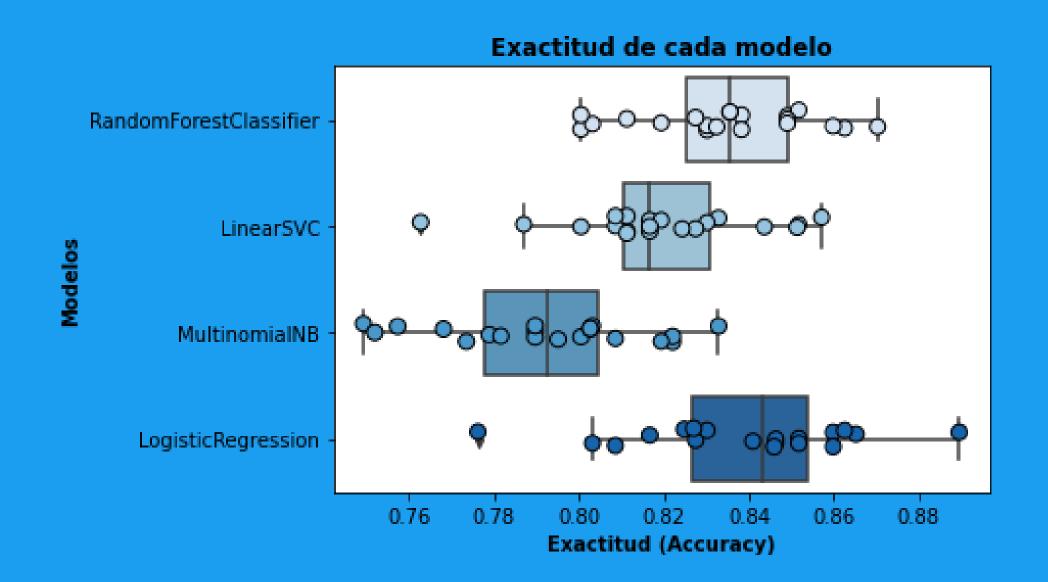
También, se utilizó una validación curzada de 20.





Resultados - Mejor modelo 9

El modelo de mejor exactitud fue el de Regresión Logística.



Modelo	Exactitud		
LinearSVC	0.819498		
LogisticRegression	0.838096		
MultinomialNB	0.791995		
RandomForestClassifier	0.834731		

Diferencia de exactitud del modelo Random Forest y Logistic Regression





Prueba Wilcoxon

H0: No existe diferencia en la diferencia de exactitud del modelo Random Forest y Logistic Regression.

H1: Existe diferencia en la diferencia de exactitud del modelo Random Forest y Logistic Regression.

p-valor	Conclusión	Intervalo de Confianza
0.3646	N.R. <i>H</i> 0	-0.013481093 0.006663832



Conclusión:

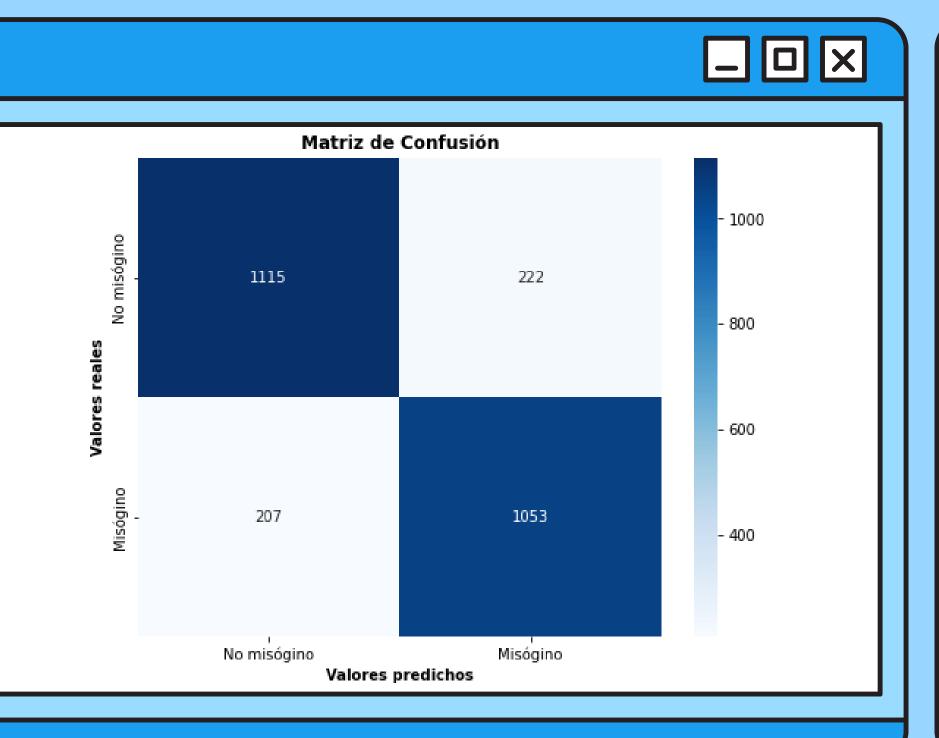
Los datos dan evidencia estadística de que No existe diferencia en la diferencia de exactitud del modelo Random Forest y Logistic Regression con una confianza del 95%.

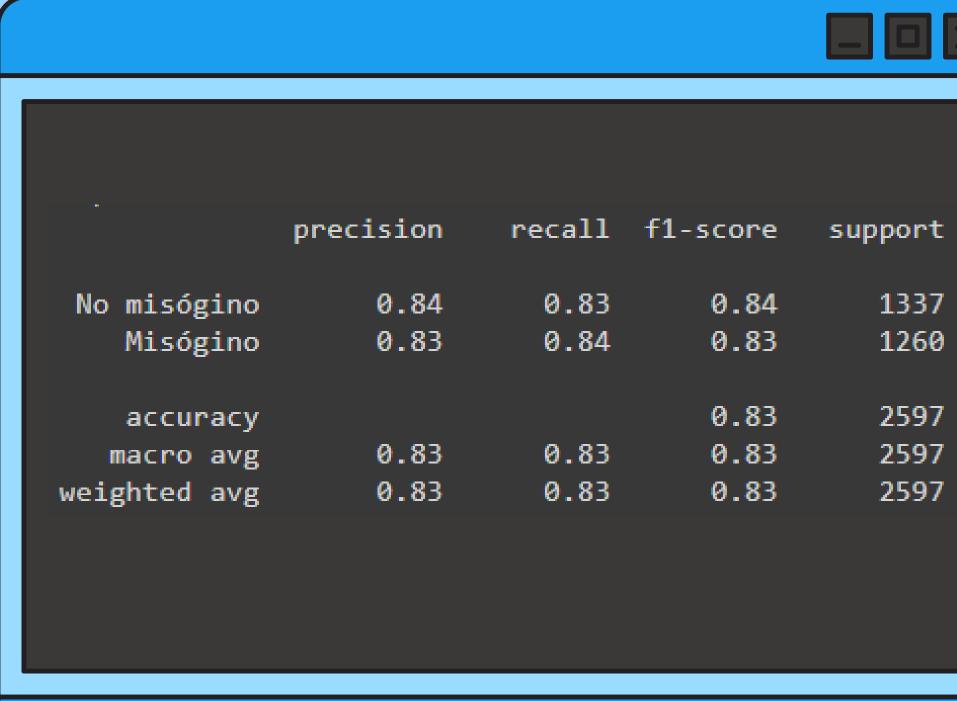
Se realizó lo mismo para una prueba parametrica, con prueba de hipotesis de varianzas y después de medias.





Resultados - Matriz de Confusión Q







Resultados - Predicción



Tweet	Tipo	Predicción
Greta la majareta nunca ha visitado ni se ha quejado en los paises q más contaminan dl mundo, una marioneta d la extrema izquierda comunista.	Misógino	Misógino
Entonces llega el momento en que todo fascista se convierte en hincha del Albacete, sigue todos los partidos de la Arandina y, sobre todo, es un entendido de Ópera, conocedor y admirador del vibrato de Plácido Domingo. De toda la vida. https://t.co/GqmQacy5k3	No Misógino	No Misógino
@monasterioR @julia_otero Tienes mucha pinta tu de agricultora Tienes mas pinta de cerda que de ganadera.	Misógino	Misógino
A mi lo de Greta me parece espectacular. ¿Qué es un monigote? Puede. ¿Qué es un producto de marketing? Quizás. ¿Qué está haciendo que todo el mundo hable de ecologismo por primera vez en la puta historia? Repito: todo el mundo. Totalmente sí.	No Misógino	No Misógino
Respirar también contamina según Greta Thunberg me cago en su vida zorra https://t.co/4Cemt5BARQ	Misógino	No Misógino



Conclusión







- Se obtienen buenos resultados (>0.8), particularmente con el modelo de regresión lineal.
- Se observo que muchas de las palabras utilizadas por tweets clasificados como misóginos son parecidas o iguales a las de no misóginos.
- Consideramos que parte de esto se debe a la red social, puesto que se basa en publicaciones (tweets) y respuestas.
- En diversas ocasiones, las respuestas incluyen las mismas palabras que las publicaciones originales.
- A pesar de esto, se observan buenos resultados.
- Se muestra que el uso de unigramas y bigramas para los modelos ofrece suficiente diferenciación entre un tipo de contenido y otro.
- Con el conjunto VARW, se obtienen de igual forma buenos resultados (>0.8).





GRACIAS!:D

