



# Índice general

<b>1. Datos: Análisis y Preproceso</b>	<b>5</b>
1.1. División entre Train y Dev . . . . .	5
1.2. Distribución de las clases en cada conjunto . . . . .	5
1.3. Descripción del preproceso . . . . .	5
1.4. Primeros resultados . . . . .	6
1.5. Descripción del Proceso de Submuestreo o Sobremuestreo . . . . .	6
1.6. Instrucciones . . . . .	6
<b>2. Algoritmos, link a la documentación y nombre de los hiperparámetros empleados</b>	<b>7</b>
2.1. Algoritmos empleados: Breve Descripción . . . . .	7
2.2. Resultados sobre el Development . . . . .	8
2.2.1. Optimizando los resultados de la clase negativa . . . . .	8
2.2.2. Discusión . . . . .	8
2.2.3. Sin optimizar ninguna clase en particular . . . . .	8
2.2.4. Discusión . . . . .	8
2.3. Conclusión . . . . .	9
<b>Bibliografía</b>	<b>9</b>
<b>3. Anexos</b>	<b>10</b>
3.1. Anexo 1: Baseline . . . . .	10
3.1.1. Resultados obtenidos . . . . .	10
3.1.2. Preproceso e hiperparámetros . . . . .	11
3.2. Anexo 2: Otros modelos probados . . . . .	15

# Índice de figuras

3.1. Resultado KNN . . . . .	10
3.2. Matriz de confusión del Modelo . . . . .	11
3.3. Preproceso de los datos . . . . .	12
3.4. Modelo seleccionado . . . . .	13
3.5. División Train/Test . . . . .	13
3.6. Columnas estudiadas . . . . .	14
3.7. Determinador del mejor modelo . . . . .	14

# Índice de cuadros

1.1. División Train y Dev . . . . .	5
1.2. Distribución Train y Dev . . . . .	5
2.1. Resultados peso de clase negativa . . . . .	8
2.2. Resultados equilibrados . . . . .	8
3.1. Combinaciones probadas . . . . .	15

# Acrónimos

- **LR**: Logistic Regression
- **XGB**: XGBoost
- **MNB**: Multinomial Naive Bayes
- **BoW**: Bag of Words
- **Tf-Idf**: Term frequency – Inverse document frequency

# 1. Datos: Análisis y Preproceso

Se puede encontrar el original en <https://es.overleaf.com/project/626b9937f5f5d5274f042ac7>

## 1.1. División entre Train y Dev

Conjunto De Datos	% de instancias	Num. de instancias
Train	80	9599
Dev	20	2400

1.1. Cuadro: División Train y Dev

## 1.2. Distribución de las clases en cada conjunto

Conjunto De Datos	Clase Neg	Clase Neutra	Clase Pos.
Train	5567	2304	1728
Dev	1392	576	432

1.2. Cuadro: Distribución Train y Dev

## 1.3. Descripción del preproceso

La aplicación de las técnicas que vamos a implementar es importante para el análisis de sentimiento. Estas nos permiten identificar patrones y tendencias en las opiniones de los clientes.

En primer lugar, hemos aplicado preprocesos como convertir los valores a unicode (útil cuando hay caracteres especiales o acentos), convertir las fechas en valores de tiempo UNIX, mapeado de valores, e imputación de valores faltantes.

En segundo lugar, hemos hecho un preprocesado del lenguaje natural para reducir la complejidad del texto y mejorar la eficacia del algoritmo. Primeramente, se ha realizado una traducción de los emoticonos a texto real. Esto se debe a que hemos pensado que los emoticonos podrían aportar información para valorar los tuits, aunque existen configuraciones en las que hemos decidido que los emoticonos se eliminen para ver si se incrementaban las figuras de mérito. Seguidamente, se han eliminado signos de puntuación y se han cambiado las mayúsculas por minúsculas (es decir, se ha normalizado el texto) para que el algoritmo no distinga entre palabras iguales con mayúsculas diferentes. Al eliminar palabras como los stop words, se reduce la dimensión del conjunto de datos y se mejora la precisión del análisis. Por último, hemos lematizado las palabras para reducir la complejidad del texto, es decir, hemos reducido las palabras para quedarnos con sus raíces.

Finalmente, hemos utilizado dos técnicas de preproceso de lenguaje natural para poder evaluar los tuits. Por un lado, hemos utilizado la técnica de bag of words que se utiliza para convertir el texto en un array de frecuencia de palabras. Esta es importante para convertir el texto en formato que puede ser

procesado por los algoritmos. Por otro lado, hemos utilizado TF-IDF para reducir el peso de las palabras frecuentes, y aumentar el peso de las palabras que son más relevantes para el análisis de sentimientos como 'bueno' y 'malo'.

## 1.4. Primeros resultados

En las primeras versiones del modelo no coordinábamos los vectores de tf-idf entre crear y probar el modelo, por lo tanto los resultados que obteníamos eran bastante mejorables, además solo se utilizaba una versión de Naive-Bayes y no existía la opción de hacer BoW ni de borrar los emoticonos, solo traducirlos. Empleando esta información el f-score obtenido rondaba entre 0.44 y 0.49.

## 1.5. Descripción del Proceso de Submuestreo o Sobremuestreo

Inicialmente nuestro algoritmo tenía problemas cuando la muestra era demasiado pequeña, por lo que a la hora de determinar si el modelo va a emplear over o under sampling este observa la cantidad de instancias. Para más de 10000 instancias se utilizará Oversampling, y para menos Undersampling, para evitar situaciones como la ocurrida con las instancias dadas para la preparación de este proyecto, con las cuales el test solo contenía una única instancia de clase neutral por lo que no era capaz de generar un modelo.

## 1.6. Instrucciones

- Ejecutar installer.sh para instalar todos los paquetes necesarios para el script.
- El fichero README.md contiene las llamadas a realizar para generar o probar el modelo respectivamente.

## 2. Algoritmos, link a la documentación y nombre de los hiperparámetros empleados

### 2.1. Algoritmos empleados: Breve Descripción

Hemos empleado los siguientes algoritmos con los siguientes hiper-parámetros:

#### ■ Gaussian Naive Bayes:

- Hiperparámetros: priors, var smoothing
- Link: Sklearn GaussianNB

En el caso de Gaussian Naive Bayes, se asume que los datos utilizan una distribución normal (o Gaussiana). Por ello, es importante discretizar los valores, porque algunos valores tienen una gran cantidad de decimales y no están uniformemente distribuidos, lo que puede resultar en un sesgo hacia unos valores en concreto.

#### ■ Mixed Naive Bayes:

- Hiperparámetros: alpha, fit prior, class prior
- Link: mixed-naive-bayes MixedNB

En el caso de Mixed Naive Bayes, se permiten diferentes distribuciones para diferentes características, por lo que se pueden utilizar diferentes funciones de densidad de probabilidad. Por lo tanto, no es necesario discretizar los datos, ya que el modelo de Mixed Naive Bayes puede manejar características continuas y discretas sin problemas.

#### ■ Multinomial Naive Bayes:

- Hiperparámetros: alpha, fit prior, class prior
- Link: Sklearn MultinomialNB

Este modelo se usa para la clasificación de texto y tiene en cuenta la frecuencia de las palabras. Intenta predecir la clase en función de esa frecuencia.

#### ■ Bernoulli Naive Bayes:

- Hiperparámetros: alpha, binarize, fit prior, class prior
- Link: Sklearn BernoulliNB

Este modelo considera la presencia o ausencia de cada característica.

#### ■ Complement Naive Bayes:

- Hiperparámetros: alpha, fit prior, class prior, norm



- Link: Sklearn ComplementNB

Es una variante del NB Multinomial y se ha creado para manejar los conjuntos de datos desequilibrados.

#### ■ Logistic Regression:

- Hiperparámetros: penalty, random state
- Link: Sklearn LogisticRegression

La Regresión Logística utiliza una función de activación para transformar la salida del modelo en una probabilidad. Es una curva que separa las instancias entre sí como una frontera, para determinar si una instancia pertenece a una clase u otra.

Hemos observado durante las pruebas que Logistic Regression tiende a clasificar más instancias como positivas que los modelos basados en Naive-Bayes. Y que de entre las distintas versiones de Naive-Bayes, la que consistentemente obtenía mejores resultados ha sido **Bernoulli Naive Bayes**.

## 2.2. Resultados sobre el Development

En esta sección se presentan los resultados obtenidos para el development.

### 2.2.1. Optimizando los resultados de la clase negativa

Algoritmo	Combinación hiperparámetros	Prec	Rec	F-sco
Logistic Regression	tf-idf y Traducir Emojis/Emoticonos	0.73	0.72	0.73

2.1. Cuadro: Resultados peso de clase negativa

### 2.2.2. Discusión

El f-score correspondiente a los comentarios negativos ha sido de 0.75. Esta composición ha sido la que ha obtenido mejores resultados de manera consistente entre las pruebas de los distintos miembros del grupo. Es la versión en la que originalmente más confiábamos ya que es la que más cantidad de información tiene en cuenta a la hora de determinar las clases.

### 2.2.3. Sin optimizar ninguna clase en particular

Algoritmo	Combinación hiperparámetros	Prec	Rec	F-sco
Logistic Regression	tf-idf y Traducir Emojis/Emoticonos	0.73	0.72	0.73

2.2. Cuadro: Resultados equilibrados

### 2.2.4. Discusión

Esta composición es igual a la mencionada en el apartado anterior por lo que no se considera necesario volver a desarrollarla. Cabe destacar que en una de las pruebas realizadas por los integrantes, la composición de hiperparámetros Logistic Regression, BoW y Borrar Emojis/Emoticonos, ha conseguido los mismos valores para Precision, Recall y F1-Score. De esto deducimos que algunos de los modelos que se mencionarán en el Anexo 2 podrían teóricamente obtener valores superiores a los que aparecen aquí.

## 2.3. Conclusión

El mismo modelo obtiene los mejores resultados en ambos apartados de este punto, pero eso no significa que siempre vaya a ser el mejor a la hora de distinguir las clases. Por esto además de entregar este modelo enviaremos también aquel que empleando Naive-Bayes obtiene el mejor resultado (Naive-Bayes, tf-idf y Traducir Emojis/Emoticonos). Aparentemente el preproceso tiene más influencia que el algoritmo a la hora de obtener resultados, aunque Logistic Regression ha sido consistentemente superior a Naive-Bayes a lo largo de todas las combinaciones.

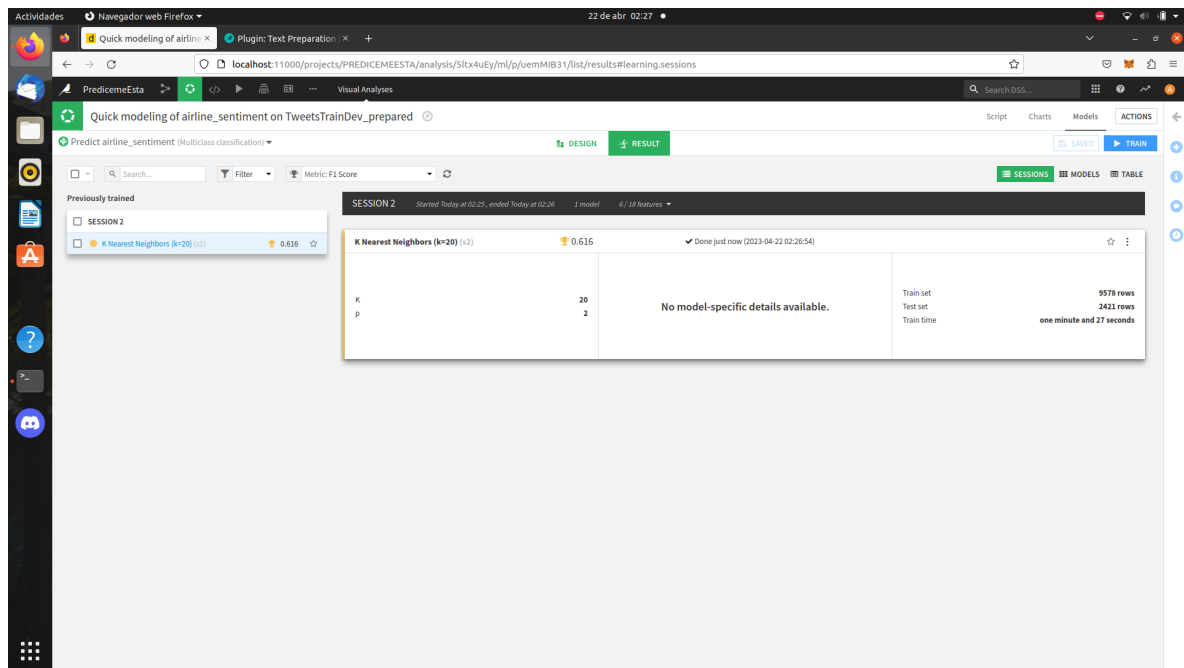
## 3. Anexos

### 3.1. Anexo 1: Baseline

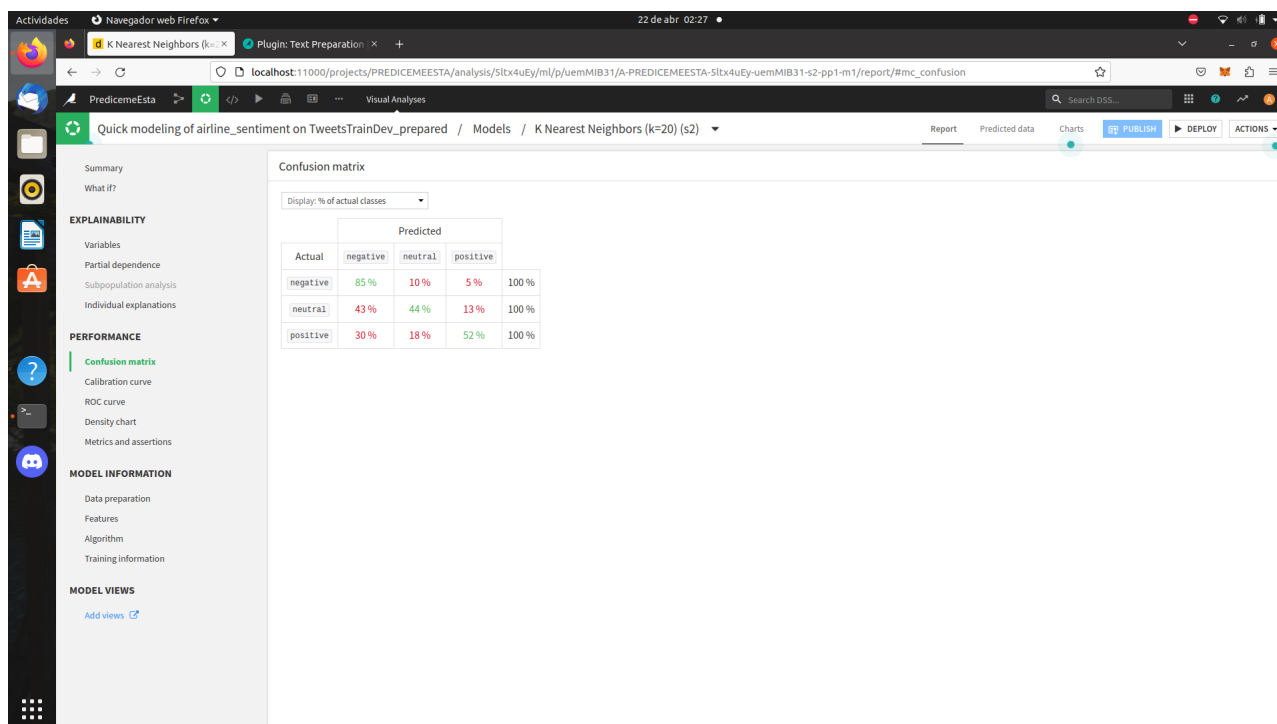
Para comprobar la efectividad de nuestro modelo frente a otros, hemos generado un Baseline empleando un modelo knn generado mediante la aplicación de Dataiku.

#### 3.1.1. Resultados obtenidos

Empleando un modelo básico KNN hemos obtenido un F1-Score del 62% aproximadamente (Fig 3.1).



3.1. Figura: Resultado KNN



3.2. Figura: Matriz de confusión del Modelo

### 3.1.2. Preproceso e hiperparámetros

Dado que la tarea no trata de optimizar el modelo KNN, el preproceso utilizado no ha sido muy riguroso, y se ha limitado a características básicas para permitir generar un BoW (Fig 3.3):

- Se simplifica el texto eliminando símbolos de puntuación.
- Se eliminan las 'stopwords' que no aportan información relevante.
- Se tokenizan las palabras y se reducen a su raíz.
- Seleccionamos algunas palabras que pueden denotar emociones positivas o negativas y las contamos.

Script output on dataset sample

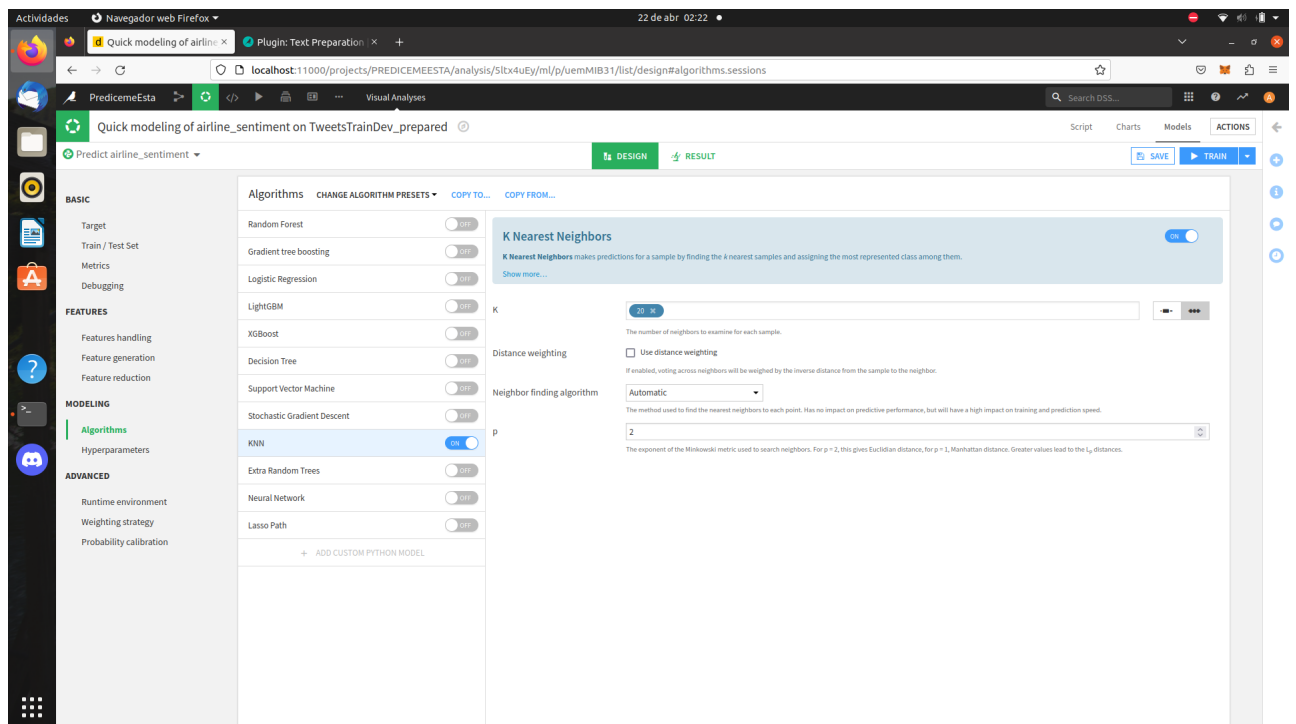
retweet_count	text	textsimpl	texttkn	badcount	thankcount	delaycount	tweet_coord	tweet_date
0	@VirginAmerica What @dhepburn said.	dhepburn said virginamerica	["dhepburn","said","virginamerica"]	0	0	0		2015-02-...
0	@VirginAmerica plus you've added commercials to ...	ad commercer experi plus tacki virginamerica	["ad","commerc","experi","plus","tacki","virginamer..."]	0	0	0		2015-02-...
0	@VirginAmerica I didn't today... Must mean I need t...	mean need today trip virginamerica	["mean","need","today","trip","virginamerica"]	0	0	0		2015-02-...
0	@VirginAmerica It's really aggressive to blast obno...	aggress amp blast entertain face guest litt obnos...	["aggress","amp","blast","entertain","face","guest",...]	0	0	0		2015-02-...
0	@VirginAmerica and it's a really big bad thing abou...	bad big reali thing virginamerica	["bad","big","reali","thing","virginamerica"]	1	0	0		2015-02-...
0	@VirginAmerica seriously would pay \$30 a flight for...	30 bad flight fl pay play reali seat serious thing va...	["30","bad","flight","fl","pay","play","reali","seat","s..."]	1	0	0		2015-02-...
0	@VirginAmerica yes, nearly every time I fly VX this "...	away ear fl go near time virginamerica vx worm yes	["away","ear","fl","go","near","time","virginamerica"...]	0	0	0		2015-02-...
0	@VirginAmerica Really missed a prime opportunity...	co hat https men miss mwpg/greep opportun pare...	["co","hat","https","men","miss","mwpg/greep","op..."]	0	0	0		2015-02-...
0	@VirginAmerica Well, I didn't...but NOW I DO! :-D	now virginamerica well	["now","virginamerica","well"]	0	0	0		2015-02-...
0	@VirginAmerica it was amazing, and arrived an ho...	amaz arriv earli good hour virginamerica	["amaz","arriv","earli","good","hour","virginamerica"]	0	0	0		2015-02-...
0	@VirginAmerica did you know that suicide is the se...	10 24 caus death know lead second suicid teen virg...	["10","24","caus","death","know","lead","second","su..."]	0	0	0		2015-02-...
0	@VirginAmerica I &lt;3 pretty graphics, so much bet...	3 better graphic iconographi lt minim much pretti ...	["3","better","graphic","iconographi","lt","minim","..."]	0	0	0		2015-02-...
0	@VirginAmerica This is such a great deal! Already L...	1st 2nd amp australia deal even gone great p think ...	["1st","2nd","amp","australia","deal","even","gone","..."]	0	0	0		2015-02-...
0	@VirginAmerica @virginmedia I'm flying your #fab...	again ahbhkhijn away co fabul fl http seduct sky s...	["again","ahbhkhijn","away","co","fabul","fl","http",...]	0	0	0		2015-02-...
0	@VirginAmerica Thanks!	thank virginamerica	["thank","virginamerica"]	0	1	0		2015-02-...
0	@VirginAmerica SFO-PDX schedule is still MIA.	mia pdx schedul sfo virginamerica	["mia","pdx","schedul","sfo","virginamerica"]	0	0	0		2015-02-...
0	@VirginAmerica So excited for my first cross countr...	29daystogo america countri cross excit first flight g...	["29daystogo","america","countri","cross","excit","fir..."]	0	0	0		2015-02-...
0	@VirginAmerica I flew from NYC to SFO last week a...	flew fulli gentleman help larg last nyc seat sfo side ...	["flew","fulli","gentleman","help","larg","last","nyc",...]	0	0	0		2015-02-...
0	♥ Flying @VirginAmerica 🙌🙌	fl virginamerica	["fl","virginamerica"]	0	0	0		2015-02-...
0	@VirginAmerica you know what would be amazing!...	amaz awesom bos fl fl know pleas virginamerica ...	["amaz","awesom","bos","fl","fl","know","pleas","vi..."]	0	0	0		2015-02-...
0	@VirginAmerica why are your first fares in May over...	avall carrier fare first seat select three time virgina...	["avall","carrier","fare","first","seat","select","three",...]	0	0	0		2015-02-...
0	@VirginAmerica I love this graphic. http://t.co/UTS...	co graphic http love ut5grwaaa virginamerica	["co","graphic","http","love","ut5grwaaa","virginam..."]	0	0	0	[40.74804263, -73...	2015-02-...

Job succeeded.

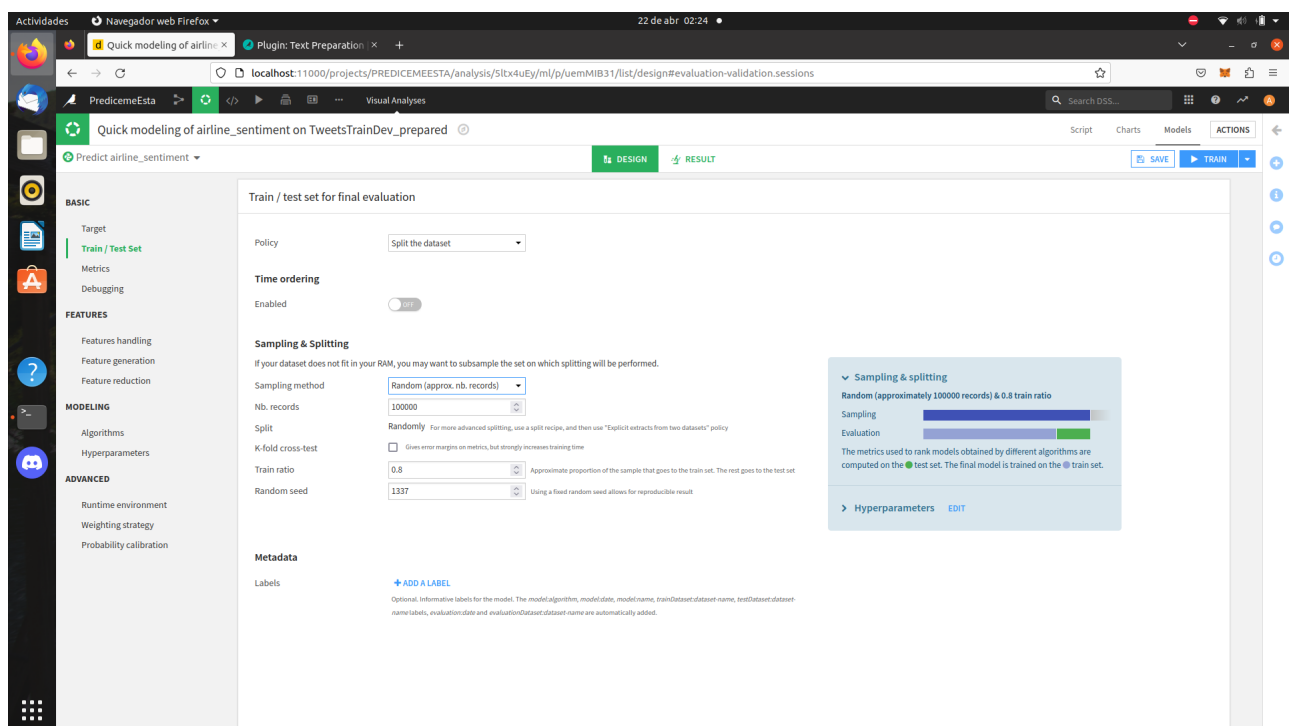
3.3. Figura: Preproceso de los datos

Además de esto la configuración del método a utilizar ha sido la siguiente:

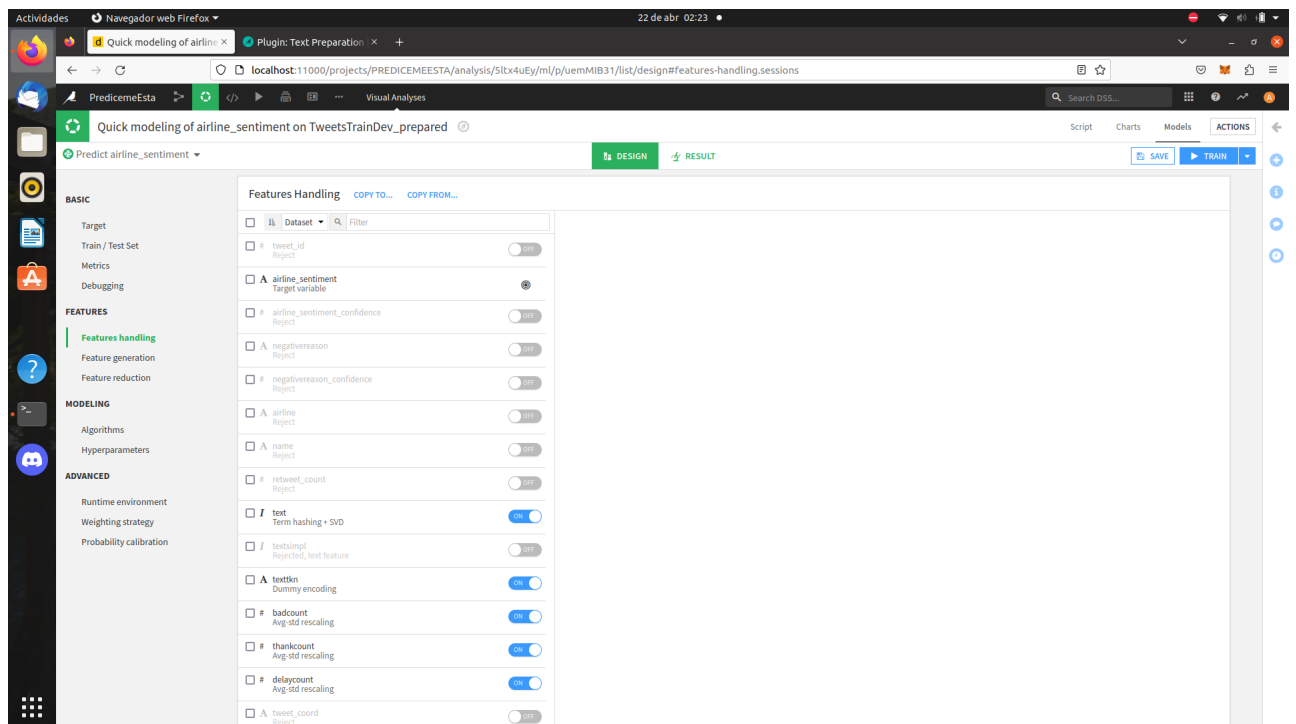
- Se han tenido en cuenta los 20 vecinos más cercanos para determinar la clase (Fig 3.4).
- La división train/test ha sido un 80/20, y se han seleccionado las filas al azar (Fig 3.5).
- Las columnas que se han tenido en cuenta han sido el texto, las palabras tokenizadas y los contadores de palabras (Fig 3.6).
- Se ha empleado el F1-Score para determinar la mejor aproximación (Fig 3.7).



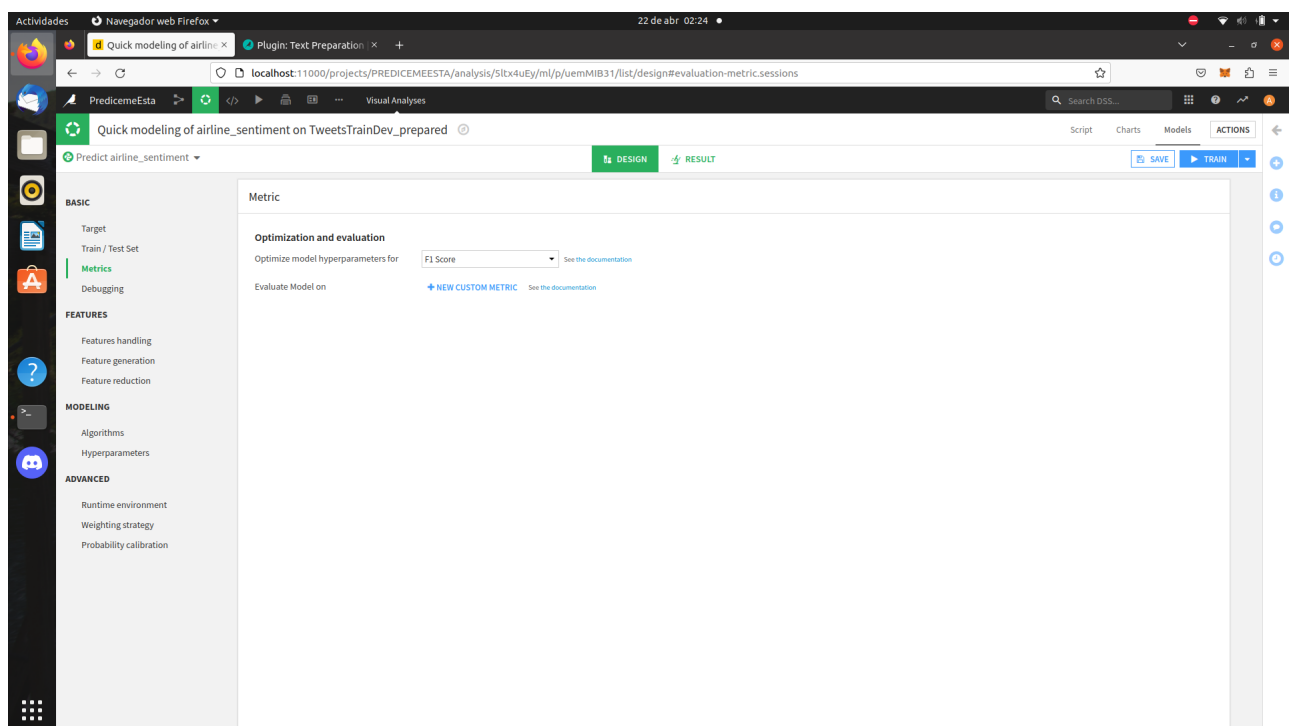
3.4. Figura: Modelo seleccionado



3.5. Figura: División Train/Test



3.6. Figura: Columnas estudiadas



3.7. Figura: Determinador del mejor modelo

## 3.2. Anexo 2: Otros modelos probados

A lo largo de la práctica hemos probado distintas variaciones en el algoritmo empleado y en el preproceso de los datos en busca de un mayor f-score. Los dos algoritmos empleados no tenían hiperparámetros personalizables. Adicionalmente, el script que genera los modelos **filtra automáticamente entre todas las versiones de Naive-Bayes** para elegir la que ofrezca un mayor f-score, por lo tanto la cantidad de tests ha sido menor que la de otros grupos.

Las combinaciones que hemos probado han sido las siguientes:

Algoritmo	Características preproceso	Prec	Rec	F-sco
Naive-Bayes	BoW, Traducir Emojis/Emotes	0,6773	0,6707	0,6618
Naive-Bayes	BoW, Borrar Emojis/Emotes	0,6949	0,6919	0,6827
Naive-Bayes	tf-idf, Traducir Emojis/Emotes	0,7017	0,6959	0,6880
Naive-Bayes	tf-idf, Borrar Emojis/Emotes	0,6982	0,6894	0,6817
Logistic Regression	BoW, Traducir Emojis/Emotes	0,7266	0,7211	0,7225
Logistic Regression	BoW, Borrar Emojis/Emotes	0,7039	0,6959	0,7225
Logistic Regression	tf-idf, Borrar Emojis/Emotes	0,7107	0,6968	0,6992

3.1. Cuadro: Combinaciones probadas

Además de estas pruebas realizadas, había una versión del script que realizaba un barrido de parámetros para Logistic Regression, pero el proceso de generar un modelo tardaba una hora aproximadamente y el resultado obtenido era muy similar al que obteníamos sin realizarlo, por lo que no hemos considerado que merezca la pena seguir adelante con esa versión.