

Tarea 1

IRVING DANIEL ESTRADA LÓPEZ
Facultad de Ciencias Físico Matemáticas
Universidad Autónoma de Nuevo León
 Nuevo León, México
 irving.estradalo@uanl.edu.mx

I. INTRODUCCIÓN

El preprocesado de datos es toda acción que se lleva a cabo antes de que empiece el proceso de análisis de datos. Se tienen que transformar los datos crudos en un nuevo conjunto de datos que tenga la información valiosa de nuestros datos reales, en la mayoría de las veces tenemos que adecuar nuestro conjunto de datos para dicho análisis. Los conjuntos de datos vienen con ruido, datos nulos, incompatibilidad entre los datos, provienen de múltiples fuentes de información, contienen datos irrelevantes para nuestro análisis, entre otras cosas. Nos referimos a información valiosa como aquella que aporte significancia a nuestro análisis.

Debemos tener en cuenta los datos con los que estemos trabajando, así como el análisis que se quiere llevar a cabo, ya que esto dependerá del preprocesado que se tiene que realizar para hacer un buen análisis.

En este documento la intención es hacer un preprocesado de texto, hacer un conteo de palabras, extraer características importantes del conjunto de datos e identificar áreas de oportunidad, para en un futuro llevar a cabo alguno de los siguientes análisis: clasificación de texto, análisis de sentimiento, análisis de tópicos o análisis de intención.

Primero comenzaremos con un análisis exploratorio de nuestras variables, de esta manera podemos perfilar hacia donde tiene que ir dirigido nuestro preprocesado. Continuando con un análisis de frecuencias para identificar características relevantes de los textos. Por último, en la parte final de este documento se identificarán áreas de oportunidad que se adecúen al conjunto de datos seleccionado. Se mencionarán posibles análisis que contrubuyan brindando información relevante del conjunto de datos. Hoy en día se genera una gran cantidad de datos diarios, de los cuales podemos extraer información valiosa, siendo el preprocesado un obstáculo para elaborar un buen análisis.

II. DATOS

El conjunto de datos con el que se trabajó en este artículo es de dominio público y fue obtenido de Kaggle. Nuestro conjunto de datos contiene 9,713 registros los cuales fueron extraídos de Flipkart utilizando Selenium y BeautifulSoup. Como podemos ver en la figura 1 el conjunto de datos contiene la valoración que le dan al producto Iphone SE, un comentario y la reseña de los clientes en el e-commerce de la India Flipkart.

Fig. 1.

	Ratings	Comment	Reviews
0	5	Super!	Great camera for pics and videos Battery life ...
1	5	Must buy!	Great device. Let me tell the Pros..1. Superb ...
2	5	Great product	Who all loves older size i.e., 4.7 inch type s...
3	5	Simply awesome	This iPhone SE is the best phone ever you get....
4	5	Classy product	This is my second iphone after iphone 4s. I've...

En la tabla 1 podemos revisar los atributos de nuestro conjunto de datos junto a una descripción, es importante destacar que la variable Ratings va de 1 a 5. Siendo la calificación más baja 1 y siendo la mejor 5.

Variable	Representación
Ratings	La calificación que los clientes le dieron al producto.
Comment	Comentario, es el título de la reseña.
Reviews	La reseña del cliente hacia el producto.

TABLE I
VARIABLES DEL CONJUNTO DE DATOS

III. METODOLOGÍA

A. Análisis Exploratorio

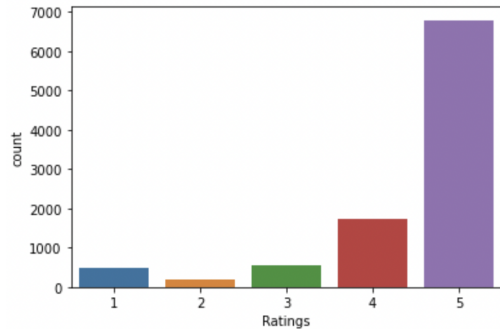
En esta sección extraeremos algunas características descriptivas de las variables y concluiremos en base a ellas. Dos de nuestras tres variables son textos, la variable Comment y la variable Reviews. Es importante mencionar que no hay datos faltantes, ya que si los hubiera deberíamos tomar la decisión si eliminar los registros o llevar a cabo un preprocesado. La variable Ratings es de tipo entero, así que podemos extraer algunas características interesantes de ella.

Fig. 2.

	Ratings
count	9713.000000
mean	4.456399
std	1.032911
min	1.000000
25%	4.000000
50%	5.000000
75%	5.000000
max	5.000000

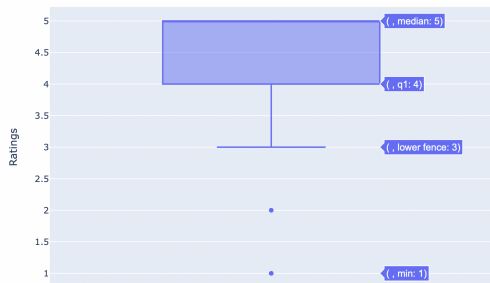
Como podemos darnos cuenta en la figura 2, el promedio de las calificaciones dadas al producto es 4.45, y viendo los cuartiles, más del 50% de nuestros datos cuentan con una calificación de 5. El 75% cuenta con una calificación de 4 o más. Esto nos indica que el producto tiene buenas valoraciones. De los 9,712 registros, 6,788 tiene 5 de calificación, casi el 70% de nuestras calificaciones.

Fig. 3.



Para verlo de una forma más ilustrativa tenemos la figura 3. Donde claramente podemos darnos cuenta de que la calificación que domina es el 5, siguiéndole la calificación 4. Los registros que cuenta con 1 como calificación son 474, menos del 5%. Indicándonos que por lo menos con los registros que tenemos los usuarios le dan en su mayoría una buena calificación.

Fig. 4.



Nos apoyamos también en la figura 4 y efectivamente el producto tiene calificaciones buenas. Incluso está tomando la calificación 1 y 2 como datos anómalos. Esto comprobando que el producto aparentemente para los usuarios es de buena calidad, al menos con los registros que tenemos.

Fig. 5.

'Great camera for pics and videos Battery life is good so far with some setting turn of which i never use and when i use i turn those on and i use it in power saving mode all t he time so a full day with light gaming of 1hr or more using camera for 1hr or more list ening music in my car on youtube and Bluetooth on for couple of hours it gives me full d ay of battery varying from 4/5 to 6/7 hours sot per dayPerformance is top notch plays eve ry game and every task with easeVery premium phone look...READ MORE'

En nuestro conjunto de datos, en los registros de la variable Reviews tenemos el siguiente caso mostrado en la figura 5. Al final de cada una de las reseñas tenemos "...READ MORE" que esto proviene directamente de cuando se extrajeron los datos de Flipkart, esta es una de las cosas que debemos de adecuar en nuestro preprocesado.

B. Preprocesado

Los procesos que se mencionarán a continuación se llevaron a cabo tanto en los registros de la columna Comment como en los de la columna Reviews.

Se llevarán acabo las siguientes acciones en el preprocesado de texto:

Limpieza de texto: Se comienza eliminando los patrones de ruido que trae consigo el texto, en algunas ocasiones por la fuente de datos en la extracción contiene patrones de ruido al inicio o al final de los textos, como mencionamos anteriormente, nuestra variable Reviews tiene en la parte final el texto "...READ MORE", el cual genera ruido dentro de nuestros elementos en nuestra variable. Utilizando la librería "re" y apoyándonos con un ciclo for podemos eliminar dicho ruido, para de esta forma se modifiquen todos nuestros registros y no hacerlo de forma individual.

Remover emojis: Hoy en día es común que las personas utilicen emojis, en este caso los removeremos y nos enfocaremos en las palabras. Hay librerías que nos pueden describir con palabras los emojis, es importante saberlo para preprocesados futuros.

Remover URL: En caso de existir alguna liga que nos lleve a un sitio web, no aporta alguna característica valiosa a los análisis posteriores.

Remover signos de puntuación: Debemos de quitar los signos de puntuación ya que en este caso no aportarán en el conteo de palabras, en otra situación como el análisis de sentimientos nos podrían brindar información valiosa. Esto se debe de hacer después de remover la URL para que no haya problemas, ya que la URL contiene signos de puntuación y todo está concatenado, de lo contrario existe la posibilidad de que se interprete como palabras.

Conversión a minúsculas: Para no tener problemas al momento de obtener la frecuencias de las palabras, convertimos todo nuestro texto en minúsculas. Si no hicieramos este paso, identificaría palabras distintas por el hecho de tener diferencia entre mayúsculas y minúsculas. En algunas ocasiones es importante no hacer este paso, debido a que las mayúsculas y minúsculas nos aportan información, como en el caso de análisis de sentimientos que se relacionan las mayúsculas como un énfasis en las palabras del usuario.

Revisión de ortografía: Este es un factor importante, tomando en cuenta que el sitio de donde fueron extraídos la información es de la India, y su primer idioma no es el inglés. Existen librerías que nos ayudan a la revisión de ortografía, ya que muchas veces nuestros conjuntos de datos contienen faltas de ortografía las cuales pueden afectar nuestros análisis, utilizamos la librería de Python "pyspellchecker" la cual utiliza la distancia de Levenshtein para hayar la palabra más parecida en caso de contener algun error ortográfico.

Tokenizar: Este paso consiste en dividir el texto en las unidades que lo conforman, entendiendo por unidad el elemento más sencillo con significado propio para el análisis en cuestión, en este caso las palabras.

Eliminación de stopwords: Las stopwords son consideradas palabras que no tienen un significado por si solas, sino acompañan a las otras palabras. Es importante removerlas debido a que abarcan una gran cantidad en los textos y al

