

Descripción del problema

En los últimos años se ha presenciado el increíble crecimiento y uso de las redes sociales, blogs y demás medios que guardan principalmente textos, estos datos no estructurados se han convertido en el mayor interés de empresas principalmente, debido a que los datos estructurados no son capaces de mostrar la riqueza o sentimiento como se describe en los comentarios de las personas que expresan su sentir respecto a una infinidad de temas. La pandemia obligó a todos los sectores a evolucionar y ofrecer mejores servicios y a la distancia de un país, poder realizar pagos, manejar saldo, consultar los movimientos de la cuenta, obtener un estado de cuenta, hacer una cda, renovar un servicio, pero principalmente a los bancos, de poder ofrecer la mayoría de los servicios sin necesidad de acudir a una sucursal, por ello, se requiere de una aplicación móvil cada vez más robusta, intuitiva, que cubra estas necesidades y que deje satisfecho al usuario. Por ello, implementamos un análisis de sentimiento y modelado de texto de las reseñas en Google Play y Apple Store orientado a la aplicación móvil de Santander para detectar las debilidades y fortalezas de dicha aplicación.

Objetivos

- 1) Realización de web scrapping a las tiendas de aplicaciones con Python y posteriormente el procesamiento para obtener un conjunto de documentos listos para poder analizarlos de manera correcta. 2) Aplicación del algoritmo de la distancia de Levenstien para corregir las palabras mal escritas por el usuario. 3) Creación de modelos de texto para detectar los puntos de dolor en las funcionalidades de las aplicaciones móviles bancarias. 4) Categorización de los tópicos generados por funcionalidad. 5) Visualización de los tópicos generados a partir de la calificación otorgada por el usuario (0 estrellas) mapeada como puntuación NPS (Detector 1-3, Promotor 5)

Proyectos relacionados

Con respecto al análisis y modelado de tópicos encontré algunos acercamientos y publicaciones, entre ellos los siguientes:

- Does the NPS reflect consumer sentiment? A qualitative examination of the NPS using a sentiment analysis approach
- Sentiment analysis and topic extraction of the twitter network of #paypalforis
- Modeling topic extraction-based sentiment analysis based on user reviews
- How can i improve my app? Classifying user reviews for software maintenance and evolution
- Covid-19 vaccine infodemic: sentiment analysis of the twitter content <ip>

Metodología

En mi experiencia estoy familiarizada con la metodología SEMMA, sin embargo, me gustará utilizar la metodología KDD (Knowledge Data Discovery) para este proyecto.

Breve descripción de la metodología:

- Selección de los datos. Los datos que usará en este proyecto se encuentran en las tiendas de aplicaciones de Google y Apple Store, la forma de extracción será a través de una técnica llamada Web Scrapping y podremos extraer datos históricos. El proceso ya lo veremos.
- Preprocesamiento de los datos. En esta parte de la metodología se realizará la limpieza de los datos, en la que se realizarán tareas como conversión a minúsculas, eliminación de acentos y caracteres especiales, stemming, stopwords, aplicación del algoritmo de Levenshtein para la corrección de las malas escrituras, dejar palabras con una longitud mínima de 5 caracteres e identificar entidades. Generar diccionarios propios de stopwords y si es posible de sinónimos.
- Interpretación de los datos. Generar y entrenar modelos de extracción de tópicos y modelos de clasificación de sentimiento para analizar los objetivos.
- Transferencia y evaluación. Para este paso se compararán los resultados obtenidos en los modelos de sentimiento con matriz de confusión o alguna otra métrica de desempeño. Mostrar los tópicos finales que tengan una interpretación más clara. Los resultados serán mostrados en un dashboard en la herramienta de Power BI.

Recopilación de la fuente de datos

** Web Scrapping Google Play Store - SuperMóvil **

```
In [20]: import pandas as pd
from google_play_scraper import Sort, reviews, reviews_all, app
from datetime import date
import datetime as dt
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize
import re

In [36]: path=C:\Users\A\Google\Documents\Videos\WASHER\DATA\TRECER SENSTRE\1\SENTENARIO DE PROYECTOS I\PROYECTO MAESTRIA\1\

In [39]: todaydate =today()
dt.today.strftime("%Y-%m-%d")
print(dt)

2022-11-23

In [41]: # obtenemos el nombre("id") de la aplicación de santander "SuperMóvil" en la google play store
# https://play.google.com/store/apps/details?id=com.bancosantander.supermovil&hl=es_MX&gl=US&pli=1
# app_santander="bv.bancosantander.supermovil"

In [18]: # apuntando a descargar los comentarios de la aplicación en android

result,cont,nation,taken=reviews(
    app_santander, #aplicación que vamos a descargar los comentarios
    location="es", #ubicación en el que se va encontrar
    country="es", #país de publicación
    sort="relevance", #ordenamiento de los datos
    count=50000, #cantidad de comentarios que vamos a recienar
    filter="score", #si existen, trae 50,000 comentarios, si no, trae todos los que hay
    filter_score=5, #filtro

)

In [17]: # observamos el número de registros obtenidos
print(len(result))
print(type(result))

50888
<class 'list'>

In [18]: # creando un dataframe de los resultados obtenidos
San_df=pd.DataFrame(result)
San_df.head(2)

Out[18]:
   reviewed      username      userimage      content      score      thumbsUpCount      reviewCreatedVersion      at      replyContent      repliedAt
0  54692838-6564-495c-80a3-9548833551b  Bengani (M)      hg.jgofguseconcent.com/a/62C9F...  Excelente App, por el momento todo bien  5  0  5.72  2022-11-22  16:10:59  None  NaN
1  3ef6b0d4-6564-495c-80a3-9548833551b  Luis Esteban  La peor app de un banco. No le permite ingresa...  1  1  5.72  2022-11-22  16:10:59  None  NaN
2  490026ac-4e8b-4a9b-80a3-9548833551b  Juan miguelchirico orozco lapez  Hay cosa que es muy difícil sacar y en sucurs...  3  0  5.72  2022-11-22  16:09:53  None  NaN
3  4734a07b-407a-4b4b-4a9b-734a054ac1a0b  Karina Morales  Bien  5  0  5.72  2022-11-22  15:56:27  None  NaN

In [19]: # observando las dimensiones del dataframe creado, vemos que tenemos 30k registros y 18 columnas
San_df.shape

Out[19]:
(50888, 18)

** Preprocesamiento y manipulación de datos google play store**

En esta parte principalmente hacemos lo siguiente:

- Eliminar stopwords
- Quitar acentos
- Convertir a minúsculas
- Quitar números y caracteres especiales
- Creación de campos auxiliares

In [20]: # Generando variables de hora y fecha
San_df["fecha"]=San_df["dt"].dt.strftime("%Y-%m-%d")
San_df["hora"]=San_df["dt"].dt.strftime("%H:%M:%S")
San_df["mes"]=San_df["dt"].dt.strftime("%Y-%m-%d")

In [21]: San_df.ues.value_counts()

Out[21]:
2022-09      8346
2022-08      6084
2022-05      5099
2022-12      5889
2022-04      5675
2022-03      5654
2022-06      5478
2022-07      5467
2022-01      5054
Name: mes, dtype: object

Observo que en promedio hay 5,000 comentarios por mes, además hay datos desde el mes de marzo con el extracto que hicimos de 50,000 registros.

In [22]: # observando algunos registros del dataframe
San_df.head(4)

Out[22]:
   reviewed      username      userimage      content      score      thumbsUpCount      reviewCreatedVersion      at      replyContent      repliedAt      fecha      hora      mes
0  54692838-6564-495c-80a3-9548833551b  Bengani (M)      hg.jgofguseconcent.com/a/62C9F...  Excelente App, por el momento todo bien  5  0  5.72  2022-11-22  16:10:59  None  NaN  2022-11-22  16:10:59  2022-11
1  3ef6b0d4-6564-495c-80a3-9548833551b  Luis Esteban  La peor app de un banco. No le permite ingresa...  1  1  5.72  2022-11-22  16:10:59  None  NaN  2022-11-22  16:10:59  2022-11
2  490026ac-4e8b-4a9b-80a3-9548833551b  Juan miguelchirico orozco lapez  Hay cosa que es muy difícil sacar y en sucurs...  3  0  5.72  2022-11-22  16:09:53  None  NaN  2022-11-22  16:09:53  2022-11
3  4734a07b-407a-4b4b-4a9b-734a054ac1a0b  Karina Morales  Bien  5  0  5.72  2022-11-22  15:56:27  None  NaN  2022-11-22  15:56:27  2022-11

In [23]: San_df.drop(columns=["userimage","repliedAt","replyContent"],inplace=True)
score_cat=["Detector","Promotor"]
San_df["score_cat"]=San_df["score"].map(score_cat)

In [24]: San_df.head(5)

Out[24]:
   reviewed      username      content      score      thumbsUpCount      reviewCreatedVersion      at      fecha      hora      mes      score_cat
0  54692838-6564-495c-80a3-9548833551b  Bengani (M)  Excelente App, por el momento todo bien  5  0  5.72  2022-11-22  16:10:59  2022-11-22  16:10:59  2022-11  Promotor
1  3ef6b0d4-6564-495c-80a3-9548833551b  Luis Esteban  La peor app de un banco. No le permite ingresa...  1  1  5.72  2022-11-22  16:10:59  2022-11-22  16:10:59  2022-11  Detector
2  490026ac-4e8b-4a9b-80a3-9548833551b  Juan miguelchirico orozco lapez  Hay cosa que es muy difícil sacar y en sucurs...  3  0  5.72  2022-11-22  16:09:53  2022-11-22  16:09:53  2022-11  Detector
3  4734a07b-407a-4b4b-4a9b-734a054ac1a0b  Karina Morales  Bien  5  0  5.72  2022-11-22  15:56:27  2022-11-22  15:56:27  2022-11  Promotor
4  92b4ac3b-d57b-4b40-8c31-4a4c742d3c19b  Evelyn Zaza  Muy fácil de usar  5  0  5.72  2022-11-22  15:52:30  2022-11-22  15:52:30  2022-11  Promotor

En el dataframe anterior observo lo siguiente:

reviewId - Identificador único del comentario
username - Nombre del usuario que dejó el comentario
content - Campo que contiene el texto escrito por el usuario
rating - Calificación (de 1 a 5 estrellas) que dio el usuario
thumbsUpCount - Número de usuarios que piensan o están de acuerdo con el comentario realizado
reviewCreatedAtVersion - Versión de la aplicación de SuperMóvil
at - Timestamp del comentario
fecha - Campo definido por mí que divide el campo "at" con el atributo de fecha
hora - Campo definido por mí que divide el campo "at" con el atributo de hora
mes - Campo definido por mí que divide el campo "at" con el atributo del mes
score_cat - Campo definido por mí para discretizar la calificación dada por el usuario y mapearla como NPS(Net Promoter Score)

In [25]: def clean_re(txt):
    try:pd.Series(txt).str.lower()
    except:pd.Series(txt).str.lower()
    # Filtrando el dataframe por el campo content que tenga una longitud de más de 1 caracter
    df_santander=df_santander[df_santander.content.str.len()>1]
    # Convirtiendo el campo de texto a minúsculas
    df_santander["texto"]=df_santander.content.str.lower()
    # Eliminando acentos
    a,b="áéíóú","aieou"
    trans=str.maketrans(a,b)
    df_santander["texto"]=df_santander.texto.str.translate(trans)
    df_santander["texto"]=df_santander["texto"].apply(clean_re)

In [27]: df_santander.head(3)

Out[27]:
   reviewed      username      content      score      thumbsUpCount      reviewCreatedVersion      at      fecha      hora      mes      score_cat      texto
0  54692838-6564-495c-80a3-9548833551b  Bengani (M)  Excelente App, por el momento todo bien  5  0  5.72  2022-11-22  16:10:59  2022-11-22  16:10:59  2022-11  Promotor  excelente app por momento todo bien
1  3ef6b0d4-6564-495c-80a3-9548833551b  Luis Esteban  La peor app de un banco. No le permite ingresa...  1  1  5.72  2022-11-22  16:10:59  2022-11-22  16:10:59  2022-11  Detector  peor app banco permite ingresar alertas dema...
2  490026ac-4e8b-4a9b-80a3-9548833551b  Juan miguelchirico orozco lapez  Hay cosa que es muy difícil sacar y en sucurs...  3  0  5.72  2022-11-22  16:09:53  2022-11-22  16:09:53  2022-11  Detector  hay cosa que muy difícil sacar sucursal ven tra...

In [28]: stopwords=stop(stopwords.words("spanish"))

In [29]: # crear un nuevo campo llamado textos que no contenga las stopwords por defecto cargadas en python
df_index,row=iter(df_santander.iterrows())
word=word_tokenize(row["texto"])
filtro=[]
for w in word:
    if w not in stopwords:
        filtro.append(w)
df_santander.loc[:,index,"texto"]='.join(filtro)

4975716 [92:37, 336.781/s]

In [30]: df_santander["listas"]=df_santander["texto"].str.split()

In [31]: df_santander.head(2)

Out[31]:
   reviewed      username      content      score      thumbsUpCount      reviewCreatedVersion      at      fecha      hora      mes      score_cat      texto
0  54692838-6564-495c-80a3-9548833551b  Bengani (M)  Excelente App, por el momento todo bien  5  0  5.72  2022-11-22  16:10:59  2022-11-22  16:10:59  2022-11  Promotor  excelente app por momento todo bien
1  3ef6b0d4-6564-495c-80a3-9548833551b  Luis Esteban  La peor app de un banco. No le permite ingresa...  1  1  5.72  2022-11-22  16:10:59  2022-11-22  16:10:59  2022-11  Detector  peor app banco permite ingresar alertas dema...

In [32]: # revisando que no haya números en el campo de texto nuevamente
df_santander=df_santander[df_santander["texto"].str.isnumeric()==False]
df_santander=df_santander[df_santander["texto"].notna()]

In [33]: df_santander.shape

Out[33]:
(49746, 14)

In [34]: # creando el conteo de palabras presentes para generar una lista propia de stopwords y eliminar ruido en los tópicos
palabras=[]
for i,row in df_santander.iterrows():
    for j,palabra in enumerate(row["listas"]):
        palabras.append(palabra)
print(len(palabras))

263827

In [36]: from collections import Counter
recuento=Counter(palabras)
df_contenido=pd.DataFrame.from_dict(recuento,orient='index').reset_index()
df_contenido.sort_values(by='count',ascending=False,inplace=True)
df_contenido.head(8)

Out[36]:
   index      count
17  aplicación  10056
42  buena  9938
1  app  8807
10  excelente  8063
3  bien  6025
26  hel  2280
93  hacer  2017
5  banco  2980

In [37]: df_contenido.shape
df_contenido.to_excel(("path\\palabras.xlsx",index=False)

In [42]: df_santander.to_excel(("path\\df_santander_(d1).xlsx",index=False)

In [39]: df_santander.to_excel(("path\\df_santander_(d1).xlsx",index=False)

In [53]: # Filtrando dataset con el que trabajaremos posteriormente
df_santander=df_santander[["reviewId","username","score","thumbsUpCount","reviewCreatedAtVersion","fecha","hora","score_cat","texto","mes"]]
df_santander.head()

Out[53]:
   reviewed      username      content      score      thumbsUpCount      reviewCreatedAtVersion      fecha      hora      mes      score_cat      texto
0  54692838-6564-495c-80a3-9548833551b  Bengani (M)  Excelente App, por el momento todo bien  5  0  5.72  2022-11-22  16:10:59  2022-11-22  16:10:59  2022-11  Promotor  excelente app por momento todo bien
1  3ef6b0d4-6564-495c-80a3-9548833551b  Luis Esteban  La peor app de un banco. No le permite ingresa...  1  1  5.72  2022-11-22  16:10:59  2022-11-22  16:10:59  2022-11  Detector  peor app banco permite ingresar alertas dema...
2  490026ac-4e8b-4a9b-80a3-9548833551b  Juan miguelchirico orozco lapez  Hay cosa que es muy difícil sacar y en sucurs...  3  0  5.72  2022-11-22  16:09:53  2022-11-22  16:09:53  2022-11  Detector  hay cosa que muy difícil sacar sucursal ven tra...

In [32]: # creando un nuevo campo llamado textos que no contenga las stopwords por defecto cargadas en python
df_index,row=iter(df_santander.iterrows())
word=word_tokenize(row["texto"])
filtro=[]
for w in word:
    if w not in stopwords:
        filtro.append(w)
df_santander.loc[:,index,"texto"]='.join(filtro)

4975716 [92:37, 336.781/s]

In [30]: df_santander["listas"]=df_santander["texto"].str.split()

In [31]: df_santander.head(2)

Out[31]:
   reviewed      username      content      score      thumbsUpCount      reviewCreatedAtVersion      at      fecha      hora      mes      score_cat      texto
0  54692838-6564-495c-80a3-9548833551b  Bengani (M)  Excelente App, por el momento todo bien  5  0  5.72  2022-11-22  16:10:59  2022-11-22  16:10:59  2022-11  Promotor  excelente app por momento todo bien
1  3ef6b0d4-6564-495c-80a3-9548833551b  Luis Esteban  La peor app de un banco. No le permite ingresa...  1  1  5.72  2022-11-22  16:10:59  2022-11-22  16:10:59  2022-11  Detector  peor app banco permite ingresar alertas dema...

In [32]: # revisando que no haya números en el campo de texto nuevamente
df_santander=df_santander[df_santander["texto"].str.isnumeric()==False]
df_santander=df_santander[df_santander["texto"].notna()]

In [34]: # creando el conteo de palabras presentes para generar una lista propia de stopwords y eliminar ruido en los tópicos
palabras=[]
for i,row in df_santander.iterrows():
    for j,palabra in enumerate(row["listas"]):
        palabras.append(palabra)
print(len(palabras))

263827

In [36]: from collections import Counter
recuento=Counter(palabras)
df_contenido=pd.DataFrame.from_dict(recuento,orient='index').reset_index()
df_contenido.sort_values(by='count',ascending=False,inplace=True)
df_contenido.head(8)

Out[36]:
   index      count
17  aplicación  10056
42  buena  9938
1  app  8807
10  excelente  8063
3  bien  6025
26  hel  2280
93  hacer  2017
5  banco  2980

In [37]: df_contenido.shape
df_contenido.to_excel(("path\\palabras.xlsx",index=False)

In [42]: df_santander.to_excel(("path\\df_santander_(d1).xlsx",index=False)

In [39]: df_santander.to_excel(("path\\df_santander_(d1).xlsx",index=False)

In [53]: # Filtrando dataset con el que trabajaremos posteriormente
df_santander=df_santander[["reviewId","username","score","thumbsUpCount","reviewCreatedAtVersion","fecha","hora","score_cat","texto","mes"]]
df_santander.head()

Out[53]:
   reviewed      username      content      score      thumbsUpCount      reviewCreatedAtVersion      fecha      hora      mes      score_cat      texto
0  54692838-6564-495c-80a3-9548833551b  Bengani (M)  Excelente App, por el momento todo bien  5  0  5.72  2022-11-22  16:10:59  2022-11-22  16:10:59  2022-11  Promotor  excelente app por momento todo bien
1  3ef6b0d4-6564-495c-80a3-9548833551b  Luis Esteban  La peor app de un banco. No le permite ingresa...  1  1  5.72  2022-11-22  16:10:59  2022-11-22  16:10:59  2022-11  Detector  peor app banco permite ingresar alertas dema...
2  490026ac-4e8b-4a9b-80a3-9548833551b  Juan miguelchirico orozco lapez  Hay cosa que es muy difícil sacar y en sucurs...  3  0  5.72  2022-11-22  16:09:53  2022-11-22  16:09:53  2022-11  Detector  hay cosa que muy difícil sacar sucursal ven tra...

In [32]: # creando un nuevo campo llamado textos que no contenga las stopwords por defecto cargadas en python
df_index,row=iter(df_santander.iterrows())
word=word_tokenize(row["texto"])
filtro=[]
for w in word:
    if w not in stopwords:
        filtro.append(w)
df_santander.loc[:,index,"texto"]='.join(filtro)

4975716 [92:37, 336.781/s]

In [30]: df_santander["listas"]=df_santander["texto"].str.split()

In [31]: df_santander.head(2)

Out[31]:
   reviewed      username      content      score      thumbsUpCount      reviewCreatedAtVersion      at      fecha      hora      mes      score_cat      texto
0  54692838-6564-495c-80a3-9548833551b  Bengani (M)  Excelente App, por el momento todo bien  5  0  5.72  2022-11-22  16:10:59  2022-11-22  16:10:59  2022-11  Promotor  excelente app por momento todo bien
1  3ef6b0d4-6564-495c-80a3-9548833551b  Luis Esteban  La peor app de un banco. No le permite ingresa...  1  1  5.72  2022-11-22  16:10:59  2022-11-22  16:10:59  2022-11  Detector  peor app banco permite ingresar alertas dema...

In [32]: # revisando que no haya números en el campo de texto nuevamente
df_santander=df_santander[df_santander["texto"].str.isnumeric()==False]
df_santander=df_santander[df_santander["texto"].notna()]

In [34]: # creando el conteo de palabras presentes para generar una lista propia de stopwords y eliminar ruido en los tópicos
palabras=[]
for i,row in df_santander.iterrows():
    for j,palabra in enumerate(row["listas"]):
        palabras.append(palabra)
print(len(palabras))

263827

In [36]: from collections import Counter
recuento=Counter(palabras)
df_contenido=pd.DataFrame.from_dict(recuento,orient='index').reset_index()
df_contenido.sort_values(by='count',ascending=False,inplace=True)
df_contenido.head(8)

Out[36]:
   index      count
17  aplicación  10056
42  buena  9938
1  app  8807
10  excelente  8063
3  bien  6025
26  hel  2280
93  hacer  2017
5  banco  2980

In [37]: df_contenido.shape
df_contenido.to_excel(("path\\palabras.xlsx",index=False)

In [42]: df_santander.to_excel(("path\\df_santander_(d1).xlsx",index=False)

In [39]: df_santander.to_excel(("path\\df_santander_(d1).xlsx",index=False)

In [53]: # Filtrando dataset con el que trabajaremos posteriormente
df_santander=df_santander[["reviewId","username","score","thumbsUpCount","reviewCreatedAtVersion","fecha","hora","score_cat","texto","mes"]]
df_santander.head()

Out[53]:
   reviewed      username      content      score      thumbsUpCount      reviewCreatedAtVersion      fecha      hora      mes      score_cat      texto
0  54692838-6564-495c-80a3-9548833551b  Bengani (M)  Excelente App, por el momento todo bien  5  0  5.72  2022-11-22  16:10:59  2022-11-22  16:10:59  2022-11  Promotor  excelente app por momento todo bien
1  3ef6b0d4-6564-495c-80a3-9548833551b  Luis Esteban  La peor app de un banco. No le permite ingresa...  1  1  5.72  2022-11-22  16:10:59  2022-11-22  16:10:59  2022-11  Detector  peor app banco permite ingresar alertas dema...
2  490026ac-4e8b-4a9b-80a3-9548833551b  Juan miguelchirico orozco lapez  Hay cosa que es muy difícil sacar y en sucurs...  3  0  5.72  2022-11-22  16:09:53  2022-11-22  16:09:53  2022-11  Detector  hay cosa que muy difícil sacar sucursal ven tra...

In [32]: # creando un nuevo campo llamado textos que no contenga las stopwords por defecto cargadas en python
df_index,row=iter(df_santander.iterrows())
word=word_tokenize(row["texto"])
filtro=[]
for w in word:
    if w not in stopwords:
        filtro.append(w)
df_santander.loc[:,index,"texto"]='.join(filtro)

4975716 [92:37, 336.781/s]

In [30]: df_santander["listas"]=df_santander["texto"].str.split()

In [31]: df_santander.head(2)

Out[31]:
   reviewed      username      content      score      thumbsUpCount      reviewCreatedAtVersion      at      fecha      hora      mes      score_cat      texto
0  54692838-6564-495c-80a3-9548833551b  Bengani (M)  Excelente App, por el momento todo bien  5  0  5.72  2022-11-22  16:10:59  2022-11-22  16:10:59  2022-11  Promotor  excelente app por momento todo bien
1  3ef6b0d4-6564-495c-80a3-9548833551b  Luis Esteban  La peor app de un banco. No le permite ingresa...  1  1  5.72  2022-11-22  16:10:59  2022-11-22  16:10:59  2022-11  Detector  peor app banco permite ingresar alertas dema...

In [32]: # revisando que no haya números en el campo de texto nuevamente
df_santander=df_santander[df_santander["texto"].str.isnumeric()==False]
df_santander=df_santander[df_santander["texto"].notna()]

In [34]: # creando el conteo de palabras presentes para generar una lista propia de stopwords y eliminar ruido en los tópicos
palabras=[]
for i,row in df_santander.iterrows():
    for j,palabra in enumerate(row["listas"]):
        palabras.append(palabra)
print(len(palabras))

263827

In [36]: from collections import Counter
recuento=Counter(palabras)
df_contenido=pd.DataFrame.from_dict(recuento,orient='index').reset_index()
df_contenido.sort_values(by='count',ascending=False,inplace=True)
df_contenido.head(8)

Out[36]:
   index      count
17  aplicación  10056
42  buena  9938
1  app  8807
10  excelente  8063
3  bien  6025
26  hel  2280
93  hacer  2017
5  banco  2980

In [37]: df_contenido.shape
df_contenido.to_excel(("path\\palabras.xlsx",index=False)

In [42]: df_santander.to_excel(("path\\df_santander_(d1).xlsx",index=False)

In [39]: df_santander.to_excel(("path\\df_santander_(d1).xlsx",index=False)

In [53]: # Filtrando dataset con el que trabajaremos posteriormente
df_santander=df_santander[["reviewId","username","score","thumbsUpCount","reviewCreatedAtVersion","fecha","hora","score_cat","texto","mes"]]
df_santander.head()

Out[53]:
   reviewed      username      content      score      thumbsUpCount      reviewCreatedAtVersion      fecha      hora      mes      score_cat      texto
0  54692838-6564-495c-80a3-9548833551b  Bengani (M)  Excelente App, por el momento todo bien  5  0  5.72  2022-11-22  16:10:59  2022-11-22  16:10:59  2022-11  Promotor  excelente app por momento todo bien
1  3ef6b0d4-6564-495c-80a3-9548833551b  Luis Esteban  La peor app de un banco. No le permite ingresa...  1  1  5.72  2022-11-22  16:10:59  2022-11-22  16:10:59  2022-11  Detector  peor app banco permite ingresar alertas dema...
2  490026ac-4e8b-4a9b-80a3-9548833551b  Juan miguelchirico orozco lapez  Hay cosa que es muy difícil sacar y en sucurs...  3  0  5.72  2022-11-22  16:09:53  2022-11-22  16:09:53  2022-11  Detector  hay cosa que muy difícil sacar sucursal ven tra...

In [32]: # creando un nuevo campo llamado textos que no contenga las stopwords por defecto cargadas en python
df_index,row=iter(df_santander.iterrows())
word=word_tokenize(row["texto"])
filtro=[]
for w in word:
    if w not in stopwords:
        filtro.append(w)
df_santander.loc[:,index,"texto"]='.join(filtro)

4975716 [92:37, 336.781/s]

In [30]: df_santander["listas"]=df_santander["texto"].str.split()

In [31]: df_santander.head(2)

Out[31]:
   reviewed      username      content      score      thumbsUpCount      reviewCreatedAtVersion      at      fecha      hora      mes      score_cat      texto
0  54692838-6564-495c-80a3-9548833551b  Bengani (M)  Excelente App, por el momento todo bien  5  0  5.72  2022-11-22  16:10:59  2022-11-22  16:10:59  2022-11  Promotor  excelente app por momento todo bien
1  3ef6b0d4-6564-495c-80a3-9548833551b  Luis Esteban  La peor app de un banco. No le permite ingresa...  1  1  5.72  2022-11-22  16:10:59  2022-11-22  16:10:59  2022-11  Detector  peor app banco permite ingresar alertas dema...

In [32]: # revisando que no haya números en el campo de texto nuevamente
df_santander=df_santander[df_santander["texto"].str.isnumeric()==False]
df_santander=df_santander[df_santander["texto"].notna()]

In [34]: # creando el conteo de palabras presentes para generar una lista propia de stopwords y eliminar ruido en los tópicos
palabras=[]
for i,row in df_santander.iterrows():
    for j,palabra in enumerate(row["listas"]):
        palabras.append(palabra)
print(len(palabras))

263827

In [36]: from collections import Counter
recuento=Counter(palabras)
df_contenido=pd.DataFrame.from_dict(recuento,orient='index').reset_index()
df_contenido.sort_values(by='count',ascending=False,inplace=True)
df_contenido.head(8)

Out[36]:
   index      count
17  aplicación  10056
42  buena  9938
1  app  8807
10  excelente  8063
3  bien  6025
26  hel  2280
93  hacer  2017
5  banco  2980

In [37]: df_contenido.shape
df_contenido.to_excel(("path\\palabras.xlsx",index=False)

In [42]: df_santander.to_excel(("path\\df_santander_(d1).xlsx",index=False)

In [39]: df_santander.to_excel(("path\\df_santander_(d1).xlsx",index=False)

In [53]: # Filtrando dataset con el que trabajaremos posteriormente
df_santander=df_santander[["reviewId","username","score","thumbsUpCount","reviewCreatedAtVersion","fecha","hora","score_cat","texto","mes"]]
df_santander.head()

Out[53]:
   reviewed      username      content      score      thumbsUpCount      reviewCreatedAtVersion      fecha      hora      mes      score_cat      texto
0  54692838-6564-495c-80a3-9548833551b  Bengani (M)  Excelente App, por el momento todo bien  5  0  5.72  2022-11-22  16:10:59  2022-11-22  16:10:59  2022-11  Promotor  excelente app por momento todo bien
1  3ef6b0d4-6564-495c-80a3-9548833551b  Luis Esteban  La peor app de un banco. No le permite ingresa...  1  1  5.72  2022-11-22  16:10:59  2022-11-22  16:10:59  2022-11  Detector  peor app banco permite ingresar alertas dema...
2  490026ac-4e8b-4a9b-80a3-9548833551b  Juan miguelchirico orozco lapez  Hay cosa que es muy difícil sacar y en sucurs...  3  0  5.72  2022-11-22  16:09:53  2022-11-22  16:09:53  2022-11  Detector  hay cosa que muy difícil sacar sucursal ven tra...

In [32]: # creando un nuevo campo llamado textos que no contenga las stopwords por defecto cargadas en python
df_index,row=iter(df_santander.iterrows())
word=word_tokenize(row["texto"])
filtro=[]
for w in word:
    if w not in stopwords:
        filtro.append(w)
df_santander.loc[:,index,"texto"]='.join(filtro)

4975716 [92:37, 336.781/s]

In [30]: df_santander["listas"]=df_santander["texto"].str.split()

In [31]: df_santander.head(2)

Out[31]:
   reviewed      username      content      score      thumbsUpCount      reviewCreatedAtVersion      at      fecha      hora      mes      score_cat      texto
0  54692838-6564-495c-80a3-9548833551b  Bengani (M)  Excelente App, por el momento todo bien  5  0  5.72  2022-11-22  16:10:59  2022-11-22  16:10:59  2022-11  Promotor  excelente app por momento todo bien
1  3ef6b0d4-6564-495c-80a3-9548833551b  Luis Esteban  La peor app de un banco. No le permite ingresa...  1  1  5.72  2022-11-22  16:10:59  2022-11-22  16:10:59  2022-11  Detector  peor app banco permite ingresar alertas dema...

In [32]: # revisando que no haya números en el campo de texto nuevamente
df_santander=df_santander[df_santander["texto"].str.isnumeric()==False]
df_santander=df_santander[df_santander["texto"].notna()]

In [34]: # creando el conteo de palabras presentes para generar una lista propia de stopwords y eliminar ruido en los tópicos
palabras=[]
for i,row in df_santander.iterrows():
    for j,palabra in enumerate(row["listas"]):
        palabras.append(palabra)
print(len(palabras))

263827

In [36]: from collections import Counter
recuento=Counter(palabras)
df_contenido=pd.DataFrame.from_dict(recuento,orient='index').reset_index()
df_contenido.sort_values(by='count',ascending=False,inplace=True)
df_contenido.head(8)

Out[36]:
   index      count
17  aplicación  10056
42  buena  9938
1  app  8807
10  excelente  8063
3  bien  6025
26  hel  2280
93  hacer  2017
5  banco  2980

In [37]: df_contenido.shape
df_contenido.to_excel(("path\\palabras.xlsx",index=False)

In [42]: df_santander.to_excel(("path\\df_santander_(d1).xlsx",index=False)

In [39]: df_santander.to_excel(("path\\df_santander_(d1).xlsx",index=False)

In [53]: # Filtrando dataset con el que trabajaremos posteriormente
df_santander=df_santander[["reviewId","username","score","thumbsUpCount","reviewCreatedAtVersion","fecha","hora","score_cat","texto","mes"]]
```