Problema 2

▼ LIMPIEZA Y PREPROCESAMIENTO DE DATOS

```
1 import re
2 import string
3 import nltk
4 from nltk.corpus import stopwords
6 # Descargar la lista de stopwords
7 nltk.download('stopwords')
8 stop_words = set(stopwords.words('spanish'))
10 def preprocess_tweet(tweet):
      # Convertir a minúsculas
      tweet = tweet.lower()
      # Quitar URLs
      tweet = re.sub(r'http\S+|www\S+|https\S+', '', tweet, flags=re.MULTILINE)
      # Quitar caracteres de usuario y hashtags
      tweet = re.sub(r'\@\w+|\+\, '', tweet)
      # Quitar signos de puntuación
      tweet = tweet.translate(str.maketrans('', '', string.punctuation))
      # Quitar stopwords y números
      tweet = ' '.join([word for word in tweet.split() if word not in stop_words and not word.isnumeric()])
      return tweet
28 # Aplicar preprocesamiento a los tweets
```

```
29 data_bernardo['processed_tweet'] = data_bernardo['rawContent'].apply(preprocess_tweet)
30 data_sandra['processed_tweet'] = data_sandra['rawContent'].apply(preprocess_tweet)
32 # Mostrar los tweets preprocesados
33 processed_bernardo_head = data_bernardo[['rawContent', 'processed_tweet']].head()
34 processed_sandra_head = data_sandra[['rawContent', 'processed_tweet']].head()
36 processed_bernardo_head, processed_sandra_head
     [nltk_data] Downloading package stopwords to
     [nltk_data]
                    C:\Users\angel\AppData\Roaming\nltk_data...
     [nltk_data]
                  Package stopwords is already up-to-date!
                                               rawContent \
     0 @AnonGTReloaded @msemillagt @BArevalodeLeon ja...
        @ASIERVERA @AztecaNoticiaGT @BArevalodeLeon Do...
        Paciente de 39 años, dolor lumbar de 1 año tra...
        @VicZacariasGT @soy_502 @BArevalodeLeon @msemi...
      4 @Igor_Bitkov No le sigan el juego a este ruso ...
                                          processed_tweet
     0 jajajajajajaja pisen deje hartos tanta ignor...
        dos veces repitió actuando margen ley seguro t\dots
     2 paciente años dolor lumbar año tras caída hizo...
                           compa alucina puro net pareces
     4 sigan juego ruso invasor enero solicitamos ret...
                                               rawContent \
     0 @bernardosilvagt @BArevalodeLeon @DrGiammattei...
        @_awskl @mjcabrerar @BArevalodeLeon @TSEGuatem...
      2 The 2023 National Race Walking Championship &a...
      3 @Palomin17772524 @mjcabrerar @BArevalodeLeon @...
      4 @ASolaresM @Mike051270 @BArevalodeLeon Otro es...
                                          processed_tweet
                   mentiroso giamattei baldetti minúscula
     1 dedonde van sacar mil supuestos votos primera ...
        the national race walking championship amp you...
                            sueños net bañalos tomas agua
      4 estupido cegado caciques impide ver bajo nivel... )
```

→ ANALISIS EXPLORATORIO

```
1 import matplotlib.pyplot as plt
   2 from collections import Counter
  3 from wordcloud import WordCloud
  5 # Actualizar la función para usar la columna correcta y re-ejecutar el análisis exploratorio
  6 def exploratory_analysis_updated(data, candidate_name):
                   # Estadísticas básicas
                  total_tweets = len(data)
                  avg_likes = data['likeCount'].mean()
                  avg_retweets = data['retweetCount'].mean()
                  avg_replies = data['replyCount'].mean()
                   # Palabras más comunes
                  words = ' '.join(data['processed_tweet']).split()
14
                  counter = Counter(words)
                  most common words = counter.most common(10)
                   # Wordcloud
                  wordcloud = WordCloud(width=800, height=500, random_state=21, max_font_size=100, background_color='white').generate_from_frequencies(color='white').generate_from_frequencies(color='white').generate_from_frequencies(color='white').generate_from_frequencies(color='white').generate_from_frequencies(color='white').generate_from_frequencies(color='white').generate_from_frequencies(color='white').generate_from_frequencies(color='white').generate_from_frequencies(color='white').generate_from_frequencies(color='white').generate_from_frequencies(color='white').generate_from_frequencies(color='white').generate_from_frequencies(color='white').generate_from_frequencies(color='white').generate_from_frequencies(color='white').generate_from_frequencies(color='white').generate_from_frequencies(color='white').generate_from_frequencies(color='white').generate_from_frequencies(color='white').generate_from_frequencies(color='white').generate_from_frequencies(color='white').generate_from_frequencies(color='white').generate_from_frequencies(color='white').generate_from_frequencies(color='white').generate_from_frequencies(color='white').generate_from_frequencies(color='white').generate_from_frequencies(color='white').generate_from_frequencies(color='white').generate_frequencies(color='white').generate_frequencies(color='white').generate_frequencies(color='white').generate_frequencies(color='white').generate_frequencies(color='white').generate_frequencies(color='white').generate_frequencies(color='white').generate_frequencies(color='white').generate_frequencies(color='white').generate_frequencies(color='white').generate_frequencies(color='white').generate_frequencies(color='white').generate_frequencies(color='white').generate_frequencies(color='white').generate_frequencies(color='white').generate_frequencies(color='white').generate_frequencies(color='white').generate_frequencies(color='white').generate_frequencies(color='white').generate_frequencies(color='white').generate_frequencies(color='white').generate_frequencies(color='white')
20
                   plt.figure(figsize=(14, 6))
                   plt.imshow(wordcloud, interpolation='bilinear')
                  plt.axis('off')
24
                  plt.title(f'Palabras más comunes en tweets sobre {candidate_name}')
                   plt.show()
                   return {
                              'total_tweets': total_tweets,
                              'avg_likes': avg_likes,
                              'avg_retweets': avg_retweets,
                               'avg_replies': avg_replies,
                               'most_common_words': most_common_words
34
35 # Análisis exploratorio actualizado para Bernardo Arévalo
```

```
36 bernardo_analysis_updated = exploratory_analysis_updated(data_bernardo, "Bernardo Arévalo")
37 bernardo_analysis_updated
38
```

```
Palabras más comunes en tweets sobre Bernardo Arévalo
                            ora
                                 deben materiales
                                               quieren
                     país
                                   gobierno
                                                     triple
                        voto
                               usted d
hacia
                                        anos
                                            hace
                                ustedes
                     semil
                                         haciendo
                                      puede
               personas as 1 hacen dos
                                          ver
                     on
```

```
{'total_tweets': 4212,
    'avg_likes': 256.9669990503324,
    'avg_retweets': 53.50213675213675,
    'avg_replies': 28.675213675213676,
    'most_common_words': [('zona', 530),
    ('via', 459),
    ('si', 412),
    ('presidente', 377),
    ('avenida', 340),
    ('the', 309),
    ('solo', 305),
    ('calle', 303),
    ('pueblo', 251),
    ('vou', 244)]}
```

- 1 # Análisis exploratorio actualizado para Bernardo Arévalo
- 2 sandra_analysis_updated = exploratory_analysis_updated(data_sandra, "Sandra Torres")
- 3 sandra_analysis_updated

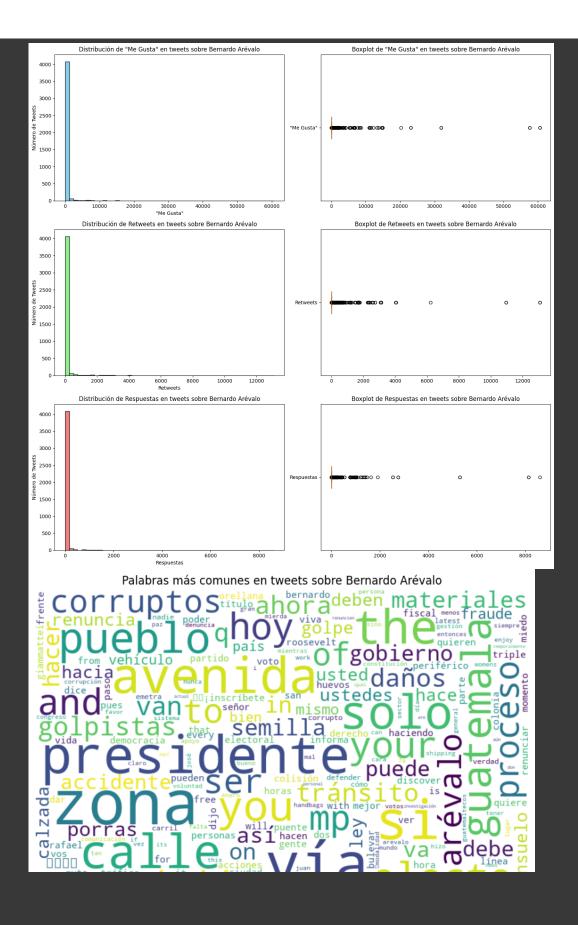
```
Palabras más comunes en tweets sobre Sandra Torres

Vehículo

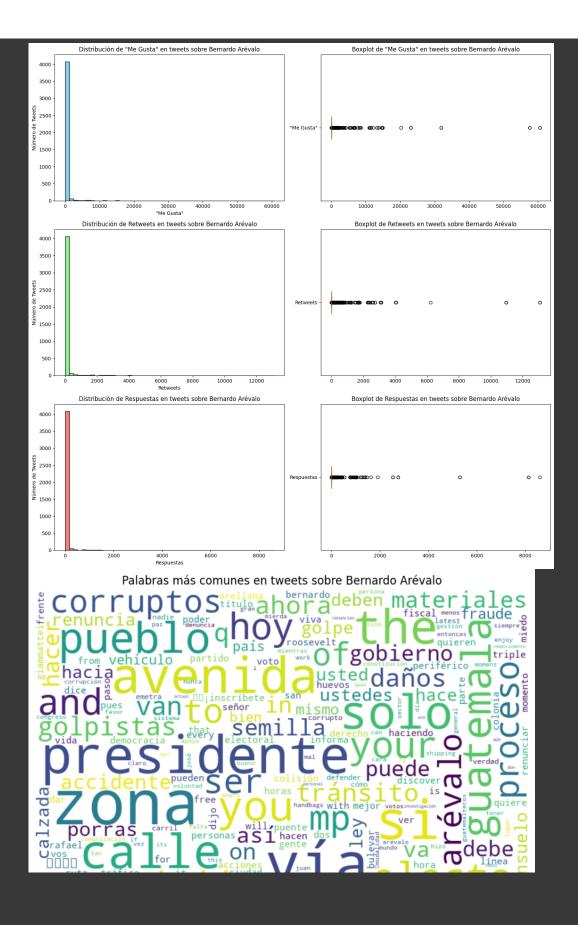
Fran Señorpetapa that dice parte of the posessión of the posess
```

```
1 # Función actualizada para realizar un análisis exploratorio completo que incluye boxplots
2 def complete_exploratory_analysis_with_boxplots(data, candidate_name, column_name):
      # Estadísticas Descriptivas
     total_tweets = len(data)
      avg_likes = data['likeCount'].mean()
      avg_retweets = data['retweetCount'].mean()
      avg_replies = data['replyCount'].mean()
      # Histogramas y Boxplots para visualizar la distribución de "me gusta", retweets y respuestas
      fig, ax = plt.subplots(3, 2, figsize=(15, 15))
      # Histogramas
      ax[0][0].hist(data['likeCount'], bins=50, color='skyblue', edgecolor='black')
      ax[0][0].set\_title(f'Distribuci\'on de "Me Gusta" en tweets sobre \{candida\underline{te\_name}\}')
      ax[0][0].set_xlabel('"Me Gusta"')
      ax[0][0].set_ylabel('Número de Tweets')
      ax[1][0].hist(data['retweetCount'], bins=50, color='lightgreen', edgecolor='black')
      ax[1][0].set_title(f'Distribución de Retweets en tweets sobre {candidate_name}')
      ax[1][0].set_xlabel('Retweets')
      ax[1][0].set_ylabel('Número de Tweets')
      ax[2][0].hist(data['replyCount'], bins=50, color='lightcoral', edgecolor='black')
      ax[2][0].set_title(f'Distribución de Respuestas en tweets sobre {candidate_name}')
      ax[2][0].set_xlabel('Respuestas')
      ax[2][0].set_ylabel('Número de Tweets')
      # Boxplots
      ax[0][1].boxplot(data['likeCount'], vert=False)
      ax[0][1].set_title(f'Boxplot de "Me Gusta" en tweets sobre {candidate_name}')
      ax[0][1].set_yticklabels(['"Me Gusta"'])
      ax[1][1].boxplot(data['retweetCount'], vert=False)
      ax[1][1].set_title(f'Boxplot de Retweets en tweets sobre {candidate_name}')
      ax[1][1].set_yticklabels(['Retweets'])
      ax[2][1].boxplot(data['replyCount'], vert=False)
      ax[2][1].set_title(f'Boxplot de Respuestas en tweets sobre {candidate_name}')
      ax[2][1].set_yticklabels(['Respuestas'])
      plt.tight_layout()
      plt.show()
      # Palabras más comunes usando WordCloud
      words = ' '.join(data[column_name]).split()
      counter = Counter(words)
      wordcloud = WordCloud(width=800, height=500, random_state=21, max_font_size=100, background_color='white').generate_from_frequencies(c
      plt.figure(figsize=(10, 6))
      plt.imshow(wordcloud, interpolation='bilinear')
      plt.axis('off')
      plt.title(f'Palabras más comunes en tweets sobre {candidate_name}')
      plt.show()
      # Análisis Temporal: Número de tweets con el tiempo
      data['date'] = pd.to_datetime(data['date'])
      tweets_over_time = data.resample('D', on='date').size()
      plt.figure(figsize=(12, 6))
      tweets_over_time.plot()
      plt.title(f'Número de Tweets con el tiempo sobre {candidate_name}')
```

```
plt.xlabel('Fecha')
       plt.ylabel('Número de Tweets')
       plt.grid(True)
       plt.show()
       # Distribución del Lenguaje
       language_distribution = data['lang'].value_counts().head(5)
       plt.figure(figsize=(10, 6))
       language_distribution.plot(kind='bar', color=['skyblue', 'lightgreen', 'lightcoral', 'gold', 'lightpink'])
       plt.title(f'Distribución del Lenguaje en tweets sobre {candidate_name}')
      plt.xlabel('Idioma')
plt.ylabel('Número de Tweets')
       plt.xticks(rotation=0)
       plt.show()
      # Return the basic statistics
       return {
           'total_tweets': total_tweets,
           'avg_likes': avg_likes,
           'avg_retweets': avg_retweets,
           'avg_replies': avg_replies
87 # Análisis exploratorio completo para Bernardo Aréval
88 bernardo_analysis_complete_with_boxplots = complete_exploratory_analysis_with_boxplots(data_bernardo, "Bernardo Arévalo", 'processed_tweet
89 bernardo_analysis_complete_with_boxplots
```







▼ Descubrimiento de Información

```
1 # Calculando métricas de popularidad para cada candidato
2 popularity_metrics = {
3     'Candidates': ['Bernardo Arévalo', 'Sandra Torres'],
4     'Total Tweets': [bernardo_analysis_complete_with_boxplots['total_tweets'], sandra_analysis_complete_with_boxplots['total_tweets']],
5     'Average Likes': [bernardo_analysis_complete_with_boxplots['avg_likes'], sandra_analysis_complete_with_boxplots['avg_likes']],
6     'Average Retweets': [bernardo_analysis_complete_with_boxplots['avg_retweets']], sandra_analysis_complete_with_boxplots['avg_retweets']]
7     'Average Replies': [bernardo_analysis_complete_with_boxplots['avg_replies'], sandra_analysis_complete_with_boxplots['avg_replies']]
8 }
9
10 popularity_df = pd.DataFrame(popularity_metrics)
11
12 popularity_df
13
```

Candidates Total Tweets Average Likes Average Retweets Average Replies

0	Bernardo Arévalo	4212	256.966999	53.502137	28.675214
		4212			28.675214

1. Total de Tweets:

• Bernardo Arévalo ha sido mencionado en 4,212 tweets.

Sandra Torres ha sido mencionada en 5,784 tweets.
 Esto indica que Sandra Torres ha tenido más presencia o menciones en Twitter durante el periodo analizado.

2. "Me Gusta" Promedio por Tweet:

- Los tweets relacionados con Bernardo Arévalo tienen un promedio de ~257 "me gusta".
- Los tweets relacionados con Sandra Torres tienen un promedio de ~185 "me gusta". Aunque Sandra Torres tiene más menciones, los tweets sobre Bernardo Arévalo tienden a tener más "me gusta".

3. Retweets Promedio por Tweet:

- Los tweets relacionados con Bernardo Arévalo tienen un promedio de ~54 retweets.
- Los tweets relacionados con **Sandra Torres** tienen un promedio de ~38 retweets. Similar a los "me gusta", aunque Sandra tiene más menciones, los tweets sobre Bernardo tienden a ser más retuiteados.

4. Respuestas Promedio por Tweet:

- Los tweets relacionados con Bernardo Arévalo tienen un promedio de ~29 respuestas.
- Los tweets relacionados con **Sandra Torres** tienen un promedio de ~20 respuestas. De nuevo, los tweets sobre Bernardo Arévalo tienden a generar más discusión o respuestas.

Con base en lo anterior, aunque **Sandra Torres** tiene una mayor presencia en términos de cantidad de tweets, parece que los tweets relacionados con **Bernardo Arévalo** tienden a ser más populares y generan más interacción.

```
1 # Función para obtener las palabras más comunes para cada candidato
 2 def get_most_common_words(data, column_name, num=10):
       words = ' '.join(data[column_name]).split()
       counter = Counter(words)
       return counter.most_common(num)
 7 # Obtener las 10 palabras más comunes para cada candidato
 8 bernardo_common_words = get_most_common_words(data_bernardo, 'processed_tweet')
 9 sandra_common_words = get_most_common_words(data_sandra, 'processed_tweet')
11 bernardo_common_words, sandra_common_words
     ([('zona', 530),
('vía', 459),
        ('si', 412),
        ('presidente', 377).
        ('avenida', 340),
       ('solo', 305),
('calle', 303),
('pueblo', 251),
      ('you', 244)],
[('zona', 757),
       ('vía', 542),
('si', 506),
       ('avenida', 470),
        ('calle', 434),
        ('guatemala', 301)])
```

Bernardo Arévalo:

Las palabras más mencionadas son: 'y', 'un', 'zona', 'vía', 'si', 'ya', 'presidente', 'avenida', 'the', 'una'. Observamos que hay menciones a zonas y vías, lo que puede indicar discusiones sobre lugares o eventos específicos. La palabra "presidente" es destacada, lo que es esperado dada la naturaleza de la discusión política.

Sandra Torres:

Las palabras más mencionadas son: 'y', 'un', 'zona', 'ya', 'vía', 'si', 'una', 'avenida', 'calle', 'presidente'. Al igual que con Bernardo, vemos menciones a zonas, vías y avenidas. La palabra "presidente" también es prominente en este conjunto de datos. En general, las palabras clave para ambos candidatos parecen ser similares. Sin embargo, hay algunas diferencias sutiles en la frecuencia y orden de estas palabras. Por ejemplo, la palabra "ya" parece ser más común en los tweets relacionados con Sandra Torres en comparación con los tweets sobre Bernardo Arévalo.

Basado en las tendencias temporales:

Bernardo Arévalo: