



Instituto Politécnico Nacional  
"La Técnica al Servicio de la Patria"



Escuela Superior de  
Cómputo

REPORTE

# Práctica Integradora

REALIZADA POR

Robles Guzmán Naomi Isabel

Ugalde Téllez Aarón

PARA LA MATERIA DE

Tecnologías de Lenguaje Natural

IMPARTIDA POR

Elizabeth Moreno Galván

GRUPO

5BM1

---

27 de junio de 2025

## Índice

<b>1. Introducción</b>	<b>2</b>
1.1. Definición de la Problemática . . . . .	2
1.1.1. Clasificación de contenido informativo . . . . .	2
1.1.2. Detección de tono y carga emocional . . . . .	2
1.1.3. Desequilibrio en los datos . . . . .	2
1.2. Propuesta de solución . . . . .	3
1.2.1. Estrategia de preprocesamiento . . . . .	3
1.2.2. Técnicas para balancear los datos . . . . .	3
1.2.3. Sistema de análisis de sentimientos por capas . . . . .	3
<b>2. Desarrollo</b>	<b>4</b>
2.1. Descripción del dataset . . . . .	4
2.2. Preprocesamiento de los datos . . . . .	6
2.2.1. Clase preprocesa . . . . .	6
2.2.2. Tokenización, Lematización y comparación de los lemas con un diccionario en español . . . . .	7
2.3. Clasificador de tipos de noticias . . . . .	7
2.3.1. Balanceo de clases y vectorización . . . . .	7
2.3.2. Entrenamiento del modelo de clasificación KNN . . . . .	10
2.4. Análisis de sentimientos . . . . .	10
2.4.1. Covariables . . . . .	11
2.4.2. Entrenamiento de modelo de K-means . . . . .	13
<b>3. Resultados</b>	<b>15</b>
3.1. Clasificador de tipos de noticias . . . . .	15
3.1.1. Métricas de rendimiento . . . . .	15
3.1.2. Matriz de confusión . . . . .	15
3.1.3. Distribución de Clusters . . . . .	17
<b>4. Conclusiones</b>	<b>18</b>
<b>Referencias</b>	<b>19</b>

## 1. Introducción

El análisis automatizado permite extraer información valiosa sobre tendencias, sentimientos y patrones de información que circulan en los medios de comunicación [1]. La capacidad de clasificar automáticamente el tipo de contenido y determinar su polaridad emocional podría significar una herramienta estratégica para organizaciones o investigadores.

Este proyecto integrador aborda dos problemáticas centrales en el procesamiento de lenguaje, la clasificación automática de tipos de noticias y el análisis de sentimientos. Para la primera problemática, se implementa un sistema de clasificación supervisada basado en el algoritmo K-Nearest Neighbors (KNN), el cual permite categorizar automáticamente las noticias según su temática específica. Para la segunda problemática, dado a que los documentos no están etiquetados según su tono emocional, se desarrolla un sistema de análisis de sentimientos que combina enfoques léxicos mediante diccionarios especializados y recursos como SenticLex, permitiendo identificar la polaridad emocional de los textos.

La metodología propuesta incluye la etapa de preprocesamiento de datos que contempla la limpieza, normalización y tokenización de textos, así como la implementación de técnicas de balanceo de clases para optimizar el rendimiento de los algoritmos de clasificación. Adicionalmente, se emplea clustering no supervisado mediante K-means para agrupar los textos según sus características sentimentales, proporcionando una perspectiva complementaria al análisis de polaridad.

Retomando el dataset utilizado a lo largo de todo el curso, el cual comprende noticias en español etiquetadas por categorías temáticas, esta característica será clave para evaluar la efectividad de las técnicas de clasificación supervisada. Los resultados esperados buscan demostrar la viabilidad y efectividad de las técnicas implementadas, estableciendo las bases para futuras investigaciones en el campo del procesamiento automatizado de textos.

### 1.1. Definición de la Problemática

En la actualidad cada día circula una gran cantidad de noticias, por tal motivo, medios de comunicación, empresas y centros de investigación enfrentan el reto de procesar y entender grandes volúmenes de contenido textual de forma rápida y eficiente. Este reto se manifiesta principalmente en tres áreas:

#### 1.1.1. Clasificación de contenido informativo

Clasificar textos como noticias en ocasiones presenta limitaciones como poca escalabilidad, criterios inconsistentes y procesamiento demasiado lento.

#### 1.1.2. Detección de tono y carga emocional

Para determinar si un texto transmite una emoción positiva, negativa o neutra se deben superar desafíos como el sarcasmo y la ironía, así como también la variación cultural y lingüística.

#### 1.1.3. Desequilibrio en los datos

El último gran reto tiene que ver con la distribución de las categorías en un conjunto de datos, frecuentemente suele ser desigual y este caso no es la excepción. El *dataset* en el que está basado este trabajo presenta una tendencia hacia clases dominantes, lo que genera problemas para generalizar y por consecuencia, peor desempeño en clases menos representadas.

## 1.2. Propuesta de solución

Se propone desarrollar un sistema que automatice el análisis de texto, con dos componentes principales: la clasificación temática y el análisis de sentimientos.

### 1.2.1. Estrategia de preprocesamiento

Se implementa una serie de pasos para limpiar y preparar los textos antes del análisis:

- **Normalización del texto:** Eliminación de signos de puntuación, uso uniforme de minúsculas y corrección de espacios.
- **Filtrado léxico:** Solo se conservan las palabras que aparecen en un diccionario en español.
- **Lematización:** Reducción de las palabras a su forma base.
- **Eliminación de palabras vacías:** Se eliminan conectores, artículos y otras palabras poco relevantes para el análisis.

### 1.2.2. Técnicas para balancear los datos

Para corregir el desbalance entre categorías, se aplica una estrategia mixta:

- **Reducción de clases dominantes** Se usa el método Near Miss para reducir de forma controlada los ejemplos de las categorías más representadas, seleccionando solo los ejemplos más relevantes para conservar diversidad.
- **Refuerzo de clases minoritarias**
  1. Se generan nuevos ejemplos mediante Word2Vec, intercambiando palabras por otras similares.
  2. Se complementa con SMOTE para crear representaciones artificiales en el espacio vectorial.

### 1.2.3. Sistema de análisis de sentimientos por capas

- **Capa básica** Diccionarios propios para identificar palabras positivas, negativas o relacionadas con emociones como alegría, tristeza, enojo y miedo.
- **Capa especializada** Se usa SenticLex para calcular un puntaje continuo que refleje la polaridad emocional del texto, obteniendo métricas agregadas para cada documento del corpus.
- **Capa de agrupamiento** Los documentos se agrupan automáticamente según sus características emocionales, identificando patrones usando el algoritmo K-means.

## 2. Desarrollo

### 2.1. Descripción del dataset

El conjunto de datos proporcionado, puede consultarse en el sitio *Kaggle* y fue construido con una herramienta de web scraping para el Dataton 2022 de Bancolombia para el entrenamiento de modelos supervisados a utilizar en una recomendación de Noticias [2].

El *dataset* tiene noticias de las categorías:

- Macroeconomía
- Sostenibilidad
- Innovación
- Regulaciones
- Alianzas
- Reputación
- Otra

Además, el documento *CSV* consta de 1217 registros y las columnas:

- **Noticia:** texto de la noticia.
- **Tipo:** Etiqueta del tipo de noticia

La distribución de instancias de noticias por categoría se aprecia en la siguiente Figura. Se observa que las clases Macroeconomía y Alianzas tienen un mayor número de muestras; mientras que el resto de clases, exceptuando reputación, oscilan entre las 100 y 200 instancias. Por su parte, Reputación es la categoría menos representada del *dataset*.

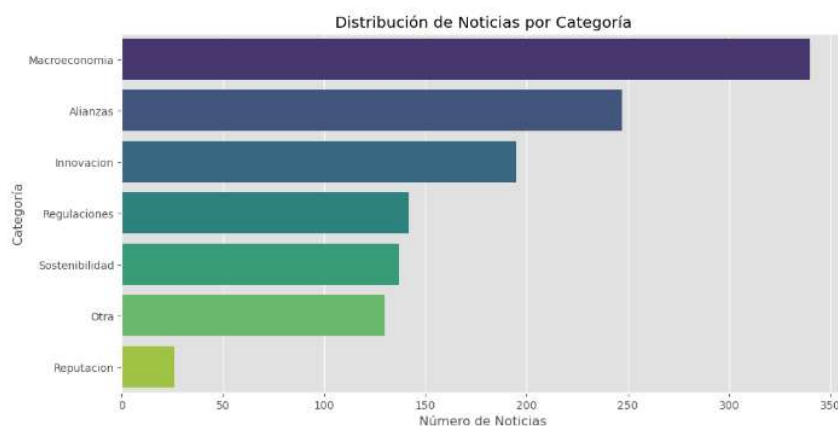


Figura 1: Gráfica de la distribución de registros de noticias por tipo.

Después de analizar el *dataset* con ayuda de Python, se encontró que la longitud promedio de palabras es de 516.55, siendo la noticia más larga de un total de 2964 palabras, y la más corta solo contenía una. En la **Figura 2** se muestra la distribución del número de palabras por noticia en cada una de las siete categorías del conjunto de datos.



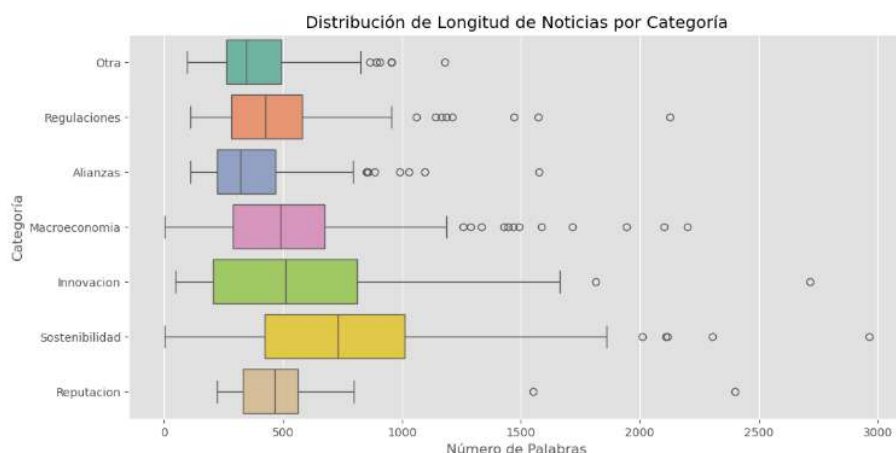


Figura 2: Gráfica de la distribución de la longitud en palabras de las noticias por tipo.

Se observa que las noticias relacionadas con Sostenibilidad e Innovación tienden a ser más extensas, mientras que las de Reputación y Alianzas son más cortas. Además, se identifican múltiples valores atípicos, lo que indica que algunas noticias presentan una longitud significativamente mayor al promedio de su categoría. Esta variabilidad puede influir en el procesamiento y análisis posterior de los textos.

Por otra parte, para conocer las palabras más frecuentes en cada categoría de noticia se hizo uso de la herramienta de *WordCloud*, debido a que permite una fácil visualización de los términos más usados.



Figura 3: Visualización de 3 de las 7 nubes de palabras generadas para cada categoría. El resto puede consultarse en la carpeta del proyecto.

Las nubes de palabras de las clases reflejan que, aunque cada categoría aborda un enfoque temático distinto, todas comparten un campo semántico relacionado con el ámbito empresarial, financiero y económico. Términos como empresa, cliente, mercado y BBVA aparecen de forma reiterada, lo que

indica una base léxica común en el corpus, posiblemente influenciada por el protagonismo de una entidad financiera específica. Sin embargo, también se observan diferencias claras entre categorías: Macroeconomía incluye vocabulario técnico relacionado con indicadores económicos; Reputación se enfoca en términos vinculados a la imagen y posicionamiento de marca; Innovación destaca palabras asociadas a tecnología y servicios digitales, mientras que Sostenibilidad introduce términos del ámbito ambiental. Esto sugiere que, si bien hay una estructura temática general compartida, el contenido de las noticias sí varía según la categoría, característica que permitirá distinguir matices y lograr un buen rendimiento para el clasificador.

## 2.2. Preprocesamiento de los datos

Se diseñó una clase `Preprocesa`: para realizar diversas operaciones de limpieza y normalización de textos en español. Apoyándose de la biblioteca `string` de Python para la lista de caracteres especiales que se eliminan del texto. Los métodos definidos por esta clase se aplican sobre el dataset mediante la siguiente línea de código.

```
1 preprocesador = Preprocesa()
2
3 dataset["news_limpio"] = dataset["news"].astype(str).apply(preprocesador.
4     remove_stopwords).apply(preprocesador.remove_punctuation).apply(preprocesador.
5     lower_words)
6
7 # Guardar el dataset procesado
8 dataset.to_csv("news.csv", index=False)
```

Además, para eliminar palabras no deseadas, se compararon las existentes con un diccionario en español y se eliminaron las que no se encontraran en este.

### 2.2.1. Clase preprocesa

```
def remove_punctuation(self, text):
```

Esta función dentro de la clase `Preprocesa` recibe un texto y define un conjunto de caracteres prohibidos que incluye signos de puntuación, caracteres especiales y secuencias no deseadas. Dicho conjunto se une al conjunto de signos de puntuación provisto por `string.punctuation`, para generar una nueva cadena que excluye dichos caracteres y finalmente estandarizar los espacios en blanco, de manera que se retorne un texto libre de puntuaciones [3].

```
def lower_words(self, text):
```

Este sencillo método se encarga de convertir todos los caracteres de la cadena minúsculas, esto para facilitar la comparación entre tokens y palabras vacías.

```
def quitarAcentos(self, s):
```

Esta función reemplaza los caracteres acentuados por su contraparte sin acento. Este proceso garantiza la uniformidad de las palabras, ya que elimina variaciones ortográficas que podrían afectar la tokenización y posteriores análisis. Si bien esta función puede ser útil en algunos contextos, para esta práctica no se aplicó al *dataset*.

```
def remove_stopwords(self, text):
```

Dentro de este método se define una lista de términos considerados poco informativos, tales como



conectores, artículos o preposiciones comunes, a continuación, cada palabra del texto se compara con dicha lista, y las que coinciden son eliminadas. Como resultado, el texto conserva únicamente aquellas palabras con mayor potencial para el análisis semántico.

### 2.2.2. Tokenización, Lematización y comparación de los lemas con un diccionario en español

Se importó un diccionario en español para la comparación con los lemas de los documentos.

```
1 with open("recursos/diccionario_limpio.txt", "r", encoding="utf-8") as f:
2 valid_words = set(word.strip().lower() for word in f.readlines())
```

Una vez leído el diccionario, se iteró sobre los documentos y se tokenizó con spacy cada uno de ellos, después de haber tokenizado los documentos, se obtuvieron los lemas del mismo y se verificó su existencia con el diccionario previamente leído.

```
1 lemmas_por_documento = []
2 for text in dataset["news_limpio"]:
3     doc = nlp(text)
4     lemmas = [
5         token.lemma_ for token in doc
6         if not token.is_punct and not token.is_space and token.lemma_.lower() in
7         valid_words
8     ]
9     lemmas_por_documento.append(lemmas)
```

## 2.3. Clasificador de tipos de noticias

Dado que nuestro dataset está etiquetado respecto al tipo de noticia, se optó por un enfoque de aprendizaje supervisado, pues se espera que instancias con características similares compartan la misma clase. KNN es un buen candidato al basarse en la vecindad en el espacio de características.

Por otro lado, dada la sensibilidad de KNN al imbalance de clases, se optó por un enfoque híbrido para balancear las clases; generando documentos similares e intercambiando ciertas palabras de las noticias por otras con un significado similar. Además, una vez vectorizados nuestros documentos, se agregaron muestras en el espacio vectorial mediante la técnica **SMOTE** [4]. Además, realizamos **Undersampling** sobre las clases con más registros para tener un balance en todas las clases.

### 2.3.1. Balanceo de clases y vectorización

Para empezar, se definió un número de registros por clase, para determinar el número de clases que se deben generar o cortar por clase. Para este proyecto, se definió un número de 200 registros por muestra. Una vez definido el número de registros deseados, se realizó undersampling sobre las clases que contienen más registros de los esperados, en este caso sobre las clases Macroeconomía y Alianzas.

```
1 clases_undersampling = ['Macroeconomia', 'Alianzas']
2 X_undersampled, y_undersampled = aplicar_undersampling(
3     documentos_originales,
4     etiquetas_originales,
5     clases_undersampling,
6     200,
7     modelo_doc2vec
8 )
```



Para realizar Undersampling sobre las clases minoritarias, se hizo uso de la librería proporcionada por **Imbalancead Learn** llamada **Near Miss** en su primera versión, que básicamente elimina las muestras mas cercanas a la clase minoritaria [5].

```

1 def aplicar_undersampling(documentos: List[List[str]], etiquetas: List[str],
2     clases_objetivo: List[str], n_muestras: int,
3     modelo: Doc2Vec) -> Tuple[np.ndarray, List[str]]:
4
5     print(f"Aplicando undersampling a clases: {clases_objetivo}")
6
7     # Filtrar documentos de las clases objetivo
8     indices_objetivo = [i for i, etiq in enumerate(etiquetas) if etiq in
9         clases_objetivo]
10    docs_filtrados = [documentos[i] for i in indices_objetivo]
11    etiquetas_filtradas = [etiquetas[i] for i in indices_objetivo]
12
13    # Vectorizar documentos
14    X = np.array([modelo.infer_vector(doc) for doc in docs_filtrados])
15
16    # Configurar estrategia de sampling
17    sampling_strategy = {clase: n_muestras for clase in clases_objetivo}
18
19    # Aplicar Near Miss
20    nm = NearMiss(version=1, sampling_strategy=sampling_strategy)
21    X_resampled, y_resampled = nm.fit_resample(X, etiquetas_filtradas)
22
23    # Convertir a lista si es necesario
24    if hasattr(y_resampled, 'tolist'):
25        y_resampled = y_resampled.tolist()
26
27    return X_resampled, y_resampled

```

Posteriormente, se realizó Oversampling sobre las clases minoritarias para alcanzar el número de registros deseados, en este caso se llevó a cabo este proceso sobre las clases Regulaciones, Sostenibilidad, Otra y Reputación.

```

1 docs_por_clase = separar_documentos_por_clase(documentos_originales,
2     etiquetas_originales)
3
4 clases_oversampling = ['Regulaciones', 'Sostenibilidad', 'Otra', 'Reputacion']
5
6 resultados_oversampling = {}
7 for clase in clases_oversampling:
8     if clase in docs_por_clase:
9         vectores_balanceados = aplicar_oversampling_hibrido(
10             docs_por_clase[clase],
11             clase,
12             200,
13             modelo_doc2vec
14         )
15         resultados_oversampling[clase] = vectores_balanceados
16     else:
17         print(f"Advertencia: Clase '{clase}' no encontrada en los datos")

```

Para aplicar el Oversampling se optó por un enfoque híbrido donde se genero la mitad de los registros faltantes mediante el remplazo de palabras por otras que estén cerca del espacio vectorial de la original y

la otra mitad mediante SMOTE. Esta técnica genera nuevas instancias sintéticas de la clase minoritaria. Mediante la identificación de cada instancia minoritaria, para después buscar sus vecinos más cercanos (generalmente cinco) dentro de la misma clase y creando nuevos puntos interpolando aleatoriamente entre una instancia y alguno de sus vecinos [4].

```

1 def aplicar_oversampling_hibrido(documentos: List[List[str]], etiqueta: str,
2                               n_muestras_objetivo: int, modelo_doc2vec: Doc2Vec) ->
  np.ndarray:
3
4     if len(documentos) >= n_muestras_objetivo:
5         return np.array([modelo_doc2vec.infer_vector(doc) for doc in documentos[:
6                               n_muestras_objetivo]])
7
8     muestras_faltantes = n_muestras_objetivo - len(documentos)
9     muestras_word2vec = muestras_faltantes // 2
10
11     # Se aplica el 1ER M todo de Oversampling
12     muestras_nuevas = generar_muestras_word2vec(documentos, muestras_word2vec)
13
14     documentos_combinados = documentos + muestras_nuevas
15
16     X_vectorizado = np.array([modelo_doc2vec.infer_vector(doc) for doc in
17                               documentos_combinados])
18     y_dummy = np.array([0] * len(documentos) + [1] * len(muestras_nuevas))
19
20     # Se aplica el 2DO M todo de Oversampling
21     muestras_smote = n_muestras_objetivo - len(X_vectorizado)
22     if muestras_smote > 0:
23         k_neighbors = min(5, len(X_vectorizado) - 1)
24         smote = SMOTE(sampling_strategy={1: len(muestras_nuevas) + muestras_smote},
25                       k_neighbors=k_neighbors, random_state=42)
26         X_resampled, _ = smote.fit_resample(X_vectorizado, y_dummy)
27         X_final = X_resampled[:n_muestras_objetivo]
28     else:
29         X_final = X_vectorizado
30
31     return X_final

```

En el código anterior se utiliza la función `generar_muestras_word2vec()` cuya definición se presenta a continuación, en esta se usa `word2vec` para obtener los vectores de las palabras y poder generar términos similares en el espacio vectorial, con una probabilidad de 0.3 de generar nuevas palabras.

```

1 def generar_muestras_word2vec(documentos: List[List[str]], n_muestras: int) -> List[
  List[str]]:
2     if n_muestras <= 0:
3         return []
4     print(f"Generando {n_muestras} muestras con Word2Vec...")
5
6     modelo_w2v = Word2Vec(documentos, window=10, min_count=1, sg=0, workers=4)
7
8     nuevas_muestras = []
9
10    for _ in range(n_muestras):
11        doc_base = documentos[np.random.randint(0, len(documentos))]
12        nueva_muestra = []
13        for palabra in doc_base:
14            if palabra in modelo_w2v.wv and np.random.random() < 0.3:

```

```

15         try:
16             similares = modelo_w2v.wv.most_similar(palabra, topn=3)
17             if similares:
18                 nueva_palabra = similares[np.random.randint(0, len(similares)
19             ][0]
20                 nueva_muestra.append(nueva_palabra)
21             else:
22                 nueva_muestra.append(palabra)
23             except:
24                 nueva_muestra.append(palabra)
25         else:
26             nueva_muestra.append(palabra)
27     nuevas_muestras.append(nueva_muestra)
28     return nuevas_muestras

```

### 2.3.2. Entrenamiento del modelo de clasificación KNN

Para determinar el valor de  $k$  que optimiza las métricas, se evaluaron valores de  $k$  entre 0 y 80, entrenando un modelo KNN con cada uno sobre el conjunto de datos previamente balanceado. Una vez identificado el  $k$  óptimo, se entrenó de nuevo el modelo con ese parámetro y se obtuvieron sus métricas junto con la matriz de confusión.

```

1 k=44
2 k_nn_model = KNeighborsClassifier(n_neighbors=k)
3
4 # Entrenamiento del modelo
5 k_nn_model.fit(X_train, y_train)
6 y_pred = k_nn_model.predict(X_test)
7
8 # Calculo de metricas y matriz de confusion
9 accuracy = accuracy_score(y_test, y_pred)
10 precision = precision_score(y_test, y_pred, average='weighted')
11 recall = recall_score(y_test, y_pred, average='weighted')
12 f1 = f1_score(y_test, y_pred, average='weighted')
13 roc_auc = roc_auc_score(y_test, k_nn_model.predict_proba(X_test), multi_class='ovr')
14 conf_matrix = confusion_matrix(y_test, y_pred)

```

El resultado de las métricas se puede consultar en la sección de resultados, donde además se hará una comparación entre la eficiencia del clasificador con clases balanceadas y sin balancear.

## 2.4. Análisis de sentimientos

En esta sección del proyecto se implementa un sistema de análisis de sentimientos en español que combina dos principales enfoques léxicos:

1. **Creación de Covariables Propias:** Se hizo uso de archivos de texto tipo diccionarios que contenían palabras que fueron clasificadas como positivas y negativas. Asimismo se manejaron listas en Python con palabras que podrían indicar un sentimiento de alegría, tristeza, enojo o miedo en una noticia. Estos archivos pueden ser consultados dentro del Drive del proyecto, en la carpeta de recursos.
2. **Puntuación con léxico especializado:** Por otra parte se hizo uso de una base de datos que asocia palabras o expresiones con valores emocionales (como positividad, negatividad, neutralidad). En este enfoque se asocia a unidades léxicas (lemas) valores de polaridad afectiva en un



intervalo continuo  $[-1, +1]$ , donde  $+1$  indica máxima positividad,  $-1$  máxima negatividad y  $0$  neutralidad [6].

### 2.4.1. Covariables

Se desarrolló la función `extraer_covariables_lemas()` que procesa las listas de lemas por documento previamente obtenidas y calcula métricas sentimentales basadas en diccionarios de palabras positivas, negativas y específicas de emociones (alegría, tristeza, enojo, miedo). Genera tanto conteos absolutos como proporciones relativas, además de frecuencias para palabras clave específicas.

```

1 def extraer_covariables_lemas(lista_lemas):
2
3     tokens = [token.lower() for token in lista_lemas]
4
5     num_pos = sum(token in palabras_positivas for token in tokens)
6     num_neg = sum(token in palabras_negativas for token in tokens)
7
8     total = len(tokens) if len(tokens) > 0 else 1 # evitar divisi n por 0
9
10    # palabras clave espec ificas
11    frecuencia_claves = {}
12    for palabra in palabras_clave:
13        frecuencia_claves[f"freq_{palabra}"] = tokens.count(palabra)
14
15    num_alegria = sum(1 for token in tokens if token in palabras_alegria)
16    num_tristeza = sum(1 for token in tokens if token in palabras_tristeza)
17    num_enojo = sum(1 for token in tokens if token in palabras_enojo)
18    num_miedo = sum(1 for token in tokens if token in palabras_miedo)
19
20    covariables = {
21        "num_palabras_positivas": num_pos,
22        "num_palabras_negativas": num_neg,
23        ... se regresan todas las listas obtenidas.
24    }
25
26    # Agregar frecuencias de palabras clave
27    covariables.update(frecuencia_claves)
28
29    return covariables

```

Las funciones `score_sentiment()` y `score_of_docs()` utilizan el recurso SentiLex (archivo XML `senti-con.es.xml`) para asignar puntuaciones numéricas de polaridad a cada documento, calculando el promedio de las puntuaciones de los lemas que aparecen en el léxico.

El archivo XML presenta una estructura jerárquica donde:

- El elemento raíz `<senticon>` especifica el idioma (“es” para español)
- Contiene capas (`<layer>`) que organizan los términos por niveles de análisis
- Cada capa agrupa lemas en categorías `<positive>` y `<negative>`.
- Cada `<lemma>` contiene texto del lema, y atributos:
  - **pos:** categoría gramatical (a=adjetivo, n=sustantivo)
  - **pol:** valor de polaridad numérico (ej. 0.708)



- **std**: desviación estándar de la anotación (medida de acuerdo inter-anotadores)

```

1  <senticon lang="es">
2    <layer level="1">
3      <positive>
4        <lemma pos="a" pol="0.708" std="0.149"> acertado </lemma>
5        <lemma pos="a" pol="0.906" std="0.125"> admirable </lemma>
6        <lemma pos="n" pol="0.45" std="0.331"> admiración </lemma>
7        <lemma pos="v" pol="0.75" std="0.177"> admirar </lemma>
8      </positive>
9      <negative>
10       <lemma pos="v" pol="-0.554" std="0.047"> abatir </lemma>
11       <lemma pos="a" pol="-0.25" std="0.0"> abochornado </lemma>
12       <lemma pos="v" pol="-0.5" std="0.072"> abochornar </lemma>
13       <lemma pos="n" pol="-0.425" std="0.056"> abominación </lemma>
14       <lemma pos="v" pol="-0.625" std="0.088"> aborrecer </lemma>
15     </negative>
16   </layer>
17   <layer level="2">
18     <positive>
19       ...
20 </senticon>

```

Figura 4: Ejemplo del contenido del archivo `senticon.xml` que se puede consultar dentro de la carpeta `recursos` del proyecto.

Para empezar a hacer uso del archivo `senticon.xml` en el código de Python, primero se lee un archivo con listas de lemas por documento y se parsea el XML de SenticLex usando `ET.parse()`.

```

1  with open('archivos/lemas_por_documento.txt', 'r', encoding='utf-8') as f:
2    listas_de_lemas = [
3      ast.literal_eval(line.strip())
4      for line in f
5      if line.strip() # descarta lineas vacias
6    ]
7
8  tree = ET.parse("recursos/senticon.es.xml")
9  root = tree.getroot()

```

Posteriormente se construye el diccionario, haciendo un mapeo sentilex donde se identifica la clave (string del lema) y el valor de la polaridad numérica (ej. 0.906). Solo se almacena la primera ocurrencia de cada lema, esto para priorizar capas superiores.

```

1  sentilex = {}
2  for layer in root.findall("layer"):
3    for sign in ("positive", "negative"):
4      for lemma in layer.find(sign).findall("lemma"):
5        text = lemma.text.strip()
6        pol = float(lemma.get("pol"))
7
8        if text not in sentilex:
9          sentilex[text] = pol

```

A continuación se manda a llamar a `score_of_docs()` para aplicar `score_sentiment` a cada documento.

```

1  scores = score_of_docs(listas_de_lemas)

```

La función `score_sentiment` calcula el sentimiento promedio de un texto mediante la búsqueda de cada lema en `sentilex` y la agregación de los valores encontrados. También normalizando por cantidad de términos encontrados.

```
1 def score_sentiment(text):
2     scores = []
3     for lemma in text:
4         if lemma in sentilex:
5             scores.append(sentilex[lemma])
6
7     return sum(scores) / len(scores)
```

Tras aplicar los cambios al conjunto de documentos lematizados, se guardan las covariables en un nuevo archivo `csv` que se usará para identificar las agrupaciones de los datos basándose en los sentimientos de las noticias, para esto se usará `k-means`.

#### 2.4.2. Entrenamiento de modelo de K-means

Se carga el archivo con las covariables y se eligen las columnas a considerar para el entrenamiento del modelo.

```
1 dataframe = pd.read_csv('archivos/sentiment_scores_combined.csv')
2
3 # Obtener las características de interés
4 feature_columns = [
5     'scores',
6     'proporcion_palabras_positivas',
7     'proporcion_palabras_negativas',
8     'sentimiento_lexico_total',
9     'proporcion_alegria',
10    'proporcion_tristeza',
11    'proporcion_enojo',
12    'proporcion_miedo'
13 ]
14
15 # Preparar los datos
16 X = dataframe[feature_columns].fillna(0)
17
18 scaler = StandardScaler()
19 X_scaled = scaler.fit_transform(X)
```

El siguiente paso consiste en implementar el algoritmo de clustering K-Means con `k=3`, debido a los grupos predefinidos (correspondientes a sentimientos positivo, negativo y neutro). Utilizando los datos estandarizados, el algoritmo agrupa automáticamente las noticias según la similitud de sus características sentimentales.

```
1 n_clusters = 3
2 kmeans = KMeans(n_clusters = n_clusters, random_state=42, n_init=10)
3 cluster_labels = kmeans.fit_predict(X_scaled)
```

Una vez entrenado K-Means se calculan métricas de calidad (`Silhouette Score` e `inercia`) para evaluar la cohesión y separación de los grupos formados. Los centroides de cada cluster se analizan para

asignar etiquetas semánticas, y se genera una representación tridimensional de los clusters que se puede visualizar en la sección de resultados.

```

1 silhouette_avg = silhouette_score(X_scaled, cluster_labels)
2 inertia = kmeans.inertia_
3
4 # Analizar centroides para interpretar clusters
5 centroids = kmeans.cluster_centers_
6 centroids_df = pd.DataFrame(centroids, columns=feature_columns)
7
8 fig = plt.figure(figsize=(10, 7))
9 ax = fig.add_subplot(111, projection='3d')
10
11 # Elegir tres características para la visualización
12 x = X_scaled[:, feature_columns.index('scores')]
13 y = X_scaled[:, feature_columns.index('proporcion_palabras_positivas')]
14 z = X_scaled[:, feature_columns.index('proporcion_palabras_negativas')]
15
16 plt.show()

```

Por último, las clases asignadas por K-Means a cada instancia, fueron mapeadas como positiva, negativa o neutra y agregadas como columna al dataset original.

proporcion_palabras_negativas	sentimientos_lectos_total	sum_palabras_alegría	sum_palabras_tristeza	sum_palabras_enojo	sum_palabras_miedo	proporcion_alegría	proporcion_tristeza	proporcion_enojo	proporcion_miedo	cluster	sentiment_label
0.0	2	0	0	0	0	0.0	0.0	0.0	0.0	1	Negativa
0.00925	-1	0	0	0	0	0.0	0.0	0.0	0.0	0	Positiva
0.0	1	1	0	0	0	0.0064516129052258	0.0	0.0	0.0	1	Negativa
0.0	2	0	0	0	0	0.0	0.0	0.0	0.0	0	Positiva
0.002832061189017	3	0	0	0	0	0.0	0.0	0.0	0.0	0	Positiva
0.0	3	0	0	0	0	0.0	0.0	0.0	0.0	1	Negativa
0.0000000000000000	2	0	0	0	0	0.0	0.0	0.0	0.0	0	Positiva
0.0	0	0	0	0	0	0.0	0.0	0.0	0.0	0	Positiva
0.0	1	0	0	0	0	0.0	0.0	0.0	0.0	0	Positiva
0.0	0	0	0	0	0	0.0	0.0	0.0	0.0	0	Positiva
0.0	1	0	0	0	0	0.0	0.0	0.0	0.0	0	Positiva
0.0	12	0	0	0	0	0.0	0.0	0.0	0.0	1	Negativa
0.0010487064870548	-1	0	0	0	0	0.0	0.0	0.0	0.0	0	Positiva
0.0	2	0	0	0	0	0.0	0.0	0.0	0.0	0	Positiva
0.0	2	0	0	0	0	0.0	0.0	0.0	0.0	0	Positiva
0.0	1	0	0	0	0	0.0	0.0	0.0	0.0	0	Positiva
0.0	0	0	0	0	0	0.0	0.0	0.0	0.0	0	Positiva
0.0	0	0	0	0	0	0.0	0.0	0.0	0.0	0	Positiva
0.0	5	0	0	0	0	0.0	0.0	0.0	0.0	1	Negativa
0.0	0	0	0	0	0	0.0	0.0	0.0	0.0	0	Positiva
0.0	1	0	0	0	0	0.0	0.0	0.0	0.0	0	Positiva
0.0	5	0	0	0	0	0.0	0.0	0.0	0.0	1	Negativa
0.0	6	0	0	0	0	0.0	0.0	0.0	0.024390243902439	1	Negativa
0.0027397260273972	4	0	0	0	0	0.0	0.0	0.0	0.0	0	Positiva
0.0042194092627004	2	0	0	0	0	0.0	0.0	0.0	0.002109704611302	0	Positiva
0.0048309178743901	2	0	0	0	0	0.0	0.0	0.0	0.0	0	Positiva

Figura 5: Fragmento del dataset que contiene las covariables y el resultado de la Clusterización.

### 3. Resultados

#### 3.1. Clasificador de tipos de noticias

##### 3.1.1. Métricas de rendimiento

En el siguiente cuadro se muestra la comparación del rendimiento del clasificador KNN medido con las métricas de *Accuracy*, *Precision*, *Recall* y *F1-Score*. La comparación entre las métricas con y sin balanceo de clases revela mejoras significativas en todos los indicadores de rendimiento

Métricas	Métricas con balanceo de clases	Métricas sin balanceo de clases
Accuracy	0.925	0.7868
Precision	0.9261	0.8213
Recall	0.925	0.7868
F1-Score	0.9251	0.7837

Cuadro 1: Comparación de métricas entre diferentes métodos de vectorización

<b>Accuracy</b>	Predictions/ Classifications	$\frac{\text{Correct}}{\text{Correct} + \text{Incorrect}}$
<b>Precision</b>	Predictions/ Classifications	$\frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}}$
<b>Recall</b>	Predictions/ Classifications	$\frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}}$
<b>F1</b>	Predictions/ Classifications	$\frac{2 * \text{True Positive}}{\text{True Positive} + 0.5 (\text{False Positive} + \text{False Negative})}$

Figura 6: Significado de las métricas de rendimiento calculadas [7]

Se hace evidente que cuando se trabaja sin balanceo de clases, el rendimiento se deteriora significativamente. La diferencia de aproximadamente 0.15 puntos en todas las métricas revela que el desbalance no solo afecta la accuracy general, sino que también compromete tanto la capacidad del modelo para identificar correctamente los casos positivos (recall) como su precisión al hacer predicciones positivas. Este patrón indica que el modelo aprende a predecir principalmente la clase mayoritaria para maximizar la accuracy, pero perdiendo capacidad discriminativa real entre las clases.

##### 3.1.2. Matriz de confusión

Al evaluar el rendimiento con datos balanceados (representado en la Figura 7), la matriz de confusión muestra que los valores más altos se concentran en la diagonal principal, lo que indica una alta tasa de clasificaciones correctas.

También es visible la clasificación perfecta en múltiples clases, como la clase 4 y la clase 5, que muestran 35 y 37 instancias correctamente clasificadas, respectivamente. Los errores de clasificación son mínimos y se distribuyen de manera dispersa, con la mayor confusión entre las clases 0 y 2, donde se



observa un solapamiento de tres instancias, lo cual puede deberse a cierta similitud semántica entre esas categorías.

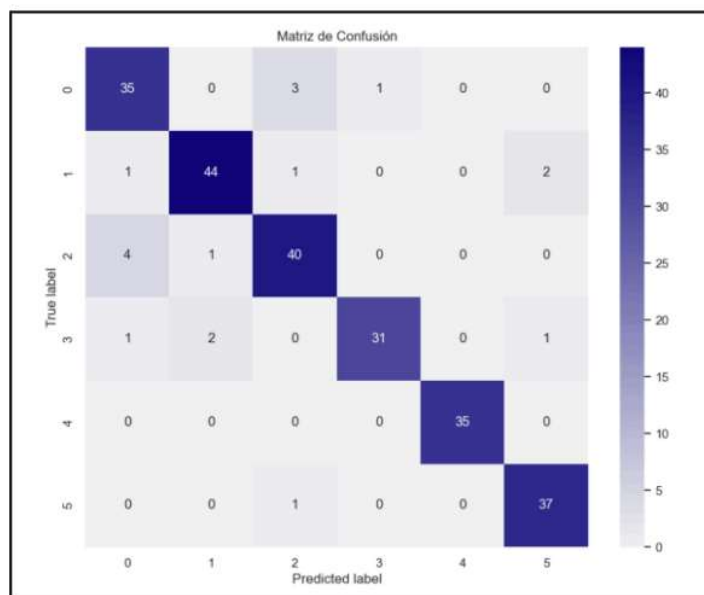


Figura 7: Matriz de confusión con los datos previamente balanceados

En contraste, al observar la matriz de confusión obtenida con datos no balanceados (Figura 8), se evidencian las limitaciones del modelo debido al sesgo hacia las clases mayoritarias. Las clases 0 y 2 presentan un número desproporcionadamente alto de instancias, con 47 y 59, respectivamente, mientras que las clases 4 y 5 están significativamente subrepresentadas, con solo 15 y 3 instancias.

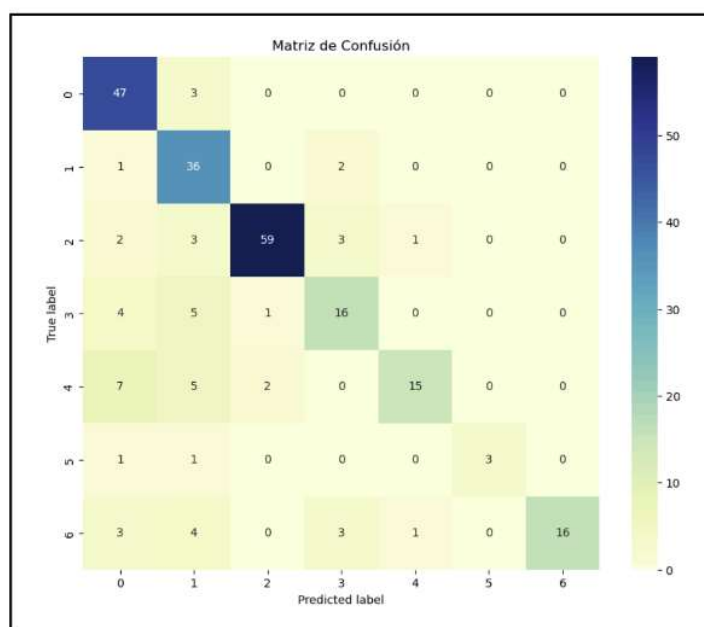


Figura 8: Matriz de confusión con los datos sin balancear

Estos resultados validan la hipótesis de que el desbalance de clases constituía una limitación significativa en el rendimiento del modelo, y que la estrategia híbrida de balanceo que combinó técnicas de undersampling como Near Miss con oversampling mediante Word2Vec y SMOTE logró mitigar efectivamente este problema.

### 3.1.3. Distribución de Clusters

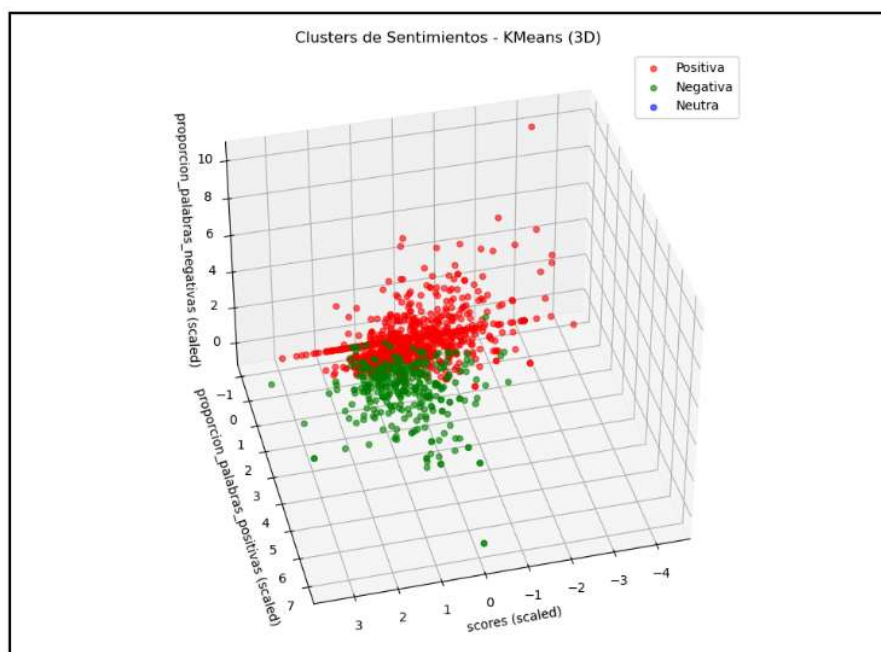


Figura 9: Gráfica 3D del agrupamiento de los datos, tomando 3 características (proporción palabras positivas, proporción palabras negativas y scores)

La visualización tridimensional de los clusters de sentimientos, presentada en la Figura 8, revela que las covariables no lograron distinguir de manera óptima las tres clases que se pretendía identificar.

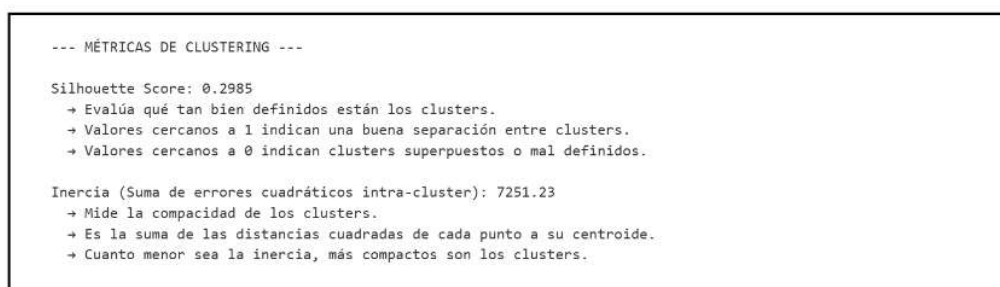


Figura 10: Resultados del rendimiento de K-Means.

El Silhouette Score de 0.2985 indica que los clusters están algo definidos, pero existe cierto solapamiento entre ellos o una separación poco clara. Por otro lado, la inercia de 7251.23 sugiere un nivel razonable de cohesión interna.

## 4. Conclusiones

En esta práctica integradora se comenzó por realizar un preprocesamiento exhaustivo, pasando por la eliminación de puntuación, la normalización y la depuración de *stopwords*. Este proceso en conjunto con la verificación de lemas frente a un diccionario en español, permitió obtener representaciones limpias y coherentes que demostraron la importancia del preprocesamiento para garantizar la calidad del texto. Gracias a este paso se obtuvo una base sólida para las etapas posteriores.

El balanceo de clases, conseguido mediante una combinación de *undersampling* con Near Miss y un enfoque híbrido de *oversampling* que incorpora generación de muestras con Word2Vec y SMOTE, resultó decisivo para reducir el sesgo hacia las categorías mayoritarias, problema que fue evidente al analizar la matriz de confusión que generaba el clasificador KNN después de entrenarlo con el *dataset* sin balancear. Gracias a este paso, el clasificador KNN mejoró su rendimiento de manera notable: la exactitud pasó de 0.7868 a 0.925 y el *F1-score* de 0.7837 a 0.9251.

Por otro lado, la implementación del análisis de sentimientos basado en diccionarios propios y en el recurso SenticLex permitió asignar puntuaciones continuas en el intervalo  $[-1, +1]$  y construir co-variables que capturan la polaridad afectiva de cada texto. Estas características sentimentales fueron lo suficientemente discriminativas como para agrupar, mediante K-Means con  $k = 3$ , las noticias en clústeres medianamente identificables como positivo, negativo y neutro. No obstante, el resultado de las métricas de *silhouette score* e inercia demostraron que el modelo tiene áreas de mejora.

Los resultados en el análisis de sentimientos son los esperados, pues no se esperaba una gran separación de los clusters debido a la naturaleza del dataset, que está conformada de noticias neutras con temáticas financieras.

En conjunto, este trabajo sienta las bases para mejoras en futuros trabajos relacionados con el procesamiento de lenguaje natural, las técnicas utilizadas demostraron que con una implementación diligente, y experimentación constante, se puede alcanzar altos niveles de rendimiento en tareas de clasificación y agrupamiento de textos.

## Referencias

- [1] S. J. Ahmed. “Natural Language Processing in News Classification: Unleashing the Power of AI in Media.” (21 de jul. de 2023), dirección: <https://medium.com/dataduniya/natural-language-processing-in-news-classification-unleashing-the-power-of-ai-in-media-34afc60c58eb> (visitado 21-07-2023).
- [2] Kaggle, *Spanish news classification*, ver. 1.0, Online; accessed 26-November-2022, Kaggle, nov. de 2022. dirección: <https://www.kaggle.com/datasets/kevinmorgado/spanish-news-classification>.
- [3] S. Vijayarani y M. J. Ilamathi, “Preprocessing techniques for text mining - An overview,” *International Journal of Computer Science & Communication Networks*, vol. 5, n.º 1, 2013. dirección: [https://www.researchgate.net/profile/Vijayarani-Mohan/publication/339529230\\_Preprocessing\\_Techniques\\_for\\_Text\\_Mining\\_-\\_An\\_Overview/links/5e57a0f7299bf1bdb83e7505/Preprocessing-Techniques-for-Text-Mining-An-Overview.pdf](https://www.researchgate.net/profile/Vijayarani-Mohan/publication/339529230_Preprocessing_Techniques_for_Text_Mining_-_An_Overview/links/5e57a0f7299bf1bdb83e7505/Preprocessing-Techniques-for-Text-Mining-An-Overview.pdf).
- [4] imbalanced-learn. “SMOTE — Version 0.13.0.” (), dirección: [https://imbalanced-learn.org/stable/references/generated/imblearn.over\\_sampling.SMOTE.html](https://imbalanced-learn.org/stable/references/generated/imblearn.over_sampling.SMOTE.html) (visitado 25-06-2025).
- [5] imbalanced-learn. “3. Under-sampling — Version 0.13.0.” (), dirección: [https://imbalanced-learn.org/stable/under\\_sampling.html#controlled-under-sampling](https://imbalanced-learn.org/stable/under_sampling.html#controlled-under-sampling) (visitado 25-06-2025).
- [6] ResearchGate. “Sentiment Analysis on Twitter Data for Portuguese Language - Scientific Figure.” (), dirección: [https://www.researchgate.net/figure/Results-without-pre-processing-and-using-the-Sentilex-lexicon\\_tbl1\\_262175717](https://www.researchgate.net/figure/Results-without-pre-processing-and-using-the-Sentilex-lexicon_tbl1_262175717) (visitado 24-06-2025).
- [7] ResearchGate, *Unlocking the potential of deep learning for marine ecology: Overview, applications, and outlook - Scientific Figure*, [Online]. Available: [https://www.researchgate.net/figure/Evaluation-metrics-accuracy-precision-recall-F-score-and-Intersection-over-Union\\_fig2\\_358029719](https://www.researchgate.net/figure/Evaluation-metrics-accuracy-precision-recall-F-score-and-Intersection-over-Union_fig2_358029719), Accessed: 2025-06-26.