



**EDUCACIÓN**

SECRETARÍA DE EDUCACIÓN PÚBLICA



TECNOLÓGICO  
NACIONAL DE MÉXICO

# Tecnológico Nacional de México

Centro Nacional de Investigación

Y Desarrollo Tecnológico

## Tesis de Maestría

Gramática libre de contexto  
para la identificación de noticias falsas  
sobre enfermedades basada  
en el análisis de literatura médica

Presentada por

**Ing. Julio Cesar Arroyo Gómez**

Como requisito para la obtención del grado de

**Maestro en Ciencias de la computación**

Director de tesis

**Dr. Noé Alejandro Castro Sánchez**

Codirector de tesis

**Dr. Juan Gabriel González Serna**

Cuernavaca, Morelos, México. Septiembre del 2024

## **Dedicatoria**

El presente trabajo de investigación lo dedico principalmente a Dios, por darme día a día la fuerza necesaria para continuar en este proceso y poder cumplir con uno de mis anhelos más deseados.

A mis padres Laura y Julio Cesar; por su amor, trabajo, apoyo y sacrificio en todos estos años, gracias por inculcar el ejemplo de esfuerzo dedicación y perseverancia ante las distintas adversidades que se me han presentado en mi camino; por estar conmigo en todo momento y brindarme palabras de aliento cuando más las necesito, les agradezco por su cariño y apoyo incondicional durante todo este proceso porque gracias a ellos he podido llegar hasta aquí y convertirme en lo que soy.

## **Agradecimiento**

Es un honor dirigirme a ustedes para expresar mi más profundo agradecimiento al Consejo Nacional de Ciencias y Tecnología (CONACYT) por el apoyo durante el desarrollo de esta tesis para la obtención del grado de Maestro en Ciencias de la Computación mediante un sistema de becas de posgrado. Deseo manifestar mi gratitud por la oportunidad que me ha sido otorgada para llevar a cabo mis estudios.

Agradezco al Centro Nacional de Investigación y Desarrollo Tecnológico (CENIDET), pertenecer al (TECNM), por brindarme la oportunidad de superarme académicamente mediante el programa de Maestrías en Ciencias y por prestarme sus instalaciones para cumplir con dicha meta.

Quiero agradecer principalmente a mi director de tesis, Dr. Noé Alejandro Castro Sánchez, por la gran oportunidad y confianza que me brindó desde un principio para desarrollar esta investigación. A mi comité revisor el Dr. Dante Mújica Vargas y Dr. Nimrod Gonzales Franco y a mi Co-Director el Dr. Juan Gabriel González Serna. Gracias por sus consejos, paciencia, tiempo y apoyo brindado durante mi estancia y formación como maestro.

## Resumen

Las noticias falsas de temas sobre la salud representan información no verificada ni comprobada como verdadera. Este tipo de noticias tiene un impacto negativo tanto en lo emocional como en lo personal, afectando a personas adultas y jóvenes; incluso, si se llegan a seguir las sugerencias que indican pueden causar la muerte, un ejemplo es lo sucedido durante la pandemia COVID-19, la pandemia de la época moderna más devastadora, donde surgieron diferentes noticias falsas que mencionan que la vacuna del COVID-19 contenía un microchip con el que el gobierno podía rastrear en tiempo real a las personas y saber su ubicación actual en todo momento. A raíz de esto, diversas personas creyeron en este tipo de noticias evitando vacunarse, y con ello se convirtieron en agentes portadores de la enfermedad y, dado que la gran mayoría no tenía inmunidad al virus, contribuyeron en la propagación y muerte de gran cantidad de personas. En este trabajo de investigación se propone un sistema de detección de noticias falsas sobre enfermedades basado en el análisis de la literatura médica a través de una gramática libre de contexto. En este contexto, el objetivo principal de este trabajo es lograr identificar hechos a través de las reglas de producción con el propósito de verificar si la noticia es falsa o verdadera. Esta metodología se presenta como una herramienta complementaria para las personas que navegan en las diferentes plataformas digitales para generar una verificación de lo que están visualizando es falso o verdadero. De acuerdo al experimento realizado se obtiene un valor de Precisión de 0.97%, Recall 0.72% y un F1 –score de 0.82%.

**Palabras claves:** Noticias falsas, salud, Procesamiento de lenguaje natural, Recuperación de información, Extracción de información.

## Abstract

Fake Health News is unverified information proven to be true. This type of news has a great social impact as well as the emotional and personal impact on adults and young people, since they are not as informed about the different events and activities that arise every day in such a way that these types of people do not manage to identify the news they see on the different platforms and social networks if it is true or false, an example is when certain events arose worldwide such as the war in Russia or events at a political level such as the electoral elections or in our case in the type of news that this research is in the area of health, highlighting the previous COVID-19 pandemic being the most devastating pandemic of modern times, emerging as a result of this different fake news where they mentioned that the COVID-19 vaccine contains a micro -chip, which contains a locator so that the government can track us in real time and know our current location at all times and as a result of this certain people, regardless of their religion, skin color or beliefs, believed in this type of news and they were not vaccinated, leading to the consequence that when they were infected with this previously mentioned disease, since they did not have a certain immunity to the virus because they did not have the COVID-19 vaccine, many people died from the current virus. This study proposes a system for detecting fake news about diseases based on the analysis of medical literature through production rules. In this context, the main objective of this work is to identify facts through production rules with the purpose of verifying whether the news is false or true. This methodology is presented as a complementary tool for people who browse different digital platforms to generate a verification of what they are viewing is true or false. According to the experiment carried out, it obtains a Precision value of 0.97%, Recall 0.72% and an F1 –score of 0.82%.

**Keywords:** Fake news, health, Natural language processing, Information retrieval, Information extraction.



Cuernavaca, Mor., 14/junio/2024

OFICIO No. DCC/074/2024  
Asunto: Aceptación de documento de tesis  
CENIDET-AC-004-M14-OFICIO


**CARLOS MANUEL ASTORGA ZARAGOZA**  
**SUBDIRECTOR ACADÉMICO**  
**PRESENTE**

Por este conducto, los integrantes de Comité Tutorial de JULIO CESAR ARROYO GÓMEZ con número de control M2ICE053, de la Maestría en Ciencias de la Computación, le informamos que hemos revisado el trabajo de tesis de grado titulado "GRAMÁTICA LIBRE DE CONTEXTO PARA LA IDENTIFICACIÓN DE NOTICIAS FALSAS SOBRE ENFERMEDADES BASADA EN EL ANÁLISIS DE LITERATURA MÉDICA" y hemos encontrado que se han atendido todas las observaciones que se le indicaron, por lo que hemos acordado aceptar el documento de tesis y le solicitamos la autorización de impresión definitiva.

  
\_\_\_\_\_  
NOÉ ALEJANDRO CASTRO SÁNCHEZ  
Director de tesis

  
\_\_\_\_\_  
JUAN GABRIEL GONZÁLEZ SERNA  
Codirector de tesis

  
\_\_\_\_\_  
DANTE MÚJICA VARGAS  
Revisor 1

  
\_\_\_\_\_  
NIMROD GONZÁLEZ FRANCO  
Revisor 2

C.c.p. Depto. Servicios Escolares.  
Expediente / Estudiante



Interior Internado Palmira S/N, Col. Palmira, C. P. 62490, Cuernavaca, Morelos  
Tel. 01 (777) 3627770, ext. 3202. e-mail: dcc@cenidet.tecnm.mx | cenidet.tecnm.mx



Cuernavaca, Mor.,  
No. de Oficio:  
Asunto:


20/junio/2024  
SAC/193/2024  
Autorización de Impresión  
de tesis

**JULIO CESAR ARROYO GÓMEZ**  
**CANDIDATO AL GRADO DE MAESTRO**  
**EN CIENCIAS DE LA COMPUTACIÓN**  
**P R E S E N T E**

Por este conducto, tengo el agrado de comunicarle que el Comité Tutorial asignado a su trabajo de tesis titulado “GRAMÁTICA LIBRE DE CONTEXTO PARA LA IDENTIFICACIÓN DE NOTICIAS FALSAS SOBRE ENFERMEDADES BASADA EN EL ANÁLISIS DE LITERATURA MÉDICA”, ha informado a esta Subdirección Académica, que están de acuerdo con el trabajo presentado. Por lo anterior, se le autoriza a que proceda con la impresión definitiva de su trabajo de tesis.

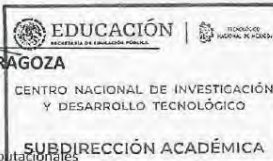
Esperanto que el logro del mismo sea acorde con sus aspiraciones profesionales, reciba un cordial saludo.

**ATENTAMENTE**  
*Excellencia en Educación Tecnológica-*  
*"Conocimiento y Tecnología al Servicio de México"*

  
**CARLOS MANUEL ASTORA ZARAGOZA**  
**SUBDIRECTOR ACADÉMICO**

C.c.p. Departamento de Ciencias Computacionales  
Departamento de Servicios Escolares

CMAZ/lmz



# Índice general

1. Introducción.....	11
1.1 Planteamiento del problema.....	12
1.2 Objetivos.....	13
1.2.1 Objetivo general.....	13
1.2.2 Objetivos Específicos.....	13
1.3 Alcances y Limitaciones.....	14
1.3.1 Alcances.....	14
1.3.2 Limitaciones.....	14
1.4 Justificación.....	15
1.5 Organización de la tesis.....	15
2. Marco conceptual.....	18
2.1 Noticias Falsas.....	18
2.2 Procesamiento de lenguaje natural.....	19
2.3 Gramática libre de contexto.....	20
2.4 Procesamiento de texto.....	20
3. Antecedentes.....	22
3.1 3HAN: A Deep Neural Network for Fake News.....	22
Detection [4].....	22
3.2 Fake News Early Detection: An Interdisciplinary Study [5].....	25
3.3 Mapeo sistemático sobre los métodos, técnicas y tecnologías orientadas a la detección de noticias [6].....	26
3.4 Detection of fake news using deep learning CNN–RNN based methods [7].....	28
3.5 La Noticias falsas y desinformación sobre el Covid-19: análisis comparativo de seis países iberoamericanos [8].....	29
3.6 Detection of fake news in a new corpus for the Spanish language [9].....	30



3.7 Computational Fast Checking from Knowledge Networks [10].....	31
3.8 IDENTIFICACIÓN AUTOMÁTICA DE NOTICIAS FALSAS EN ESPAÑOL UTILIZANDO TÉCNICAS DE MINERÍA DE DATOS Y PROCESAMIENTO DE LENGUAJE NATURAL [11].....	31
3.2 Trabajos representativos... ..	33
4. Método de Solución .....	38
4.1 Propuesta de Solución.....	38
FASE 1: Creación de Gramática.....	39
FASE 2: Pre procesamiento de noticia.....	49
FASE 3: Recuperación de Información.....	51
FASE 4: Identificación de contradicciones.....	66
FASE 5: Diseño de Interfaz.....	69
5. Resultados.....	75
5.1 Experimentación.....	75
5.2 Resultados.....	78
6. Conclusiones.....	82
6.1 Objetivos y Alcances Logrados.....	82
6.2 Resultados del trabajo de investigación.....	84
6.2.1 Productos.....	84
6.3 Conclusiones.....	84
6.4 Trabajos a Futuro.....	85
Referencias.....	86
Anexo.....	88
Anexo A cambios de interfaz.....	88
Anexo B Producción científica.....	89

# Índice de Figuras

Figura 1 Arquitectura de la red neuronal utilizada para el modelo 3HAN .....	23
Figura 2 Los nodos rojos y azules corresponden a miembros del congreso, los nodos grises a Ideologías y nodos blancos a vértices de cualquier otro tipo .....	30
Figura 3 Metodología de solución .....	38
Figura 4 Fase 1 de la metodología de solución.....	39
Figura 5 Ejemplo de noticia que hace referencia a la noticia fuente.....	40
Figura 6 Ejemplo de noticias que no hacen referencia a la noticia fuente.....	41
Figura 7 Ejemplo de corpus generado .....	42
Figura 8 Ejemplo de campos semántico .....	43
Figura 9 Gramática.CFG.....	48
Figura 10 Fase 2 de la metodología de solución.....	50
Figura 11 Fase 3 de la metodología de solución.....	51
Figura 12 Diagrama de planificación para la conexión al API de Google.....	53
Figura 13 Nombre de Api a utilizar Custom Search API.....	54
Figura 14 Apis y servicios habilitados .....	55
Figura 15 Nombre de tipos de credenciales .....	55
Figura 16 Nombre de clave de Api de credencial .....	56
Figura 17 Motor de búsqueda programable .....	56
Figura 18 Crear un buscador .....	57
Figura 19 Todos los buscadores.....	57
Figura 20 Datos Básicos .....	58
Figura 21 Código añadido a la herramienta en Python .....	59
Figura 22 Fórmula BM-25 .....	61
Figura 23 Fase 4 de la metodología de solución.....	66
Figura 24 Árbol sintáctico de afirmaciones positivas .....	67
Figura 25 Árbol sintáctico de afirmaciones negativas... ..	67
Figura 26 Fase 5 de la metodología de solución.....	69
Figura 27 Vista de inicio.....	70
Figura 28 Analizador... ..	70
Figura 29 Artículo Capturados .....	71
Figura 30 Hechos detectados .....	71
Figura 31 Veredictos.....	72
Figura 32 Vista de noticias populares .....	72
Figura 33 Acerca de .....	73
Figura 34 Matriz de confusión .....	76
Figura 35 Fórmula de Precisión .....	77
Figura 36 Fórmula de Recall .....	77

Figura 37 Fórmula de F1 Score .....	77
Figura 38 Rediseño de vista de noticias populares.....	90
Figura 39 Constancia de ponencia.....	91
Figura 40 Constancia de ponencia.....	92

## Lista de Tablas

Tabla 1 Resumen del estado del arte.....	33
Tabla 2 Clasificación y características del corpus.....	42
Tabla 3 Ejemplo de campos semánticos de noticias falsas alusivos a su color de categoría.....	44
Tabla 4 Categoría con sus palabras en común.....	45
Tabla 5 Acciones positivas y negativas de cada noticia falsa.....	46
Tabla 6 Reglas identificadas.....	47
Tabla 7 Método de pre procesamiento de noticia.....	50
Tabla 8 Pruebas de Api BM-25.....	63
Tabla 9 Pruebas Api Fast BM – 25.....	64
Tabla 10 Pruebas de Api Sparce – Retriver.....	65
Tabla 11 Resultados de matriz de confusión.....	79
Tabla 12 Resultados de las métricas de evaluación utilizadas.....	79
Tabla 13 Comparación de métodos utilizados para la identificación de noticias falsas.....	80
Tabla 14 Objetivos cumplido.....	84
Tabla 15 Alcances realizados.....	84

# Capítulo I

Introducción

## 1. Introducción

Las noticias falsas son noticias ficticias creadas principalmente para manipular e influenciar las ideologías de las personas llegando a tener el riesgo de afectar su salud mental y física por los temas delicados a los que pueden hacer referencia, este tipo de noticias han tenido una rápida y fácil difusión a través de las redes sociales.

Anteriormente, para publicar cualquier comunicado de forma masiva era necesario ser partícipe de una empresa televisiva, periodística, entre etc., hoy en día solo con tener acceso a internet y en tan solo unos minutos se pueden realizar publicaciones de forma masiva sin supervisar la veracidad y posteriormente ser publicadas. Se ha reconocido la importancia que ha tenido el lograr combatir la demanda de divulgación y propagación de noticias falsas, identificando que se trata de un problema de una gran complejidad logrando formular la pregunta ¿cómo saber si una noticia es falsa o verdadera? Ya que, para lograr conocer su veracidad, se necesita conocer qué fragmento de texto de esa noticia es falso o verdadero, esto motiva la realización de esta investigación y la propuesta de un método poco abordado en los trabajos que se han realizado en torno a este tema, que enfoca su interés en contrastar el contenido de las noticias con lo que medios confiables y especializados comentan respecto a los hechos mencionados en las noticias.

## **1.1 Planteamiento del problema**

Las noticias falsas pueden llegar a ser muy perjudiciales dependiendo de su contenido, en particular las relacionadas con temas de salud, pues creer en ellas pone en peligro la salud física y afecta de manera directa la vida de las personas.

La detección automática de noticias es un problema abierto pues enfrenta varios desafíos siendo uno de ellos la falta de datos que se necesitan para entrenar modelos de clasificación automática, principalmente en idiomas diferentes al inglés, lo que obliga a buscar métodos alternos que permitan lograr su identificación. En este trabajo se propone un método basado en la creación de una gramática libre de contexto para identificar hechos en noticias y buscar su confirmación o contradicción en información médica confiable.

## 1.2 Objetivos

### 1.2.1 Objetivo general

Reconocer noticias falsas relacionadas con temas de salud comparando su contenido con información de notas médicas obtenidas de fuentes confiables a través de una gramática libre de contexto.

### 1.2.2 Objetivos Específicos

Los objetivos específicos planteados son los siguientes:

- 1) Identificar y recolectar noticias falsas y notas que traten temas de salud.
- 2) Conocer los patrones de uso del lenguaje utilizado en noticias falsas y notas médicas a través de su análisis manual.
- 3) Determinar la confirmación o contradicción de hechos contenidos en las noticias a través de una gramática libre de contexto.
- 4) Proponer un proceso automático de recuperación de notas relacionadas con noticias falsas.
- 5) Evaluar las reglas definidas de acuerdo con la precisión y cobertura que muestran en la detección de noticias falsas.

## **1.3 Alcances y Limitaciones**

### 1.3.1 Alcances

- 1) El método tomará como fuente confiable más de una fuente médica para obtener información relacionada con la noticia a analizar.
- 2) Se identificarán los fragmentos de información del contenido de las noticias como falsos o verdaderos.

### 1.3.2 Limitaciones

- 1) El modelo para la identificación de noticias falsas solo funcionará con temas de salud.
- 2) Las noticias y fuentes de información que se procesen por el modelo deberán estar escritas en el idioma español.



## **1.4 Justificación**

Debido al alto crecimiento de los medios de comunicación se han publicado masivamente noticias falsas desarrolladas por individuos para dañar o desacreditar a personas o instituciones con o sin fines de lucro. Las noticias falsas de temas sobre la salud tienen una gran repercusión en las personas afectando sus ideologías, emociones e incluso de manera extrema, este tipo de noticias puede traer graves consecuencias como la muerte si se siguen este tipo de noticias.

Así, el presente trabajo permitirá que las personas tengan un mejor control e identificación de conocer si una noticia es falsa o verdadera a través de los diferentes métodos de identificación, ya que hablar de noticias de la salud puede traer graves consecuencias.

## **1.5 Organización de la tesis**

Este documento de tesis está compuesto por seis capítulos, anexos y referencias.

### **Capítulo I**

**Introducción:** En este capítulo se describe el origen de la investigación, el problema que se trata de solucionar y la justificación para elaborar esta tesis, así como sus objetivos, alcances y limitaciones.

### **Capítulo II**

**Marco Conceptual:** En este capítulo se definen los principales conceptos para el tema de tesis desarrollado.

### **Capítulo III**

**Estado del Arte:** En este capítulo se proporciona información de los trabajos de investigación que se han realizado con relación al tema de investigación de esta tesis.

## **Capítulo IV**

**Metodología de Solución:** En este capítulo se describe a detalle el método propuesto para identificar comentarios con contenido engañoso; el cual se compone por cinco fases:

Creación de Gramática, Procesamiento de Noticia, Recuperación de información, Identificación de contradicciones e Interfaz.

## **Capítulo V**

**Pruebas y Resultados:** En este capítulo se analizan los resultados obtenidos a través de las pruebas experimentales que se realizaron durante el trabajo de investigación.

## **Capítulo VII**

**Conclusiones:** En este capítulo se exponen las conclusiones que se derivan de esta investigación.

# Capítulo II

Marco Conceptual

## 2. Marco conceptual

En el presente capítulo se presentan los conceptos teóricos relacionados con las temáticas tratadas en esta investigación, con la finalidad de proporcionar una definición concreta que permita abordar la lectura del documento con una base teórica.

### 2.1 Noticias Falsas

El término noticias falsas, se define como informaciones publicadas deliberadamente en medios digitales que no han sido comprobadas ni verificadas que crecen de fuentes identificadas y que no cuentan con la información de un editor. Actualmente, las noticias falsas han logrado un alto grado de notoriedad debido a su utilización en campañas políticas en Brasil y en Estados Unidos, pero su uso y estudio es tan antiguo como la prensa misma. Las noticias falsas no son un fenómeno nuevo, por el contrario, siempre han existido en el periodismo, la diferencia es que en la actualidad estas son más visibles debido a la popularidad de las redes sociales que es donde se pueden hacer virales y logran alcanzar un mayor número de receptores. En este sentido, el uso de las noticias falsas al menos a nivel conceptual no es un fenómeno reciente y su popularidad tiene que ver más con los procesos de masificación de la información desarrollados con la llegada del internet. [1]

La relación noticias falsas-internet cobra relevancia en los últimos años a partir de cambios en los comportamientos informacionales y su incidencia en agendas mediáticas, ya que esto ocurre porque el internet y las redes sociales posibilitan que cualquier persona produzca contenido y se viabilice la interacción, configurándose como una especie de micrófono con el que antes no se contaba y brindando una ilusión de poder en el manejo y consumo de la información. [1]

Se plantea que esto ocurre en internet y en las redes sociales posibilitando que cualquier persona produzca contenido y se viabilice la interacción, configurándose como una especie de micrófono con el que antes no se contaba y brindando una ilusión de poder en el manejo y consumo de la información, logrando reconocer que cualquier persona puede crear contenido y difundirlo. [1]

## **2.2 Procesamiento de lenguaje natural**

El Procesamiento del Lenguaje Natural (PLN) es una subdisciplina de la Inteligencia Artificial y rama de la ingeniería Lingüística Computacional: ahora bien, la razón principal del PLN es construir sistemas y mecanismos que permitan la comunicación entre personas y máquinas por medio de lenguaje naturales, tal como que una computadora aprenda a interpretar el lenguaje natural debido a dos caminos, uno epistemológico y otro heurístico:

- 1) Epistemológico: Define el espacio de conceptos que el programa puede aprender.
  
- 2) El heurístico: Define los algoritmos para el aprendizaje.

El primer avance obtenido en el PLN se dio en el área del acceso a las bases de datos con el sistema lunar (1973) construido en la NASA por William Woods.

El PLN busca poder crear programas que puedan analizar, entender; generar lenguajes que los humanos utilizan habitualmente de manera que el usuario pueda llegar a comunicarse con la máquina o computador de la misma forma que lo haría con un ser humano. [1]

## **2.3 Gramática libre de contexto**

La gramática de tipo 2 o gramática independientes del contexto, son las que generan los lenguajes libres o independientes del contexto, los lenguajes libres del contexto son aquellos que pueden ser reconocidos por un autómata de pila determinístico o no determinístico. Según [2]

- Una Gramática libre de contexto (GLC), describe un lenguaje libre de contexto.
- Son útiles para describir bloques anidados en lenguajes de programación ya que describen su sintaxis.

## 2.4 Procesamiento de texto

El procesamiento de textos hace referencia a una etapa inicial que contempla la labor de estructurar el texto en un formato simplificado, eliminando caracteres no relevantes y palabras no representativas, así como signos de puntuación y otros elementos no esenciales para el análisis del texto. [3]

## 2.5 Campo semántico

Un campo semántico de acuerdo a [4] es un grupo de palabras que comparten uno o varios significados, cabe recalcar que para que varias palabras conformen el campo semántico deben de tener el mismo significado en común entre ellas.

A continuación, se presenta un ejemplo de campo semántico de palabras que comparten el mismo significado en común en la Figura 1.

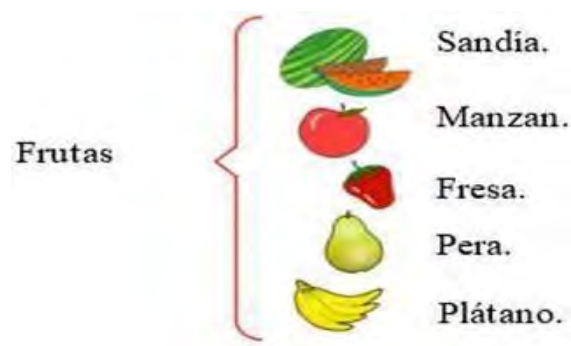


Figura 1 Ejemplo de campos semántico

# Capítulo III

Estado del arte

En esta sección se presentan las investigaciones relacionadas a este trabajo de tesis más relevantes en la literatura, en cada una de ellas se mencionan los objetivos y las herramientas que utilizaron para su desarrollo, teniendo un total de artículos recopilados de 23.

### 3. Estado del arte

#### Antecedentes

Algunas de las investigaciones que anteceden a la actual propuesta de tesis son:

#### 3.1 3HAN: A Deep Neural Network for Fake News Detection [4]

Habla sobre la investigación de donde se presenta un detector automatizado basado en Deep learning a través de una red neuronal de atención jerárquica de tres niveles (3HAN): palabras, oraciones y encabezados, para la detección precisa de noticias falsas. Su modelo se basa en la representación de un artículo de noticia como un vector de noticia, el que es utilizado para la clasificar un artículo asignando una probabilidad de ser falso, 3HAN proporciona una puntuación de importancia para cada palabra y oración de un artículo en función de su

Relevancia para retornar dicha probabilidad. La arquitectura de 3HAN se muestra en la Figura 2.

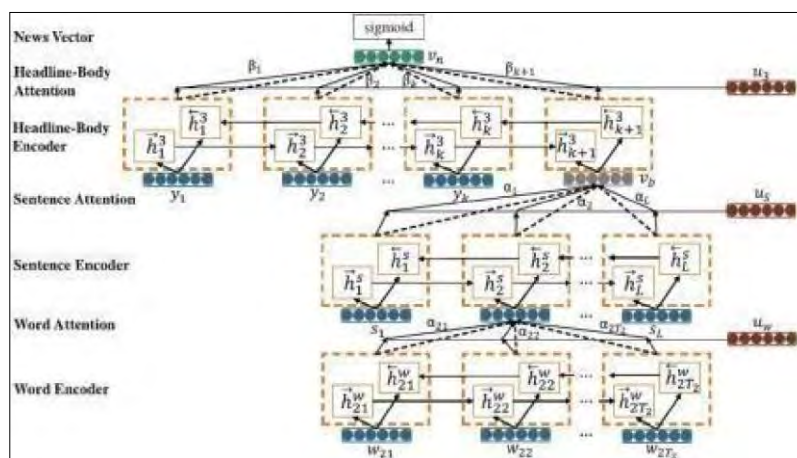


Figura 2 Arquitectura de la red neuronal utilizada para el modelo 3HAN



Los datos utilizados para el entrenamiento de este modelo son de carácter lingüístico a un nivel de palabras y oraciones, extraídos de la plataforma PolitiFact para el conjunto de noticias falsas y del listado de sitios verificados populares de EE. UU proporcionado por Forbes para el conjunto de noticias reales. Los conjuntos de noticias pertenecen al periodo de elecciones presidenciales en EE. UU del año 2016. 3HAN en base a sus experimentos se logra una excelente precisión del 96,77 % superando otros modelos basados en redes neuronales.

### **3.2 Fake News Early Detection: An Interdisciplinary Study [5]**

Este artículo hace mención a la difusión de noticias falsas que tiene como objetivo detectar noticias erróneas explorando cómo se propagan en las redes sociales, ya que este trabajo explora patrones potenciales de noticias falsas para mejorar la interpretabilidad en la ingeniería de características de noticias ficticias y estudia las relaciones entre las noticias, el engaño, la desinformación y los clickbaits.

En este artículo, se lleva a cabo un estudio interdisciplinario para la detección temprana de noticias falsas explicables para predecir noticias falsas antes de que comiencen a propagarse en las redes sociales, nuestro trabajo estudia y representa exhaustivamente el contenido de las noticias en cuatro niveles lingüísticos: nivel léxico, nivel de sintaxis, nivel semántico; nivel de discurso. Esta representación está inspirada en teorías bien establecidas de la psicología social y forense.

Los resultados indican que la identificación de artículos de noticias falsas, respectivamente, en función de la calidad, el sentimiento, la cantidad y la especificidad de su contenido funciona de manera similar, con una precisión del 60% al 70% y una puntuación F1 utilizando los datos de PolitiFact, y una precisión del 50% al 60% y una puntuación F1 utilizando los datos de BuzzFeed. Combinar todos los atributos para detectar noticias falsas funciona mejor que usar cada tipo de atributo por separado, lo que puede lograr una precisión del 70 % al 80 % y una puntuación F1 en los datos de Politifact, y una precisión del 60 % al 70 % y una puntuación F1 en los datos de BuzzFeed.

### **3.3 Mapeo sistemático sobre los métodos, técnicas y tecnologías orientadas a la detección de noticias [6]**

El objetivo de este estudio es investigar los principales métodos, técnicas y tecnologías orientadas a detectar noticias falsas en redes sociales, que emplean técnicas de IA, mediante la ejecución de un proceso de revisión sistemática de la literatura.

Para el logro de los objetivos planteados en la investigación se empleó el método inductivo que parte de lo particular a lo general, con el propósito de determinar cuál eran los principales métodos usados por otros investigadores para detectar noticias falsas en las redes sociales y en función de ello determinar las técnicas de IA más comunes para ello, ayudando a definir el aprendizaje automático, aprendizaje profundo y la minería de datos como las tres técnicas más usuales. Usando la IA.

De manera complementaria, se aplicó el método deductivo, que parte del general a lo particular, para determinar el alcance de las técnicas de IA que se han empleado en la detección de noticias falsas, y así determinar aquellos modelos que han tenido mayor precisión en la detección de noticias falsas.

También, el método ayudó a determinar los principales factores que determinan el buen o mal funcionamiento de clasificadores automáticos en la detección de noticias falsas. Además, se usó el método experimental puesto que en el estudio se llevó a cabo la manipulación de varias técnicas de aprendizaje en el cual se realizó pruebas modificando argumentos de los clasificadores lo cual ayudó en la ejecución de los métodos investigados y que, además, permitieron comparar los modelos estudiados. Finalmente, es importante señalar que también se utilizó el método descriptivo con el objetivo de describir las tres técnicas estudiadas para la detección de noticias falsas: aprendizaje automático, aprendizaje profundo, minería de datos y sus modelos y algoritmos de clasificación, algoritmos de aprendizaje supervisado y redes neuronales.

Se concluye que se está trabajando en técnicas de aprendizaje automático y profundo para crear modelos que permitan la detección de noticias falsas. Los algoritmos que mejor se desempeñaron fue uno de los pocos usados en los trabajos el cual fue *passive aggressive*, dicho algoritmo tuvo porcentajes muy alentadores de 100%. Otro algoritmo que cabe mencionar que es muy usado para la detección de noticias falsas es el algoritmo de máquina de vectores de soporte, por parte del aprendizaje profundo el que mejor se desempeña es un modelo de red neuronal con 99.90 % de exactitud. Con esto los clasificadores se puede utilizar para crear sistemas que permita detectar noticias falsas, además se han creados modelos para que los usuarios puedan probarlos haciendo uso de la URL y que puedan verificar si la noticia es verdad, pero integrarlo en las redes sociales todavía es un reto debido a que para entrenar los modelos solo usan datos extraídos de dos redes sociales o sitios web.

Con la investigación se obtuvo que si se aplican técnicas de minería de texto los modelos aplican eliminación de palabras debido a que no aportan información importante y lo que ocasiona es que el rendimiento del clasificador sea bajo. Entre estas técnicas de eliminación de palabras estuvieron las eliminaciones de palabras vacías, derivación y *Tokenización*.

Por otra parte, está la extracción de característica TF-IDF, que es la más utilizada, pero algo interesante fue que la herramienta LIWC 2015 puede dar buenos resultados ya que permite extraer diferentes características textuales de los artículos de noticias, esta permitió obtener los mayores porcentajes de exactitud en diferentes algoritmos por lo que usarlo puede beneficiar a ser más precisos en la detección de noticias falsas.

### **3.4 Detection of fake news using deep learning CNN–RNN based methods [7]**

Este estudio analiza exhaustivamente el rendimiento de varios métodos de aprendizaje profundo combinados con la incorporación de palabras de última generación en conjuntos de datos de referencia de noticias falsas, utilizando incrustaciones de palabras previamente entrenadas, como *Word2Vec*, *Glove* y *fastText*, estas incrustaciones de palabras previamente entrenadas son incrustaciones de palabras pre entrenadas populares y se eligieron porque se entrenaron utilizando un corpus masivo para producir suficiente vocabulario, cada incrustación de palabras previamente entrenada se combina con los respectivos métodos de aprendizaje profundo, a saber, CNN, LSTM bidireccional y ResNet, para determinar su rendimiento en la detección de noticias falsas. Como resultados se identificó que el modelo bidireccional LSTM + *fastText* nuevamente supera a otros modelos con un rendimiento de 99,24 % de precisión, 99,19 % de precisión, 99,26 % de recuperación y 99,23 % de puntuación F1, Los resultados de la prueba también muestran que ResNet + GloVe ocupa el segundo lugar con un rendimiento de 98,99% de precisión, 99,05% de precisión, 98,89% de recuperación y 98,97% de puntuación F1. CNN + GloVe con un resultado con un valor de precisión del 98,24%, 98,32% de precisión, 98,09% de recuperación y una puntuación F1 del 98,2%.

### **3.5 La Noticias falsas y desinformación sobre el Covid-19: análisis comparativo de seis países iberoamericanos [8]**

El objetivo principal de este estudio habla sobre las noticias falsas sobre Covid-19 que circularon en 6 países iberoamericanos, muestran un análisis comparativo en torno a los temas recurrentes, la relación entre el tipo de noticias falsas y la técnica de engaño que emplean los desinformadores y la identificación de sus posibles internacionalidades.

Se implementó esta investigación la presente metodología fundamentada en aprendizaje estadístico que detecta de manera precisa y automática noticias falsas basándose en la forma en que están escritas, por lo que se utilizó procesamiento de lenguaje natural, algoritmos clásicos de clasificación supervisada y el lenguaje de programación Python.

El estudio implementado demostró la capacidad de los algoritmos de clasificación supervisada para lograr identificar con más del 90% de precisión noticias falsas que circulan diariamente en internet, especialmente en redes sociales por parte de páginas de sátira política enfocada en problemas ecuatorianos.

### **3.6 Detection of fake news in a new corpus for the Spanish language [9]**

En este artículo se presenta un modelo para analizar y detectar información engañosa, en gran cantidad de sitios web en idioma español. Para la creación del modelo se usó un conjunto de datos de noticias recopiladas manualmente de distintos sitios web para crear un nuevo corpus de noticias etiquetadas como noticias falsas y reales con el objetivo de su detección automática, también se proporciona relacionando a la temática de la noticia: ciencia, deporte, economía, educación, entretenimiento, política, salud seguridad y sociedad.

El entrenamiento se realizó siguiendo el proceso para un modelo de aprendizaje supervisado, particionado el conjunto de datos en entrenamiento y test, 70% y 30% respectivamente, y utilizando los siguientes algoritmos de clasificación: máquina de soporte vectorial, regresión logística, bosques aleatorios y Boosting, evaluando el desempeño para dos escenarios: removiendo palabras repetitivas y considerando palabras repetitivas. La extracción de información y su representación se realizó a través de características lingüísticas obtenidas de tres técnicas: bag-of-words, n-grams y POS tags n-grams. En base a los experimentos realizados los algoritmos llegan a niveles de precisión cercanos a 75%.

### 3.7 Computational Fast Checking from Knowledge Networks [10]

Se presenta en este artículo de investigación un modelo basado en redes o grafos en el que la verificación de hechos realizadas en el proceso de fast-checking manual, puede aproximarse bastante bien al encontrar el camino más corto entre nodos de conceptos bajo métricas de proximidad semántica adecuadamente definidas sobre grafos de conocimiento, este enfoque es visible con técnicas computacionales eficientes identificado en la Figura 3.

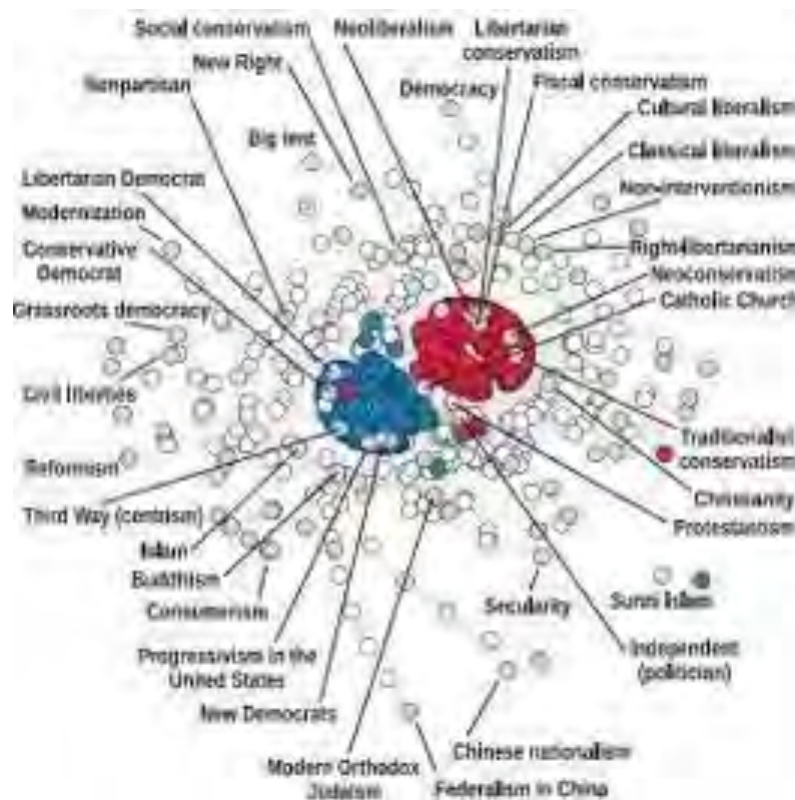


Figura 3 Los nodos rojos y azules corresponden a miembros del congreso, los nodos grises a ideologías y nodos blancos a vértices de cualquier otro tipo

En este modelo un grafo de conocimiento se produce de la unión de los nodos que denotan afirmaciones sujetos u objetos de declaraciones y las aristas que denotan conectores lingüísticos. Para verificar si una afirmación es cierta, se espera que exista una arista del grafo de conocimiento que conecte con la entidad, o si existe el camino más corto entre entidades del grafo de conocimiento. Por otro lado, si la afirmación es falsa, no debe de haber aristas ni caminos cortos que conecten a dicha afirmación.

En la creación del modelo los datos utilizados fueron extraídos de Wikipedia en los que constan todas las declaraciones fácticas extraídas de las cajas de información de Wikipedia, creándose así un grafo de conocimiento con 3 millones de entidades unidas por 23 millones de aristas aproximadamente.

En base a los experimentos realizados. El mejor resultado obtenido es para una representación con grafos no dirigidos evaluando las clasificaciones mediante dos algoritmos de clasificación: vecinos más cercanos y árboles aleatorios, llegando a niveles de precisión que sobrepasan el 95% aproximadamente.



### **3.8 IDENTIFICACIÓN AUTOMÁTICA DE NOTICIAS FALSAS EN ESPAÑOL UTILIZANDO TÉCNICAS DE MINERÍA DE DATOS Y PROCESAMIENTO DE LENGUAJE NATURAL [11]**

El objetivo principal de este tema de investigación es crear un modelo que permita identificar de manera precisa y automática noticias falsas en español de páginas de noticias y sátiras ecuatorianas. Reconociendo la circulación de noticias falsas, en español las de sátira política, por medios de comunicación digitales ha afectado a la mayoría de la población ecuatoriana por su masiva propagación y la falta de regulación, por lo que en este trabajo se presenta una metodología fundamentada en aprendizaje estadístico que detecta de manera precisa y automática noticias falsas basándose en la forma en que están escritas.

El estudio realizado muestra la capacidad de los algoritmos de clasificación supervisada para lograr identificar con más del 90% de precisión noticias falsas que circulan diariamente en internet, especialmente en redes sociales por parte de páginas de sátira política enfocadas en problemáticas ecuatorianas, que, si bien su objetivo no es malicioso si no humorístico, pero la información que se crea y es compartida puede ser mal interpretada fuera de contexto desinformando a las personas.

### **3.9 Detección automática de noticias falsas usando representaciones textuales tradicionales y soluciones basadas en aprendizaje profundo [23]**

En este trabajo de tesis se estudia el impacto del uso de diversas representaciones de características del contenido de la noticia para la detección de noticias falsas en español con técnicas de aprendizaje automático, incluyendo arquitecturas profundas. Métodos basados en texto con técnicas de aprendizaje profundo: Las técnicas de aprendizaje profundo también han sido utilizadas para detectar automáticamente noticias falsas desde arquitecturas simples hasta las complejas que fusionan diferentes modelos. Entre las más usuales se encuentran las Redes Neuronales Convolucionales (CNN) y Recurrentes (RNN). También ha sobresalido el uso de redes recurrentes especiales, como son GRU (Gated Recurrent Unit) y LSTM (Long Short-Term Memory). Recientemente, el uso de BERT (Bidirectional Encoder Representations from Transformers) ha sido frecuentemente probado en este dominio.

Este trabajo de tesis es sumamente un competidor directo con respecto a mi tema, ya que aborda la identificación de noticias falsas en español basados en el contenido de la noticia. La metodología combina técnicas de aprendizaje.

En este trabajo se estudian diversas representaciones textuales para la detección automática de noticias falsas en español. Los modelos presentados en este trabajo solo consideran el contenido de la noticia sin usar metadatos. Los resultados muestran que el desempeño de los modelos, depende en gran medida de las características usadas en la representación del texto, así como de los algoritmos de aprendizaje automático aplicados.

### **3.10 Inteligencia Artificial como alternativa en la detección de noticias falsas [24]**

El objetivo del presente artículo es analizar cómo la Inteligencia Artificial ha servido como herramienta para la detección de noticias falsas. Se trata de un artículo de revisión bibliográfica, descriptivo, exploratorio cuya unidad de análisis está representada por artículos publicados en bases de datos tales como: Scielo, Dialnet, Science Direct y Scopus. Se obtuvo como resultado que existen distintos mecanismos de Inteligencia Artificial en el subcampo del Aprendizaje Automático, el Aprendizaje Profundo y Procesamiento de Lenguaje Natural (NLP), como el uso de máquinas de soporte vectorial, el clasificador Naive Bayes y distintos tipos de modelos de Redes de Neuronas Artificiales; como las Redes Neuronales Recurrentes (RNN), las Redes Neuronales Convolucionales (CNN), que son bastante populares en el Procesamiento de Lenguaje Natural, especialmente en el uso de Redes de Memoria a Corto y Largo Plazo (LSTM), las cuales permiten establecer mecanismos confiables y precisos para detectar patrones dentro de un objetivo con contenido textual. El artículo desarrollado es de tipo revisión donde la unidad de análisis está representada por artículos publicados en bases de datos tales como: Scielo, Dialnet, Science Direct, latindex, Scopus, los cuales abordan la temática de la detección de noticias falsas utilizando la Inteligencia Artificial. Se utilizó una metodología descriptiva-exploratoria, para identificar cómo se detectan las noticias falsas utilizando Inteligencia Artificial, específicamente en el campo del Aprendizaje Automático.

### **3.11 Enmascaramiento de la Información para la Detección Automática de Noticias Falsas [25]**

Como Objetivo General de este estudio se planteó: Detectar noticias falsas siguiendo un enfoque basado en estilo, mediante técnicas de enmascaramiento para el idioma inglés y español en notas periodísticas con resultados semejantes al estado del arte. La idea principal con esta tesis es proponer un método para la detección de noticias falsas mediante el enmascarado de los textos de entrada, de tal manera que se mantenga la estructura textual, relacionada con el estilo de las noticias, mientras se enmascaran las ocurrencias de las palabras menos frecuentes, correspondientes a información temática. Para efectos de evaluación y comparación también se considera el escenario inverso, se enmascaran las ocurrencias de las palabras más frecuentes y se conservan las palabras con información temática, relacionadas al contenido de la noticia.

### **3.12 A New Application of Social Impact in Social Media for Overcoming Fake News in Health [26]**

El Objetivo del presente trabajo es analizar los retos que hay en la actualidad que es afrontar las fake news (información falsa) en salud por su potencial impacto en la vida de las personas. Este artículo contribuye a una nueva aplicación de la metodología del impacto social en las redes sociales (SISM). Este estudio se centra en el impacto social de la investigación para identificar qué tipo de información de salud es falsa y qué tipo de información es evidencia del impacto social compartido en las redes sociales. El análisis de las redes sociales incluye Reddit, Facebook y Twitter.

En este artículo se concluyó que el impacto social en las redes sociales (SISM) constituye una metodología novedosa tanto en el análisis de las redes sociales como en la evaluación del impacto social de la investigación. Este artículo aplica la metodología SISM al caso específico de las fake news en salud para identificar el tipo de interacciones relacionadas con la información compartida en las redes sociales.

El método utilizado es la metodología de impacto social en redes sociales (SISM), que combina análisis de contenido cuantitativo y cualitativo de la muestra seleccionada considerando los aportes del impacto social de la investigación.

Este artículo demuestra que SISM es una metodología replicable que se ha aplicado con éxito en el análisis de redes sociales en relación con la salud y las noticias falsas, contribuyendo a una mayor exploración de las posibilidades de esta metodología. Este estudio ofrece la posibilidad de identificar, por un lado, evidencias de impacto social compartidas en redes sociales y, por otro lado, desinformación o información falsa relacionada con la salud.

### **3.13 COVID-19-related Fake News in Social Media [27]**

El propósito de este artículo, se centra en las noticias falsas de las redes sociales indias relacionadas con COVID-19, analizando datos digitales y analógicos para responder a las consultas de investigación el cual analiza 120 noticias destacadas relacionadas con la actual pandemia difundida en redes sociales del 29 de enero al 11 de abril de 2020. El cual habla de los 5 parámetros del análisis que son temas, tipo de contenido, fuentes, cobertura e intenciones, ya que los temas principales de noticias falsas son la salud, religión, política, crimen y entretenimiento.

En este estudio se realizó un análisis de contenido cuantitativo, en el cual se implementó la codificación deductiva de los datos recopilados. Dos codificadores capacitados codificaron los datos recopilados. Busca explorar la naturaleza dinámica de las noticias falsas en las redes sociales relacionadas con COVID-19. Se preparó una lista de noticias falsas en las redes sociales del 29 de enero al 11 de abril de 2020.

Después de filtrar, relacionado con COVID-19. Se recopilaron 125 noticias falsas y se recopilaron datos e información relevantes relacionados con cada noticia logrando descubrir que las noticias falsas no son consistentes a lo largo del tiempo, más bien fluctúan, y se nota una disminución lenta y sutil en los casos diarios.

### 3.2 Trabajos representativos

A continuación, en la Tabla 1 se presentan algunos de los trabajos más representativos encontrados en la literatura, los cuales fueron seleccionados con el propósito de identificar las actividades desarrolladas en el análisis de detección de noticias falsas en temas de salud.

#### Tabla comparativa de trabajos relacionados con la detección de Noticias Falsas

Se realizó una comparativa de las herramientas y resultados obtenidos en la literatura revisada para el desarrollo de este trabajo de investigación.

Tabla 1 Resumen del estado del arte

Título	Objetivo	Metodología	Resultado / Conclusión	Año
Identificación automática de noticias falsas en español utilizando técnicas de minería de datos y procesamiento de lenguaje natural. [11]	En este estudio se creará un modelo que permita identificar de manera precisa y automática noticias falsas en español de páginas de noticias y sátiras ecuatorianas.	Se utilizó aprendizaje estadístico que detecte de manera precisa y automática noticias falsas basándose en la forma en que están escritas la noticia. Utilizando procesamiento de lenguaje natural, algoritmos clásicos de clasificación supervisada.	El trabajo realizado muestra la capacidad de los algoritmos de clasificación supervisada para lograr identificar con más del 90 % de precisión para la detección de noticias falsas que circulan diariamente en internet. Se identificó que el modelo entrenado presenta una buena capacidad para identificar una noticia falsa de una real, el cual tiene un área bajo la curva ROC de 97.27%.	2021
Fake news early detection: an interdisciplinary study. [5]	Este tema de investigación estudia cómo Detectar noticias erróneas explorando cómo se propagan en las redes sociales.	El método investiga el contenido de las noticias en varios 4 niveles: nivel de léxico, nivel de sintaxis, nivel semántico y nivel de discurso, que detecta noticias falsas	Los resultados indican que la identificación de artículos de noticias falsas respectivamente en función de la calidad y la especificidad de su contenido funciona de manera similar con una precisión del 60% al 70% y F-1 Score con una puntuación	2018

		que se llevan a cabo dentro de un marco de aprendizaje automático supervisado.	utilizando datos PolitiFact de 50% a 60% de precisión	
Fake news detection using machine learning approaches. [12]	Generar análisis de la investigación relacionada con la detección de noticias falsas y explorar los modelos tradicionales de aprendizaje automático para elegir el mejor con el fin de crear un modelo de producto con algoritmo de aprendizaje automático supervisado que puede clasificar noticias falsas como verdaderas o falsas mediante el uso de herramientas como Python, PNL para analizar texto.	Este artículo propone una metodología para crear un modelo que detecte si un artículo es auténtico o falso en función de sus palabras, frases, fuentes y títulos, mediante la aplicación de algoritmos de aprendizaje automático supervisado en un conjunto de datos anotados (etiquetados), que se clasifican y garantizan manualmente.	Se obtuvo como resultado utilizando el clasificador Naive Bayes para detectar noticias falsas de diferentes fuentes con un resultado de precisión del 74%. Se utilizó un algoritmo ML que depende de un umbral de probabilidad poco confiable con 58-91%, los datos de Kabble con una precisión de 74.5%, Naive Bayes con una precisión de 70% a 71,2%	2021
3HAN: A Deep Neural Network for Fake News Detection. [4]	Se presenta una red de atención jerárquica de tres niveles, que crea una representación efectiva de un artículo noticioso llamado vector de noticias	Se empleó un detector automático basado en aprendizaje profundo a través de una red de atención jerárquica de tres niveles (3HAN) para una detección rápida y precisa de noticias falsas	Se demostró que la jerárquica de 3 niveles 3HAN otorgó una importante diferencia a las partes de un artículo debido a sus tres niveles de atención. Mediante los experimentos en un gran conjunto de datos del mundo real, se observó la efectividad de 3HAN con una precisión de 97.77%.	2017
Detección de noticias falsas en redes sociales basadas en aprendizaje automático y profundo: una breve revisión sistemática. [13]	Se busca desarrollar una revisión de la literatura en la que se determine como el aprendizaje automático y aprendizaje profundo, han aprobado para el desarrollo de clasificadores de noticias falsas en redes sociales	Se implementó una revisión sistemática de la literatura que consiste en identificar, evaluar e implementar los estudios más relevantes de detección de noticias falsas.	Los resultados mostraron que los modelos de aprendizaje de IA han sido ampliamente empleados para la creación de sistemas de detección automática de noticias falsas. Con alto y bajos porcentajes de exactitud, con una exactitud (accuracy) muestran que el modelo de red neuronal que mejor se adaptó fue una red neuronal genérica de la cual no se dio detalles con 99,90% (S33),	2020

			luego, las redes convolucionales con el 96,00% (S13) y finalmente, el modelo de memoria a largo plazo con el 95,30% (S9).	
Noticias falsas y desinformación sobre el Covid-19: análisis comparativo de seis países iberoamericanos. [8]	hablar sobre las noticias falsas sobre Covid-19 que circularon en 6 países iberoamericanos, mostrando un análisis comparativo en torno a los temas recurrentes, e identificar la relación entre el tipo de noticias falsas y la técnica de engaño que emplean los Desinformadores y la identificación de sus posibles intencionalidades.	Se realizó un análisis de contenido de 371 noticias falsas, previamente verificadas por fact-checkers. Se procedió a clasificar los bulos según su tipo de internacionalidad, tema principal abordado, las redes en que circularon, las técnicas de engaño, el país de origen, su carácter transnacional.	Se obtuvo que en primer lugar se pudo evidenciar que Facebook fue la red más empleada para divulgar noticias falsas, ya que se usó en el 32.9 % de los 371 analizados. El segundo puesto fue para los bulos que se difundieron en dos o más redes los cuales corresponden al 31.9%.	2020
La red sanitaria y su participación en la difusión o contención de las fake news y bulos relacionados con COVID-19: el caso de Lima, Perú. [14]	Identificar qué tan frecuente los profesionales de la salud recibieron noticias falsas y bulos a través de las redes sociales	Este es un estudio de enfoque cuantitativo, con un diseño no experimental transversal. Se estudian las participaciones de la red sanitaria de Lima-Perú en la producción, consumo, difusión, así como el aprendizaje sobre las <i>fake news</i> y los bulos.	Durante la pandemia por el COVID 19, en Lima, el 90 % de los 60 profesionales de la salud que fueron encuestados recibió noticias falsas y bulos a través de las redes sociales. Las redes sociales por las cuales recibieron estos contenidos fueron WhatsApp 53,7 % y Facebook 40,7 %.	2020
El impacto de las fake news en la investigación en ciencias sociales revisión bibliográfica sistematizada. [15]	Es desarrollar una revisión sistematizada de los estudios realizados hasta la fecha sobre las noticias falsas, partiendo de las principales bases de datos (Web of Science, Dilanet ) y con el fin de verificar qué temas han traído la atención	Se realizará una revisión bibliográfica sistematizada descriptiva sobre noticias falsas y ofrecer un panorama del estado del arte a partir de la producción científica nacional e internacional en los últimos 20 años	Los resultados obtenidos permiten validar la hipótesis de la investigación que señalaba que la principal producción científica sobre noticias falsas procede del Ámbito científico. En la misma línea, el estudio ha permitido verificar que este conjunto de trabajos se centra en el ámbito de la Comunicación y que su eclosión se ha producido en los dos últimos años a raíz,	2019



	de la comunidad científica al respecto, qué técnicas de investigación se han utilizado en estos estudios, donde se concentran geográficamente y qué campos quedan aún por cubrir en torno a este fenómeno		especialmente, de los acontecimientos políticos que tuvieron lugar en el año 2016.	
Approaches to Identify Fake News: A Systematic Literature Review [28]	En este artículo habla de la importancia de combatir las noticias falsas, porque queda claramente ilustrada que durante la actual pandemia COVID-19. Las redes sociales están intensificando el uso de herramientas digitales de detección de noticias falsas y la educación del público para detectarlas.	El propósito de este documento es categorizar los enfoques utilizados para identificar noticias falsas. Para ello, se realizó una revisión sistemática de la literatura para el reconocimiento de las noticias falsas.	En este artículo se analiza la prevalencia de las noticias falsas, y como ha cambiado la tecnología en los últimos años permitiendo conocer nuevas herramientas que se puedan utilizar en la lucha contra las noticias falsas. Contra la batalla actual llamada Covid-19, identificando los enfoques principales para reconocer noticias falsas.	2020
A New Application of Social Impact in Social Media for Overcoming Fake News in Health [29]	Uno de los retos de este artículo es afrontar la información falsa en la salud por su potencial impacto en la vida de las personas. De igual manera se centra en el impacto social de la investigación para identificar qué tipo de información de salud es falsa y que tipo de información es evidencia del impacto social compartida en las redes sociales.	El método utilizado es la metodología de impacto social en redes sociales (SISM) que combina análisis de contenido cuantitativo de la muestra seleccionada considerando los aportes del impacto social de la investigación.	Este artículo demuestra que SISM es una metodología replicable que se ha aplicado con éxito en el análisis de redes sociales en relación con la salud y las noticias falsas. Este estudio ofrece la posibilidad de identificar, por un lado, evidencias de impacto social compartidas en redes sociales y por otro lado desinformación o información falsa.	2020
COVID-19-related Fake News in Social Media [27]	Este artículo se centra en las noticias falsas de las redes sociales indias relacionadas con COVID-19, el cual analiza 120 noticias destacadas relacionadas con la actual	En este estudio se realizó un análisis de contenido cuantitativo, en el cual se implementó la codificación deductiva de los datos recopilados. Dos	Este estudio busca explorar la naturaleza dinámica de las noticias falsas en las redes sociales relacionadas con COVID-19. En el cual se realizó una lista de noticias falsas en las redes sociales del 29 de enero al 11 de abril del 2020 se recopilaron 125 noticias	2020

	<p>pandemia difundida en redes sociales del 29 de enero al 11 de abril de 2020. El cual habla de los 5 parámetros del análisis que son temas, tipo de contenido, fuentes, cobertura e intenciones, ya que los temas principales de noticias falsas son la salud, religión, política, crimen y entretenimiento.</p>	<p>codificadores capacitados codificaron los datos recopilados.</p>	<p>falsas y se recopilaron datos de información relevantes relacionadas con cada noticia.</p>	
<p>CSI: A Hybrid Deep Model for Fake News Detection [30]</p>	<p>Este estudio hace mención a las noticias falsas en las redes sociales, ya que la detección automática de noticias erróneas es un problema importante, pero desafiante, que aún no se comprende en su totalidad. Ya que este trabajo se ha centrado en gran medida en adaptar las soluciones a una característica particular que ha ilimitado su éxito y generalidad.</p>	<p>Este trabajo, propone un modelo que combina las tres características para una predicción más precisa y automatizada. En el cual se propuso un modelo llamado CSI que se compone de 3 módulos: Captura, Puntuación e integrar. El primer módulo se basa en la respuesta y texto utilizando una red neuronal para capturar el patrón temporal de la actividad del usuario, el segundo módulo aprende la fuente la característica basada en el comportamiento del usuario y los 2 se integran con el 3 para clasificar un artículo como falso o verdadero.</p>	<p>En este trabajo se estudió el problema oportuno de la detección de noticias falsas. Enfocándose en el texto, la respuesta que recibe un artículo, o los usuarios que reciben la noticia.</p>	<p>2017</p>
<p>Neural User Response Generator: Fake News Detection with Collective User Intelligence [31]</p>	<p>Estudia el problema de la detección temprana de noticias falsas bajo el supuesto de que el texto del artículo es la única</p>	<p>En este estudio se aborda la arquitectura del modelo de detección temprana de noticias falsas propuesto TCNN-URG se compone de dos partes: TCNN representa</p>	<p>Los trabajos existentes no se pueden aplicar al problema de la detección temprana de noticias falsas porque la mayoría de ellos se basan principalmente en la respuesta del usuario que no está</p>	<p>2018</p>

	información disponible en el momento de la detección	cada artículo en dos niveles y es capaz de aplicar una clasificación de texto puro en artículos de noticias y URG está capacitado para aprender cómo los usuarios responden a los artículos de noticias y puede generar respuestas de usuario para ayudar a TCNN con la sabiduría del usuario cuando la respuesta del usuario no está disponible.	disponible para la detección temprana de noticias falsas.	
Facing Fake News: the case of the students of the University of the Basque Country [32]	La presente investigación fue desarrollada con el objetivo de investigar los métodos, técnicas y herramientas para la detección de noticia falsas que emplean técnicas de Inteligencia artificial (IA). Para llevar a cabo esta investigación se realizó un mapeo sistemático de literatura considerando criterios de inclusión y criterios de exclusión,	Se empleó el método inductivo, que parte de lo particular a lo general, con el propósito de determinar cuál eran los principales métodos usados por otros investigadores para detectar noticias falsas en las redes sociales y en función de ello determinar las técnicas de IA más comunes, para ello ayudó a definir el aprendizaje automático, aprendizaje profundo y la minería de datos como las tres técnicas más usuales. Usando la IA. De manera complementaria, se aplicó el método deductivo, que parte de lo general a lo particular, para determinar el alcance de las técnicas de IA que se han empleado en la detección de noticias falsas, y así determinar aquellos modelos que han tenido	Se concluye que se está trabajando en técnicas de aprendizaje automático y profundo para crear modelos que permitan la detección de noticias falsas.	2021

		mayor precisión en la detección de noticias falsas.		
Procesamiento de Lenguaje Natural: una solución para detectar noticias falsas sobre la 4T en México [33]	El objetivo de este trabajo es analizar las noticias verdaderas/falsas de la 4T en Twitter	Para el análisis de los datos se determinaron características divididas en características de usuario y del tweet. Se utilizó la librería NLKT que proporciona más de 50 recursos corporales y léxicos y ofrece un conjunto de bibliotecas de procesamiento de texto para clasificación, tokenización, derivación, etiquetado, análisis y razonamiento semántico para bibliotecas PLN	En particular, Twitter es un medio que permite escribir y compartir información sin filtros, por tanto, es casi inmediata la publicación del tweet. El usuario en su entusiasmo por escribir en Twitter, de manera incluso irreflexiva, se expresa rápidamente, con descuidos frecuentes sobre la fuente de la información que comparte y, en consecuencia, el usuario asume que la información que publica es verdadera, otorgándole al tweet una ingenua credibilidad, de tal manera se implementó un análisis de datos para analizar las noticias falsas y verdaderas.	2020
MAPEO SISTEMÁTICO SOBRE LOS MÉTODOS, TÉCNICAS Y TECNOLOGÍAS ORIENTADAS A LA DETECCIÓN DE NOTICIAS FALSAS [34]	La presente investigación fue desarrollada con el objetivo de investigar los métodos, técnicas y herramientas para la detección de noticia falsas que emplean técnicas de Inteligencia artificial (IA). Para llevar a cabo esta investigación se realizó un mapeo sistemático de literatura considerando criterios de inclusión y criterios de exclusión,	Se empleó el método inductivo, que parte de lo particular a lo general, con el propósito de determinar cuál eran los principales métodos usados por otros investigadores para detectar noticias falsas en las redes sociales y en función de ello determinar las técnicas de IA más comunes, para ello ayudó a definir el aprendizaje automático, aprendizaje profundo y la minería de datos como las tres técnicas más usuales. Usando la IA. De manera complementaria, se aplicó el método deductivo, que parte de lo general a lo particular, para determinar el alcance de las técnicas de IA que se han empleado en la	Se concluye que se está trabajando en técnicas de aprendizaje automático y profundo para crear modelos que permitan la detección de noticias falsas.	2021

		detección de noticias falsas, y así determinar aquellos modelos que han tenido mayor precisión en la detección de noticias falsas.		
--	--	--	--	--

# Capítulo IV

Método de solución

## 4. Método de Solución

### 4.1 Propuesta de Solución

Para el desarrollo de este trabajo de investigación se implementó una metodología de solución compuesta por 5 fases: Fase 1: Creación de Gramática, Fase 2: Procesamiento de Noticia, Fase 3: Recuperación de Información, Fase 4: Identificación de Contradicciones y Fase 5: Diseño de Interfaz. Dichas fases se muestran en la siguiente Figura 4.

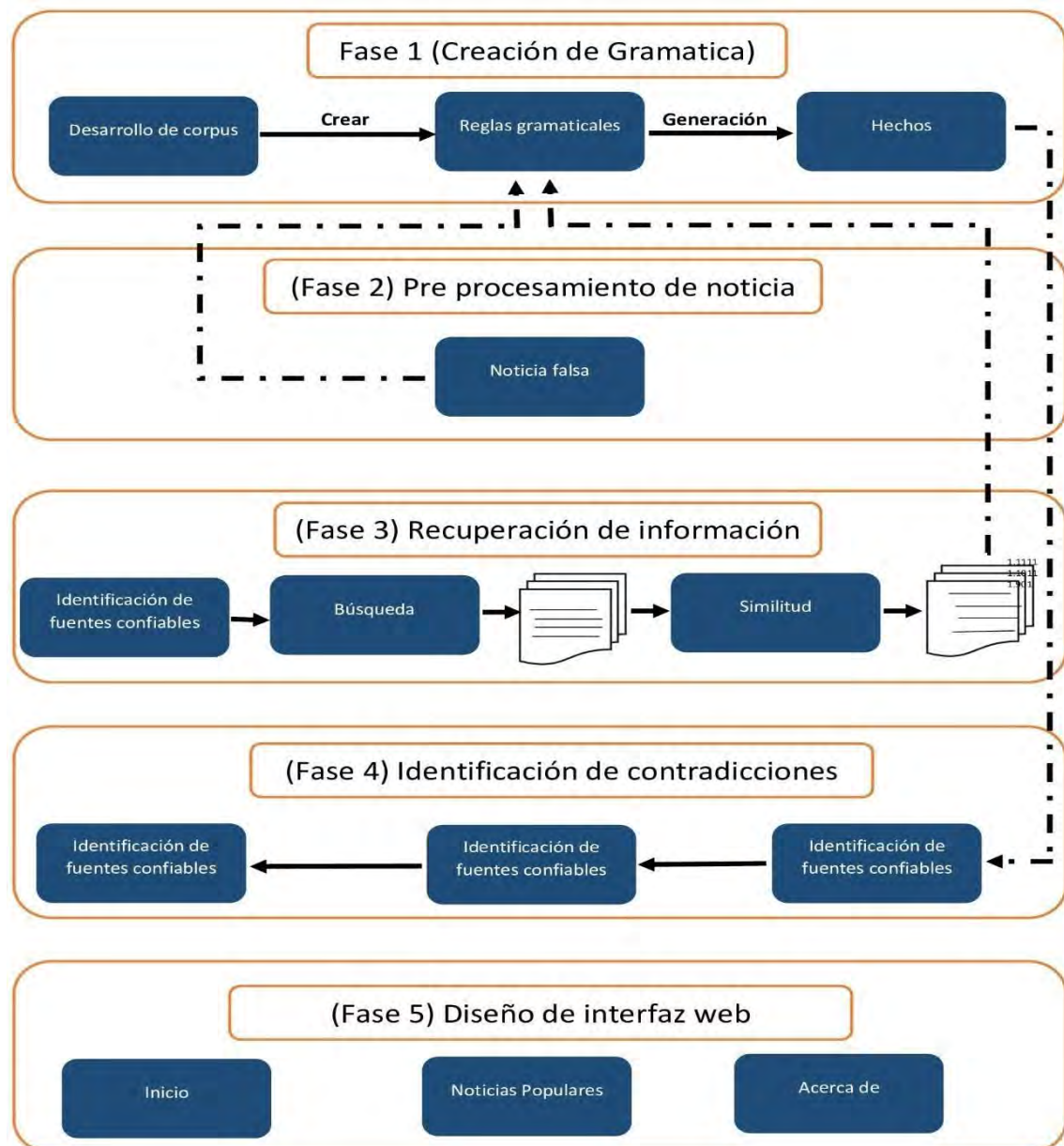
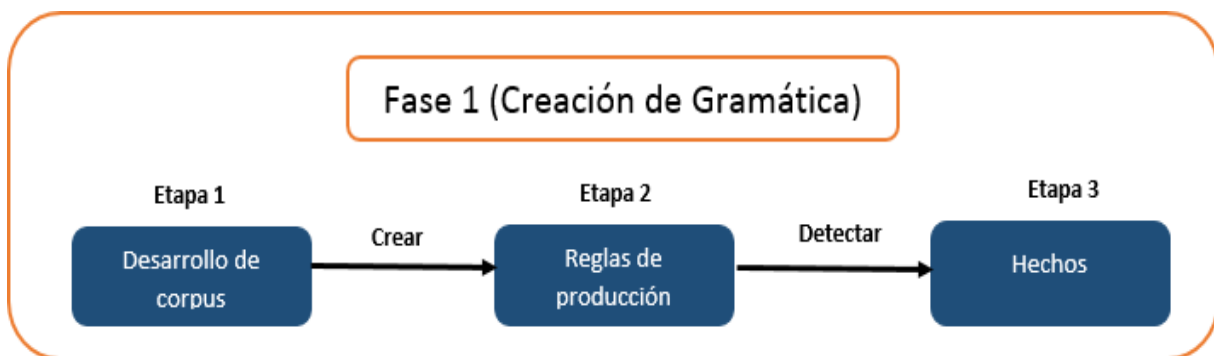


Figura 4 Metodología de solución

## FASE 1: Creación de Gramática

La primera fase de la metodología está conformada por tres etapas: Etapa 1: Identificar y recolectar noticias falsas de temas de la salud, Etapa 2: Crear las reglas gramaticales que existan en las noticias recolectadas y almacenadas en el Corpus, Etapa 3: Identificar los hechos tanto en la noticia con la que se genera la búsqueda tanto en las notas recolectadas a través de la búsqueda.

A continuación, se muestra la Fase 1 de la metodología de solución también llamada Creación de Gramática dentro de la Figura 5.



*Figura 5 Fase 1 de la metodología de solución*

### Etapa 1: Desarrollo de Corpus

Se generó un corpus de noticias falsas de temas sobre la salud extraídas de diferentes redes sociales y fuentes de información falsas, el corpus consta de un total de 300 noticias las noticias falsas fueron recolectadas de las publicaciones que hicieron diferentes personas y bots.

### Extracción y recolección de los textos

Se procedió a realizar una búsqueda manual en las diferentes redes sociales y fuentes de información de noticias falsas de temas sobre la salud.

En la Figura 6 se presenta el ejemplo de una noticia falsa intentando hacerla pasar como que está haciendo referencia a la noticia fuente original, aunque este texto que se muestra es parte de la noticia original carece de datos de la estructura de la noticia tales como:



Título, cuerpo de noticia, así como su fuente de información, el cual se muestra un tipo de noticia que no presenta el formato como tal pues carece de título, y por lo tanto ese tipo de publicaciones no se iban a procesar.

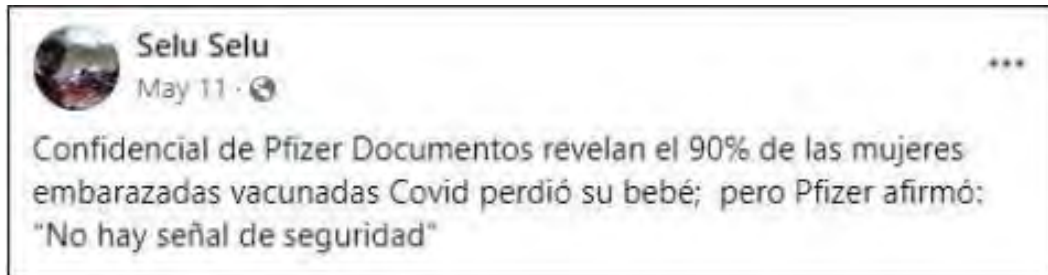


Figura 6 Ejemplo de noticia que hace referencia a la noticia fuente

A continuación, se presenta en la Figura 7 una publicación que presenta el formato de noticia y, por lo tanto, se considera para este trabajo de investigación.



## Beber agua caliente, sola o con limón, no cura el cáncer

Un mensaje en redes sociales compartido a través de varias publicaciones más de un millón de veces desde 2013 afirma que el agua caliente, sola o con limón, ayuda a curar el cáncer, entre otras enfermedades. Otras entradas dicen que el agua con limón es 10.000 veces más poderosa que la quimioterapia. Sin embargo, expertos en salud definen esta teoría como "absurda" y alertan sobre sus riesgos para los pacientes.

Figura 7 Ejemplo de noticias que no hacen referencia a la noticia fuente

## Corpus generado

El Corpus desarrollado consta de un total de 300 noticias falsas de temas sobre la salud, las cuales hacen referencia a varias enfermedades letales, tales como Cáncer, Sida, Covid-19 y demás. En la Figura 8 se muestra un ejemplo del Corpus generado.

2	La Organización Mundial de la Salud ha avalado un estudio que dice que vacunarse <b>contra</b> el COVID-19 sea un 339% más riesgoso que pasar la enfermedad	"un estudio de la Organización Mundial de la Salud (OMS) concluye que el riesgo de <b>sufrir</b> daños serios como consecuencia de la vacuna para el COVID-19 es 339% más alto que el de ser hospitalizado por COVID-19".	<a href="https://factchequeado.com/verificaciones/20220713/oms-no-avala-estudio-rse-339-riesgoso/">https://factchequeado.com/verificaciones/20220713/oms-no-avala-estudio-rse-339-riesgoso/</a>
3	Vacuna contra COVID-19 reduce el conteo de espermatozoides de los hombres	la vacuna de Pfizer reduce el conteo de espermatozoides en hombres	<a href="https://factchequeado.com/verificaciones/20220713/oms-no-avala-estudio-rse-339-riesgoso/">https://factchequeado.com/verificaciones/20220713/oms-no-avala-estudio-rse-339-riesgoso/</a>
4	Viruela del mono es consecuencia de la vacunación anticovid	La viruela de Mono? No!!! El principal efecto de un veneno causado "(vacuna) Covid"	<a href="https://archive.ph/Wz33f">https://archive.ph/Wz33f</a>
5	La vacuna contra el coronavirus también sirve para combatir la viruela	La vacuna contra el coronavirus también sirve para combatir la viruela del mono.	<a href="https://www.europapress.es/verificaciones/noticia-sabe-supuesta-doble-fijidad-vacuna-covid-19-podria-combatir-viruela-mono-20220608111935.htm">https://www.europapress.es/verificaciones/noticia-sabe-supuesta-doble-fijidad-vacuna-covid-19-podria-combatir-viruela-mono-20220608111935.htm</a>
	Los documentos de Pfizer revelan que el 90% de las mujeres embarazadas vacunadas contra la COVID-19 perdió su bebé.	Confidencial de Pfizer Documentos revelan el 90% de las mujeres embarazadas vacunadas Covid perdió su bebé; pero Pfizer afirmó: "No hay señal de seguridad"  Los documentos confidenciales de Pfizer que la FDA se vio obligada a publicar por orden judicial revelan que del 82 % al 97 % de las mujeres que fueron expuestas por error a la inyección de mRNA Covid-19 sufrieron un aborto espontáneo o tuvieron que presenciar la muerte de su hijo recién nacido al dar a luz.	<a href="https://www.facebook.com/erzeli/posts/3131598267110086">https://www.facebook.com/erzeli/posts/3131598267110086</a>

Figura 8 Ejemplo de corpus generado

En la Tabla 2 se muestran las características del corpus, además muestra la descripción de este, a través de la clasificación implementada por Listerry y Toourella, (1999), artículo titulado "Diseño de corpus textual y oral". Los datos obtenidos durante la construcción del corpus se incluyen a continuación.

Tabla 2 Clasificación y características del corpus

Clasificación	
Origen	Corpus Textual
Codificación y Anotado	Anotado
Tiempo	Diacrónico
Cantidad de Texto	Pequeño
Características	
Número de Textos	300
Idioma	Español
Datos Incluidos	Título, Texto, Link

## Etapa 2: Desarrollo de reglas gramaticales

Se desarrolló un total de 9 reglas de producción afirmativas y 8 negativas contando con una variante total de 32 reglas negativas y 9 afirmativas con el fin de categorizar por título el tipo de noticia extraída de las diferentes plataformas y redes sociales, almacenadas dentro del corpus de texto donde se tuvo que identificar por categoría cada grupo de noticias, con el propósito de generar reglas de producción para lograr conocer si la noticia que estamos identificando en ese momento es falsa o verdadera, para compararla con las notas extraídas desde el método propuesto.

En la siguiente Tabla 3, se observa el desarrollo del campo semántico identificado a partir de cada palabra de cada noticia falsa, a través de la agrupación de palabras en verbos, sustantivos, etc., junto con sus sinónimos de cada una de ellas.

- ❖ **Amarillo / Medicamento:** Se relaciona con Medicina, Remedio, Fármaco, Tratamiento. Sirve para curar o prevenir enfermedades.
- ❖ **Verde / Enfermedad:** Se relaciona con Dolores, Sufrimientos y padecimiento. Es cuando existe un virus o bacteria dentro de los cuerpos afectando la salud del mismo.
- ❖ **Gris / Objeto:** Se refiere a cualquier instrumento, dispositivo, equipo o herramienta utilizada en el diagnóstico, tratamiento, prevención o atención de enfermedades y condiciones médicas. Está relacionado con Cosa, Elemento, Cuerpo, Ente, Pieza.
- ❖ **Rojo / Acciones:** Se refiere a la actividad que hace el sujeto sobre una cosa relacionándola con un Objeto, Enfermedad, Medicamento.

*Tabla 3 Ejemplo de campos semánticos de noticias falsas alusivos a su color de categoría*

MEDICAMENTO	ENFERMEDAD	OBJETOS
Vacuna contra COVID-19	ELA	Documentos de Pfizer
Vacuna contra el coronavirus	Coronavirus	Mascarillas
Vacuna covi-19	Cáncer	Test PCR
Vacuna de ARN	Células de cáncer	Aire acondicionado de los autos
Contener la respiración	Infartos	Hisopos de PCR
Jugo	Ataques al corazón	Preservativo
Zumo de remolacha	Derrames cerebrales	
Apio	Cáncer de útero	
Zanahoria	Esterilidad	
Verduras	Sida	
Hojas de guanábana		
Agua antes de ir a la cama		
Agua cada 15 minutos		
Agua caliente con limón		
Vacunas contra el COVID-19		
Agua helada		
Agua con gas		
Coco		
Pepino		
Champú durante la menstruación		
Vacunas		
Proteínas		
Espigas tóxicas		
Vitaminas		

A continuación, se presenta la Tabla 4 con ejemplos de noticias falsas de temas de la salud, identificando los campos semánticos de cada una de ellas alusivas a su color por categoría.

Tabla 4 Categoría con sus palabras en común

Ejemplo
El jugo o zumo de remolacha, apio, zanahoria y verdura cura el cáncer
Las vacunas causa autismo
Usar uñas acrílicas causa cáncer

En la Tabla 5, se presentan las acciones Positivas y Negativas identificadas de cada una de las noticias dentro del corpus por categorías a través de los verbos en Infinitivo, con el propósito de conocer las estructuras gramaticales de las noticias para identificar si se está reconociendo una contradicción entre la noticia y la nota extraída de las fuentes de información médicas.

Tabla 5 Acciones positivas y negativas de cada noticia falsa

ACCION (+)	ACCIONES(-)
Aumentar, Causar, Curar, Beber, Evitar, Tomar, Acabar, Provocar, Inocular, Generar, Comer, Usar, Consumir,	No Aumentar, No Curar, No Beber, No Evitar, No Tomar, No Acabar, No Provocar, No Inocular, No Generar, No Comer, No Usar, No consumir

En la siguiente Tabla 6, se presenta algunas de las reglas de producción identificadas a través de las categorías y acciones encontradas para conocer qué patrones se pueden identificar de cada noticia falsa, con la nota a comparar extraída del método propuesto, recalando que esa anotación fue hecha solo para ejemplificar las reglas, pero esa no es la sintaxis de las reglas.

En color rojo se destaca como ejemplo un verbo de tipo acción y que presenta un matiz semántico dependiendo de la clase donde se menciona:

- ❖ Amarillo: Se enlaza a categoría de medicamento.
- ❖ Verde: Hace mención a la categoría de enfermedad.
- ❖ Gris: Se comunica con los objetos.

Tabla 6 Reglas identificadas

Categoría	MEDICAMENTO
Regla	<MEDICAMENTO> <ACCION><ENFERMEDAD>
Ejemplo	El jugo o zumo de remolacha, apio, zanahoria y verdura cura el cáncer. El ajo cura el cáncer. El Agua caliente con limón cura el cáncer. La dieta alcalina cura el cáncer. La vacuna contra la COVI-19 induce el sida. La hidroxiclороquina cura la viruela del mono. La vacuna australiana contra la COVID-19 contagia de VIH.
Categoría	ENFERMEDAD
Regla	<ENFERMEDAD><ACCION><MEDICAMENTO>
Ejemplo	El coronavirus se cura con gárgara caliente de bicarbonato de sodio. La COVID-19 se cura con jengibre, limón, aspirina, canela, miel, ajo y cebolla roja.
Categoría	OBJETO
Regla	<OBJETO><ACCION><ENFERMEDAD>
Ejemplo	El test PCR daña la barrera hematoencefálica del cerebro. Termómetro eléctrico infrarrojo, pistola con radiación genera cáncer cerebral. La mascarilla provoca cáncer y enfermedades autoinmunes. Las mascarillas incuban cáncer.

## Etapa 2: Implementación de reglas de producción

En esta sección se presentan las reglas de producción que fueron creadas dentro de un archivo .CFG, teniendo 9 reglas de producción afirmativas y 8 negativas contando con una variante de 31 reglas negativas sumando un total de 40 reglas de producción utilizadas para identificar hechos tanto de la noticia posiblemente falsa y las notas recolectadas por el método propuesto al generar la búsqueda para asignar un veredicto, a partir de una gramática de tipo 2 la cual se utiliza para definir la sintaxis de los lenguajes naturales y de programación y demás lenguajes formales, en función del conjunto T de símbolos terminales, el conjunto N de símbolos no terminales, el axioma S y el conjunto P de reglas de producción (N,T,P,S).

Se presenta a continuación la definición de la gramática libre de contexto o tipo dos para la gramática asignada para el método propuesto.

- ❖ Axioma S = Es el símbolo inicial desde el cual se empieza a aplicar las reglas de producción para generar todas las posibles cadenas de símbolos terminales pertenecientes al lenguaje descrito por la gramática.
- ❖ No terminales: son los componentes intermedios que se utilizan en las reglas, todo símbolo no terminal debe ser definido como una secuencia de otros símbolos. En nuestro caso, los símbolos no terminales van a ser las categorías gramaticales.
- ❖ Los símbolos terminales: son los componentes finales reconocidos por la Gramática. Los terminales van a ser las palabras de las oraciones que queremos analizar sintácticamente.
- ❖ Reglas de producción: Son utilizadas para definir patrones gramaticales o semánticos que ayudan a analizar y comprender el lenguaje humano.

A continuación, se presenta la gramática libre de contexto para el método propuesto. Donde se observa el axioma S denominado GRAMATICA, a partir del cual se define la regla más general.

<GRAMATICA>=<PATRON1>|<PATRON2>|<PATRON3>|<PATRON4>|<PATRON5>|<PATRON6>|<PATRON7>|<PATRON8>|<PATRON9>|<PATRON1\_NEGADA>|<PATRON2\_NEGADA>|<PATRON3\_NEGADA>|<PATRON4\_NEGADA>|<PATRON5\_NEGADA>|<PATRON6\_NEGADA>|<PATRON7\_NEGADA>|<PATRON8\_NEGADA>

El conjunto de los símbolos no terminales  $N_{lerrmm}$  se define de la siguiente manera:

$N_{lerrmm}$ =<ENFERMEDAD>|<ACCION>|<MEDICAMENTO>|<OBJETO>|<SIN\_EVIDENCIA>  
<FALSO>|<NO\_ACCION>|<INFECCION>|<SISTEMA\_NERVIOSO>|<GASTROINTESTINAL>|<AUTOINMUNE>|<\_COVID>|<\_CANCER>|<CANCER\_CUERPO>|<DIABETES>|<CIRCULATORIO>|<INFARTO>|<HIPERTENCION>|<OCULAR>|<CEGUERA>|<DERMIS>|<MENTAL>|<AUDITIVO>|<TOMAR>|<\_TOMAR>|<\_COMER>|<\_USAR>|<\_CURAR>|<PREVEINIR>|<CONTRARRESTAR>|<\_CAUSAR>|<\_EMPEORAR>|<PROPICIAR>|<CUERPO>|<AGUA>|<LICOR>|<REFRESCO>|<BEBIDA>|<\_LENTILLAS>|<\_LENTES>|<SMARTPHONE>|<\_PC>|<DINERO>|<CUBREBOCAS>|<MEDICINA>|<CANNABIS>.

El conjunto de símbolos terminales  $T_{\text{lermm}}$  se define como:

$T_{\text{lermm}} = \langle \text{no} \rangle | \langle \text{tampoco} \rangle | \langle \text{impedir} \rangle | \langle \text{prevenir} \rangle | \langle \text{donar} \rangle | \langle \text{nadar} \rangle | \langle \text{ingerir} \rangle | \langle \text{consumir} \rangle | \langle \text{administrar} \rangle | \langle \text{fumar} \rangle | \langle \text{inyectar} \rangle | \langle \text{usar} \rangle | \langle \text{emplear} \rangle | \langle \text{utilizar} \rangle | \langle \text{curar} \rangle | \langle \text{eliminar} \rangle | \langle \text{detener} \rangle | \langle \text{combatir} \rangle | \langle \text{matar} \rangle | \langle \text{evitar} \rangle | \langle \text{disminuir} \rangle | \langle \text{reducir} \rangle | \langle \text{debilitar} \rangle | \langle \text{provocar} \rangle | \langle \text{generar} \rangle | \langle \text{producir} \rangle | \langle \text{acelerar} \rangle | \langle \text{empeorar} \rangle | \langle \text{propagar} \rangle | \langle \text{contagiar} \rangle | \langle \text{infectar} \rangle | \langle \text{alzheimer} \rangle | \langle \text{parkinson} \rangle | \langle \text{cerebrovascular} \rangle | \langle \text{artritis} \rangle | \langle \text{cáncer} \rangle | \langle \text{cancerígeno} \rangle | \langle \text{infarto} \rangle | \langle \text{ceguera} \rangle | \langle \text{hematoma} \rangle | \langle \text{sordera} \rangle | \langle \text{corazón} \rangle | \langle \text{colon} \rangle | \langle \text{útero} \rangle | \langle \text{próstata} \rangle | \langle \text{apio} \rangle | \langle \text{espinaca} \rangle | \langle \text{lechuga} \rangle | \langle \text{tomate} \rangle | \langle \text{sal} \rangle | \langle \text{azúcar} \rangle | \langle \text{agua} \rangle | \langle \text{alcohol} \rangle | \langle \text{cerveza} \rangle | \langle \text{soda} \rangle | \langle \text{lente} \rangle | \langle \text{pupilete} \rangle | \langle \text{celular} \rangle | \langle \text{smartphone} \rangle | \langle \text{teléfono} \rangle | \langle \text{computadora} \rangle | \langle \text{mascarilla} \rangle | \langle \text{cubrebocas} \rangle | \langle \text{aspirina} \rangle | \langle \text{paracetamol} \rangle | \langle \text{metformina} \rangle | \langle \text{insulina} \rangle | \langle \text{desinfectante} \rangle | \langle \text{radiografía} \rangle | \langle \text{tomografía} \rangle | \langle \text{ultrasonido} \rangle.$

Por último, el conjunto de las reglas de producción  $P_{\text{lermm}}$  se cómo:

$P_{\text{lermm}} =$

1.  $\langle \text{ENFERMEDAD} \rangle ::= \langle \text{ACCION} \rangle \langle \text{MEDICAMENTO} \rangle$
2.  $\langle \text{MEDICAMENTO} \rangle ::= \langle \text{ACCION} \rangle \langle \text{ENFERMEDAD} \rangle$
3.  $\langle \text{OBJETO} \rangle ::= \langle \text{ACCION} \rangle \langle \text{ENFERMEDAD} \rangle$
4.  $\langle \text{ENFERMEDAD} \rangle ::= \langle \text{ACCION} \rangle \langle \text{OBJETO} \rangle$
5.  $\langle \text{ACCION} \rangle ::= \langle \text{MEDICAMENTO} \rangle ::= \langle \text{ACCION} \rangle \langle \text{ENFERMEDAD} \rangle$
6.  $\langle \text{ACCION} \rangle ::= \langle \text{OBJETO} \rangle \langle \text{ACCION} \rangle \langle \text{ENFERMEDAD} \rangle$
7.  $\langle \text{ENFERMEDAD} \rangle ::= \langle \text{ACCION} \rangle \langle \text{ACCION} \rangle \langle \text{MEDICAMENTO} \rangle$
8.  $\langle \text{MEDICAMENTO} \rangle ::= \langle \text{ACCION} \rangle \langle \text{ACCION} \rangle \langle \text{OBJETO} \rangle$
9.  $\langle \text{SIN\_EVIDENCIA} \rangle ::= \langle \text{ENFERMEDAD} \rangle \langle \text{ACCION} \rangle \langle \text{MEDICAMENTO} \rangle$
10.  $\langle \text{SIN\_EVIDENCIA} \rangle ::= \langle \text{MEDICAMENTO} \rangle \langle \text{ACCION} \rangle \langle \text{ENFERMEDAD} \rangle$
11.  $\langle \text{SIN\_EVIDENCIA} \rangle ::= \langle \text{OBJETO} \rangle \langle \text{ACCION} \rangle \langle \text{ENFERMEDAD} \rangle$
12.  $\langle \text{SIN\_EVIDENCIA} \rangle ::= \langle \text{ENFERMEDAD} \rangle \langle \text{ACCION} \rangle \langle \text{OBJETO} \rangle$
13.  $\langle \text{SIN\_EVIDENCIA} \rangle ::= \langle \text{ACCION} \rangle \langle \text{MEDICAMENTO} \rangle \langle \text{ACCION} \rangle \langle \text{ENFERMEDAD} \rangle$
14.  $\langle \text{SIN\_EVIDENCIA} \rangle ::= \langle \text{ACCION} \rangle \langle \text{OBJETO} \rangle \langle \text{ACCION} \rangle \langle \text{ENFERMEDAD} \rangle$
15.  $\langle \text{SIN\_EVIDENCIA} \rangle ::= \langle \text{ENFERMEDAD} \rangle \langle \text{ACCION} \rangle \langle \text{ACCION} \rangle \langle \text{MEDICAMENTO} \rangle$



16. <SIN\_EVIDENCIA>:=<ENFERMEDAD><ACCION><ACCION><OBJETO>
17. <FALSO>:=<ENFERMEDAD><ACCION><MEDICAMENTO>
18. <FALSO>:=<MEDICAMENTO><ACCION><ENFERMEDAD>
19. <FALSO>:=<OBJETO><ACCION><ENFERMEDAD>
20. <FALSO>:=<ENFERMEDAD><ACCION><OBJETO>
21. <FALSO>:=<ACCION><MEDICAMENTO><ACCION><ENFERMEDAD>
22. <FALSO>:=<ACCION><OBJETO><ACCION><ENFERMEDAD>
23. <FALSO>:=<ENFERMEDAD><ACCION><ACCION><MEDICAMENTO>
24. <FALSO>:=<ENFERMEDAD><ACCION><ACCION><OBJETO>
25. <ENFERMEDAD>:=<ACCION><MEDICAMENTO><FALSO>
26. <MEDICAMENTO>:=<ACCION><ENFERMEDAD><FALSO>
27. <OBJETO>:=<ACCION><ENFERMEDAD><FALSO>
28. <ENFERMEDAD>:=<ACCION><OBJETO><FALSO>
29. <ACCION>:=<MEDICAMENTO><ACCION><ENFERMEDAD><FALSO>
30. <ACCION>:=<OBJETO><ACCION><ENFERMEDAD><FALSO>
31. <ENFERMEDAD>:=<ACCION><ACCION><MEDICAMENTO><FALSO>
32. <ENFERMEDAD>:=<ACCION><ACCION><OBJETO><FALSO>
33. <ENFERMEDAD>:=<NO\_ACCION><MEDICAMENTO>
34. <MEDICAMENTO>:=<NO\_ACCION><ENFERMEDAD>
35. <OBJETO>:=<NO\_ACCION><ENFERMEDAD>
36. <ENFERMEDAD>:=<NO\_ACCION><OBJETO>
37. <ACCION>:=<MEDICAMENTO><NO\_ACCION><ENFERMEDAD>
38. <ACCION>:=<OBJETO><NO\_ACCION><ENFERMEDAD>
39. <ENFERMEDAD>:=<NO\_ACCION><ACCION><MEDICAMENTO>
40. <ENFERMEDAD>:=<NO\_ACCION><ACCION> <OBJETO>

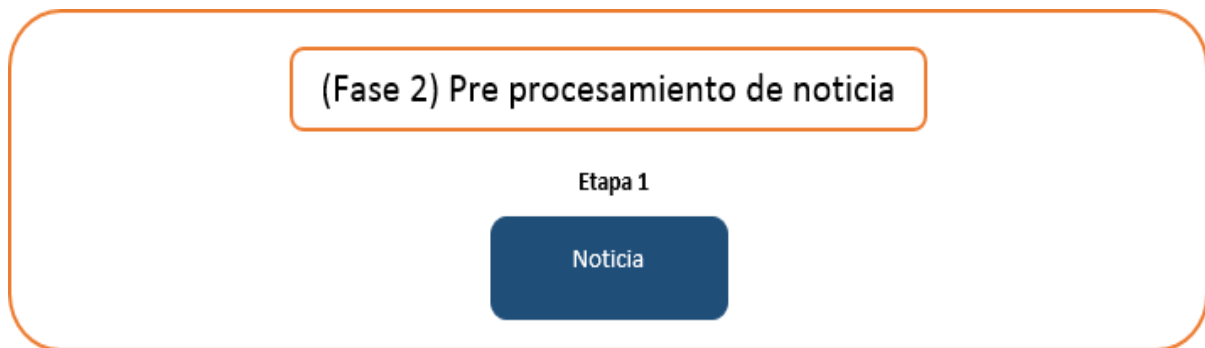
### Etapa 3: Obtención de hechos

El proceso de identificación de hechos se explica a fondo en la Etapa 4, también llamada Identificación de contradicciones.

### FASE 2: Pre procesamiento de noticia

La segunda fase de la metodología está conformada por una etapa, con la finalidad de generar el pre procesamiento de noticia, esta es una etapa fundamental en el análisis de texto que implica la limpieza y transformación de los datos.

A continuación, se observa la Fase 2 de la metodología de solución también llamada Pre procesamiento en la Figura 9.



*Figura 9 Fase 2 de la metodología de solución*

### Etapa 1: Noticia Falsa

Para generar la búsqueda y recolectar las notas, en primer lugar, se realiza un pre procesamiento donde se realiza la eliminación de *Stopwords* y la Tokenización del título de la noticia candidata a ser falsa, para lo cual se utiliza la librería de Python NLTK, ampliamente utilizada en el Procesamiento de Lenguaje Natural.

A continuación, se presenta la Tabla 7 un ejemplo de una búsqueda donde se eliminan las *Stopwords* y se genera la Tokenización de palabras del título de la noticia a buscar.

Tabla 7 Método de pre procesamiento de noticia

Tarea	Salida
Extracción de título	Beber agua caliente de limón cura el cáncer
Eliminación de <i>stopwords</i>	Beber agua caliente limón cura cáncer
<i>Tokenización</i>	Beber agua calient elimón cura cáncer
Lematización	Beber agua caliente limón curar cáncer
Infinitivo	Beber agua caliente limón curar cáncer

### FASE 3: Recuperación de Información

La tercera Fase de la metodología está conformada por tres etapas: Etapa 1: identificar y recolectar las fuentes de información médicas que se usarán durante el proceso de identificación de las notas médicas, Etapa 2: Proceso general de datos necesarios para generar una búsqueda, Etapa 3: Generar el proceso de búsqueda para extraer las notas relevantes a partir de la Query del título de la noticia a buscar, Etapa 4: Aplicar un método de similitud que nos ayude filtrar notas que no tienen ninguna o poca relevancia con la búsqueda.

A continuación, se muestra la Fase 3 de la metodología de solución también llamada Recuperación de Información de la Figura 10.

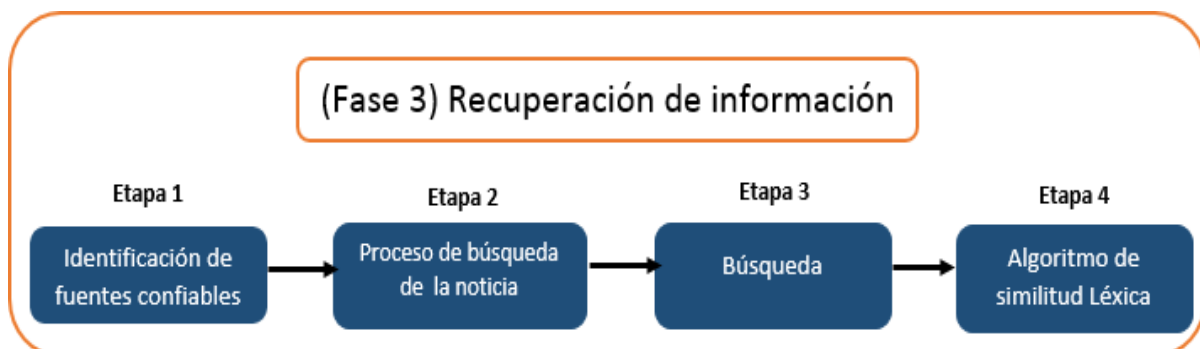


Figura 10 Fase 3 de la metodología de solución

## Etapa 1: Identificación de fuentes de información médicas

Se generó el proceso de identificación de fuentes de información médica las cuales se utilizaron para poder generar búsquedas y almacenar notas médicas que hablen sobre diferentes áreas de la salud.

A continuación, se presentan las fuentes de información médicas utilizadas para el proceso de búsquedas.

1. MedlinePlus: Es un servidor de información en línea provisto por la Biblioteca Nacional de Medicina de los Estados Unidos. Brinda información sobre salud de forma gratuita, en inglés y español. [16]
2. PubMed: Es un motor de búsqueda de libre acceso que permite consultar principal y mayoritariamente los contenidos de la base de datos. [17]
3. Newtral: Es un medio de comunicación español constituido como empresa dedicada a la comprobación de hechos y producción de programas de televisión, es un Startup fundado por Ana Pastor. El cual es dedicado a Fast-checking dedicada a desmentir noticias falsas. [18]
4. AFP Factual: Los fast-checkers de la AFP investigan afirmaciones sospechosas que son virales, tienen impacto en la sociedad y son potencialmente dañinas para el público. [19]

Las afirmaciones que verifica pueden aparecer en diversos formatos como Redes Sociales, Blogs, Sitios Web, Aplicaciones de mensajería y otros foros de la esfera pública.

## Etapa 2: Proceso de búsqueda de la noticia

Esta etapa se implementó con el fin de dar a conocer que datos se necesitan principalmente para el proceso de generación de búsqueda, el cual nos lleva a la siguiente pregunta.

Los datos necesarios utilizados para el proceso de búsqueda son los siguientes:

1. Título: Para generar la búsqueda se necesita un título de una noticia posiblemente falsa.
2. Pre procesamiento: En esta etapa se pasa la noticia posiblemente falsa a través de un proceso de filtrado del texto del título de la noticia tales como Eliminación de StopWords, Tokenización y lematización de palabras.
3. Hechos: Para generar el proceso de identificación de hechos existentes tanto en el título de la noticia como en las notas recolectadas a partir de la búsqueda a través de la Query del título, se crean ciertas reglas de producción alimentadas a partir de una Gramática las cuales ayudarán a identificar los hechos existentes.

### Etapa 3: Búsqueda (Query)

Para generar la búsqueda se tiene que generar la consulta a través de la Query del Título de la noticia posiblemente falsa a través del método propuesto, para conocer la información más relevante de este título de noticia extraída de las diferentes fuentes de información médicas.

Posteriormente antes de pasar a la búsqueda se procede a realizar una conexión a través de una API que nos brinde apoyo para generar una Búsqueda. A continuación, se presenta el proceso de conexión.

### **Conexión**

Este módulo permite la conexión a través de la API de *Google Custom Search*, con el propósito de generar una búsqueda a través de una consulta en más fuentes de información médica obteniendo como resultado ciertas notas médicas. (Vea Figura 11).

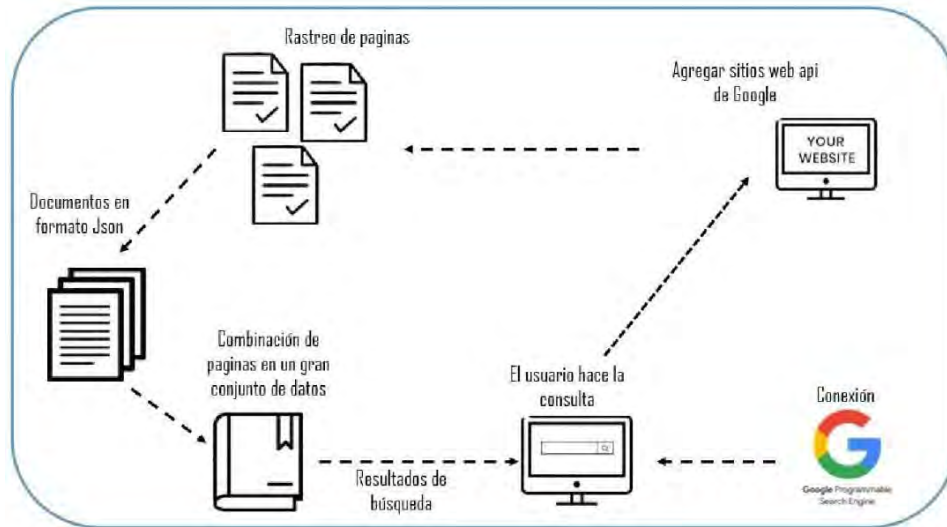


Figura 11 Diagrama de planificación para la conexión al API de Google

A continuación, se presenta el proceso de etapas para implementar la conexión a la API.

En este apartado se presenta la Figura 12 con el nombre del API a utilizar llamada Custom Search API, con el propósito de crear la clave del API para generar la conexión, donde como primer proceso se tiene que dar clic en el apartado de habilitar para entrar a la página de “Apis y servicios habilitados” con la finalidad de poder crearla.

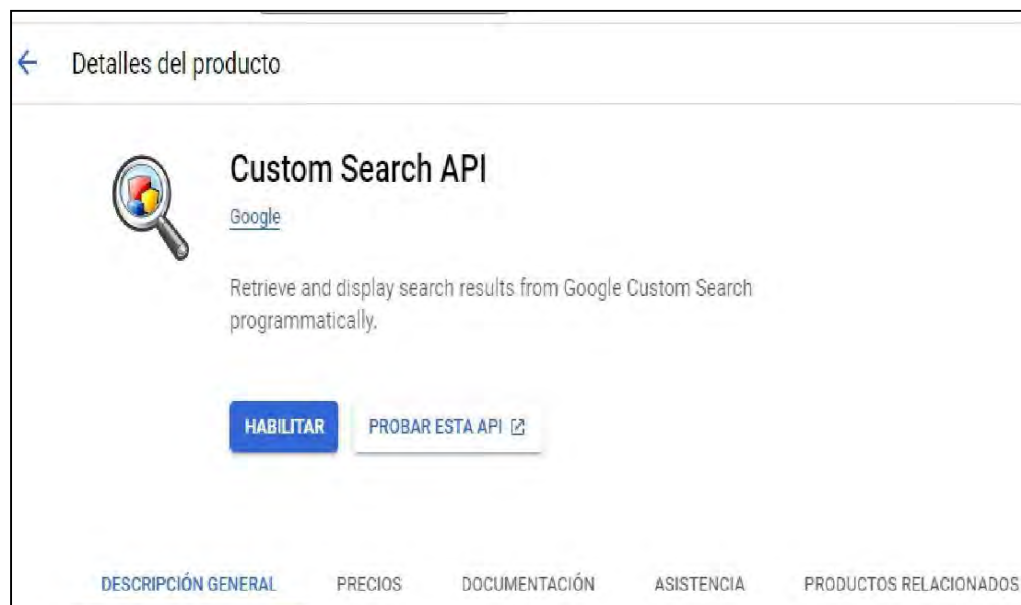


Figura 12 Nombre de Api a utilizar Custom Search API

Ya que se ingresó a la pestaña de “Apis y servicios habilitados” se tiene que ingresar en el apartado de credenciales. (Vea Figura 13)

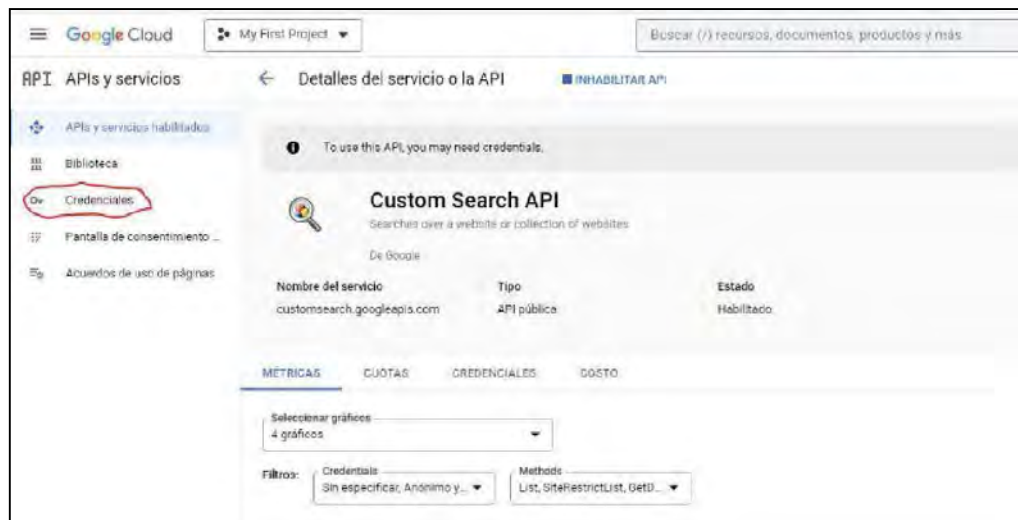


Figura 13 Apis y servicios habilitados

A continuación, ya que se seleccionó el apartado de credenciales aparecerá un recuadro con los tipos de cada una de ellas para seleccionar, (Vea Figura 14) donde se seleccionó el apartado de “Clave de API” para identificar mi proyecto a través de una clave de API simple para verificar la cuota y el acceso proporcionado.

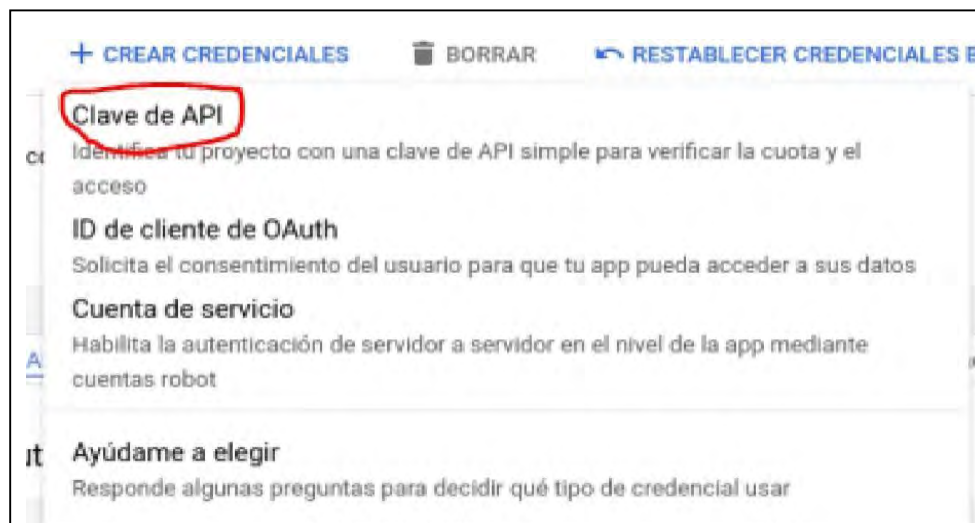


Figura 14 Nombre de tipos de credenciales

Ya que se seleccionó el apartado de “Clave de API” te direcciona a la pestaña donde se crea la clave, esta clave se utilizará más tarde en nuestra aplicación. Como se presenta en la (Vea Figura 15).

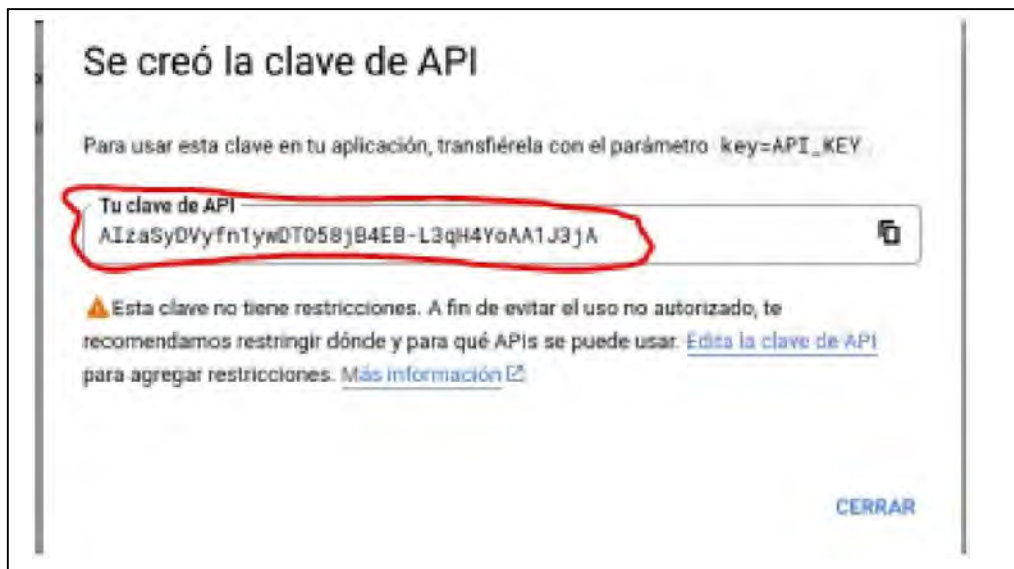


Figura 15 Nombre de clave de Api de credencial

Ya que se cuenta con la clave del API, se tiene que crear el “Motor de búsqueda programable” en la siguiente pestaña como se muestra la Figura 16 seleccionando el botón de Añadir para crearlo.

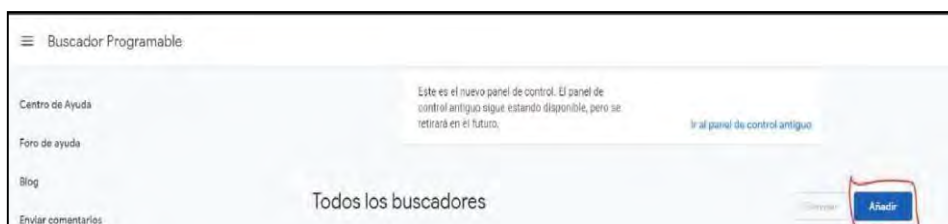


Figura 16 Motor de búsqueda programable



Una vez seleccionado el botón de “añadir”, te aparecerá el siguiente recuadro donde tendrás que asignar el nombre para tu buscador y el tipo de configuración de búsqueda. (Vea Figura 17).

Crear un nuevo buscador

Para empezar, proporciona información básica sobre tu buscador. Podrás personalizar más las configuraciones del buscador (idiomas, regiones, etc.) después de crearlo. [Más información](#)

Asigna un nombre a tu buscador

¿Dónde buscar?  Buscar sitios o páginas específicos

Puedes añadir cualquiera de los elementos siguientes:

Páginas web: `www.ejemplo.com/página.html`  
Todo un sitio: `www.ejemplo.com/*`  
Partes de un sitio: `www.ejemplo.com/documento/*` o `www.ejemplo.com/documento/*`  
Todo un dominio: `*ejemplo.com`

Buscar en toda la Web

Configuración de búsqueda  Búsqueda por imágenes  
 Búsqueda Segura

No soy un robot

Al hacer clic en Crear, aceptas los [Términos del Servicio](#).

Figura 17 Crear un buscador

Posteriormente se hace clic en el nombre del buscador que recién se añadió (Vea Figura 18) llamado “Noticias falsas de temas sobre la salud”.

<input type="checkbox"/>	Nombre	Rol	URL pública	Hora de la última actualización
<input checked="" type="checkbox"/>	Noticias falsas de temas sobre la salud	Propietario		17 abr 2023, 11:06

Filas por página: 10 1-1 de 1

Figura 18 Todos los buscadores

A continuación, se muestran los datos el nombre de tu buscador, la Figura 19 te muestra los datos básicos donde se observa el Nombre del buscador, Descripción, Código, Id del buscador y URL pública.

The image shows a web interface titled "Vista general" (General View). Inside, there is a section titled "Datos básicos" (Basic Data) which contains five rows of information:

Field	Value	Icon
Nombre del buscador	Noticias falsas de temas sobre la salud	✎
Descripción	Añadir descripción	
Código	Obtener código	
ID de buscador	f0272d58b7f354512	📄
URL pública	<a href="https://cse.google.com/cse?cx=f0272d58b7f354512">https://cse.google.com/cse?cx=f0272d58b7f354512</a>	

*Figura 19 Datos Básicos*

Para concluir, ya que se tiene la "Clave de API" y "Id de buscador" se añade a la herramienta en Python como se muestra el ejemplo en la Figura 20 para posteriormente visualizar nuestro proyecto programable.

```

class WebSearchGoogle:
    def __init__(self):
        # Estableciendo credenciales y la url de búsqueda
        self.api_key = ' AizaSyBjA46_n8pWZb0tNjZdHRe0p*****'
        self.engine_id = 'e5ded972e9f4049da'
        self.url = 'https://www.googleapis.com/customsearch/v1'
        self.server = ''
        self.file = ''
        self.term = ''
        self.source = ''

```

*Figura 20 Código añadido a la herramienta en Python*

**Agregar sitios web al API de Google:** En esta etapa se ingresan la URL de los sitios web que se encuentran validados como factibles para conocer información relacionada con temas sobre la salud, para generar la consulta dentro de los sitios web proporcionados.

**Rastreo de páginas:** Se genera el rastreo de las fuentes de información médicas a través de la Query del título de la noticia que se quiere mandar a traer la información para que nos arroje información relacionada.

**Documentos en formato Json:** El tipo de formato que se quiere es un formato ligero de intercambio de datos y fácil de leer.

**Combinación de fuentes de temas sobre la salud:** Al generar el rastreo de las páginas, devuelve los nombres de los documentos de las fuentes de información que fueron asignadas para generar la búsqueda de forma agrupada.

#### Etapa 4: Aplicación de algoritmo de Similitud Léxica

Se implementó un API de algoritmo de similitud léxica en las notas recolectadas anteriormente a través de la búsqueda realizada sobre la Query del título de la noticia, para poder conocer las similitudes de cada una de las notas extraídas comparadas con la Query del título de la noticia, con el propósito de conocer las notas que tengan una mejor relación y significado de palabra con la noticia que están ingresando los usuarios.

A continuación, se presentan los algoritmos utilizados y sus conceptos con el propósito de conocer cuál de ellos fue el óptimo para utilizar que nos brindara mayor accesibilidad para conocer a través de la consulta un conjunto de documentos y devolver los más relevantes para la consulta.

A continuación, se presentan los algoritmos utilizados:

BM-25: Es un algoritmo utilizado para generar una búsqueda a través de la Query de un título para la extracción de documentos que contengan similitudes semejantes entre ellos. [20]

A continuación, se presenta en la Figura 21 la fórmula de la API del algoritmo de BM25 utilizada:

$$\text{score}(D, Q) = \sum_{i=1}^n \text{IDF}(q_i) \cdot \frac{f(q_i, D) \cdot (k_1 + 1)}{f(q_i, D) + k_1 \cdot \left(1 - b + b \cdot \frac{|D|}{\text{avgdl}}\right)},$$

Figura 21 Fórmula BM-25

Donde

- ❖  $(D, Q)$ : Valor de salida para calcular la puntuación de un documento (D) dada una consulta (Q).
- ❖  $\sum_{i=1}^n$ : Es la suma de una secuencia de términos donde el índice  $i$  toma valores desde
- ❖  $1=$  hasta  $n$ , dependiendo del término que se sume estando a la derecha.
- ❖  $IDF(q_1)$ : Se refiere a calcular el valor (IDF) siendo la medida que indica que tan relevante es una palabra ( $q_i$ ) para un documento.
- ❖  $F(q_i, D)$ : Es la frecuencia de palabras ( $q_i$ ) que aparecen dentro del documento (D) de los términos que aparecen en la consulta inicial.
- ❖  $K_1, b$ : Es el parámetro de como la frecuencia de una palabra en un documento influye en la puntuación inicial. Valores de medida  $k_1=2.0$  y  $b=0.75$ .
- ❖  $|D|$ : Es la longitud del documento (D) en termino de numero de palabras.
- ❖  $Avgdl$ : Es la medida de la longitud de los documentos donde estamos realizando la búsqueda.

Fast-BM25: Es una optimización de BM25 el cual utilizan (BestMacht25) para la recuperación de información recalando que aplican el mismo algoritmo de similitud, el cual se centra en mejorar la eficiencia computacional del cálculo de la puntuación de relevancia. Esta optimización es importante en sistemas de búsqueda de gran escala donde se necesita procesar grandes volúmenes de documentos y consultas de forma eficiente. [21]

En términos generales, el algoritmo de FAST BM-25 calcula la similitud entre un documento y una consulta utilizando una fórmula que tiene en cuenta factores como la frecuencia de los términos en el documento y en la consulta, así como la longitud del documento y la longitud promedio de los documentos en la colección.

*Sparse-Retriver*: Es una técnica utilizada en la recuperación de información y en sistemas de búsqueda para encontrar los documentos más relevantes para una consulta dada. El objetivo principal de este algoritmo de recuperación de informaciones que los documentos relevantes a menudo contienen solo unos pocos términos claves de la consulta. Por lo tanto, en lugar de comparar la consulta con todos los documentos en la colección, el algoritmo *Sparse-Retriver* se enfoca en identificar y recuperar sólo un subconjunto seleccionado de documentos que tienen una alta probabilidad de ser relevantes. [22]

## **Resultados de la experimentación de las búsquedas implementadas a través de APIS de similitud léxica**

En esta sección se presentan los resultados obtenidos a través de las pruebas realizadas de los diferentes algoritmos implementados con el propósito de identificar cuál de ellos está logrando encontrar la mayor relevancia de similitud entre la búsqueda generada y los documentos devueltos identificados extraídos de las fuentes de información médicas.

En total se realizaron 10 pruebas por cada uno de las APIS implementadas dando un total de 30 resultados A continuación, se presentan las pruebas aplicadas a cada algoritmo.

Tabla 8 Pruebas de Api BM-25

Usar cubrebocas provoca neumonía		
N	TÍTULO	SIMILITUD
1	Las mascarillas limpias y bien manipuladas no provocan neumonía   Factual	3.31
2	Verificaciones de la AFP en español sobre el nuevo coronavirus   Factual	2.82
3	Es falso que el tapabocas cultive hongos y bacterias que bloquean el oxígeno y provocan neumonía - Chequeado	1.55
4	Neumonía en adultos, adquirida en la comunidad: MedlinePlus enciclopedia médica	1.39
5	Uso de guantes en el hospital: MedlinePlus enciclopedia médica	0.85
6	COVID-19, una emergencia de salud pública mundial - PMC	0.76
7	Facebook	0.55

Se puede observar que en la Tabla 8 a través de las pruebas aplicada a la API de BM-25 está identificando ciertas similitudes léxicas en las notas recolectadas a través de la búsqueda, identificando que esta Api es eficiente si la comparamos con el título de la búsqueda y las notas devueltas se puede observar que sí hay similitudes entre la primera nota con mayor rango de similitud.

Tabla 9 Pruebas Api Fast BM - 25

Usar cubrebocas provoca neumonía		
N	TÍTULO	SIMILITUD
1	Uso de guantes en el hospital: MedlinePlus enciclopedia médica	3.32
2	COVID-19, una emergencia de salud pública mundial - PMC	2.82
3	Neumonía en adultos, adquirida en la comunidad: MedlinePlus enciclopedia médica	1.55
4	Verificaciones de la AFP en español sobre el nuevo coronavirus   Factual	1.48
5	Las mascarillas limpias y bien manipuladas no provocan neumonía   Factual	0.86
6	Facebook	0.77
7	Es falso que el tapabocas cultive hongos y bacterias que bloquean el oxígeno y provocan neumonía - Chequeado	0.56

Los resultados de la búsqueda de la Tabla 9 se identificó que las similitudes léxicas identificadas con mayor porcentaje de semejanza no tienen nada que ver con el título de la noticia buscada, y a su vez se reconoció que las que cumplen con menor porcentaje de similitud tiene una mayor semejanza a la noticia buscada.

Ya que se pasó por el filtro las notas recolectadas y se almacenaron las notas que contienen mayor similitud con la Query del título de la noticia a través de la búsqueda, se pasan a través de la Gramática desarrollada para encontrar hechos en las notas que nos brinden oraciones que nos ayuden a verificar si la noticia es falsa o verdadera.



Tabla 10 Pruebas de Api Sparce - Retrriver

Usar cubrebocas provoca neumonía		
N	TÍTULO	SIMILITUD
1	Uso de guantes en el hospital: MedlinePlus enciclopedia médica	42.03
2	COVID-19, una emergencia de salud pública mundial - PMC	4.00
3	Neumonía en adultos, adquirida en la comunidad: MedlinePlus enciclopedia médica	31.45
4	Verificaciones de la AFP en español sobre el nuevo coronavirus   Factual	25.09
5	Las mascarillas limpias y bien manipuladas no provocan neumonía   Factual	14.26
6	Facebook	10.30
7	Es falso que el tapabocas cultive hongos y bacterias que bloquean el oxígeno y provocan neumonía - Chequeado	0.81

A través de las pruebas realizada de la Tabla 10 se identificó que las similitudes léxicas identificadas con mayor porcentaje de semejanza no tienen nada que ver con el título de la noticia buscada, y a su vez se reconoció que las que cumplen con menor porcentaje de similitud tiene una mayor semejanza a la noticia buscada

## Conclusión

En resumen, se reconoció que a partir de las pruebas realizadas de cada una de las Apis utilizadas en el reconocimiento de similitudes léxicas, el denominado BM-25 está logrando identificar y reconocer mejor las notas que contienen una mayor similitud con la Query del título de la noticia que se está buscando.

## Diferencias entre BM-25 y FAST BM-25.

Para notar las diferencias que existen primero hay que mencionar que ambos algoritmos poseen a *BestMacht25* que se utiliza comúnmente en la recuperación de información y la búsqueda de texto.

Las diferencias entre estos algoritmos pueden surgir en los cálculos de los parámetros de similitud como la frecuencia de los términos y la longitud del documento, que pueden afectar la relevancia asignada a un documento en relación con una consulta de búsqueda.

## FAST- BM25

Se trata de una alteración en el algoritmo para agilizar peso, e incrementar rendimiento, a cambio de afinidad, por ejemplo:

En un Corpus de texto, las palabras más comunes (*stopwords*) suelen ser las menos informativas. Al excluirlos de la consulta y buscar únicamente documentos que contengan al menos una palabra de la consulta, BM25 gana mucha velocidad y pierde muy poca precisión.

En resumen, la diferencia esencial y característica entre ambos siempre será la exactitud y la rapidez del resultado, cuando BM25 trabaja de forma habitual, Fast BM25 tiene alteraciones para funcionar 10 veces más rápido, a coste de su exactitud.

### **Ejemplo**

Algoritmo = [BM25]

Consulta = [Usar cubrebocas provoca neumonía]

Similitud = [3.31]

Resultado = [Las mascarillas limpias y bien manipuladas no provocan neumonía

| Factual (3.31)]

Algoritmo = [Fast BM25]

Consulta = [Usar cubrebocas provoca neumonía]

Similitud = [0.83]

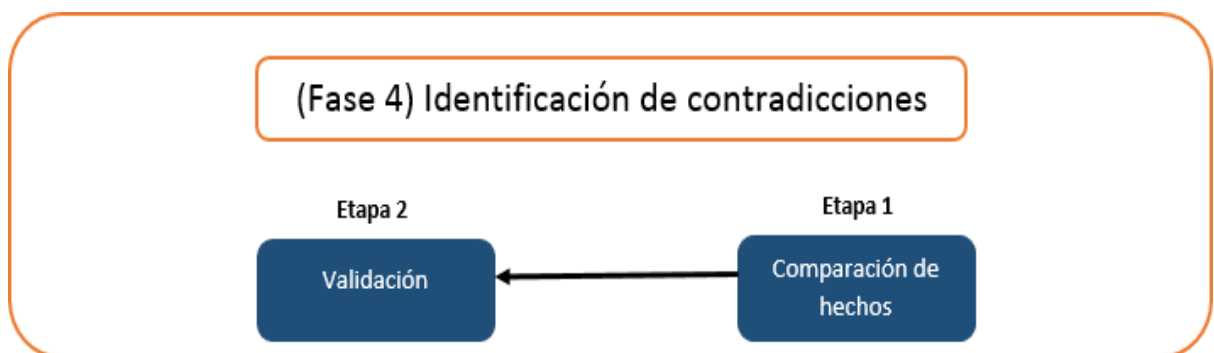
Resultado = [Las mascarillas limpias y bien manipuladas no provocan neumonía

| Factual (0.83)]

## FASE 4: Identificación de contradicciones

La cuarta fase de la metodología está conformada por dos etapas. Con la finalidad de conocer los hechos existentes en las noticias para dar un veredicto de si es falsa o verdadera: Etapa 1: Comparación entre el hecho identificado de la noticia con la que se generó la búsqueda y compararla con los hechos de las notas recolectadas a través de la búsqueda, Etapa 2: Validación para generar un veredicto que indique si es falsa o verdadera la noticia identificada.

A continuación, se observa la Fase 4 de la metodología de solución también llamada Identificación de contradicciones dentro de la Figura 22.



*Figura 22 Fase 4 de la metodología de solución*

### Etapa 1: Comparación de hechos

A continuación, se presentan las pruebas aplicadas a través de las reglas de producción asignadas dentro de la Gramática en las Figura 23 y 24, se identifican los árboles sintácticos generando una representación gráfica de un hecho identificando dentro del título de la noticia y la nota recolectada el cual nos sirve para realizar la comparación de hechos entre las dos y generar un veredicto.

### Ejemplo

#### **Beber agua de limón cura el cáncer**

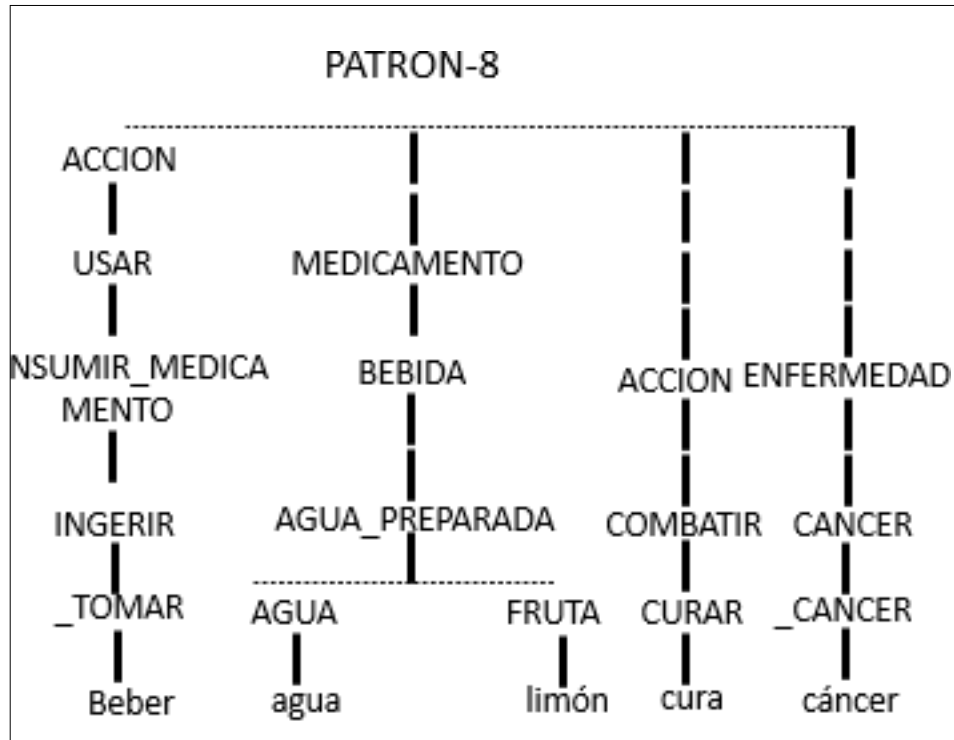


Figura 23 Árbol sintáctico de afirmaciones positivas

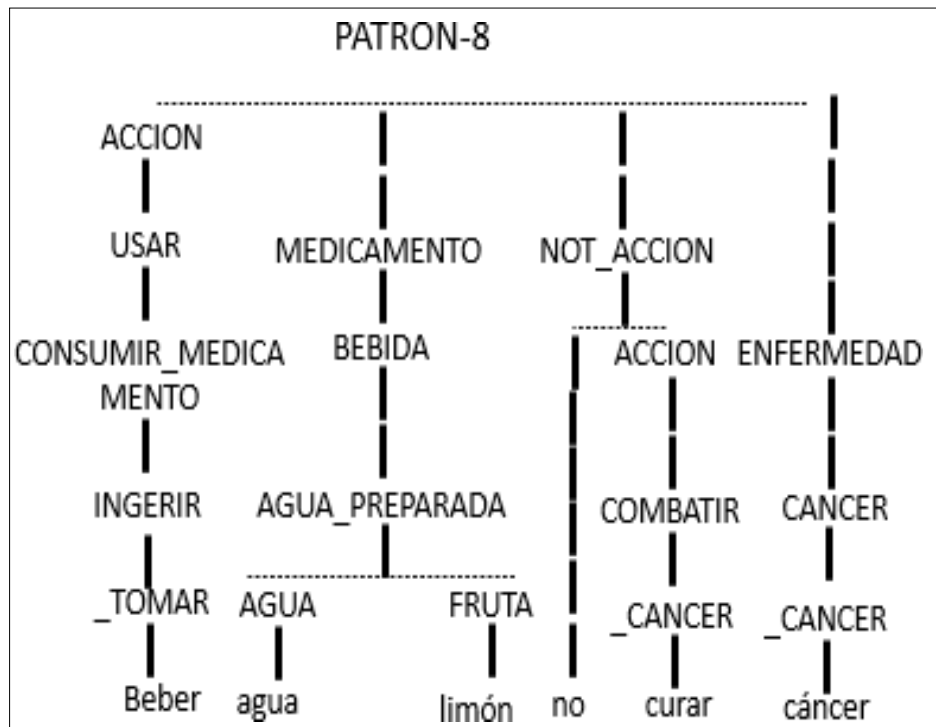


Figura 24 Árbol sintáctico de afirmaciones negativas

## Etapa 2: Comparación y validación

Ya que se tiene identificados los hechos tanto de la noticia al buscar y de las notas anteriormente de la búsqueda realizada se pasa al siguiente proceso el cual se implementó con el propósito de lograr conocer si existe una contradicción entre los 2 patrones devueltos por el sistema, con el fin de conocer si la noticia es falsa o verdadera a través de los tipos de veracidad.

Tipos de veracidad:

- ❖ Si existe contradicción: El sistema devolverá el patrón identificado como contradictorio y la nota completa donde se encontró ese fragmento de texto, junto con su liga de acceso de la fuente de información de donde se está validando que esa noticia es falsa.
- ❖ Si no existe contradicción: El sistema devolverá un mensaje en donde haga mención de que no se encontró información necesaria para procesar esa noticia.

## FASE 5: Diseño de Interfaz

La quinta Fase de la metodología está conformada por 3 Etapas. Con la finalidad de conocer el diseño de cada vista que complementan el sistema de VERACIA. Etapa 1: La vista de inicio representa la sección donde se genera el proceso de generación de búsqueda y validación para reconocer el veredicto de una noticia, Etapa 2: Hace mención a la vista donde se observan las noticias más populares, Etapa 3: Se presenta la vista de acerca de donde se presenta un manual informativo de cómo funciona el sistema

A continuación, se observa la Fase 5 de la metodología de solución también llamada Diseño de interfaz web dentro de la Figura 25.

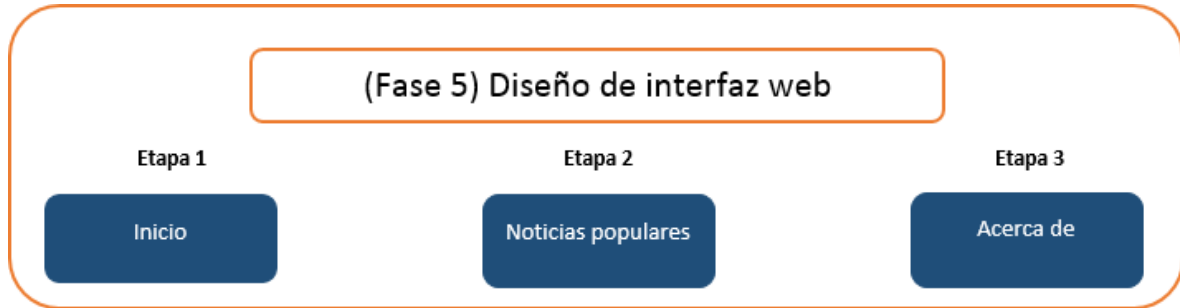


Figura 25 Fase 5 de la metodología de solución

Esté módulo fue desarrollado utilizando HTML, CSS, *Boostrap* y el *framework* de *Django* para poder hacer la interconexión del sistema de detección de noticias falsas de temas sobre la salud con la interfaz web para la detección de noticias y la iteración con el usuario.

A continuación, la interfaz se divide en 3 vistas distintas, que se explican a continuación.

### Etapa 1: Vista de Inicio

La vista de inicio brinda banners publicitarios haciendo mención a que no te dejes engañar o sorprender por noticias falsas y que las reportes para su prevención, contando con el apartado del analizador de noticias falsas, además la página contiene el logotipo del tecnológico nacional de México y el del centro nacional de investigación y desarrollo tecnológico (CENIDET) y el logotipo de la secretaría de educación y del nombre del sistema llamado VERACIA en la Figura 26 se muestra como es la vista de inicio. Cabe recalcar que para el diseño de la página se incluyó un residente para mejorar la calidad de las vistas de cada módulo de la misma.



Figura 26 Vista de inicio

A continuación, se presenta la vista del analizador, que se encuentra dentro del módulo de inicio al generar una búsqueda realizada por un usuario como se presenta en la Figura 27 el siguiente título “Beber agua caliente de limón cura el cáncer”.

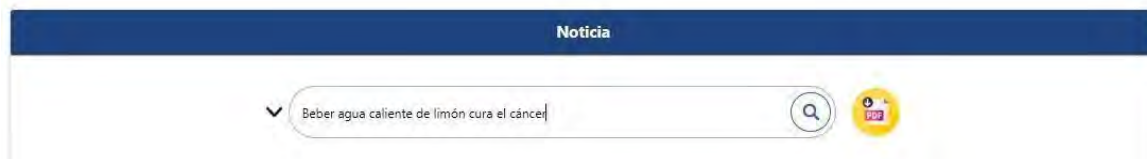


Figura 27 Analizador

Como ya se mencionó una vez generada la búsqueda el sistema ingresará a una serie de fuentes de información médicas estipuladas las cuales recolectan ciertas notas que tengan similitudes con la Query del título de la noticia que se ha ingresado, logrando identificar las notas recolectadas en el apartado de ARTÍCULOS CAPTURADOS. (Vea Figura 28)



Figura 28 Artículo Capturados

A continuación, se presentan los resultados obtenidos a través de la búsqueda realizada dentro del sistema, (Vea Figura 29) presentando una tabla llamada HECHOS DETECTADOS identificando el apartado que presenta las oraciones identificadas a través de las reglas de producción que sirven con el propósito de identificar contradicciones dentro de las notas identificadas y estipular un veredicto tanto como FALSO/VERDADERO.

Hechos detectados			
Hecho en la noticia	Nota	Hechos en las notas relacionadas	Resultado
beber agua caliente de limón cura el cáncer	1	Beber agua caliente, sola o con limón, no cura el cáncer   Factual	Contradicción

Figura 29 Hechos detectados

Ya identificando si la noticia ingresada por un usuario existe contradicción, se encuentra la tabla de VEREDICTO que sirve para reafirmar a través de una leyenda si es FALSA O VERDADERA. (Vea Figura 30)

Veredicto
• La noticia 1 contiene información Falsa.

Figura 30 Veredictos

## Etapa 2: Vista de Noticias Populares

Esta vista permite encontrar las posibles noticias que fueron más populares o las que tuvieron gran impacto social en las personas, (Vea Figura 31) ya que se identificó que en la época que se dio a conocer la pandemia del covid-19 circularon diferentes noticias en las redes sociales donde como ejemplo mencionan que la (vacuna covid-19 causa cáncer), a consecuencia de la noticia muchas personas no se vacunaban sobre el actual virus llamado covid-19 para tener cierta inmunidad a ella, y como consecuencia las personas se enfermaban al contraerlo y en ciertos casos fallecían por creer en noticias que en su momento fueron propagadas y divulgadas de manera masiva en diferentes plataformas y redes sociales haciéndolas pasar como noticias verdaderas.





Figura 31 Vista de noticias populares

### Etapa 3: Vista de Acerca De

En esta vista se muestra a detalle una breve descripción de qué es el sistema de detección de noticias falsas de temas sobre la salud llamado Veracia e identificando un pequeño desglose de funcionalidad del mismo, para que se les facilite a los usuarios como usar el sistema (Vea Figura 32).



Figura 32 Acerca de

# Capítulo V

Resultados

## 5. Resultados

En este capítulo se presentan las pruebas con el objetivo de comprobar y validar la metodología, así como también analizar los resultados.

### 5.1 Experimentación

Se realizó la evaluación de la Gramática porque siempre es fundamental generar el análisis, ya que no solo sirve para medir su funcionamiento sino para mejorarlo, e incluso compárarlo o complementarlo o sustituirlo por otras metodologías. Para generar los resultados se cuenta con dos clases de términos: Errores y Aciertos, con los que se comparan los resultados los cuales se clasifican en distintos grupos.

Esta información puede interpretarse bajo las métricas de Precisión, Cobertura y F1-Scores, muy utilizadas en el ámbito de recuperación de información. Para ello es necesario introducir cuatro variables: Verdaderos Positivos (VP), Verdaderos Negativos (VN), Falsos Negativos (FN) y Falsos Positivos (FP).

- ❖ **Verdaderos positivos (VP):** Estos son los casos en los que el modelo predijo correctamente que una noticia es verdadera, y en realidad lo son.
- ❖ **Falso negativo (FN):** Estos son los casos en los que el modelo predijo correctamente que una noticia es falsa, y en realidad lo es.

Los dos valores adicionales de la matriz de confusión son los falso positivos y falsos negativos, se obtienen cuando existe una contradicción entre el valor de referencia y el valor obtenido por el modelo.

- ❖ **Falsos positivos (FP):** Estos son los casos en los que el modelo predijo incorrectamente que una noticia es verdadera, pero en realidad es falsa.
- ❖ **Falsos Negativos (FN):** Estos son los casos en los que el modelo predijo incorrectamente que una noticia es falsa, pero en realidad es verdadera.

Para validar que los datos obtenidos con el sistema de clasificación son correctos se utilizó la matriz de confusión y sus métricas correspondiente, la cual es una herramienta utilizada en el análisis de clasificación de datos para evaluar el rendimiento de un modelo predictivo

En la Figura 33 se muestra la matriz que resume toda esta información.

		Predicción	
		Positivos	Negativos
Observación	Positivos	Verdaderos Positivos (VP)	Falsos Negativos (FN)
	Negativos	Falsos Positivos (FP)	Verdaderos Negativos (VN)

Figura 33 Matriz de confusión

Precisión: Es la proporción de verdaderos positivos sobre la suma de verdaderos positivos y falsos positivos. Esta métrica se utilizó para generar la medida de cuán precisa son las predicciones del sistema en clasificar noticias como falsas o verdaderas (Vea Figura 34) la fórmula de precisión.

$$Precisión = \frac{VP}{VP+FP} \quad (34)$$

*Recall*: Es la proporción de verdaderos positivos sobre la suma de verdaderos positivos y falsos negativos. Se utilizó con el fin de proporcionar una medida de cuán efectivo es el sistema para capturar y detectar todas las noticias falsas presentes en el conjunto de datos (Vea Figura 35) la fórmula de *Recall*.

$$Recall = \frac{VP}{VP+FN} \quad (35)$$

*F1 Score*: Es la estadística media de la Precisión y Recall. Se utilizó con el propósito de proporcionar una medida equilibrada del rendimiento del sistema (Vea Figura 36) la fórmula de *Recall*.

$$F1 - Score = 2 * \frac{Precisión*Recall}{Precisión+Recall} \quad (36)$$

## 5.2 Resultados

Se realizaron pruebas con 200 noticias con el propósito de evaluar la Gramática desarrollada, contando con un total de 100 noticias falsas y 100 noticias verdaderas. Las noticias fueron recolectadas manualmente de diferentes plataformas y redes sociales.

El proceso de recolección de resultados arrojados por el sistema de VERACIA se almacena de manera manual, dentro de un archivo. XML donde se crearon las siguientes columnas llamadas Noticia, Valor real y Valor esperado. El cual significa:

- ❖ **Noticias**: Se guarda la noticia con la que se está generando la búsqueda.
- ❖ **Valor Real**: Es el veredicto real donde estipula si la noticia es falsa o verdadera.
- ❖ **Valor Esperado**: Es el veredicto que actualmente está arrojando el sistema.

A continuación, se presentan los resultados arrojados durante la prueba realizada identificadas a través de la matriz de confusión en la Tabla 11.

*Tabla 11 Resultados de matriz de confusión*

	<b>POSITIVO</b>	<b>NEGATIVO</b>
<b>POSITIVO</b>	VP=85	FP=8
<b>NEGATIVO</b>	FN=12	VN=95

Los resultados del cálculo de las métricas antes mencionadas, se muestran en la Tabla 12.

*Tabla 12 Resultados de las métricas de evaluación utilizadas*

Resultados de las métricas de evaluación	
Métricas	Resultado
Precisión	0.91
Cobertura	0.87
F1-Score	0.93

### **5.3 Experimentación.**

Comparación de Resultados con métodos del Estado del Arte

Para este experimento se realizó la comparación de los resultados obtenidos con trabajos del estado del arte respecto a las métricas: *Precisión*, *Recall*, *Accuracy* y *F1- Score*. De acuerdo a la Tabla 13.

*Tabla 13 Comparación de métodos utilizados para la identificación de noticias falsas*

Metodología	BD	Precisión	Recall	Accuracy	F1-Score
CNN(Redes convolucionales(Pablo A et al., 2021))	Propia	88.78%	n. a	99.90%	n. a
Método basado en aprendizaje automático, características basadas en estilo	LIAR	76.27%	74.24%	n. a	76%
Modelo basado en estilo y las palabras más frecuentes del corpus(Jennifer et al.,2022)	Propia	80.85%	80.85%	80.85%	80.83%
Método de probabilidad en base a un historial de publicaciones (Luis et al., 2022)	Propia	80.85%	80.81%	n.a	81.85%
Método basado en aprendizaje automático(Luis et al.,2022)	GossipCop	0.94%	0.93%	0.94%	0.95%
Método basado en reglas de producción (esta propuesta)	Propia	0.91%	0.87%	n. a	0.93%

De acuerdo a los resultados de la tabla anterior se puede ver como el método de solución basado en reglas de producción comparado con el estado del arte, se obtuvo el siguiente resultado con una precisión 0.91%, Recall 0.87% y F1-Score de 0.93% identificando que el método que más se acercó fue el que aplicó redes convolucionales con una precisión de 88.78% y el método que ha superado es el método propuesto es el método basado en aprendizaje automático superando nuestros resultados.

# Capítulo VI

Conclusión



## **6. Conclusiones**

En este último capítulo, después de haber realizado la evaluación correspondiente, posteriormente se genera un análisis de las conclusiones de la investigación a partir de los resultados obtenidos reconociendo los objetivos planteados. Describiendo los productos generados y las aportaciones. Reconociendo por último los trabajos a futuro.

### **6.1 Objetivos y Alcances Logrados**

A partir de este tema de investigación propuesto en este trabajo de tesis Gramática libre de contexto para la identificación de noticias falsas sobre enfermedades basada en el análisis de la literatura médica. A continuación, se detallan los logros de los objetivos en la Tabla 14 y los alcances esperados de este tema en la Tabla 15.

Tabla 14 Objetivos cumplido

Objetivo	Actividad
Identificar y recolectar noticias falsas y notas que traten temas de salud.	Se realizó el proceso de búsqueda y recolección de manera manual en las diferentes plataformas y redes sociales, identificando noticias falsas que tengan obvias razones para reconocerlas como falsas con tal motivo de no necesitar un experto en el área.
Analizar manualmente el contenido de noticias falsas y notas médicas para detectar patrones de uso del lenguaje.	Se realizó el análisis manualmente de cada una de las noticias recolectadas, identificando la sintaxis de la noticia con el motivo de categorizar cada oración de cada una de las noticias recolectadas
Determinar la confirmación o contradicción de hechos contenidos en las noticias a través de una gramática libre de contexto.	Se definieron las reglas de producción a partir del estudio de Noam Chomsky también llamada Gramática libre de contexto. El cual se utiliza para identificar los hechos que existan en la noticia y las notas recolectadas sin importar el orden estructural que exista con el fin de generar un veredicto
Proponer un proceso automático de recuperación de notas relacionadas con noticias falsas.	Se aplicó un algoritmo de recuperación de información para filtrar las notas que no tienen similitud con la búsqueda generada, recuperando las notas que si tienen cierta semejanza entre la búsqueda.
Evaluar las reglas definidas de acuerdo con la precisión y cobertura que muestran en la detección de noticias falsas.	Se evaluó la Gramática con las métricas propuestas y se hizo una comparación de los resultados con 7 metodologías del estado del arte

Tabla 15 Alcances realizados

<b>Alcances</b>	<b>Actividad</b>
El método tomará como fuente confiable más de una fuente médica para obtener información relacionada con la noticia a analizar	Se generó la búsqueda de manera manual de fuentes de información médica con el motivo de contener más de una fuente de información que se pueda acceder para generar la búsqueda contando con 4 fuentes.
Se identificarán los fragmentos de información del contenido de las noticias como falsas o verdaderas.	Se desarrolló un sistema web llamado VERACIA el cual se utiliza para identificar a través de una búsqueda de una noticia, compararla con las notas recolectadas a partir de la búsqueda y generar un veredicto para comprobar si es falsa o verdadera la noticia identificada.

## 6.2 Resultados del trabajo de investigación

### 6.2.1 Productos

A continuación, se presentan los entregables como producto final los siguientes:

- ❖ Reporte del estado del arte.
- ❖ Implementación del sistema propuesto.
- ❖ Presentación de Resultados.
- ❖ Reporte de resultados.
- ❖ Artículo de congreso.
- ❖ Constancia de ponencia.
- ❖ Documento de tesis.

### 6.3 Conclusiones

- ❖ En este trabajo de tesis se generó un corpus lingüístico conformado por 300 noticias pertenecientes a noticias extraídas y recolectadas de redes sociales y plataformas digitales, el cual se desarrolló un agramatical con el propósito de identificar hechos dentro de las noticias, desarrollando una interfaz web para poder identificar y asignar un veredicto de si una noticia es falsa o verdadera.
- ❖ Se generó un bot, el cual es un programa informático diseñado para realizar diferentes tareas y repetitivas en internet tales como dentro de la interfaz de forma independiente del flujo del sistema, para identificar noticias falsas a través de una fuente de información que están destacadas como más populares o las más recientes publicaciones.
- ❖ En general, los resultados obtenidos sugieren que el modelo tiene un buen rendimiento en la identificación de las instancias positivas, con una alta proporción de instancias positivas correctamente clasificadas.
- ❖ Se identificó que las noticias que clasificó erróneamente el sistema se debe a que las reglas gramaticales implementadas desconocían la estructura de la oración analizada en la noticia. La solución a este hecho consiste en analizar más noticias para descubrir nuevas reglas gramaticales.

## 6.4 Trabajos a Futuro

- ❖ Al analizar los resultados obtenidos a partir de las pruebas que se realizaron, se considera que los siguientes trabajos futuros pueden complementar el desarrollo y la investigación que se llevaron a cabo en este proyecto de tesis
- ❖ Aumentar el corpus no solo de noticias falsas de temas sobre la salud sino de todo tipo de noticias falsas, Extraer fuentes de información de diferentes temas de las cuales se pueda comparar las noticias falsas de otras áreas e identificar su veracidad.
- ❖ Aumentar la gramática con nuevas palabras y nuevas categorías que ayuden a una mejora para la identificación de hechos, a través de nuevas noticias falsas sobre otros temas de investigación.

## Referencias

- [1] D. F. E. C. y. A. D. Q. Cortés, Modelo prototipo de inteligencia artificial basado en procesamiento de lenguaje natural y redes de neuronales artificiales para la detección de noticias falsas en español, Bogotá, Colombia, 2021.
- [2] E. A. F. Martínez, Gramáticas libres de contexto, p. 18.
- [3] M. I. O. Velázquez, Procesamiento de Lenguaje Natural para el Análisis de la Escritura Emocional Autorreflexiva (EEA) en mujeres víctimas de violencia, Cuernavaca, Morelos,, 2022.
- [4] N. F. & S. R. Sneha Singhanía, 3HAN: A Deep Neural Network for Fake News Detection, 2017.
- [5] A. J. V. V. P. a. R. Z. XINYI ZHOU, Fake News Early Detection: An Interdisciplinary Study, 2020.
- [6] A. D. N. VALERIA, MAPEO SISTEMÁTICO SOBRE LOS MÉTODOS, TÉCNICAS Y TECNOLOGÍAS ORIENTADAS A LA DETECCIÓN DE NOTICIAS FALSAS, 2021.
- [7] I. B. D. M. S. A. I. Kadek Sastrawan, Detection of fake news using deep learning CNN–RNN based methods, 2022.
- [8] L. Gutiérrez-Coba, P. Coba-Gutiérrez y J. A. Gómez-Díaz, La Noticias falsas y desinformación sobre el Covid-19 análisis comparativo de seis países iberoamericanos, 2020.
- [9] H. G. A. G. S. J. M. Juan Pablo Francisco Posadas Durán, Detection of fake news corpus for the Spanish language, 2019.
- [10] P. S. M. R. B. M. F. Giovanni Luca Ciampaglia, Computational Fact Checking from Knowledge Networks, 2015.
- [11] NICOLÁS JESÚS MAFLACHECA, IDENTIFICACIÓN AUTOMÁTICA DE NOTICIAS FALSAS EN ESPAÑOL UTILIZANDO TÉCNICAS DE MINERÍA DE DATOS Y PROCESAMIENTO DE LENGUAJE NATURAL, 2021.
- [12] B. N. A. H. S. a. M. R. Z Khanam<sup>1</sup>, Fake News Detection Using Machine Learning Approaches, 2021.
- [13] D. d. N. F. e. R. S. B. e. A. A. y. P. U. B. R. Sistemática, Nathaly Álvarez-Daza<sup>1</sup>, Pablo Pico-Valencia<sup>1, 2</sup>, Juan A. Holgado-Terriza<sup>2</sup>, 2020.
- [14] M. A. L. C. d. C. A. T. María Mercedes PALOMINO GONZALES, La red sanitaria y su participación de las fake news y bulos relacionados con la COVID-19: el caso de Lima Perú, 2020.
- [15] C. G. G. S. T. C. Ignacio Blanco Alfonso, El impacto de las fake news en la investigación en ciencias sociales revisión bibliográfica sistematizada, 2019.

- [16] «MedlinePlus,» [En línea]. Available: <https://medlineplus.gov/spanish/>.
- [17] PubMed. [En línea]. Available: <https://pubmed.ncbi.nlm.nih.gov/>.
- [18] Newtral. [En línea]. Available: <https://www.newtral.es/>.
- [19] «Factual,» [En línea]. Available: <https://factual.afp.com/>.
- [20] «Learn Microsoft,» [En línea]. Available: <https://learn.microsoft.com/es-es/azure/search/index-ranking-similarity>.
- [21] Witiko, «github,» [En línea]. Available: [https://github.com/dorianbrown/rank\\_bm25](https://github.com/dorianbrown/rank_bm25).
- [22] AmenRa, «github,» [En línea]. Available: [https://github.com/AmenRa/retriv/blob/main/docs/sparse\\_retriever.md](https://github.com/AmenRa/retriv/blob/main/docs/sparse_retriever.md).

## Anexo

### Anexo A cambios de interfaz

En esta sección se muestra el módulo de la interfaz nombrado Noticias Populares, donde se almacenaron noticias que en su momento fueron propagadas y divulgadas en las diferentes plataformas y redes sociales las cuales fueron redireccionadas a la interfaz en el módulo ya antes mencionado tratando de almacenar las noticias falsas que probablemente fueron las más populares en un determinado momento. El cambio en el apartado del módulo de Noticias Populares se implementó con el objetivo de lograr almacenar las noticias más populares a través de un bot que se conecte a una fuente de información médica y extraiga las notas más populares directamente desde su sitio web desde el apartado de salud, reflejándose en el módulo ya antes mencionado.

A continuación, se presenta la fuente de información utilizada, con el propósito de que el bot desarrollado logre conectarse y extraiga las notas más populares, ya que se estará actualizando siempre que se quiera ingresar al sistema de VERACIA estará extrayendo las más populares de ese momento.

Afp-Factual: Los fast-checkers de la AFP investigan afirmaciones sospechosas que son virales, tienen impacto en la sociedad y son potencialmente dañinas para el público. Las afirmaciones que verificamos pueden aparecer en diversos formatos como redes sociales, blogs, sitios web, aplicaciones de mensajería y otros foros de la esfera pública.

A continuación, se presenta el diseño del módulo de Noticias Populares implementando el uso de un bot que facilite la identificación de notas publicadas como noticias falsas, siendo las más populares. (Vea Figura 37)



## NOTICIAS POPULARES



### Tomar alcohol mata el coronavirus

Pandemia COVID-19



### El sida se contagia tocando una persona infectada

VIH sida



### ¿Es posible contagiarse de VIH al compartir un rastrillo?

VIH sida



### Usar uñas postizas provoca cáncer

Cosmética, estilo de vida



### Cáncer es deficiencia de vitamina b17

Nutrición, oncología



### Las vacunas causan autismo

Prevención, vacunas



### Hacer crujir los nudillos causa artritis

Mitos, articulaciones



### Sentarse demasiado cerca de la TV empeora la vista

Vista, miopía

Figura 37 Rediseño de vista de noticias populares

## Anexo B Producción científica

Se realizó una presentación “Método para la identificación de noticias falsas sobre enfermedades basado en el análisis de la literatura médica” en el XI Coloquio de Lingüística Computacional CoLiCo 2022 en la ciudad de México, 12 de septiembre de 2023 por la Universidad Nacional Autónoma de México (UNAM). La Figura 38 muestra la constancia de participación.



*Figura 38 Constancia de ponencia*

Se realizó una presentación denominado “Método para la identificación de noticias falsas sobre enfermedades basado en el análisis de la literatura médica” en el 1er Coloquio sobre procesamiento del lenguaje natural en México, en Cuernavaca Morelos el 25 de noviembre de 2022 por el Tecnológico nacional de México. La Figura 39 muestra la constancia de participación.



Figura 39 Constancia de ponencia