

TRABAJO FIN DE MÁSTER

Análisis de textos médicos mediante NLP

Autor:

Jesús Enrique Cartas Rascón

Tutora:

Rocío Romero Zaliz

Resumen

En el ámbito de la medicina se almacena una gran cantidad de información relevante: desde valores numéricos correspondientes a signos vitales hasta texto plano que realiza un especialista para completar un informe. Muchas veces los datos guardados en el historial médico de un paciente, que no tiene una estructura determinada, son ignorados. Este proyecto propone recuperar texto médico sin formato utilizando técnicas de Procesamiento del Lenguaje Natural y modelos generativos de lenguaje, para extraer nuevo conocimiento que pueda utilizarse para complementar la información estructurada y mejorar en la clasificación y tratamiento de los pacientes, así como para generar potenciales conjuntos de datos infinitos de forma aleatoria para los desarrolladores.

Abstract

Medical reports hold very important information, ranging from numerical and categorical values to plain text written by the professional in charge. These excerpts are often ignored because of their unstructured nature. This project aims to focus on these pieces of text, using Natural Language Processing techniques, as well as generative language models, in order to recover the information and knowledge they hold, and present them in an easier to interpret, easier to index way. This process ultimately will create new structured data that will complement and enrich the original report, allowing for better patient treatment and data management, as well as generating potentially infinite random datasets for developers to use.

ÍNDICE GENERAL

1. Introducción	7
1.1. Documentos médicos	7
1.2. Datos delicados y escasos	8
1.3. Objetivos	9
1.4. Esquema general del proyecto	9
2. Fundamentos de la minería de texto	11
2.1. Minería de datos	11
2.2. Minería de texto	12
2.2.1. Términos	12
2.2.2. Técnicas	16
2.3. Estado del arte	18
2.3.1. Terminología médica	19
2.4. Reconocimiento de entidades en cuerpos de texto	20
2.5. Transformers	21
2.5.1. Redes neuronales	21
2.5.2. Redes Neuronales Recurrentes	22
2.5.3. Arquitectura	23
2.5.4. Atención	24
2.5.5. Modelos basados en transformers	25

3. Metodología	27
3.1. Datos de entrada y preprocesamiento	27
3.1.1. Codificación de los tokens	28
3.2. Ajuste de los pesos	29
3.2.1. Guardado del estado del modelo	30
3.3. Generación de comentarios	30
3.3.1. Carga del modelo	31
3.3.2. Sampling y otros métodos de generación de texto	31
4. Experimentación: modelos, entrenamiento y generación	37
4.1. Datos: fuente y forma	37
4.1.1. Dataset: <i>Medical Text</i>	37
4.1.2. Dataset: <i>Medical Transcriptions</i>	38
4.2. Preprocesamiento	40
4.2.1. Medical Text	40
4.2.2. Medical Transcriptions	41
4.3. Resultados	43
4.3.1. Métricas y calidad de los comentarios	47
5. Análisis de textos médicos	51
5.1. Despliegue de MEDGEN	51
5.2. Evaluación de las herramientas	52
5.3. Instalación y modo de uso	53
6. Conclusiones	54

6.1. Problemática inicial	54
6.2. Objetivos propuestos	54
6.3. Metodología y resultados	55
6.4. Posibles trabajos futuros	56
A. Apénice	57
A.1. Comentarios	57
Bibliografía	64

ÍNDICE DE FIGURAS

1.1. Diagrama de flujo general de todo el proyecto	10
2.1. Ejemplo del Med7 en acción. Vemos cómo las entidades se han reconocido en diferentes colores, haciendo la extracción de información mucho más eficiente. <i>bid</i> viene del latín <i>bis in die</i> , que se traduce por dos veces al día. <i>PO</i> viene de <i>Per os</i> , vía oral. <i>Suspension</i> hace referencia a una disolución en agua.	21
2.2. Ilustración de una red neuronal clásica, donde se pueden apreciar las neuronas, interconectadas.	22
2.3. Ilustración de una red neuronal recurrente.	23
2.4. Ilustración de una neurona en una red LSTM	23
2.5. Diagrama de un transformer [extraído de [36]]	24
3.1. Diagrama resumen del preprocesamiento efectuado en los datos.	27
3.2. Ilustración del proceso de codificación – decodificación necesario	28
3.3. Diagrama resumen del ajuste de los pesos del GPT-2	29
3.4. Diagrama resumen de la generación de comentarios	31
3.5. Ilustración de los pasos que daría el algoritmo greedy, tomando siempre la posibilidad más alta. Se muestra en rojo el camino que el algoritmo tomaría en este particular ejemplo, ya que <i>nice</i> supera a los otros dos términos en la probabilidad de ser elegido.	32
3.6. Ilustración demostrando el funcionamiento del <i>beam search</i> , con 2 beams en este caso. Vemos como se toman dos caminos, en línea continua el camino definitivo y en línea discontinua la otra alternativa. La probabilidad conjunta del camino elegido es $0.4 \times 0.9 = 0.36$, en contraste con el camino alternativo cuya probabilidad es de $0.5 \times 0.4 = 0.2$. Se toma un camino que a priori no es tan probable pero que a largo plazo, sí.	33

3.7.	Ilustración de un ejemplo de muestreo. Se toman palabras aleatoriamente, aunque se pondera en función de su probabilidad.	34
3.8.	Ilustración del proceso de cálculo en muestreo con un valor de $k = 4$	35
4.1.	Visualización de la distribución de nuestro conjunto de entrenamiento . . .	39
4.2.	Visualización de la distribución de nuestro conjunto de evaluación	40
4.3.	Visualización del dataset Medical Transcriptions	41
4.4.	Wordcloud de todo el conjunto de datos	43
4.5.	Matriz de comparativa de distancia de los 100 comentarios generados en todas las posibles combinaciones.	50
5.1.	Captura de pantalla de la aplicación elaborada. A la izquierda podemos generar comentarios y a la derecha, analizarlos.	52
5.2.	Captura de pantalla del análisis que ofrece la aplicación acerca de un determinado comentario.	53
A.1.	Pipeline para el procesamiento de los comentarios de Medical Transcriptions	59

1. INTRODUCCIÓN

En este capítulo introduciremos los principales problemas existentes en el contexto de minería de datos en texto médico y propondremos una solución que se desarrollará a lo largo del documento.

1.1. Documentos médicos

Toda atención médica dispone de documentos que recogen toda la información relacionada con el paciente, su historial de enfermedades, así como los recursos a utilizar. Los documentos se suelen organizar en un formato de campo-valor, que relaciona los diferentes datos a guardar con su valor esperado. Por ejemplo, nombre, apellidos, o referencias de códigos que se utilicen, como medicamentos o tratamientos.

En estos documentos suele haber una sección en la que el o la profesional en cuestión describe en texto libre el estado del paciente, así como otros matices que no estuvieran considerados en los campos anteriores. Precisamente por esta razón existe dicha sección en el documento.

La medicina personalizada trata de acercar los tratamientos al paciente lo máximo posible, de forma que los profesionales sanitarios puedan hacer un seguimiento en profundidad de los pacientes sin tener que ello conllevar una enorme carga cognitiva de recordar cada paso en el tratamiento. Esto, entre otras cosas, involucra bases de datos en las que guardar toda esta información, así como el análisis de la taxonomía de dichos documentos. Esto facilita su posterior búsqueda para poder retomar el punto en el que se acabó anteriormente.

Dados estos datos, se pueden efectuar una serie de análisis y minería de datos que nos provean con muchos tipos de información. Podemos encontrar patrones de distintas enfermedades dados los síntomas, llevando un seguimiento de todos los pacientes y acometiendo técnicas

El objetivo es centrar nuestra atención en esas secciones de texto sin formato anteriormente mencionadas, con objeto de obtener la mayor cantidad de información posible y anexarla, ahora con formato, al documento del que provienen, enriqueciendo el informe

y habilitando nuevas claves de búsqueda, así como mejorando el indexado de los documentos.

Para facilitar esta tarea, se han creado una colección de términos fácilmente procesables por un ordenador, siendo uno de los ejemplos más destacables el denominado Systemized Nomenclature of Medicine - Clinical Terms (SNOMED) [1]. Este conjunto de términos fácilmente indexables ayudan en la informatización y el procesamiento automático de los documentos médicos.

1.2. Datos delicados y escasos

Uno de los principales problemas a los que nos enfrentamos es la falta de *datasets* o conjuntos de datos en los que estos documentos estén presentes.

Gran parte de este problema es la delicada naturaleza de estos datos. No estamos hablando de el ancho y el alto de los pétalos de una flor, o de las acciones en bolsa de una determinada empresa. Hablamos de datos profundamente íntimos de personas reales con problemas reales. Por ley, concretamente la Ley Orgánica de Protección de Datos (LOPD) [2], las entidades deben proteger dichos datos y garantizar la privacidad de las personas involucradas.

Una de las posibles soluciones a esto es la anonimización de los datos. Puede parecer una tarea simple pero es muy costosa, proporcional al número de datos que poseamos. Debemos no solo garantizar la anonimidad de los pacientes, sino también la de todos los profesionales involucrados. Esto es información que no aporta nada al modelo, en cualquier caso.

Afortunadamente, ya se ha trabajado en esto y existen bases de datos públicas, precisamente para esto. Uno de nuestros principales objetivos es el de buscar y fusionar tantas fuentes de datos como sea posible, unificarlas y crear una herramienta que aproveche todos los datos disponibles públicamente para generar datos nuevos sintéticos, pero suficientemente realistas.

De ser nuestro proyecto exitoso, podríamos obtener ingentes cantidades de información estructurada de la secciones de texto mencionadas anteriormente, lo que posibilitaría una nueva dimensión en el tratamiento de datos médico y habilitaría a tratamientos mucho más precisos y cercanos al paciente.

1.3. Objetivos

Dado este marco, describiremos en esta sección los objetivos de nuestro trabajo:

- Recopilar todas las fuentes de información públicas que nos provean con datos de comentarios médicos listos para su minería y análisis.
- Evaluar las herramientas que ya existen en el estado del arte tanto para clasificar texto como para generarlo. Haremos una revisión de cómo se utilizan y del rendimiento de dichas herramientas.
- Crear un modelo que sea capaz de generar tantos comentarios médicos como sea necesario. La idea es suplir la carencia de datos con un modelo generativo, de forma que no se tenga que lidiar con aspectos de privacidad o licencia, ya que todos los comentarios serían generados de forma sintética. Si bien los comentarios son sintéticos, deben ser lo suficientemente convincentes como para que la evaluación de las herramientas sea fiel y rigurosa. Esto ofrece una herramienta esencial para los desarrolladores de los sistemas que habilita a un mejor y más fructífero desarrollo, ya que se dispone de una cantidad *idealmente infinita* de comentarios sobre los que testear el modelo.

1.4. Esquema general del proyecto

La Figura 1.1 representa el esquema general del proyecto, como un diagrama de flujo.

En cada uno de los capítulos siguientes ahondaremos en cada uno de las fases, respectivamente, explicando en profundidad en qué consisten, cual es su entrada, su salida y su objetivo.

Diagrama de flujo general del proyecto

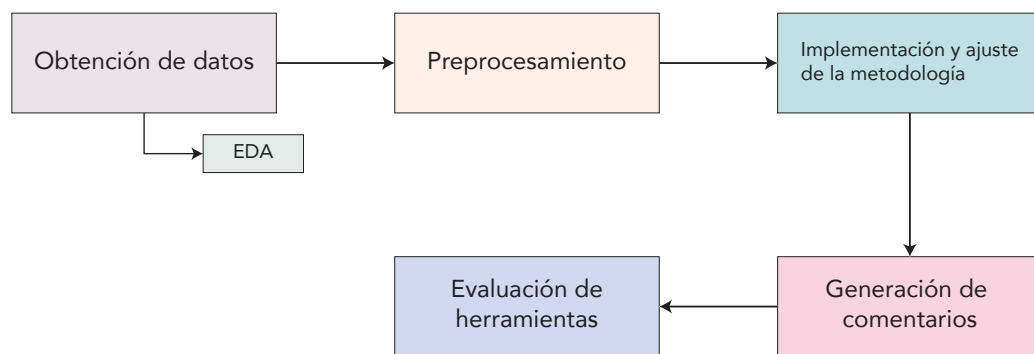


Figura 1.1: Diagrama de flujo general de todo el proyecto

2. FUNDAMENTOS DE LA MINERÍA DE TEXTO

En este capítulo discutiremos algunas de las técnicas y términos más importantes a la hora de hablar de minería de texto, así como minería de datos en general, con objeto de que todas las consideraciones realizadas posteriormente queden claras.

En *Text Mining Applied to Electronic Medical Records: A Literature Review* [3] se hace una revisión de los diferentes aspectos a tener en cuenta durante el procesamiento de textos médicos. Nos apoyaremos en gran medida en la estructura, contenidos y referencias de este artículo, que resume muy bien todo lo que necesitamos saber para resolver nuestro problema.

2.1. Minería de datos

La minería de datos es una rama de la informática que se dedica a encontrar tendencias y patrones en grandes volúmenes de información. Estas tendencias y patrones crean *conocimiento* a partir de los datos, es decir: información estructurada desde los datos no estructurados. Esta información es muy valiosa y contribuye en las decisiones que se vayan a tomar o a monitorizar algunos aspectos que sean de vital importancia para el interesado.

La minería de datos puede dividirse en un número de técnicas que funcionan de forma diferente en función del tipo de datos que tengamos y la información que busquemos.

1. **Asociación:** esta técnica se centra en encontrar relaciones entre las distintas variables de nuestros datos, con objeto de encontrar muestras que sean estadísticamente dependientes. Una de las técnicas más utilizadas son las reglas de asociación, cuya salida tras el cálculo son un conjunto de reglas con antecedentes y consecuentes, muy fácilmente interpretables por cualquier persona, familiarizada o no con la ciencia de datos. [4]
2. **Clasificación:** el proceso de clasificación trata de asignar una categoría a un conjunto de elementos que tengan algún aspecto en común. La clasificación en la minería

de datos es una de las técnicas más utilizadas, ya que la naturaleza de gran parte de los datos responden bien a este método. [5]

3. **Agrupamiento:** también denominado *clustering* trata de agrupar muestras que tengan características similares. A diferencia de la clasificación, aquí no tenemos una etiqueta o categoría a la que asignar las muestras, sino que las agrupamos *a ciegas*, simplemente basándonos en alguna métrica para evaluar la distancia que haya entre un determinado par de muestras. [6]
4. **Predicción:** la predicción nos ayuda a encontrar tendencias entre variables, generalmente en datos con una componente temporal fuerte. [7] Es común poder predecir si un paciente sufrirá una determinada enfermedad conociendo su historial médico, por ejemplo.
5. **Identificación de patrones secuenciales:** Al igual que la predicción, se trabaja sobre datos con una componente temporal marcada. En este caso, se buscan patrones, es decir, conjuntos o cadenas de muestras que aparecen de forma frecuente en un orden concreto.

2.2. Minería de texto

En esta sección, discutiremos los diferentes aspectos a tener en cuenta en la minería de textos en concreto, tras haber abordado el concepto de minería de datos en un ámbito más general.

2.2.1. Términos

Definiremos algunos de los términos más utilizados en esta disciplina, guiándonos principalmente por el trabajo de Kamran Kowsari, *Text Classification Algorithms: A Survey* [8].

Tokens

El término más esencial en minería de textos es *token*. Un token es la mínima unidad en la que dividiremos un cuerpo de texto a la hora de analizarlo. Este elemento suele corresponderse con una palabra, que en el contexto de la mayoría de los idiomas corresponde con un conjunto de letras separado por espacios anterior y posteriormente. Esto da lugar a la creación de *Tokenizers*, algoritmos que toman un cuerpo de texto como una

cadena de caracteres muy larga, y devuelven un vector de palabras. Estos *tokenizers* no han de tomar el espacio en blanco necesariamente ni exclusivamente como criterio divisor, aunque suele ser lo más común. Algunos de los *tokenizers* más famosos son:

- **Tokenizers de palabras**

- **Standard Tokenizer:** El Standard Tokenizer divide el texto en términos siguiendo los límites de las palabras según están definidos en el algoritmo *Unicode Text Segmentation* [9].
- **Letter Tokenizer:** Divide el texto en términos cada vez que encuentra un carácter que no es una letra.
- **Whitespace Tokenizer:** Toma como criterio divisor el espacio en blanco.
- **Language Tokenizer:** Otros tipos de tokenizers adaptados a diferentes idiomas, como el inglés, que es el idioma más estudiado con diferencia, pero también otros idiomas con caracteres y reglas diferentes a aquellos basados en reglas occidentales, como el tailandés o el chino.

- **Tokenizers de palabras parciales**

- **N-Gram Tokenizer:** Este tokenizador incluye un parámetro adicional. Primero divide el texto con alguna de las reglas mencionadas anteriormente, y posteriormente, divide cada término del vector resultante en una ventana deslizable de n elementos, (de ahí *N-Gram*). Por ejemplo: *quick fox* devolvería [qu, ui, ic, ck], [fo, ox], dado un $n = 2$. Estos tokenizers también pueden utilizarse a nivel de párrafo, por lo que se devolverían pares de palabras, algo que puede ser muy útil para el análisis de *dichos* o expresiones.

- **Tokenizers de texto estructurado**

- **Pattern Tokenizer:** este tokenizer utiliza el patrón provisto como parámetro para la división de texto, utilizando expresiones regulares.
- **Simple Pattern Tokenizer:** este tokenizer utiliza el patrón provisto como parámetro para la división de texto, utilizando expresiones optimizadas para el patrón dado, lo que hace que funcione generalmente más rápido pero también será más específico.

En resumen, un tokenizer es un algoritmo que divide el texto provisto siguiendo los criterios definidos por el usuario, devolviendo un vector con los elementos del texto divididos atendiendo a dichos criterios. Es una de las herramientas esenciales en la minería de texto, ya que permite generar la mínima unidad de información a partir de la que se extraerá conocimiento.

Palabras vacías

Las palabras vacías o *stopwords* son términos presentes en un idioma que sirven de apoyo para la formulación de oraciones pero que no poseen información en sí. Nos referimos a los artículos, determinantes, preposiciones, etc.

Estos términos son considerados como *ruido* en el procesamiento de texto, por lo que lo más usual es disponer de un diccionario de términos vacíos y filtrar el texto original, eliminando dichos términos. De esta forma, nos quedamos con las palabras más importantes. Los símbolos de puntuación también se suelen considerar como ruido; si bien son esenciales para la comprensión y estructuración de texto para los humanos, suponen un detrimento para algoritmos de clasificación.

Sin embargo, esta operación es delicada y no siempre ofrecerá buenos resultados. Por ejemplo, si tratamos de inferir la intención de la oración *No me gusta el fútbol* y pasamos previamente un filtro de palabras vacías, el texto resultante sería *gusta fútbol*. Dados estos términos, se infiere que se está opinando de forma positiva acerca del tema *fútbol*, cuando no es así.

Este caso particular está descrito en la literatura como *negation handling*, en trabajos como [10] o [11]. Aún así, hay muchos factores que se deben tener en cuenta antes de eliminar términos de una oración.

Stemming y Lematización

Stemming hace referencia a la gestión de palabras con prefijos o sufijos para su integración en una frase, como plurales (casa, casas). Se trata de eliminar los posibles complementos añadidos con objeto de normalizar las palabras y que todas tengan la misma forma. En este caso también han de tenerse en cuenta las negaciones (típico, atípico).

La lematización va un paso más allá y trata de encontrar la raíz de las palabras, obteniendo una normalización más estricta. Un buen ejemplo son la conjugación de los verbos: de *estudiando*, *estudiante* o *estudio* obtenemos *estudi-*. [12]

Frecuencias: TF, IDF

Uno de los datos más importantes a obtener de un texto es la frecuencia de palabras. Esta operación es tan simple como suena: contar cuántas veces aparece cada palabra y anotarlo en una estructura similar a un diccionario. Este término se conoce como *Term*

Frequency o TF. Estos valores suelen representarse en una escala logarítmica, con objeto de que las palabras muy dominantes no eclipsen a las menos frecuentes.

Del campo de teoría de la información [13] conocemos que aquellos términos que aparezcan con una frecuencia muy alta poseerán menos información que aquellos que aparezcan menos. Como vimos en la Sección 2.2.1, eliminamos las palabras vacías porque aparecían mucho. Es decir, un artículo como *el* o una preposición como *de* tendrían una frecuencia desproporcionada, cuando en realidad no aportan ninguna información.

De forma similar, el valor *Inverse Document Frequency* [14] trata de abarcar esta frecuencia pero en un conjunto de documentos, añadiendo la inversa de la frecuencia por documento. Esta métrica se utiliza mucho en conjunción con la TF, resultando en la TF-IDF, que trata de medir la relevancia de un término en un conjunto de documentos. Esto resulta en el cálculo:

$$W(d, t) = TF(d, t) * \log\left(\frac{N}{df(t)}\right) \quad (2.1)$$

donde d es un documento del conjunto de documentos con cardinalidad N , t es el término en concreto y $df(t)$ es el número de documentos que contienen el término t .

Bolsas de palabras

Conociendo el concepto de frecuencias de palabras, una de las aplicaciones directas son las bolsas de palabras, que recogen en una estructura con forma de diccionario cada término y su frecuencia.

De esta forma, tenemos un *ranking* para cada término. Esto se utiliza extensamente en sistemas de recomendación, en donde una consulta provista por un usuario se compara con la bolsa de palabras del posible conjunto de documentos, y este conjunto va afinándose conforme se van comparando los conjuntos de palabras. El resultado es una búsqueda más refinada que devuelve documentos más relevantes con respecto a la consulta realizada.

Word Clouds

Las nubes de palabras o *Word Clouds* son un tipo de visualización especializada en conjuntos de datos de texto. Consisten en representar el conjunto de palabras más relevantes del texto que analicemos, disponiendo dicho conjunto de forma distribuida por la imagen. Por lo general, se suelen utilizar códigos de color o tamaño para aludir a factores como la frecuencia o la relevancia de dicha palabra en el texto.

Dada una bolsa de palabras, podemos tomar los n términos más presentes y escribirlos en una imagen, haciendo la fuente tanto más grande cuanto más aparezca dicha palabra. Esto ofrece una manera rápida e intuitiva de averiguar el tema del cuerpo del texto y sus términos más relevantes.

2.2.2. Técnicas

En esta sección discutiremos algunas de las técnicas avanzadas más utilizadas en procesamiento y análisis de texto que además utilizaremos en nuestra implementación directa o indirectamente.

Word Embeddings

Esta técnica esencialmente trata de convertir los diferentes términos en vectores de números reales, ya que esto los convierte en objetos matemáticos fáciles de comparar y procesar. Matemáticamente hablando corresponde con una representación de un espacio n -dimensional a un espacio vectorial continuo de menor tamaño, donde n es el número total de términos presentes en todos los documentos.

Esta técnica ha sido estudiada en profundidad en varios proyectos:

- **Word2Vec**: esta técnica trata de representar las palabras como vectores utilizando una red neuronal con dos capas, haciendo uso de una bolsa de palabras continua (CBOW) y el modelo del Skip Gram. [15]
- **GloVe**: acrónimo de *Global Vectors for Word Representation*, es una técnica muy similar a la *Word2Vec*, con la particularidad de estar preentrenada en grandes corpus de texto, basados en Wikipedia y Gigaword. [16]
- **FastText**: es una técnica desarrollada por Facebook. Esta técnica hace uso de la técnica de los n -grams para su entrenamiento, obteniendo una representación de los términos mucho más granular. [17]
- **Contextualized Word Representations**: esta técnica hace uso del contexto de las palabras para tratar de encontrar una representación y relación entre ellas. Esta técnica basa su funcionamiento en el uso de Long-Short Term Memory, un tipo de red neuronal recurrente muy utilizada en procesamiento de texto, en la que ahondaremos más en profundidad en secciones posteriores de este documento. [18]

Reducción de dimensionalidad

La reducción de dimensionalidad es una técnica que permite proyectar el espacio en el que se hallan nuestros datos en un subespacio de menor dimensionalidad, con objeto de facilitar el cálculo de las propiedades de dichos datos sin tener que utilizar todas sus características. Es común encontrar conjuntos de datos con un número de dimensiones muy alto, que hace inviable su estudio.

Las principales técnicas desarrolladas para reducir la dimensionalidad incluyen:

- **Principal Component Analysis (PCA):** PCA o análisis de componentes principales trata de encontrar un subespacio latente que represente a los datos encontrando aquellas variables que estén menos relacionadas y que maximicen la varianza, para conservar la mayor cantidad de variabilidad posible. [19]
- **Independent Component Analysis (ICA):** es una técnica similar que trata de expresar los datos con transformaciones lineales. [20]
- **Linear Discriminant Analysis (LDA):** es otro método muy utilizado cuando los datos son de carácter categórico y no tienen una proporción uniforme intracase. [21]

Toda esta familia de algoritmos resulta muy conveniente para *comprimir* datos de alta dimensionalidad y extraer solo las **características principales** de los mismos. Se suelen usar en etapas de preprocesamiento, donde los datos resultantes se pasan a los algoritmos a entrenar, consiguiendo un mejor resultado y rendimiento en comparación con los datos sin preprocesar.

Clasificación de texto

Por último, abordaremos las principales técnicas para clasificar texto, ámbito importante en nuestro proyecto, así como una breve explicación de las mismas. Al igual que en las secciones anteriores, no es nuestro objetivo estudiarlas en profundidad, pero más bien ofrecer una vista general del panorama en cuanto a esta tecnología con objeto de que el lector se familiarice con los términos.

Los principales algoritmos de clasificación de texto, entre otros, son los siguientes:

- **Naïve Bayes:** estos modelos utilizan el poder de inferencia del cálculo de probabilidades condicionales para la categorización, más en concreto, la categorización

de texto. Esta familia de modelos se ha demostrado que funcionan particularmente bien para esta tarea, además de ser muy rápidos. [22]

- **Regresión logística:** la regresión logística es uno de los métodos de aprendizaje más simples, junto con la regresión lineal. La regresión logística es una especialización de la regresión lineal, de forma que se utiliza una función logística para predecir categorías discretas, no continuas. [23]
- **Máquinas de Soporte Vectorial:** Las Support Vector Machines (SVM) son modelos muy conocidos de aprendizaje supervisado para clasificación y regresión. Aunque originalmente se concibieron como modelos de separación lineal entre las clases a clasificar, se pueden modificar sus *kernels* para obtener clasificadores no lineales. El funcionamiento de estos modelos en clasificación de texto es muy bueno como se sugiere en [24].
- **Redes neuronales recurrentes:** las redes neuronales recurrentes son una especialización de las redes neuronales en las que un subconjunto de las neuronas reciben su salida como una entrada, generándose ciclos de retroalimentación o *feedback*. Estas redes funcionan especialmente bien con datos con patrones y componentes temporales, características especialmente destacables del lenguaje humano, así como de la música o vídeo. [25]. Profundizaremos en este tema en las secciones siguientes.
- **Boosting y bagging:** boosting y bagging son dos métodos basados en lo que se denomina en la literatura como *ensemble learning*. Estos métodos no son algoritmos de clasificación en sí, sino que son métodos aplicables a cualquier modelo. Esta técnica utiliza un gran número de modelos que funcionan muy bien para casos muy específicos pero no generalizan correctamente. La idea es que la respuesta conjunta de todos los modelos provea la respuesta correcta de forma más probable. [26]

Existen otras técnicas en clasificación de texto, como *K Nearest Neighbours* [27], árboles de decisión [28] o *Random Forests* [29], entre muchas otras.

2.3. Estado del arte

Por último, hablaremos de los diferentes lenguajes y vocabularios médicos disponibles en el estado del arte que ayudan a extraer conocimiento de la terminología médica.

Además, mencionaremos brevemente los tres proyectos más grandes de lenguaje generativo hasta la fecha, con objeto de entender qué es lo que hacen para poder integrar dichas técnicas en nuestra metodología.

2.3.1. Terminología médica

Debido al vocabulario y términos tan específicos con los que se tratan en el ámbito de la medicina, si deseamos automatizar cualquier proceso de extracción de conocimiento debemos tener en cuenta este factor. Es por ello que en 1999 se creó el Systemized Nomenclature of Medicine - Clinical Terms (SNOMED) [1].

El SNOMED es una colección de términos sistemáticamente organizada, que provee códigos, términos, sinónimos y definiciones utilizadas en el ámbito de la medicina. Podemos ver un pequeño ejemplo de las estadísticas que nos ofrece la web de SNOMED en la Tabla 2.1. Se muestran los códigos generales, el nombre colectivo de todos los términos que caen bajo dicha categoría y cuántos términos exactamente pertenecen a dicha categoría. Esta tabla es solo un extracto, se pueden consultar el resto de términos en la [web del SNOMED](#).

Semantic Type ID	Semantic Type Name	Count
T047	Disease or Syndrome	42225
T033	Finding	41837
T061	Therapeutic or Preventive Procedure	41496
T037	Injury or Poisoning	27635
T023	Body Part, Organ, or Organ Component	25149
T121	Pharmacologic Substance	19557
T074	Medical Device	14109

Cuadro 2.1: Ejemplo de los tipos de entidades más comunes en el SNOMED junto con el número de términos en total.

A raíz del SNOMED aparecen vocabularios más específicos que matizan aspectos que pudieran no haberse considerado en el sistema inicial. Todas estas colecciones se aglomeran en el Unified Medical Language System (UMLS) [30], que incluye al SNOMED.

Estos vocabularios son de vital importancia para modelos de lenguaje médico, ya que proveen un sistema rápido y accesible para el etiquetado y clasificación, resumen o extracción de palabras clave de extractos médicos o similares. Los **reconocedores de entidades** hacen un especial uso de este sistema, como veremos en capítulos posteriores.

Como vemos, el SNOMED trata de unificar y sistematizar el lenguaje médico, por lo que la cuestión de cómo esto afecta al personal sanitario real empieza a ser relevante. Como en muchas otras disciplinas, se habla de que las *máquinas sustituyen a las personas*, cuando el trabajo en cuestión es altamente automatizable. Esto es solo parcialmente cierto, ya que sí se requiere de un personal cualificado capaz de comprender y utilizar este sistema y más aún si queremos extraer conocimiento automáticamente del mismo. Se requiere

de personal más especializado y formado, pero este sistema no va a sustituir a dichos expertos, al menos en los años venideros.

2.4. Reconocimiento de entidades en cuerpos de texto

Para el análisis y minería de texto no estructurado una de las técnicas más útiles es el **Named Entity Recognition (NER) Tagging**, es decir, el etiquetado de entidades. La idea es, dado un cuerpo de texto, extraer aquellas entidades que tengan un significado propio o relevante por sí mismas. En nuestro contexto, esto corresponde, por ejemplo, a nombres de enfermedades, partes del cuerpo, medicamentos, entre otras. Un NER Tagger entrenado sería capaz de detectar dichas entidades en el texto y devolver una lista de las mismas, automatizando la extracción de información.

La librería SciSpacy [31] nos provee con una selección de NER Taggers preentrenados en diferentes corpus de carácter biomédico y los hace disponibles mediante la famosa librería Spacy [32], que se ha convertido en el estándar de la industria para el análisis del lenguaje natural.

Estos son los taggers que utilizaremos en nuestra aplicación, aunque resulta sencillo incluir algún otro más que pudiera ser de interés.

- **Med7** [33]: El Med7 es un NER entrenado en datos clínicos capaz de detectar 7 entidades diferentes: nombres de medicamentos (Aspirina, Advil), vía de administración del medicamento (oral, intravenosa, respiratoria), frecuencia (cada 8 horas), dosis (mg o ml), cantidad (número de pastillas), formato (pastilla, polvos, etc), y tiempo total de medicación (semanas, meses).
- **BC5CDR** [34]: Este tagger recibe su nombre del corpus en el que fue entrenado. Dicho corpus consta de 1500 artículos *PubMed* con 4409 sustancias químicas etiquetadas, 5818 enfermedades y 3116 interacciones enfermedad-medicamento.
- **BIONLP13CG** [35]: De igual forma, este tagger recibe su nombre del corpus que lo origina. Este tagger se especializa más en términos genéticos, cánceres, órganos y compuestos químicos como aminoácidos o proteínas, ligeramente mejor adecuado para prescripciones quirúrgicas.

En la Figura 2.1 vemos un ejemplo de cómo una de estas herramientas puede ayudar al análisis del texto, destacando con diferentes colores las entidades encontradas. Esto

A patient was prescribed **Magnesium hydroxide** **DRUG** **400mg/5ml** **STRENGTH** **suspension** **FORM**
PO **ROUTE** of total **30ml** **DOSAGE** **bid** **FREQUENCY** **for the next 5 days** **DURATION** .

Figura 2.1: Ejemplo del Med7 en acción. Vemos cómo las entidades se han reconocido en diferentes colores, haciendo la extracción de información mucho más eficiente. *bid* viene del latín *bis in die*, que se traduce por dos veces al día. *PO* viene de *Per os*, vía oral. *Suspension* hace referencia a una disolución en agua.

es, sin embargo, solo una aplicación de demostración. El verdadero poder reside en que dichas entidades están internamente representadas en forma de diccionario, asociando cada término encontrado en el texto con la entidad correspondiente.

Estos diccionarios pueden anexarse y guardarse en la base de datos correspondiente, haciendo su búsqueda y análisis muchísimo más rápidos, además de habilitando comparativas que no serían posibles de otra forma.

2.5. Transformers

En esta sección hablaremos del Transformer, un tipo de arquitectura de red neuronal creada por ingenieros de Google [36], en más profundidad, sus características principales en contraste con las demás tecnologías y de la configuración escogida para nuestro problema.

2.5.1. Redes neuronales

Para entender lo que es un Transformer, debemos primero comprender lo que es una red neuronal y lo que es una red neuronal recurrente.

Una red neuronal es un modelo algorítmico y matemático que trata de modelar el funcionamiento del cerebro en tareas de aprendizaje, de forma que el objetivo principal de estos modelos es ser capaz de aprender automáticamente la solución a la tarea en cuestión simplemente exponiéndose a ejemplos de dicha tarea.

Estos modelos se estructuran de forma muy similar a la de un cerebro real, constando de **neuronas** unidas y agrupadas en capas que desempeñan diferentes roles en esta tarea de aprendizaje. Podemos apreciar un ejemplo de una red neuronal en la Figura 2.2.

Aunque a primera vista parezca ser una tecnología muy nueva, ya que es en los últimos años cuando más hemos visto la integración de estos sistemas en problemas reales, es un campo que lleva en desarrollo desde los años 40. Podemos atribuir los inicios más determinantes a McCulloch y Pitts, que crearon el primer modelo computacional de una

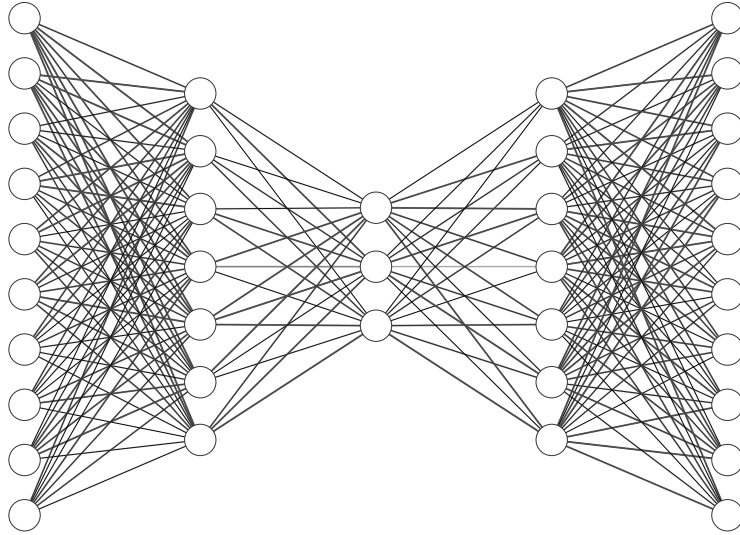


Figura 2.2: Ilustración de una red neuronal clásica, donde se pueden apreciar las neuronas, interconectadas.

red neuronal en [37], y a Donald Hebb, considerado el padre de las redes neuronales tras desarrollar el concepto de *aprendizaje hebbiano* [38], uno de las claves para el desarrollo de estos sistemas.

2.5.2. Redes Neuronales Recurrentes

El auge en la investigación de este campo llevó a crear diferentes categorías de redes neuronales. Dichas categorías se correspondían con la arquitectura de la misma, que, en gran medida, define el propósito para el que estaba diseñada. Esto dio lugar a las redes neuronales convolucionales, cuyo principal elemento era la capa de convolución, que se utiliza ampliamente en tareas de procesamiento de imágenes, o los autoencoders, que se utilizan principalmente en tareas de detección de anomalías.

Una de estas arquitecturas es la red neuronal recurrente, cuya particularidad es que las distintas neuronas no solo están conectadas con neuronas adyacentes, sino consigo mismas. Es decir, su salida depende directamente de su entrada. De esta forma, estos modelos se suelen concebir con una componente temporal, y siempre hablamos de la salida en el momento n , que dependerá de la entrada en el momento n y de la salida en el momento $n - 1$. Podemos ver este comportamiento en la Figura 2.3, donde la neurona inicial se *desenrolla* en el tiempo para visualizar las salidas y las entradas en cada instante.

Esta arquitectura cíclica, recurrente, convierte a esta familia de redes en sistemas muy potentes de procesamiento del lenguaje, vídeo o música, tipos de datos que tienen una componente temporal muy importante. Los modelos más relevantes son las redes

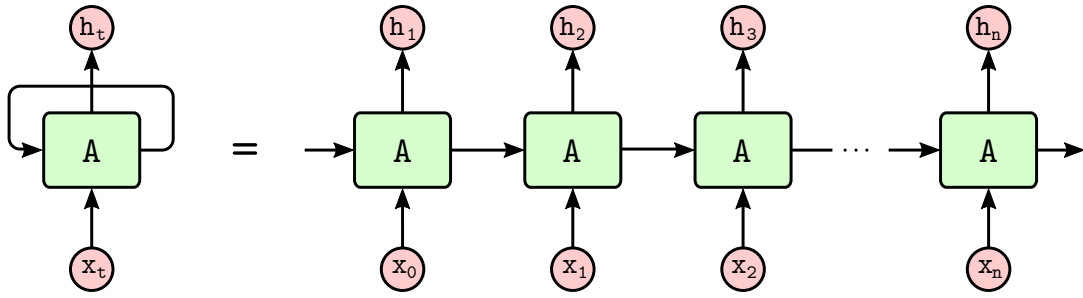


Figura 2.3: Ilustración de una red neuronal recurrente.

recurrentes (RNN), las Long-Short Term Memory (LSTM), ambas estudiadas en [39] y, finalmente, los muy nuevos Transformers.

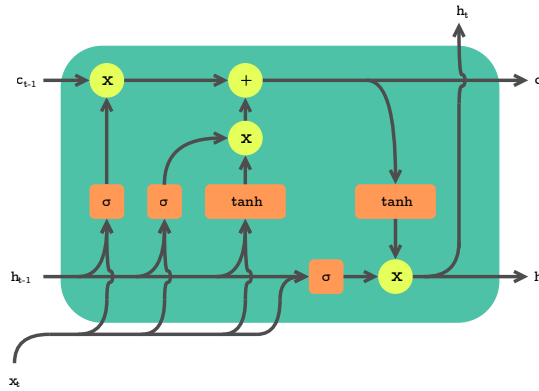


Figura 2.4: Ilustración de una neurona en una red LSTM

2.5.3. Arquitectura

El Transformer corresponde con una de las técnicas más modernas de procesamiento de lenguaje natural, y una de las más relevantes del momento. Se la presenta muchas publicaciones como *estado del arte en procesamiento del lenguaje natural*, constituyendo un modelo mucho más potente y **considerado** que las anteriores LSTM.

El transformer está principalmente basado en redes recurrentes, redes cuya entrada está conectada a la salida y, generalmente, se las construye con un contexto temporal en mente. Esto quiere decir que dada una entrada en el momento n , predeciremos la salida en función de la salida que se obtuvo en el momento $n - 1$, además de la entrada del momento n en sí.

Esto es en sí el funcionamiento de una red recurrente estándar. El transformer añade varios conceptos clave a su implementación que lo hacen particularmente poderoso que serán detallados a continuación.

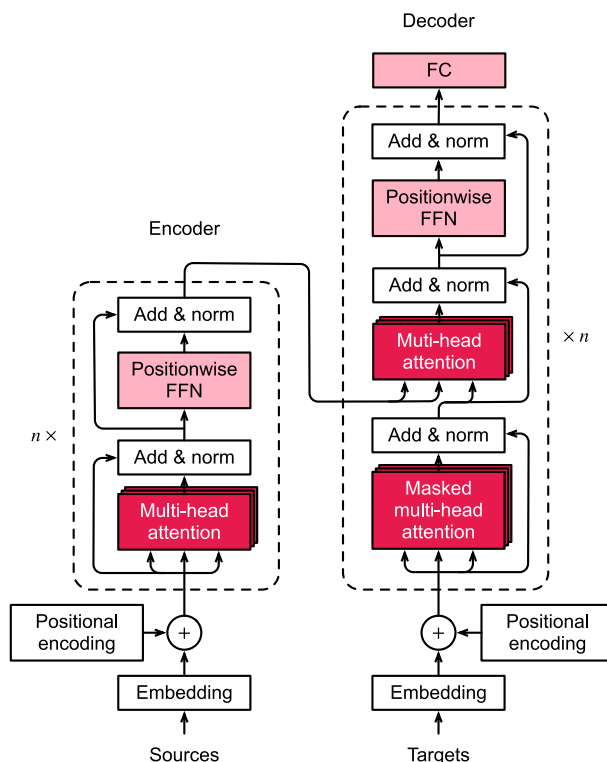


Figura 2.5: Diagrama de un transformer [extraído de [36]]

2.5.4. Atención

Uno de los principales conceptos que añade el transformer es de la **atención**. Este concepto trata de emular el concepto de atención con el que todos estamos familiarizados: el de ponderar y distribuir los recursos cognitivos de forma que aquellos estímulos más importantes reciban más recursos de cómputo por parte del sistema, de la misma forma que nuestro cerebro procesa de forma más potente aquello en lo que estamos concentrados, ignorando aquellos estímulos que sean menos importantes en un determinado instante.

Uno de los principales problemas de las redes neuronales recurrentes clásicas es la imposibilidad de la paralelización de su ejecución, y el crecimiento de los términos a considerar conforme más avanzado es el análisis de la secuencia en concreto.

La técnica de *multi-head attention*, que podemos apreciar en color en la Figura 2.5, soluciona estos problemas. En primer lugar, hace la paralelización del modelo no solo posible, sino también muy fácil. Cada módulo calcula por separado la atención que le corresponda y posteriormente se concatenan y se transforman linealmente en la salida de la dimensión que se espera.

Por otro lado, los modelos tradicionales sufren de no considerar dependencias entre elementos si estos están distantes en la secuencia. Los diferentes módulos calculan la aten-

ción en diferentes puntos de la secuencia de forma paralela. Las diferentes zonas pueden efectivamente considerar zonas muy distantes en tiempo constante, gracias a la paralelización. Esto efectivamente soluciona el problema de olvidar características importantes de puntos distantes, así como evitando tener que crear caminos de cálculo cada vez más largos conforme avanzamos en el análisis de la secuencia.

2.5.5. Modelos basados en transformers

Habiendo descrito brevemente la arquitectura de un transformer en general, veamos qué se ha desarrollado con esta nueva tecnología.

GPT

GPT son las siglas de *Generative Pre-trained Transformer* [40]. Este modelo trata de abarcar los problemas que otros transformers solían tener, como es la habilidad de un humano para continuar una tarea lingüística dado muy poco contexto o instrucciones. Escalando los modelos se logra mejorar significativamente la generalización del mismo y se descubre que no es necesario afinar los parámetros de los modelos, sino que esto puede corresponderse más con un problema de meta-aprendizaje.

Existen varias versiones del GPT, como GPT-2 o GPT-3, así como varios tamaños dentro de cada versión. En el caso del GPT-3, se han creado versiones que empiezan con 125 millones de parámetros entrenables en la versión que denominan pequeña, hasta 175 mil millones, que es la versión que los investigadores denominan como “GPT-3” a secas. Jared Kaplan sugiere en [41] que la función de pérdida en la etapa de validación debería corresponderse con la ley de potencia, (también conocida como el principio de Pareto) en función del tamaño de dicha red, de ahí que se probaran distintas configuraciones y tamaños.

BERT

BERT son las siglas en inglés de *Bidirectional Encoder Representations from Transformers* [42]. Es un modelo de lenguaje generativo basado en un encoder bidireccional como su nombre indica. Este modelo fue uno de los primeros en realmente conseguir una fluidez comunicativa convincente. Su arquitectura preentrenada permite, con una sola capa extra, crear modelos para tareas en casi cualquier ámbito, como responder preguntas o inferencia del lenguaje, sin necesidad de modificar particularmente su arquitectura interna.

LaMDA

LaMDA corresponde con las siglas de *Language Model for Dialogue Applications* [43], un proyecto de Google que compite de forma directa con los modelos antes mencionados. Este proyecto se creó con una aplicación en concreto: un *chatbot* automático lo más natural posible, al que llamaron *Meena*. Según las estadísticas de Google, Meena tiene casi el doble de capacidad de predicción e inferencia que el antiguo GPT-2, y se entrenó en 8 veces más datos. Es el buque insignia de la empresa.

Para nuestro proyecto, hemos escogido el modelo GPT-2, de OpenAI. Concretamente usaremos la versión pequeña, que consta de 117M de parámetros. Este modelo fue entrenado en una base de datos de texto tomada de Wikipedia, constando de más de 60GB de datos de texto. Utilizaremos el modelo preentrenado, que ya de por sí es capaz de hablar de forma genérica, y lo entrenaremos con nuestra base de datos para que aprenda a hablar tal y como lo haría un médico.

3. METODOLOGÍA

En este capítulo ahondaremos en los conceptos más importantes que hemos de entender de cara al funcionamiento de nuestro modelo generativo, tanto en el entrenamiento como en la generación de comentarios.

3.1. Datos de entrada y preprocesamiento

En primer lugar, debemos ofrecer los datos de entrada al modelo. Estos datos deben estar previamente formateados y unificados. Veamos cómo se ha acometido dicho preprocesamiento.

En la sección 4.2 expondremos los diferentes métodos y técnicas utilizadas de preprocesamiento de texto para darle una forma apropiada y unificada al conjunto de datos, antes de poder introducirlos en el modelo.

En la Figura 3.1 se ilustra un esquema de los pasos a seguir en dicha fase del proyecto. Explicaremos cada paso desde el punto de vista de los dos conjuntos de datos a unificar, ya que cada uno necesitará de un tratamiento particular en función de su estado inicial.

Diagrama de flujo: Preprocesamiento

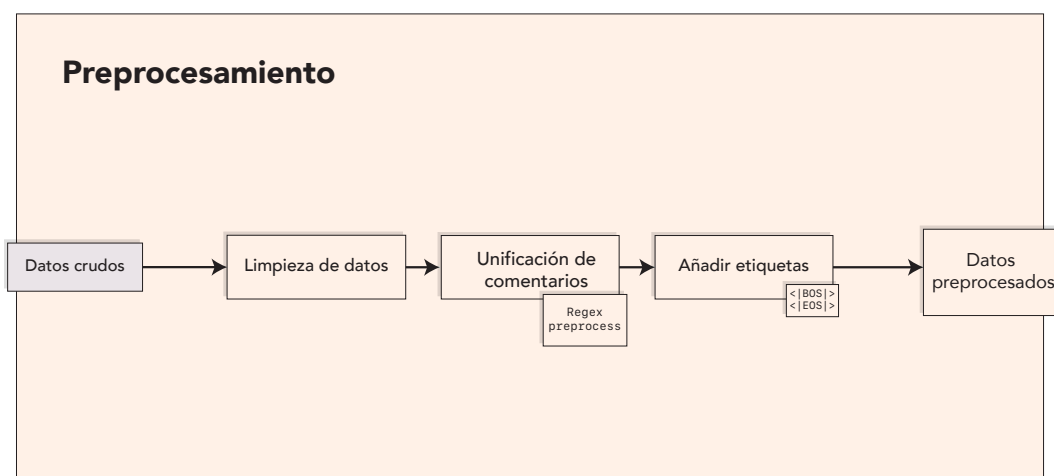


Figura 3.1: Diagrama resumen del preprocesamiento efectuado en los datos.

Es importante la adición de las etiquetas `<|BOS|>` y `<|EOS|>`. Estas etiquetas son las siglas de *begin of sentence* y *end of sentence*. Son marcadores que indican el inicio y fin de una oración. Estas etiquetas son necesarias para que el modelo entienda cuál es el criterio a partir del cual empieza y acaba una oración.

3.1.1. Codificación de los tokens

Antes de pasar los datos al modelo, debemos codificarlos. Codificar significa obtener un código único en función del token que estemos considerando. La codificación es esencial para estos modelos, ya que ofrece una representación matemática de las palabras. La codificación de los tokens se efectúa mediante un diccionario preentrenado, de forma que la codificación y decodificación sean consistentes, dado un modelo determinado.

Este proceso es relativamente simple: para cada palabra que encontremos en el texto, le asociamos un índice y guardamos la entrada en un diccionario. De esta forma, cada vez que queramos que la red procese un extracto de texto, sustituimos cada término por su correspondiente número en el diccionario. De esta tarea también se suele encargar el Tokenizer.

Como vimos en la Sección 2.2.1, estos algoritmos solo se encargan de transformar el texto en vectores, pero al ser una tarea perteneciente a la fase de preprocesamiento, las librerías suelen ofrecer esta funcionalidad de forma conjunta. Esto también garantiza que el tokenizer utilizado para un determinado modelo sea el mismo de una ejecución a otra, ya que de otra forma habría problemas con la codificación.

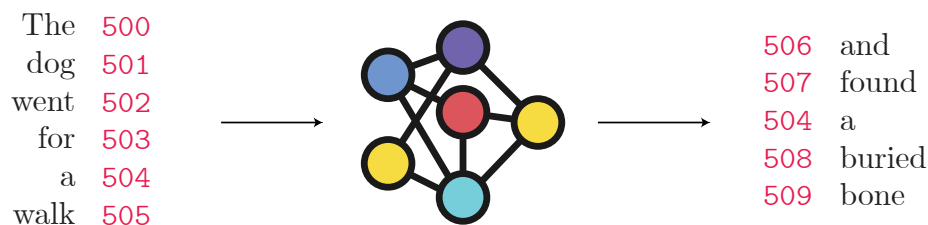


Figura 3.2: Ilustración del proceso de codificación – decodificación necesario

Una vez procesado, el modelo nos devolverá una secuencia de números, de los que podemos volver a obtener el texto subyacente deshaciendo la operación.

3.2. Ajuste de los pesos

Diagrama de flujo: Ajuste de la red neuronal

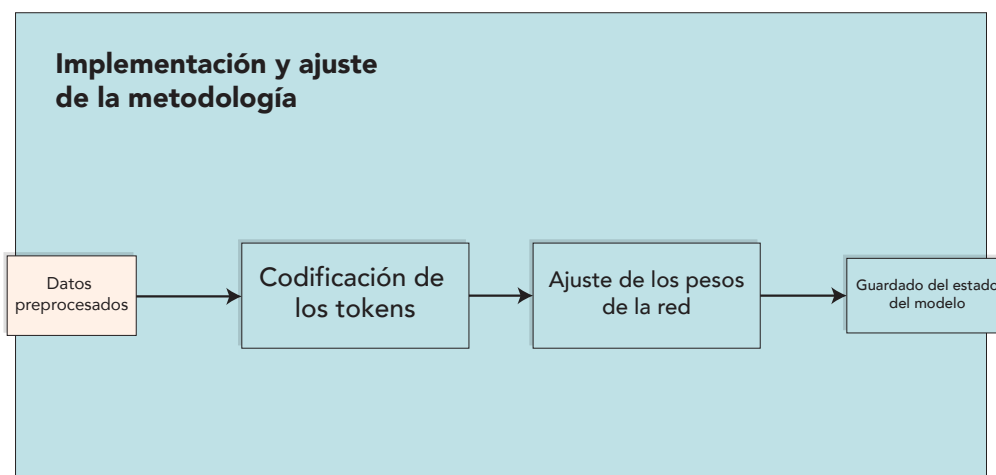


Figura 3.3: Diagrama resumen del ajuste de los pesos del GPT-2

Nuestro modelo ya está construido. Esto es, los diseñadores ya han decidido las distintas capas y el orden de estas según hemos visto en las secciones anteriores. Este modelo se nos ofrece preentrenado de forma bastante simple. Tal y como viene en el paquete, es capaz de generar frases con relativo sentido, es decir, de algún modo *sabe hablar*.

El ajuste se efectúa como un proceso de aprendizaje no supervisado, en la que exponemos a la red a un conjunto de comentarios de forma que ésta podrá generar nuevos comentarios que caigan bajo la función de distribución de los comentarios de entrada.

En nuestro caso, esta función de distribución la define nuestro conjunto de datos de informes médicos. Este proceso efectivamente ajusta los pesos de la red de forma que se familiarice al modelo con el vocabulario y expresiones comunes encontradas en los extractos médicos presentados.

Este proceso es, computacionalmente, extraordinariamente costoso. Debemos calcular el peso de millones de parámetros (117 millones en nuestro caso particular), debido a la magnitud del modelo. Para ello, hemos de tener disponible un equipo con, como mínimo, una tarjeta gráfica decente que nos permita hacer cálculos matriciales en paralelo, operaciones muy comunes en el entrenamiento de las redes neuronales.

Dichos equipos pueden ser muy caros. Para ello, hicimos uso del servicio de clústeres del Instituto DaSCI. En dicho servicio se nos puede asignar una máquina dependiendo

de la carga que tengan las demás. Por referencia, **Selene**, uno de los equipos disponibles, consta de:

- **Procesador:** Dual 20-Core Intel Xeon E5-2698 v4 2.2 GHz
- **RAM:** 512 GB 2,133 MHz DDR4 RAM
- **Disco:** 4×1.92 TB SSD RAID 0
- **Gráfica:** $8 \times$ GPU NVIDIA Tesla V100 32GB

Mediante `ssh`, nos conectamos a la máquina que el equipo nos haya asignado y enviamos nuestros archivos usando `scp`. Una vez nuestros archivos están en la máquina de destino, entrenamos el modelo, obtenemos los pesos y los descargamos de vuelta de la misma forma.

Gracias a `pytorch` [44] y `transformers` [45], podemos reentrenar nuestro modelo GPT-2 con nuestros datos, como se puede ver en [Fine-tuning with custom datasets](#).

3.2.1. Guardado del estado del modelo

Finalmente, una vez entrenado el modelo, lo más importante es guardar su estado. Mediante la función `torch.save()`, podemos determinar el lugar en el que el estado del modelo se guardará. Este estado se guarda en un archivo, que, generalmente, consta de un diccionario en el que las claves corresponden a los nodos de las capas en sí, es decir, a sus parámetros entrenables, y cuyo valor es el peso de dicho nodo.

De la misma forma, mediante `torch.load()`, podemos cargar un modelo *vacío* con los pesos del archivo que especifiquemos, sobrescribiendo los valores por defecto y efectivamente recuperando el estado en el que dejamos el modelo tras el entrenamiento.

3.3. Generación de comentarios

En esta sección hablaremos de la generación de comentarios, una vez nuestro modelo está entrenado y listo para funcionar.

Diagrama de flujo: Generación de comentarios

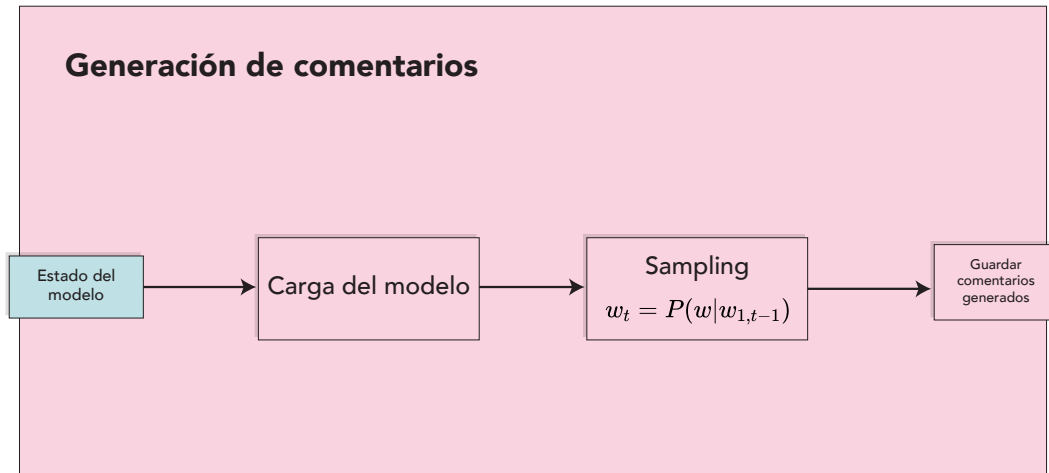


Figura 3.4: Diagrama resumen de la generación de comentarios

3.3.1. Carga del modelo

Con nuestro archivo con los pesos localizado, podemos cargar de vuelta los pesos relevantes en los parámetros correspondientes de la red, recuperando el estado original. Este estado es el que nos permitirá generar comentarios de forma automática.

3.3.2. Sampling y otros métodos de generación de texto

Una vez entrenado el modelo y ajustados los pesos, ya tenemos un marco de trabajo sobre el que poder generar comentarios de texto libre. Aun así, el modelo no es capaz de ofrecernos inferencias como en modelos clásicos de predicción o clasificación. Para poder acometer el proceso de generación de palabras, debemos efectuar una *decodificación* de las mismas, y existen varias formas para acometer esto.

Existen varias maneras de generar lo que en la jerga se denomina *texto abierto*, es decir, generar texto de forma relativamente libre y sin restricciones. Muchos otros modelos son capaces de hacerlo, aunque nosotros nos centraremos en los métodos que conciernen a nuestro GPT-2, que es un modelo autorregresivo.

Los modelos autorregresivos asumen que la siguiente palabra a generar se calcula como una función de probabilidad de todas las palabras anteriores, dada una palabra de contexto inicial. Dicho esto, existen varias maneras de *decodificar* texto de un modelo de lenguaje. Veamos las más relevantes.

Los métodos de decodificación aquí mencionados aparecen en la publicación [How to generate text](#) de Patrick von Platen. Es uno de los integrantes de Hugging Face, una de las empresa de código abierto que se dedica a la creación de los diferentes modelos mencionados anteriormente, desarrolladores de la librería `transformers`, entre otras.

Greedy search y beam search

La búsqueda voraz obtiene la palabra con mayor probabilidad de todas las opciones dada una palabra inicial.

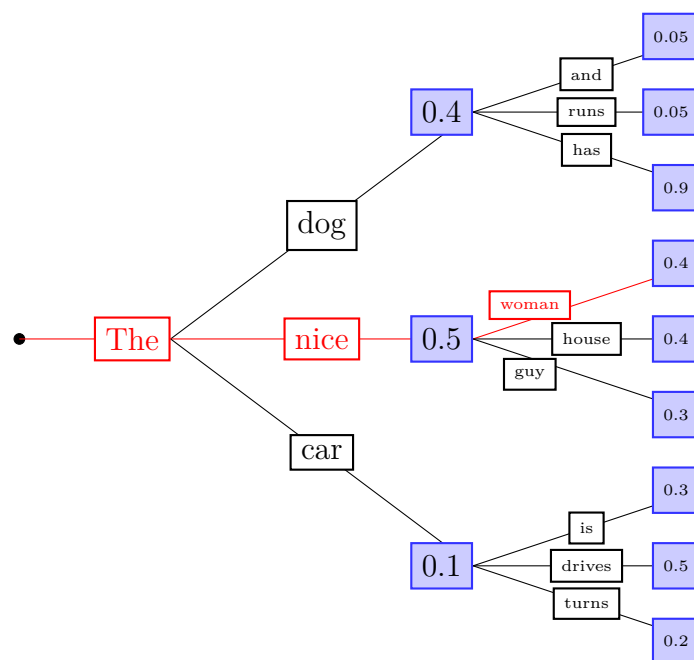


Figura 3.5: Ilustración de los pasos que daría el algoritmo greedy, tomando siempre la posibilidad más alta. Se muestra en rojo el camino que el algoritmo tomaría en este particular ejemplo, ya que *nice* supera a los otros dos términos en la probabilidad de ser elegido.

Dada una palabra inicial que tomamos del usuario, por ejemplo, siempre tomaremos aquella palabra de todas las posibles opciones que más probabilidad tenga de aparecer, dada la palabra anterior. Dichas probabilidades se calculan en el entrenamiento del modelo. Podemos ver un ejemplo de este comportamiento en la Figura 3.5.

El algoritmo toma siempre la palabra más probable y esto funcionará bien, aunque este algoritmo peca de empezar a repetirse bastante pronto. Esto es, generará secuencias que contienen palabras finales e iniciales similares, por lo que el ciclo empieza de nuevo al escogerse siempre la palabra más probable.

Una solución planteada es el *beam search*, que podemos visualizar en la Figura 3.6, algo así como una búsqueda distribuida. Dadas las probabilidades de cada palabra, el algoritmo calcula varios pasos en avanzadilla para varias alternativas, y devuelve aquella

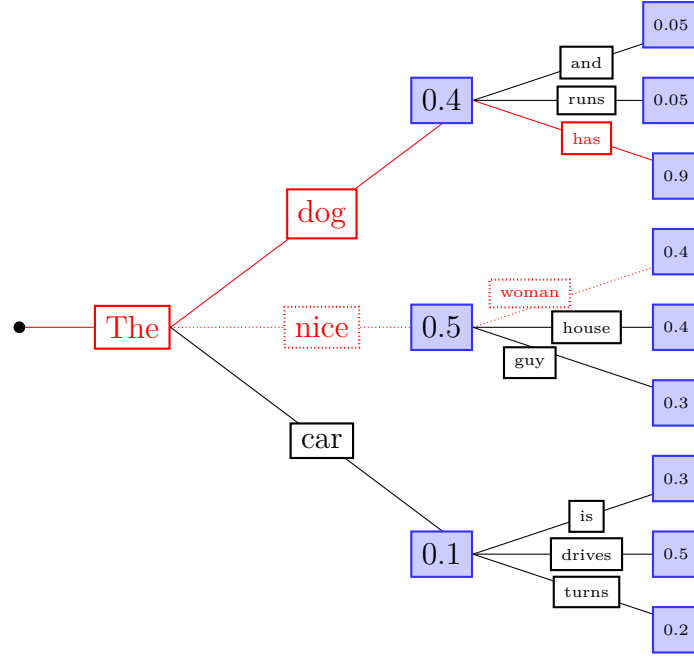


Figura 3.6: Ilustración demostrando el funcionamiento del *beam search*, con 2 beams en este caso. Vemos como se toman dos caminos, en línea continua el camino definitivo y en línea discontinua la otra alternativa. La probabilidad conjunta del camino elegido es $0.4 \times 0.9 = 0.36$, en contraste con el camino alternativo cuya probabilidad es de $0.5 \times 0.4 = 0.2$. Se toma un camino que a priori no es tan probable pero que a largo plazo, sí.

secuencia que más probabilidad general tenga. Esto soluciona algunos problemas, pero no termina de lidiar con el factor de repetitividad anterior.

Por otro lado, Ari Holtzmann sugiere en [46] que la generación del lenguaje humano no se guía por la elección de las palabras más probables. El lenguaje humano real es mucho menos predecible, ya que de otra forma sería *aburrido* de leer o escuchar.

Es por ello que los métodos alternativos han de evitar las técnicas **deterministas** de generación de texto, y apostar por los métodos con gran influencia aleatoria.

Sampling

Sampling es otra técnica de generación de lenguaje autorregresiva. En la forma más básica, el muestreo simplemente toma una palabra de forma aleatoria, ponderándolas con una distribución de probabilidad tal que:

$$w_t = P(w|w_{1,t-1}) \quad (3.1)$$

es decir, la probabilidad de escoger la palabra w_t depende de todas las palabras anteriores.

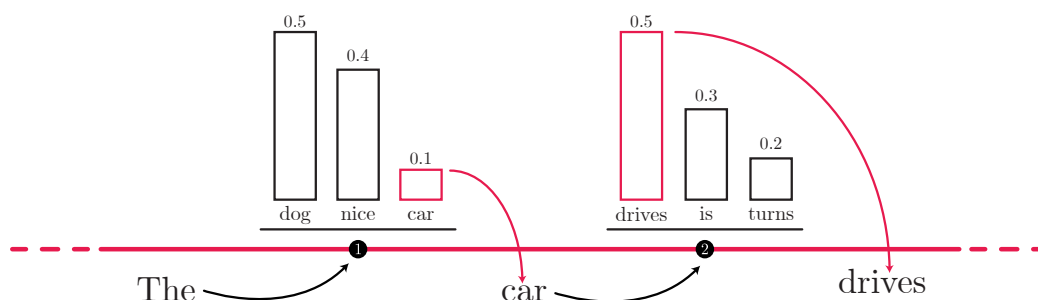


Figura 3.7: Ilustración de un ejemplo de muestreo. Se toman palabras aleatoriamente, aunque se pondera en función de su probabilidad.

El hecho de que una palabra presente una gran probabilidad de ser escogida no corresponde con que al final se la acabe escogiendo, como vemos en el paso 1 de la Figura 3.7. En este caso, se escogió *car* y posteriormente sí se escogió la palabra más probable. Este proceso nos ayuda a explorar más el espacio de búsqueda subyacente compuesto de todas las combinaciones de palabras posibles.

Temperatura

Podemos graduar la ponderación de la que hablamos con lo que los autores denominan la *temperatura* de la función de activación *softmax* de la última capa del modelo.

Al bajar la temperatura por debajo de 1 (siendo 1 un ajuste nulo), las palabras más probables se saturan, es decir, se hacen más probables, y las menos probables se comprimen, haciéndolas menos probables. En el otro extremo, al ajustar la temperatura a límites muy cercanos a 0, nos encontraríamos con la búsqueda voraz de la que hablamos antes.

Este factor nos ofrece una especie de regulador entre **máximo determinista** o **máximo aleatorio**.

La adición de esta técnica ayuda a que la generación no sea *extremadamente* aleatoria. Si bien en la sección anterior comentábamos que debíamos añadir aleatoriedad, todo en exceso es malo. Una generación completamente aleatoria es, en su mayoría, poco coherente. La bajada sutil de la temperatura del modelo provoca que se elijan, normalmente, palabras bastante probables, pero que de vez en cuando se escojan palabras poco probables. Dichas palabras ahora abren nuevos caminos por los que continuar generando, que no hubieran sido considerados de otra forma, pero continuamos eligiendo, por lo general, combinaciones de palabras más probables, para garantizar la coherencia del texto.

En esencia, el proceso es un gran compromiso entre coherencia y predecibilidad, además de los demás grandes problemas presentes como los bucles infinitos.

Top K Sampling

Dadas las premisas anteriores, las siguientes técnicas son mejoras incrementales que tratan de solucionar algunos otros problemas existentes.

El muestreo de los K mejores, tal y como su nombre indica, solo considera las K mejores opciones de toda la batería de palabras posibles. Dado este subconjunto K, se muestrea una palabra como en el Sampling original, tomándola aleatoriamente de forma ponderada. Puede verse en la Figura 3.8

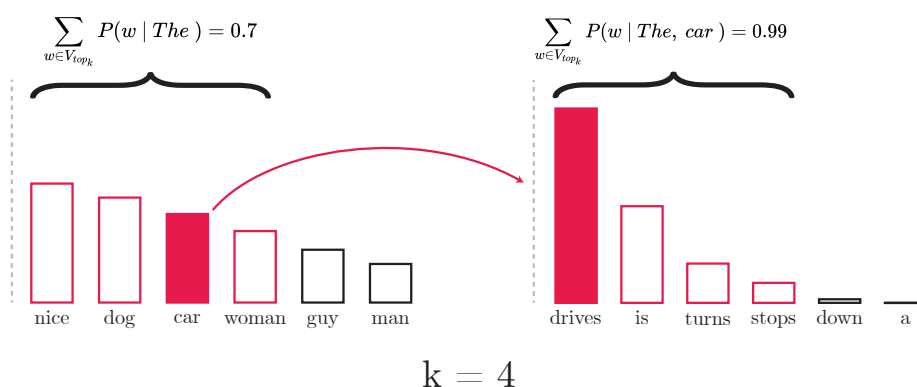


Figura 3.8: Ilustración del proceso de cálculo en muestreo con un valor de $k = 4$.

Esta técnica elimina posibles elecciones bastante malas, lo que Holtzmann denomina como "*the unreliable tail*", que podemos acuñar como *la cola traicionera*, que puede observarse en la ilustración. Holtzmann afirma que existe una enorme cantidad de palabras con muy baja probabilidad en cada elección que, en realidad, están sobrerrepresentadas cuando hablamos de términos agregados. Es decir, la **probabilidad agregada** de todas esas palabras con una bajísima probabilidad puede corresponder a **grandes proporciones de la distribución de probabilidad**, lo que, efectivamente, va en contra de nuestro objetivo.

Tomando solo aquellas K mejores, se elimina dicha cola y se redistribuye la probabilidad de los términos elegidos de forma oportuna, creando un modelo de decodificación mucho más natural que los otros.

Top P nucleus

Finalmente, esta última técnica se denomina Top P nucleus, refiriéndose a un valor diferente al K anterior.

En este caso, en lugar de tomar los K mejores términos, tomamos el conjunto más pequeño de términos cuya probabilidad acumulada supere el p establecido por el usuario. Esto es, dado un $p = 0.9$, tomamos las n mejores palabras que, dada la suma de su probabilidad, se supere el p determinado. El proceso es básicamente idéntico al ilustrado en la Figura 3.8, solo que consideramos la probabilidad acumulada, no un número de términos.

Esto provoca que en lugar de tomar un subconjunto de palabras de cardinalidad constante durante todo el proceso, ahora el tamaño varía. Se ha demostrado que esta flexibilidad ayuda y contribuye a una generación de texto más natural.

Estos dos últimos métodos son los más utilizados en el estado del arte, y este último es el que **utilizamos nosotros** en nuestro proyecto a la hora de generar palabras. Como vemos, aun disponiendo de una red neuronal muy potente, no es trivial extraer información coherente y de calidad.

4. EXPERIMENTACIÓN: MODELOS, ENTRENAMIENTO Y GENERACIÓN

En este capítulo usaremos todos los conceptos vistos en los capítulos anteriores para justificar las elecciones de los diferentes modelos, hablaremos de las características principales de los mismos y finalmente entraremos en la fase del entrenamiento y los resultados de dichos modelos.

4.1. Datos: fuente y forma

Los datos escogidos provienen del dataset [Medical Text](#) publicado por Chaitanya Krishna Kasaraneni, y de [Medical Transcriptions](#), publicado por Tara Boyle. La naturaleza de los mismos es ligeramente diferente así que explicaremos el proceso de preprocesamiento y unificación posteriormente.

En la Figura 3.1 podemos ver un pequeño resumen del preprocesamiento que se acometerá a los datos. En la siguiente sección se detallan cada uno de estos pasos.

4.1.1. Dataset: *Medical Text*

El dataset tiene formato `.dat`, estructurado como un `.tsv` (Tab Separated Values). La primera columna corresponde con una categoría determinada –ya que el dataset estaba diseñado para clasificación– y la segunda columna contiene fragmentos de documentos médicos.

El dataset es en inglés, está anonimizado y los comentarios principalmente consisten en descripciones quirúrgicas o relacionadas con operaciones complejas. Podemos ver los dos primeros de nuestro conjunto de datos en los Comentarios 1 y 2.

Comentario 1 *Excision of limbal dermoids. We reviewed the clinical files of 10 patients who had undergone excision of unilateral epibulbar limbal dermoids. Preoperatively, all of the affected eyes had worse visual acuity (P less than .02) and more astigmatism (P less than .01) than the contralateral eyes. Postoperatively, every patient was cosmetically improved. Of the eight patients for whom both preoperative and postoperative visual acuity measurements had been obtained, in six it had changed minimally (less than or equal to 1 line), and in two it had improved (less than or equal to 2 lines). Surgical complications included persistent epithelial defects (40 %) and peripheral corneal vascularization and opacity (70 %). These complications do not outweigh the cosmetic and visual benefits of dermoid excision in selected patients.*

Comentario 2 *Retained endobronchial foreign body removal facilitated by steroid therapy of an obstructing, inflammatory polyp. Oral and topical steroids were used to induce regression in an inflammatory, obstructing endobronchial polyp caused by a retained foreign body. The FB (a peanut half), which had been present for over six months, was then able to be easily and bloodlessly retrieved with fiberoptic bronchoscopy.*

El dataset está descompuesto en un archivo `train.dat` y otro `test.dat`. El archivo de entrenamiento contiene 14438 comentarios, y el de evaluación, 14442. En total, disponemos de 28880 comentarios.

Como podemos observar en las Figuras 4.1a y 4.1b, la distribución de los distintos elementos de nuestro dataset de entrenamiento está muy normalmente distribuída.

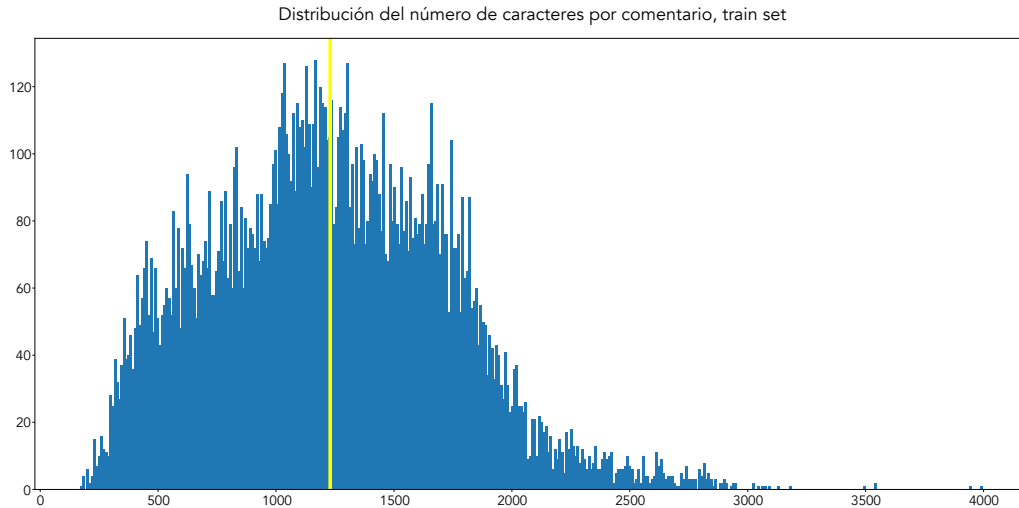
El número medio de caracteres por comentario es de 1230, y el número medio de tokens por comentario es de unos 180, correspondiéndose con las líneas amarillas en las figuras.

Se pueden apreciar, aún así, algunos valores atípicos de comentarios particularmente largos. Esto, sin embargo, no es necesariamente malo en nuestro caso. En definitiva, cuanto más texto tengamos a nuestra disposición, mejor para el modelo.

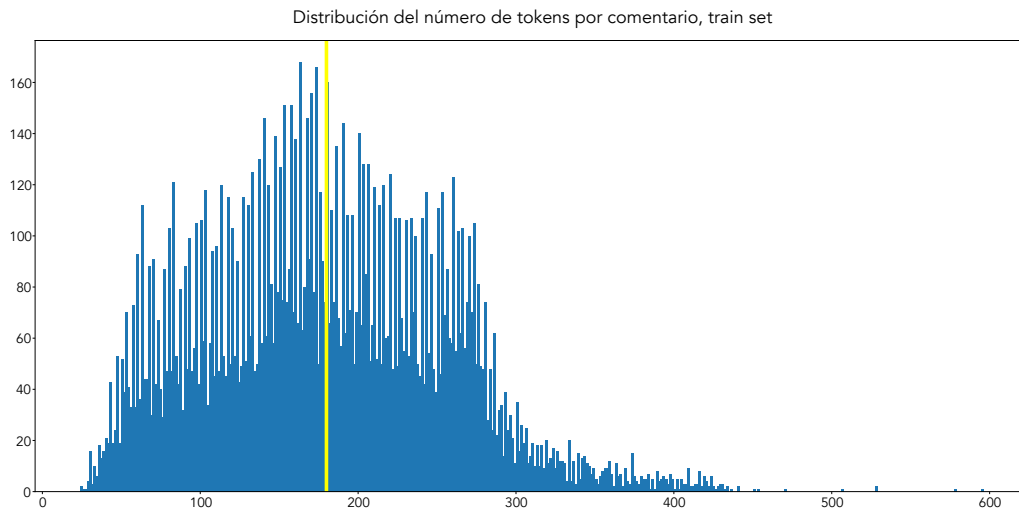
Una media de casi 200 palabras por comentario con comentarios alcanzando las 500 corresponde con comentarios relativamente largos. Esto nos vendrá bien de cara al entrenamiento de nuestro modelo, para poder formar oraciones con más sentido.

4.1.2. Dataset: *Medical Transcriptions*

Este dataset es en realidad una extracción de la página web mmsamples.com, donde se halla una respetable cantidad de transcripciones médicas. La autora extrajo todos



(a) Distribución del número de caracteres por comentario, en el conjunto de entrenamiento



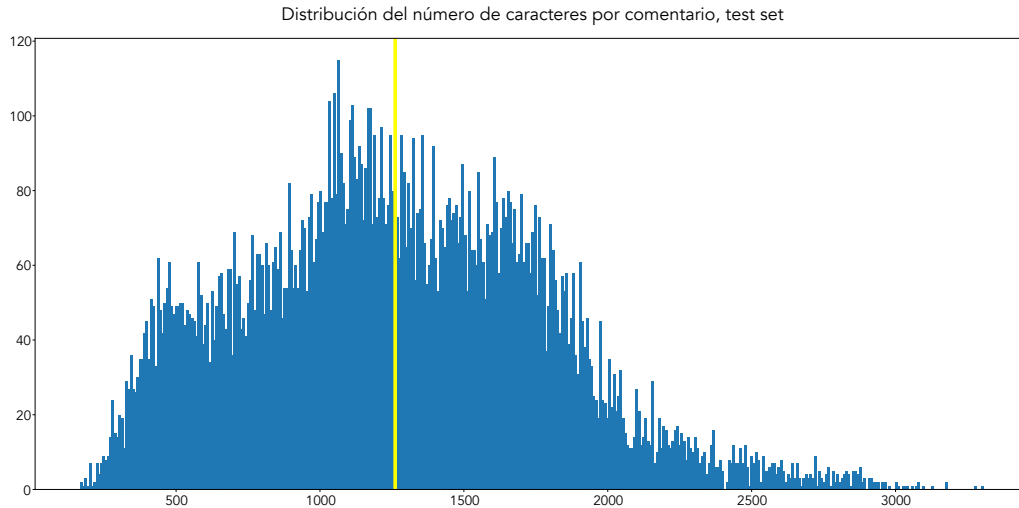
(b) Distribución del número de tokens por comentario en el conjunto de entrenamiento

Figura 4.1: Visualización de la distribución de nuestro conjunto de entrenamiento

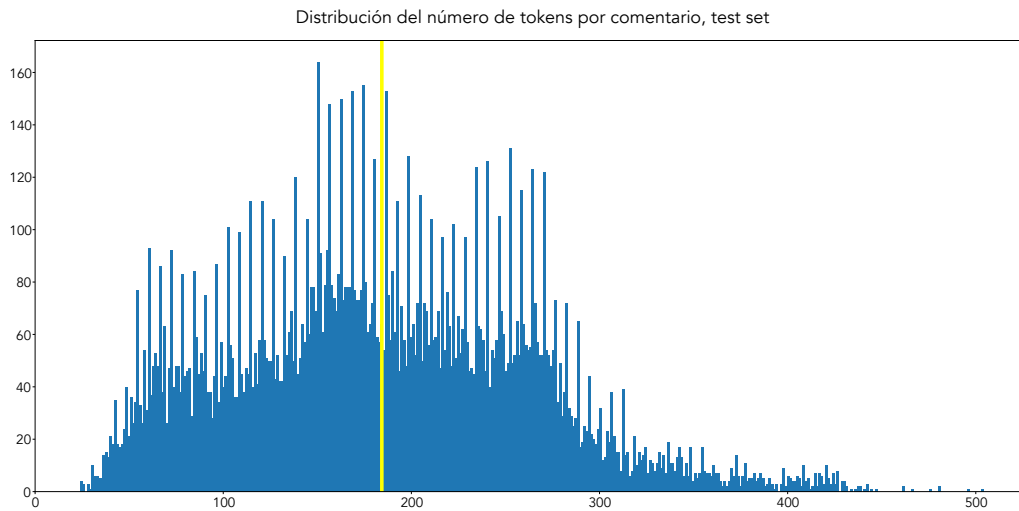
los comentarios, así como los diferentes metadatos que los acompañaban mediante *web scraping* y los provee en la columna `transcription`.

Al igual que el anterior, este conjunto de datos se corresponde con informes de operaciones quirúrgicas en inglés, y de igual forma anonimizado.

En este caso, como podemos apreciar en la Figura 4.3, las distribuciones son ligeramente asimétricas, predominando comentarios más cortos. Aún así, disponemos de comentarios excepcionalmente largos, con alrededor de 18000 caracteres.



(a) Distribución del número de caracteres por comentario, en el conjunto de evaluación



(b) Distribución del número de tokens por comentario en el conjunto de evaluación

Figura 4.2: Visualización de la distribución de nuestro conjunto de evaluación

4.2. Preprocesamiento

En esta sección, describiremos el preprocesamiento acometido en cada uno de los datasets. Proviene de fuentes diferentes así que cada uno recibirá un trato diferente, con objeto de normalizar y unificar el formato de todos de cara al entrenamiento.

4.2.1. Medical Text

Este conjunto de datos, siendo específicamente texto, el formato, ortografía y en general formato de los archivos es muy bueno. Simplemente hemos de eliminar las categorías adjuntas a cada comentario, para obtener una lista de comentarios crudos en sí. Por lo demás, los comentarios carecen de problemas de formato, codificación o cualquier otra

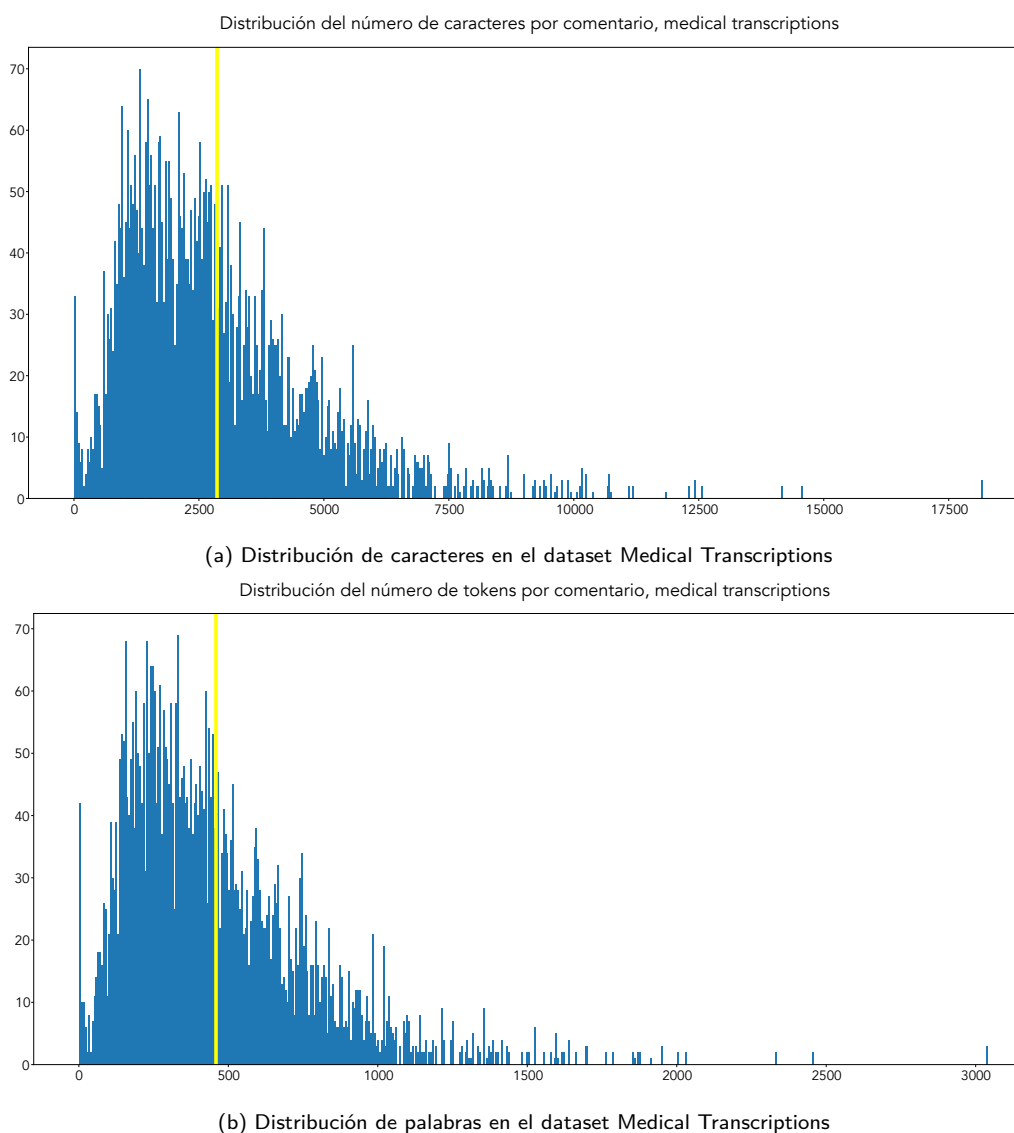


Figura 4.3: Visualización del dataset Medical Transcriptions

cosa que pudiera interferir con el proceso de entrenamiento. Probablemente el autor ya hiciera esto por nosotros antes de publicarlo.

4.2.2. Medical Transcriptions

En el caso de las transcripciones médicas, la tarea es considerablemente más compleja. El conjunto de datos proviene de la página web metsamples.com, como especificamos anteriormente. La autora efectuó un proceso de *web scraping* para obtener toda la información y recogerla en el archivo `.CSV`.

Esto facilita las cosas, pero desde luego los comentarios deben ser tratados en profundidad antes de poder pasarlos a cualquier modelo. Los trazos de formato en HTML se

dejan entrever en los comentarios con signos de puntuación o tabulaciones fuera de lugar, apreciables en el Comentario 3, así que debemos arreglarlo previo entrenamiento.

Para ello, se ha hecho un fuerte uso de expresiones regulares, y se ha creado un *pipeline* para procesar todo el texto a la vez.

El pipeline elimina todas las posibles trazas o residuos que hubieran quedado del *web scraping*. Podemos ver el pipeline diseñado en la Figura A.1 del Apéndice.

El resumen del proceso es eliminar signos de puntuación mal colocados, eliminar títulos o cabeceras de secciones de la página web, sustituir múltiples espacios por uno solo o eliminar los números de listas enumeradas (1., 2., etc). Finalmente, se añaden las etiquetas que vemos en la Figura 3.1. Se explicará su funcionamiento en la sección de la experimentación.

El resultado es un texto muy limpio y claro, mucho más apto para la fase de entrenamiento.

Podemos ver una comparativa del antes (Comentario 3) y el después (Comentario 4) del preprocesamiento de un determinado comentario de nuestro dataset. Los comentarios han sido truncados debido a su longitud.

Comentario 3 *SUBJECTIVE; This 23-year-old white female presents with complaint of allergies. She used to have allergies when she lived in Seattle but she thinks they are worse here. In the past, she has tried Claritin, and Zyrtec. Both worked for short time but then seemed to lose effectiveness. She has used Allegra also. She used that last summer and she began using it again two weeks ago. [...]*

Comentario 4 *< | BOS | > This 23-year-old white female presents with complaint of allergies. She used to have allergies when she lived in Seattle but she thinks they are worse here. In the past, she has tried Claritin, and Zyrtec. Both worked for short time but then seemed to lose effectiveness. She has used Allegra also. She used that last summer and she began using it again two weeks ago. [...] < | EOS | >*

Tras haber preprocesado y cribado todos los elementos que pudieran suponer un problema a la hora de entrenar nuestro modelo, lo que obtenemos es un dataset unificado con un total de 33846 comentarios, sumando un total de 7537697 palabras con una media de 222 tokens y 1482 caracteres por comentario.

Vemos además, en la Figura 4.4 una nube de palabras de todo el conjunto de datos.

Se aprecia que las palabras más comunes son *patient*, *case*, o *treatment*, junto con *study* o *performed*. Esta nube de palabras se ha calculado en el texto preprocesado, pero

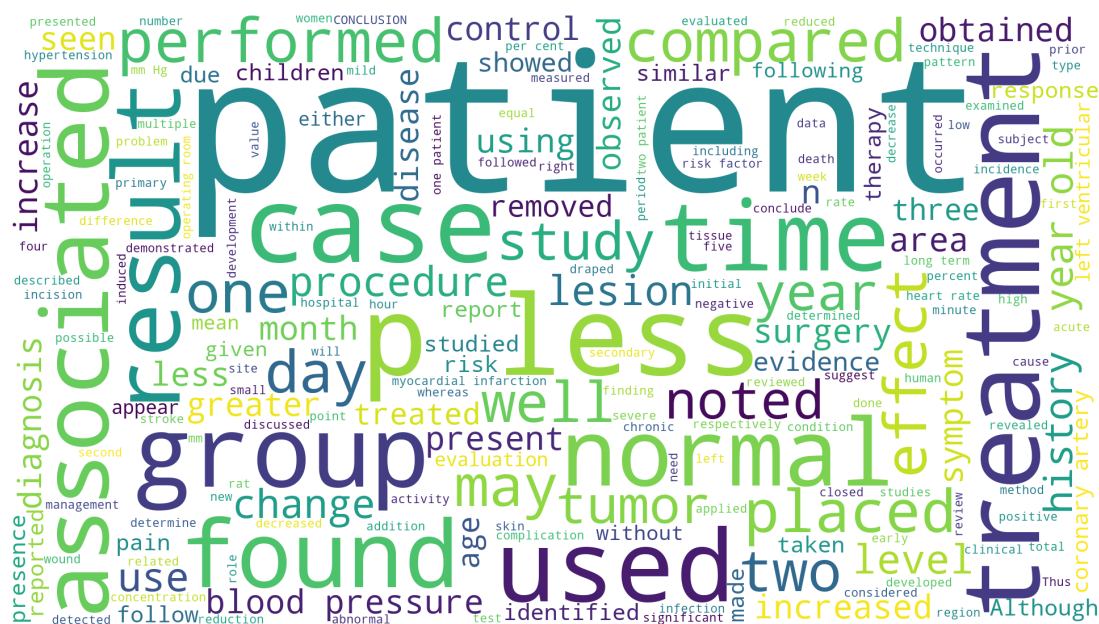


Figura 4.4: Wordcloud de todo el conjunto de datos

el preprocesamiento solo ha tenido en cuenta cuestiones de formato, mayoritariamente. De ahí que términos como *used*, *using* o *use* aparezcan a la vez. No hemos hecho stemming ni lemmatization porque nos interesa el contenido de las oraciones tal y como está, que es como el modelo fue diseñado para ser entrenado.

Disponiendo de este dataset, estamos listos para poder entrenar y ajustar nuestro modelo para que genere comentarios muy similares a los presentes en nuestro conjunto de datos.

4.3. Resultados

Habiendo preprocesado los datos, entrenado el modelo, guardado su estado y decidido el método de decodificación, estamos listos para poder generar comentarios.

La generación, como vimos en las secciones anteriores, viene determinada por dos parámetros clave en nuestro caso: la temperatura y el valor de p para el Top P Nucleus. Veamos cómo estos valores afectan en la generación de los comentarios. Generaremos un par de comentarios por cada combinación de valores y estudiaremos lo que ocurre.

Temp: 0.1, Top p: 0.2

Comentario 5 *The effect of a high-dose of aspirin on the incidence of coronary artery disease in a prospective study of the elderly. A prospective study of the incidence of coronary artery disease in a prospective study of the elderly was conducted in a large population of elderly men and women. The incidence of coronary artery disease was studied in a large population of elderly men and women. The incidence of coronary artery disease was studied in a large population of elderly men and women. The incidence of coronary artery disease was studied in a large*

Comentario 6 *The effect of a high-dose of aspirin on the incidence of coronary artery disease in a prospective study of the elderly. A prospective study of the incidence of coronary artery disease in a prospective study of the elderly was conducted in a large population of elderly men and women. The incidence of coronary artery disease was studied in a large population of elderly men and women. The incidence of coronary artery disease was studied in a large population of elderly men and women. The incidence of coronary artery disease was studied in a large*

Vemos que en este caso, se han generado dos comentarios exactamente iguales. Esto es debido al alto grado determinista del generador en este caso. En azul se ha indicado la frase que inicia el ciclo, y en rojo el bucle infinito que obtendríamos de seguir generando dicho comentario. Esta combinación de parámetros no es trabajable.

Temp: 0.4, Top p: 0.4

Comentario 7 *The effect of a high-dose of aspirin on the incidence of coronary artery disease in a prospective study of the elderly. A prospective study of the incidence of coronary artery disease in a prospective study of the elderly was conducted in a large population of elderly men and women. The incidence of coronary artery disease was studied in a cohort of 5,000 men and women aged 65 to 74 years. The incidence of coronary artery disease was studied in a cohort of 5,000 men and women aged 65 to 74 years*

Comentario 8 *The effect of a high-dose of aspirin on the incidence of coronary artery disease in a prospective study of the elderly. A prospective study of the incidence of coronary artery disease in a prospective study of the elderly was conducted in a large population of patients with acute myocardial infarction. The incidence of acute myocardial infarction was significantly higher in patients with a single dose of a selective serotonin reuptake inhibitor than in patients with a single dose of a selective serotonin reuptake inhibitor.*

En este caso, vemos que el inicio, marcado en verde, es exactamente igual que la generación anterior. En este caso, al aumentar la temperatura, obtenemos unos caminos ligeramente distintos al final del comentario, pero continúan apareciendo ciclos.

Temp: 0.6, Top p: 0.6

Comentario 9 *Treatment of chronic pain in patients with chronic pain. A prospective study. A retrospective study of patients with chronic pain in whom treatment was discontinued or discontinued for chronic pain was conducted. Patients were randomly assigned to receive either a placebo (n = 14) or a sham (n = 14) treatment. Patients were treated with either a placebo (n = 14) or sham (n = 14) for 2 months. Patients were treated with either a placebo (n = 14) or sham (n = 14)*

Comentario 10 *The effects of a variety of drugs on the central nervous system.*
We report the effects of a variety of drugs on the central nervous system. We
report the effects of a variety of drugs on the central nervous system. We report
the effects of a variety of drugs on the central nervous system. We report the
effects of a variety of drugs on the central nervous system. We report the effects
of a variety of drugs on the central nervous system. We report the effects of a
variety of drugs on the central

En este caso, a pesar de haber aumentado la temperatura, hemos tenido un caso particular en el que los ciclos infinitos empezaron particularmente pronto.

Temp: 0.9, Top p: 0.8

Comentario 11 *Optimal effect of systemic and local opioid administration on serum glucose levels in patients with preoperative seizures. Fourteen patients (8 %) had an abnormal serum glucose level (mean \pm SEM) measured by intra-peritoneal glucose measurements. The mean serum glucose level of the patients with preoperative seizures was markedly elevated ($P = 0.002$), compared with baseline values ($P = 0.003$) (Figure 2 A). Glucose plasma concentrations decreased with treatment and did not change after the initial systemic*

Comentario 12 *A case of cardiac arrhythmia in 14 patients with coronary artery disease. The risk of arrhythmia was evaluated by using clinical hemodynamic and pulmonary blood pressure measurements. Patients underwent standardized coronary artery bypass grafts and intraoperative radiculograms (ICGs). The history of the disease and the follow-up of patients with an early arrhythmia were compared with those with a later arrhythmia. Arterial pressure was significantly higher in the more severe patients.*

Conforme aumentamos el valor de la temperatura y estrechamos el rango de palabras que tendemos a escoger, los resultados se aproximan mucho más a lo que esperábamos del modelo: comentarios que no se repiten y con una coherencia mucho más visible. Es más complicado sacar defectos de estos comentarios, lo cual es buena señal.

Temp: 1, Top p: 0.8

Comentario 13 *Paediatric invasive immunosuppression (IPIS) in the treatment of recurrent kidney disease. Epidemiologic evidence suggests a association between immunosuppression and recurrence of recurrent kidney disease. Acute renal injury is a major risk factor for recurrent kidney disease. As the patient develops a family history of recurrence, his status will likely be monitored to determine whether he may benefit from continued immunosuppression.*

Comentario 14 *Predominant asymptomatic single cell carcinoma, melanoma, and cerebrovascular myopathy. The diagnosis of melanoma, cerebrovascular myopathy, and primary carcinoma requires extensive examination. Among these cancers, recurrent myopathy is a common occurrence. From 1980 to 1989, the incidence of melanoma, cerebrovascular myopathy, and secondary carcinoma increased from 8.5 % to 12.6 %.*

En la Sección A.1 del Apéndice se exponen algunos más, por si su inspección resultase de interés.

Se puede apreciar que, a veces, los comentarios son muy cortos, como el comentario 17. En contadas ocasiones el generador ha devuelto una línea vacía.

En ocasiones, los comentarios pueden no ser rigurosos médicamente hablando, pero esto, de cara al análisis de las herramientas disponibles para el procesamiento de los comentarios, no es estrictamente necesario. Necesitamos que los comentarios incluyan conceptos que sean importantes y destacables, como tipos de enfermedades, medicaciones, etc. De esta forma, podremos analizar cómo las herramientas detectan diferentes tipos de información y cómo de fiables son.

Es por ello que poseer un generador de comentarios nos resulta conveniente, ya que se pueden considerar muchos más casos, prácticamente de forma ilimitada. Dados los resultados, el valor finalmente elegido para los parámetros del GPT-2 es de **0.9 para la temperatura y 0.8 para el valor de p**.

4.3.1. Métricas y calidad de los comentarios

Uno de los grandes retos de la creación de modelos de lenguaje, a parte de en sí ser el desarrollo de los mismos extremadamente complejo, es la creación de métricas de calidad de los resultados obtenidos.

Con modelos clásicos que trabajan con números, o incluso imágenes, existe una serie de herramientas que permiten ofrecer, de forma objetiva, rigurosa y consistente, un valor cuantitativo que mide la calidad de la salida de dicho modelo.

En los modelos de lenguaje, y más aún, en los modelos generativos, obtener un valor cuantitativo que evalúe la calidad de los comentarios generados es, más que complicado, subjetivo.

Las preguntas que nos hacemos realmente son: ¿cómo de buenos son los comentarios generados? ¿Son comentarios coherentes? ¿Son comentarios rigurosos? ¿Cómo podemos evaluar de forma automática un conjunto de comentarios suficientemente representativo?

Para solucionar estos problemas, se han inventado algunas técnicas, aunque las más relevantes son la métrica BLEU [47] y la métrica Rouge [48].

BLEU

$$BLEU = \min \left(1, \exp \left(1 - \frac{ref_{length}}{output_{length}} \right) \right) \left(\prod_{i=1}^4 precision_i \right)^{\frac{1}{4}} \quad (4.1)$$

La métrica **Bi**Lingual **E**valuation **U**nderstudy (BLEU) se ideó originalmente con objeto de evaluar traducciones automáticas. Dado un texto en un idioma y la traducción automática a otro, poder ofrecer un valor objetivo de cómo de bien se corresponden dichas oraciones. Esta tarea, como decíamos anteriormente, se compone de un importante carácter subjetivo, y es la razón por la que existen intérpretes y traductores que se encargan de este trabajo. La traducción de texto no es un mapeo directo de un conjunto de palabras de un idioma al de otro, ya que han de considerarse expresiones y valores culturales.

Dicho eso, disponer de, al menos, una medida que pueda servir de forma orientativa para evaluar la calidad de un traductor automático es, cuanto menos, conveniente.

El algoritmo en cuestión toma un conjunto de oraciones de referencia y otro conjunto de oraciones generadas que deseemos evaluar, y calcula un *score* en función de cómo de probable es que una de las frases generadas pertenezca al conjunto de oraciones de referencia. Esta métrica también se utiliza extensamente en evaluación de modelos generativos de lenguaje, como el nuestro, ya que la naturaleza de la evaluación se corresponde directamente con la nuestro caso de uso.

Para evaluar el modelo, tomamos 100 oraciones aleatorias de nuestro conjunto de datos, y generamos 100 oraciones con nuestro modelo. El valor obtenido es de 0.0036.

Por referencia, la misma métrica nos ofrece un valor entre 0 y 1, donde un valor de 0 indica un nulo solapamiento del texto generado con el valor de referencia, y un valor cercano a 1 indica una buena traducción del texto dadas las referencias obtenidas.

Hay que tener en cuenta que es una medida de **solapamiento** especialmente diseñada para la evaluación de la traducción automática, por lo que nuestro valor cercano a 0 no necesariamente indica una mala calidad de la generación de los comentarios.

Rouge

La métrica Recall-Oriented Understudy for Gisting Evaluation (Rouge) funciona de una manera muy similar, calculando el solapamiento de las palabras en los textos a comparar y haciendo especial énfasis en las diferencias entre el número de tokens, ya que es

una métrica especialmente diseñada para la evaluación de resúmenes. En esta métrica se obtuvo un valor de 0.0669, indicando un supuestamente mal resumen del texto.

Como vemos, no son medidas muy útiles para nuestro propósito, pero lo cierto es que no existen muchas otras alternativas para la evaluación automática del lenguaje generativo libre, ya que es una tarea tan intrínsecamente subjetiva, como mencionábamos antes. Existe una opción en [49], pero el código no ha sido publicado a día de hoy ya que continúa en desarrollo.

Las medidas utilizadas anteriormente puede que no nos ofrezcan un valor objetivo de la calidad de los comentarios, pero sí nos ofrecen información acerca de su solapamiento.

Similitud

Esto nos lleva a otra métrica a utilizar: la distancia en los comentarios. Un factor que sí podemos determinar de forma sencilla es si los comentarios son muy similares entre sí. Idealmente, se desearía que el generador pudiera proporcionar comentarios que no estuvieran relacionados entre sí, ya que no se exploraría lo suficiente en el espacio de todas las posibles combinaciones de entidades que pudiésemos encontrar. En esencia, deseamos simular casos específicos y poco relacionados entre sí.

Para ello, podemos utilizar la medida de la distancia del coseno del valor TF-IDF mencionado en la Sección 2.2.1. Calculando el valor TF-IDF de todos los comentarios generados y multiplicándolo por su traspuesta, generamos una matriz que nos define cómo de *cerca* se encuentra cada comentario de todos los demás en el espacio latente. Podemos apreciar dicha matriz en la Figura 4.5.

Se aprecia que la diagonal principal tiene un valor de 1, ya al comparar un comentario consigo mismo, estamos lo más cerca que se puede estar, o lo que es lo mismo, la distancia es 0.

Por lo general, todos los comentarios ofrecen valores de similitud muy bajos, que se intuye por la gran predominancia del color morado en la imagen, que corresponde con valores cercanos a 0.

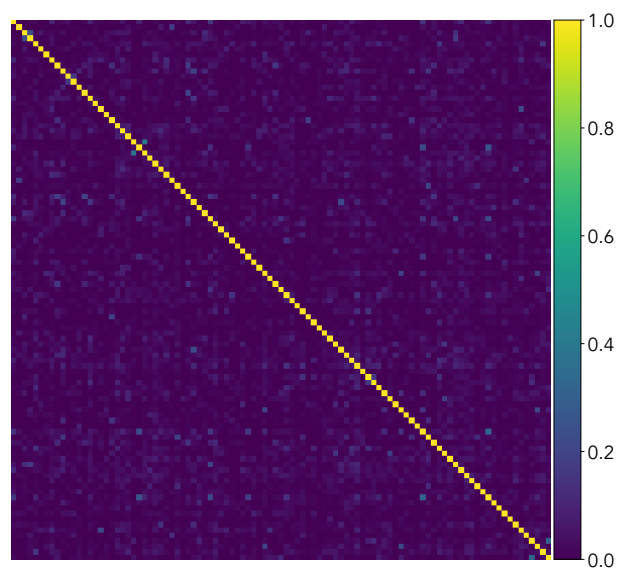


Figura 4.5: Matriz de comparativa de distancia de los 100 comentarios generados en todas las posibles combinaciones.

5. ANÁLISIS DE TEXTOS MÉDICOS

En este capítulo discutiremos cómo hemos llevado a cabo la evaluación de las herramientas disponibles en el estado del arte para la extracción de información en texto no estructurado de carácter médico.

Para ello se ha creado una web app disponible de forma muy accesible para que sea fácil ver el poder combinado de todas estas herramientas. Dicha herramienta se denomina [MEDGEN](#).

5.1. Despliegue de MEDGEN

Para hacer accesible todo el trabajo que se ha acometido, se ha elaborado una aplicación que permite hacer uso de todas las técnicas y herramientas discutidas a lo largo del documento.

El framework utilizado ha sido [Streamlit](#). Streamlit es un framework para Python diseñado para el prototipado ágil de aplicaciones web. Los creadores idearon esta herramienta para disminuir el esfuerzo que hay que acometer a la hora de desplegar una aplicación basada en inteligencia artificial, es decir, una aplicación que consta de sistemas de inferencia o similares, sistemas que no se encuentran en aplicaciones web usuales.

De esta forma, es muy accesible y sencillo diseñar una aplicación en Python, donde se puede integrar cualquier sistema de IA en desarrollo y habilitar un espacio para que personas ajenas al desarrollo experimenten de primera mano cómo funciona, ofreciendo la oportunidad de crear casos de uso reales.

La aplicación permite generar comentarios utilizando el modelo generativo, y permite analizar dichos comentarios, así como importar desde distintas fuentes. La aplicación efectúa un pequeño análisis del texto y utiliza el tagger que el usuario seleccione.

5.2. Evaluación de las herramientas

Como vemos en las Figuras 5.1 y 5.2, podemos efectuar nuestro análisis de diferentes formas. La aplicación nos permite utilizar los NER Taggers indicados en la Sección 2.4: Med7, BC5CDR y BIONLP13CG.

Podemos, en caso de no disponer de un conjunto de datos, generar comentarios usando nuestro modelo generativo. Una vez generados, podemos analizarlos con el NER Tagger que deseemos. Aquí, el desarrollador estaría programando el suyo propio y podría integrarlo al sistema fácilmente para testearlo de forma apropiada.

El análisis nos ofrece una versión etiquetada del texto original, así como unas estadísticas del comentario en concreto. Además de ello, podemos importar nuestros propios comentarios desde un archivo o copiarlos y pegarlos en el campo de texto, con tal de maximizar la comodidad de uso de la aplicación.

En función del texto de entrada y el Tagger seleccionado obtendremos mejores o peores resultados. El generador de comentarios está entrenado para generar comentarios de carácter quirúrgico o posoperatorios. Para ello, el BIONLP13CG es una buena elección, pero el Med7 rara vez será capaz de ofrecer resultados, ya que está especializado en recetas médicas.

La ventaja de esto es el gran abanico de posibilidades que tenemos de forma muy accesible, siendo posible analizar casi cualquier tipo de texto médico que sea de interés. De no disponer de un Tagger apropiado, siempre podemos incluirlo en la aplicación, como mencionábamos en la Sección 2.4.



Figura 5.1: Captura de pantalla de la aplicación elaborada. A la izquierda podemos generar comentarios y a la derecha, analizarlos.



Figura 5.2: Captura de pantalla del análisis que ofrece la aplicación acerca de un determinado comentario.

5.3. Instalación y modo de uso

La aplicación puede compilarse desde el código fuente siguiendo el [enlace al repositorio](https://github.com/jesi-rgb/medical-text-analysis), descargándolo y ejecutando los siguientes comandos:

```
git clone https://github.com/jesi-rgb/medical-text-analysis
```

Activamos el entorno que más nos guste, ya sea de python o conda.

```
cd medical-text-analysis
```

```
pip install -r requirements.txt
```

Una vez finalizado,

```
streamlit run src/streamlit_gen_test.py
```

se nos abrirá una ventana en el navegador y la aplicación estará lista para funcionar.

La primera ejecución tarda un poco más, ya que ha de descargar el modelo de Internet y cargarlo en memoria. Tras eso, los modelos se guardan en caché y la ejecución es mucho más rápida.

6. CONCLUSIONES

En este último capítulo abordaremos todos los problemas y objetivos propuestos a lo largo de este proyecto con objeto de refrescarlos, y finalmente se ofrecerá una evaluación acerca de cómo estos objetivos se han cumplido. Finalmente se reflexionará acerca de cómo podría extrapolarse este proyecto de cara al futuro.

6.1. Problemática inicial

Para empezar, hablaremos del problema inicial que nos ha motivado en la elaboración de esta herramienta.

Los informes médicos ofrecen grandes cantidades de información clave para los profesionales de cara al tratamiento de los pacientes, y gran parte de esta información se almacena en cuerpos de texto libre. Dichos extractos usualmente no se utilizan de cara a la extracción de información automática debido a la naturaleza no estructurada de los mismos, que hace su análisis muy complicado.

Aún así, en estas secciones se puede encontrar mucha información muy valiosa y relevante, por lo que tratar de encontrar un sistema automático que pueda ofrecer información estructurada dado un conjunto de información no estructurada puede ser de gran utilidad en el ámbito médico profesional.

6.2. Objetivos propuestos

Dada la problemática inicial, el objetivo es crear una herramienta que pueda ofrecer información estructurada. Una de las herramientas más útiles en este caso es un reconocedor de entidades, o, en inglés, un NER Tagger. Dado un texto, es capaz de reconocer entidades importantes automáticamente.

En el contexto médico, debemos atender a las entidades más importantes como los medicamentos, enfermedades, partes del cuerpo, etc. Para ello, existe un vocabulario unificado, denominado SNOMED. Gracias a esto, podemos averiguar con precisión dónde se hallan los datos más importantes.

Si bien se han creado varios reconocedores de entidades, cada uno construido para diferentes tareas dentro del ámbito médico (desde recetas de medicamentos hasta informes quirúrgicos), se pretendía crear un sistema con el que el desarrollo de nuevos reconocedores fuera más fácil. Esto implica disponer de ingentes cantidades de datos que, por su delicada corte, pueden no estar fácilmente disponibles.

Es por ello, que en pos de alcanzar este primer objetivo de facilitar el desarrollo de reconocedores de entidades (o cualquier otro tipo de herramienta relacionada), se plantea un segundo objetivo: un generador de comentarios automático.

Un generador de comentarios automático elimina la necesidad de inspeccionar en busca de conjuntos de datos, ya que se podrán generar cuantos se deseen, de cara al desarrollo de las herramientas. Esto es además viable gracias al destacable avance acometido en los campos de inteligencia artificial, donde se han creado modelos de lenguaje muy competentes.

6.3. Metodología y resultados

Es por ello que tenemos varios objetivos:

1. Recopilar todas las fuentes de información públicas que nos provean con datos de comentarios médicos listos para su análisis.
2. Crear un modelo de lenguaje generativo que sea capaz de generar comentarios de forma automática.
3. Disponiendo de un conjunto de datos *infinito*, habilitar un sencillo desarrollo de herramientas de extracción de conocimiento de texto.

Para ello, se han utilizado modelos generativos de lenguaje preentrenados, como el GPT-2. El GPT-2 es un *transformer*, un tipo de arquitectura de red neuronal especializada en procesamiento de texto de forma paralela. El GPT-2 es un modelo preentrenado y abierto, lo que nos habilita a poder crear una herramienta que genere comentarios de forma automática.

Se ha observado que el modelo es capaz de crear comentarios de forma bastante competente con un entrenamiento previo en un conjunto de datos preprocesado y unificado manualmente, lo que los creadores denominan como *fine-tuning* de la red.

Esto hace el desarrollo de dichos modelos más fácil para todos, contribuyendo a una mayor y mejor producción de herramientas de asistencia médica, clave para cualquier

complejo hospitalario medianamente grande, donde se manejan volúmenes de datos, a menudo, insostenibles.

6.4. Posibles trabajos futuros

Partiendo del punto en el que nos encontramos, podemos dirigir el proyecto en varias direcciones.

Podemos centrarnos en desarrollar un sistema de reconocimiento de entidades que unifique todos los existentes y los mejore, con ayuda del SNOMED y de los datos públicos disponibles.

Por otro lado, podemos centrarnos en el ámbito de la red neuronal que hemos construido. El GPT-2 no solo puede generar comentarios, puede resumir y puede clasificar texto. Esto permite que, con la correcta configuración, obtengamos, por ejemplo, un resumen de un informe médico de forma que no tengamos que leer todo el contenido para obtener toda la información.

Con ayuda del clasificador, se puede crear un sistema de recomendación y búsqueda para una base de datos muy poderoso. Buscando términos relativamente ambiguos, tal y como un humano preguntaría de forma natural, el sistema es capaz de devolver todos aquellos documentos relacionados, de forma que el acceso a los documentos en la base de datos es ahora mucho más fácil y natural.

En resumen, la asistencia que cualquiera de estas herramientas ofrecen puede suponer una mejora en la calidad de la atención que cada paciente recibe, además de la reducción de carga cognitiva que los correspondientes profesionales deben soportar, mejorando la calidad de vida de ambas partes. Con todo ello, se apunta a una mejora del sistema médico del que disponemos, del que ya de por sí podemos estar orgullosos siendo uno de los mejores en el mundo.

A. APÉNICE

Incluiremos en este apéndice todos los bloques de código o cualquier otro tipo de contenido necesario para comprender la estructura del documento, pero evitando que interfiera con la lectura del mismo.

A.1. Comentarios

Comentario 15 *Isolation of subcutaneous growth factor-alpha (SFL-alpha) in the stroma of lung carcinoma of the basolateral part of the rat lung and related tumors. We performed a prospective study to identify SFL-alpha-like growth factor-alpha (SFL-alpha) expression in the tissue and the body of lung carcinoma of the basolateral part of the rat lung and related tumors.*

Comentario 16 *Fibrous elastobrachial septal pressure syndrome in adults with delayed progressive disease. The prevalence of fibrous elastobrachial septal pressure syndrome in adult patients with delayed progressive disease is greater than 5 %. The duration of the disease is related to genetic factors, operative time, and duration of cure. The existence of this syndrome was examined in 103 children with the rare disease and in 17 patients with the rare disease.*

Comentario 17 *Intraperitoneal extracorporeal tone transplantation in children with renal echocardiography.*

Comentario 18 *Biology of anaphylaxis in cardiomyopathy. Two case reports. Case reports of patients treated with biofilms for intracranial pressure syndrome and histologic abnormalities were reviewed. Anaphylaxis in patients with cardiomyopathy was described, and histologic abnormalities were found in six patients treated with biofilms for intracranial pressure syndrome. Treatment of cardiomyopathy requires strict care with regard to the pharmacological and hormonal effects of the biofilms.*

Comentario 19 *Antibody of the prostate with abnormal lumen density. Ultrasound measurements were performed with respect to 30 age-matched patients on 11-year follow-up and for all clinical variables. Sixteen patients were identified by clinical examination as having abnormal lumen density; one patient was selected for transplant and one was selected for subsequent elective bone marrow transplantation. A total of 45 elective bone marrow transplantations were made.*

```

1 def regex_processing(text):
2     # Remove capital letters surrounded by 0 or more ` ` and a colon,
3     #   ↳ i.e. the titles
4     no_caps = re.sub(r'*,*([A-Z\s]+):', '', text)
5
6     # Remove weirdly positioned commas. Find commas that dont have any
7     #   ↳ letter before and some space after them.
8     weird_commas = re.sub(r'(?<!\w),\s+', '', no_caps)
9
10    # Remove commas that dont have spaces around them. (Commas should
11    #   ↳ always have a trailing space after them)
12    more_commas = re.sub(r'(?<!\s),(?! \s)', ' ', weird_commas)
13
14    # Remove digits adjacent to dots or commas, as in enumerated lists.
15    no_digits = re.sub(r'[\.,]*\d[\.,,]+', ' ', more_commas)
16
17    # Remove any other commas left behind the process. Particularly these
18    #   ↳ cases: Hello. ,How are you?
19    trailing_commas = re.sub(r'\s,(?=[A-Z\d])', ' ', no_digits)
20
21    # Substitute any number of spaces for 1 single space.
22    no_double_spaces = re.sub(r'\s+', ' ', trailing_commas)
23
24    # Solve these problems: Hello .How are you? => Hello. How are you?
25    final_text = re.sub(r'(?<!\s)\.(?! \s)', '. ', no_double_spaces)
26
27    # Finally, strip the text from any trailing commas or white spaces.
28    # The result is hopefully a clean version of the text, ready to be
29    #   ↳ tokenized
30    # and passed to the models.
31    return final_text.strip(' , ')

```

Figura A.1: Pipeline para el procesamiento de los comentarios de Medical Transcriptions

BIBLIOGRAFÍA

- [1] Guillermo Reynoso, Ernesto Martin-Jacod, María Carolina Berra, Olga Burlak, Patricia Houghton, and María Cecilia Vallese. SNOMED: la nomenclatura sistematizada de medicina del College of American Pathologists. *Panacea: boletín de medicina y traducción*, 4, 2003.
- [2] Ley Orgánica de Protección de Datos. *BOE 298*, pages 43088 – 43099, Diciembre 1999.
- [3] Luis Pereira, Rui Rijo, Catarina Silva, and Ricardo Martinho. Text mining applied to electronic medical records: A literature review. *International Journal of E-Health and Medical Communications (IJEHMC)*, 6:1–18, 07 2015.
- [4] G. Piatetsky-Shapiro. Discovery, analysis, and presentation of strong rules. In *Knowledge Discovery in Databases*, 1991.
- [5] R. Kumar and Rajesh Verma. Classification algorithms for data mining: A survey. *International Journal of Innovations in Engineering and Technology (IJIET)*, 1:7–14, 2012.
- [6] Anil K. Jain, M. N. Murty, and P. Flynn. Data clustering: a review. *ACM Comput. Surv.*, 31:264–323, 1999.
- [7] Jiawei Han, Micheline Kamber, and Jian Pei. 6 - mining frequent patterns, associations, and correlations: Basic concepts and methods. In Jiawei Han, Micheline Kamber, and Jian Pei, editors, *Data Mining (Third Edition)*, The Morgan Kaufmann Series in Data Management Systems, pages 243–278. Morgan Kaufmann, Boston, third edition edition, 2012.
- [8] Kamran Kowsari, K. Meimandi, Mojtaba Heidarysafa, Sanjana Mendu, Laura E. Barnes, and D. Brown. Text classification algorithms: A survey. *Inf.*, 10:150, 2019.
- [9] Mark Davis. Unicode text segmentation. Technical Report 29, 10 2010.
- [10] Umar Farooq, Hasan Mansoor, A. Nongaillard, Y. Ouzrout, and M. Qadir. Negation handling in sentiment analysis at sentence level. *J. Comput.*, 12:470–478, 2017.
- [11] Omar Ali, A. Gegov, Ella Haig, and R. Khusainov. Conventional and structure based sentiment analysis: A survey. *2020 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8, 2020.

- [12] Vimala Balakrishnan and Lloyd-Yemoh Ethel. Stemming and lemmatization: A comparison of retrieval performances. *Lecture Notes on Software Engineering*, 2:262–267, 01 2014.
- [13] Thomas Cover and Joy Thomas. *Elements of Information Theory*, volume 36, pages i – xxiii. 10 2001.
- [14] K. Jones. A statistical interpretation of term specificity and its application in retrieval. *J. Documentation*, 60:493–502, 2004.
- [15] Tomas Mikolov, Kai Chen, G. Corrado, and J. Dean. Efficient estimation of word representations in vector space. In *ICLR*, 2013.
- [16] Jeffrey Pennington, R. Socher, and Christopher D. Manning. Glove: Global vectors for word representation. In *EMNLP*, 2014.
- [17] Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics*, 5:135–146, 2017.
- [18] Oren Melamud, J. Goldberger, and Ido Dagan. context2vec: Learning generic context embedding with bidirectional lstm. In *CoNLL*, 2016.
- [19] I.T. Jolliffe. *Principal Component Analysis*. Springer Verlag, 1986.
- [20] A. Hyvärinen. Topographic independent component analysis. In *Encyclopedia of Computational Neuroscience*, 2014.
- [21] Stan Z. Li and Anil Jain, editors. *LDA (Linear Discriminant Analysis)*, pages 899–899. Springer US, Boston, MA, 2009.
- [22] Eibe Frank and Remco R. Bouckaert. Naive bayes for text classification with unbalanced classes. In *Proceedings of the 10th European Conference on Principles and Practice of Knowledge Discovery in Databases, ECMLPKDD’06*, page 503510, Berlin, Heidelberg, 2006. Springer-Verlag.
- [23] Georgiana Ifrim, Gökhan Bakir, and Gerhard Weikum. Fast logistic regression for text categorization with variable-length n-grams. In *Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, page 354362, New York, NY, USA, 2008. Association for Computing Machinery.
- [24] Thorsten Joachims. Text categorization with support vector machines: learning with many relevant features. In Claire Nédellec and Céline Rouveirol, editors, *Proceedings of ECML-98, 10th European Conference on Machine Learning*, number 1398, pages 137–142, Chemnitz, DE, 1998. Springer Verlag, Heidelberg, DE.

- [25] Chunting Zhou, Chonglin Sun, Zhiyuan Liu, and Francis C. M. Lau. A c-lstm neural network for text classification. *CoRR*, abs/1511.08630, 2015.
- [26] Eric Bauer and R. Kohavi. An empirical comparison of voting classification algorithms: Bagging, boosting, and variants. *Machine Learning*, 36:105–139, 2004.
- [27] Xiao-Peng Yu and Xiao-Gao Yu. Novel text classification based on k-nearest neighbor. In *2007 International Conference on Machine Learning and Cybernetics*, volume 6, pages 3425–3430, 2007.
- [28] Wan Noormanshah, Puteri Nohuddin, and Zuraini Zainol. Document categorization using decision tree: Preliminary study. *International Journal of Engineering and Technology*, 7:437–440, 12 2018.
- [29] Leo Breiman. Random forests. *Machine Learning*, 45(1):5–32, 2001.
- [30] Olivier Bodenreider. The unified medical language system (umls): integrating biomedical terminology. *Nucleic Acids Res.*, 32(Database-Issue):267–270, 2004.
- [31] Mark Neumann, Daniel King, Iz Beltagy, and Waleed Ammar. ScispaCy: Fast and Robust Models for Biomedical Natural Language Processing. In *Proceedings of the 18th BioNLP Workshop and Shared Task*, pages 319–327, Florence, Italy, August 2019. Association for Computational Linguistics.
- [32] Matthew Honnibal, Ines Montani, Sofie Van Landeghem, and Adriane Boyd. spaCy: Industrial-strength Natural Language Processing in Python, 2020.
- [33] Andrey Kormilitzin, Nemanja Vaci, Qiang Liu, and Alejo Nevado-Holgado. Med7: a transferable clinical natural language processing model for electronic health records, 2020.
- [34] Jiao Li, Yueping Sun, Robin J. Johnson, Daniela Sciaky, Chih-Hsuan Wei, Robert Leaman, Allan Peter Davis, Carolyn J. Mattingly, Thomas C. Wieggers, and Zhiyong Lu. BioCreative V CDR task corpus: a resource for chemical disease relation extraction. *Database*, 2016, 05 2016.
- [35] Pratyay Banerjee, Kuntal Kumar Pal, Murthy Devarakonda, and Chitta Baral. Knowledge guided named entity recognition for biomedical text, 2020.
- [36] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. *CoRR*, abs/1706.03762, 2017.
- [37] Warren Mcculloch and Walter Pitts. A logical calculus of ideas immanent in nervous activity. *Bulletin of Mathematical Biophysics*, 5:127–147, 1943.

- [38] Donald O. Hebb. *The organization of behavior: A neuropsychological theory*. Wiley, New York, June 1949.
- [39] Alex Sherstinsky. Fundamentals of recurrent neural network (rnn) and long short-term memory (lstm) network, 2018.
- [40] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. *CoRR*, abs/2005.14165, 2020.
- [41] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models, 2020.
- [42] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics.
- [43] Daniel Adiwardana, Minh-Thang Luong, David R. So, Jamie Hall, Noah Fiedel, Romal Thoppilan, Zi Yang, Apoorv Kulshreshtha, Gaurav Nemade, Yifeng Lu, and Quoc V. Le. Towards a human-like open-domain chatbot. *CoRR*, abs/2001.09977, 2020.
- [44] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019.
- [45] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen

- Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. Transformers: State-of-the-art natural language processing. pages 38–45, Online, October 2020. Association for Computational Linguistics.
- [46] Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration. 2020.
- [47] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: A method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*, page 311318, USA, 2002. Association for Computational Linguistics.
- [48] Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81, 2004.
- [49] Weizhe Yuan, Graham Neubig, and Pengfei Liu. Bartscore: Evaluating generated text as text generation, 2021.