



UNIVERSITAT OBERTA DE CATALUNYA (UOC)
MÁSTER UNIVERSITARIO EN CIENCIA DE DATOS (*Data Science*)

TRABAJO FINAL DE MÁSTER

ÁREA: 2

CLASIFICADOR DOCUMENTOS MÉDICOS HOPE

Autor: Rubén Vasallo Gonzalez

Profesor: Jordi Casas Roma

Tutor: Carlos Luis Sanchez Bocanegra

Co-Tutor: Rafael Pastor Vargas

Barcelona, 2 de diciembre de 2020

Créditos/Copyright



Esta obra está sujeta a una licencia de Reconocimiento - NoComercial - SinObraDerivada
3.0 España de Creative Commons.

FICHA DEL TRABAJO FINAL

Título del trabajo:	CLASIFICADOR DOCUMENTOS MÉDICOS HOPE
Nombre del autor:	Rubén Vasallo Gonzalez
Nombre del colaborador/a docente:	Carlos Luis Sanchez Bocanegra, Rafael Pastor Vargas
Nombre del PRA:	Jordi Casas Roma
Fecha de entrega (mm/aaaa):	01/2021
Titulación o programa:	Máster Universitario en Ciencia de Datos
Área del Trabajo Final:	M2.979 - TFM
Idioma del trabajo:	Español
Palabras clave	hope, clasificador, medicina

Dedicatoria/Cita

Quiero dedicarle este trabajo a mis mentores *Carlos Luis Sanchez Bocanegra* y *Rafael Pastor Vargas* que me han ayudado, apoyado y guiado en todo momento para conseguir los objetivos del máster.

Agradecimientos

Quiero agradecer a *Carlos Luis Sanchez Bocanegra* por invitarme participar en el Proyecto HOPE y poder aportar mi granito de arena a este gran proyecto.

También quiero dar la gracias a todos los miembros del proyecto HOPE que me han dado la bienvenida al grupo y me han facilitado la vida en unas circunstancias en las que, cuando entre en este, no eran las más idóneas. La mayoría de integrantes del grupo son médicos y la saturación de trabajo que había por el COVID-19 era enorme. Tengo claro que sin ellos y sin la ayuda en especial de Carlos y *Nicolas Passadore* que ha estado luchando para conseguir un dataset con más observaciones, este Máster no habría sido posible.

Abstract

This Final Master's Thesis *FMT* was born from the need to be able to have in a simple, up-to-date and immediate way, medical bibliographic references cataloged according to the information and the patient's symptoms, being able to make a ranking of more or less interest function of the feedback provided by health professionals on these bibliographic references.

Resumen:

Este Trabajo de Final de Máster (*TFM*) nace de la necesidad de poder disponer de una manera sencilla, actualizada e inmediata, referencias bibliográficas médicas catalogadas según la información y los síntomas que tienen del paciente, pudiendo hacer una clasificación (*ranking*) de más o menos interés en función de la valoración (*feedback*) aportada por los profesionales sanitarios sobre estas referencias bibliográficas.

Palabras clave: clasificador, artículos, médicos, PCA, KNN, Regresión Logística, Random Forest, SVM

Índice general

Abstract	IX
Índice	XI
Listado de Figuras	XIII
Listado de Tablas	1
1. Introducción	3
1.1. Descripción general del problema	3
1.2. Motivación personal	4
2. Objetivos del Máster	7
2.1. Objetivo principal	7
2.2. Objetivos secundarios	7
3. Metodología	9
3.1. Reuniones con el cliente	9
3.2. Limitaciones detectadas	9
3.3. Análisis de datos	10
3.4. Análisis de componentes principales	11
3.5. Enriquecimiento de los datos.	13
4. Planificación	15
4.1. Fechas importantes	15
4.2. Diagrama Gantt	15
5. Estado del arte	17
5.1. Recomendadores actuales	17
5.1.1. Recomendadores en el ámbito de la medicina	18

5.1.2. Recomendadores en el ámbito de la salud	19
5.2. Recomendadores en el futuro	22
6. Resultados	23
6.1. Análisis de componentes principales	23
6.1.1. Conjunto de datos completo	23
6.1.2. Conjunto de datos solo con atributo utilidad definido	24
6.1.3. Conjunto de datos solo con atributo utilidad definido, añadiendo el mes y año del artículo, eliminando los atributos <i>gender</i> y artículo y expandiendo el atributo <i>respuesta.pubmed_keys</i>	25
6.2. Enriquecimiento de los datos. Aproximación por Vecinos más próximos (K-NN)	26
6.2.1. Características de los conjuntos de datos a analizar	26
6.2.2. (Caso de estudio 1) Resultados del entrenamiento con el Conjunto de datos solo con atributo utilidad definido	27
6.2.3. (Caso de estudio 2) Conjunto de datos solo con atributo utilidad definido, añadiendo el mes y año del artículo, eliminando los atributos <i>gender</i> y artículo y expandiendo el atributo <i>respuesta.pubmed_keys</i>	28
6.2.4. Conclusiones del algoritmo K-NN	29
7. Conclusiones	31
7.1. Objetivos Secundarios	31
8. Anexos	33
8.1. Procesado de datos	33
8.2. Análisis de componentes principales	33
8.2.1. Iteración 1	33
8.2.2. Iteración 2	33
8.2.3. Iteración 3	33
8.3. Enriquecimiento de los datos	33
8.3.1. Iteración 1	33
8.3.2. Iteración 2	33
Bibliografía	34

Índice de figuras

4.1. Gantt del proyecto.	15
5.1. Web de Pubmed.	19
5.2. Vista principal de ClinicalTrials.	20
5.3. Resultados de búsqueda de ClinicalTrials.	20
5.4. Web de MedlinePlus.	21
6.1. Resultado del <i>PCA</i>	23
6.2. Proyección de los resultados del <i>PCA</i> en una gráfica.	24
6.3. Proyección de los resultados del <i>PCA</i> sobre los atributos del conjunto de datos.	24
6.4. Resultado del <i>PCA</i> de la segunda ejecución.	24
6.5. Proyección de los resultados del <i>PCA</i> sobre los atributos del conjunto de datos para la segunda ejecución.	25
6.6. Resultado del <i>PCA</i> de la tercera ejecución.	25
6.7. Proyección de los resultados del <i>PCA</i> sobre los atributos del conjunto de datos para la tercera ejecución.	26
6.8. Distribución del atributo a predecir (utilidad) en el conjunto de entrenamiento.	27
6.9. Distribución del atributo <i>pubmed_keys</i> en el conjunto de entrenamiento.	27
6.10. Modelo K-NN en el caso de estudio 1.	28
6.11. Modelo K-NN en el caso de estudio 2.	29

Índice de cuadros

4.1. <i>Tabla de fechas clave del proyecto.</i>	15
---	----

Capítulo 1

Introducción

1.1. Descripción general del problema

El Trabajo de Final de Máster (*TFM*) que aquí se presenta nace de la necesidad por parte de los profesionales sanitarios de poder tener la información más exacta posible sobre las mejores referencias bibliográficas actuales sobre tratamientos a aplicar a un paciente, dado unos síntomas concretos. Actualmente existe infinidad de referencias medicas que los profesionales sanitarios pueden consultar, pero esta información es tan abundante que acaba siendo engorrosa de consultar[1]. Esto hace que, muchas veces sea complicado encontrar la información sobre estas referencias bibliográficas para tratar algunas enfermedades. En el ámbito de la medicina el tiempo perdido puede costar vidas y es un precio demasiado elevado a pagar, tanto a nivel económico como emocional.

Una de las necesidades que nos transmite el cliente es que los profesionales sanitarios necesitan poder disponer de una plataforma que sea capaz de poder facilitarles estas referencias bibliográficas actuales lo más exactas y personalizadas que sean posible, ajustándose a la búsqueda de la información que poseen de sus pacientes, ya que ahora mismo, tienen que utilizar múltiples herramientas y buscadores para poder obtener esas referencias[2].

Bajo esa premisa nace el proyecto HOPE[3] (que significa *Health Operations for Personalized Evidence* en ingles) con el objetivo de ayudar a estos profesionales sanitarios a encontrar estas referencias bibliográficas que necesitan de la manera más rápida y fácil posible. Actualmente existen bases de datos de confianza en donde los profesionales sanitarios y el publico en general puede buscar estas referencias de ensayos y estudios clínicos con información fiable y desarrollados con anterioridad, pero no siempre es fácil o rápido encontrar estos resultados[4].

El proyecto HOPE es un sistema basado en inteligencia artificial para identificar la información de casos clínicos registrados en la Historia Clínica Electrónica, en base a los cuales realiza una búsqueda única por paciente para proporcionar al profesional sanitario recomendaciones de referencias bibliográficas donde constan tratamientos, estudios de investigación e información para ayudar al paciente. Todo en base a registros de fuentes científicas validas de información. En este proyecto, los profesionales sanitarios de todo el mundo puede consultar en una base de datos estas referencias y ver que otros tratamientos relacionados con los síntomas de sus pacientes han dado resultado. Todo y con eso, el sistema no siempre devuelve las referencias actuales más relevantes por lo que, no siempre la información consultada es útil[5].

Ademas, actualmente los profesionales sanitarios pueden valorar si la información recibida ha sido útil o no respecto a la búsqueda que han realizado, por lo que con esta valoración (*feedback*), se pretende mejorar el sistema actual complementándolo con un modelo clasificador capaz de ayudar al actual a entregar realmente las referencias que son más útiles basándose en la valoración que los profesionales sanitarios dan al sistema.

Este Trabajo de final de máster (*TFM*) pretende ayudar al proyecto HOPE a mejorar su algoritmo de Inteligencia Artificial para que los resultados se ajusten a las necesidades de información que requieren los profesionales sanitarios, en base a las búsquedas personalizadas que puedan hacer respecto a la información que tienen de sus pacientes. Para hacer esto se realizara un estudio de aproximación a conocer cual es el mejor modelo predictivo que puede ayudar a devolver esa información lo más exacta, fiable y actualizada posible.

1.2. Motivación personal

En los duros tiempos en los que estamos viviendo actualmente, tanto económica como emocionalmente, es grato ver como la humanidad es capaz de dejar a un lado sus diferencias y unirse para afrontar problemas comunes. En el caso del proyecto HOPE, lo que más me atrajo fue la oportunidad de poder ayudar a encontrar soluciones a enfermedades que ya están en el último paso (o como médicamente se le describe, en cuidados paliativos).

Es evidente que este proyecto no va a dar una solución para curar cualquier enfermedad, pero si puede ayudar a los profesionales sanitarios a poder encontrar posibles soluciones a enfermedades complejas, ayudando a estos profesionales a buscar en la infinidad de documentación medica que existe, el tratamiento que más pueda ayudar a paliar o, quien sabe, curar una enfermedad

que ya se daba por incurable. Esto mismo es lo que me motiva y mucho el poder ayudar ante estas situaciones.

También me ha motivado muchísimo el conocer a gente profesional que, independientemente del país al que pertenece o la profesión que tiene, se una al proyecto HOPE para participar y ayudar con sus conocimientos a hacer de este mundo, un lugar un poco mejor en el que vivir. Esta experiencia me esta enriqueciendo muchísimo personalmente y espero poder seguir contribuyendo al proyecto cuando este máster acabe.

Capítulo 2

Objetivos del Máster

2.1. Objetivo principal

OP - Poder recomendar al profesional sanitario las referencias bibliográficas actuales más exactas útiles y personalizadas que sean posibles, que pueden ayudar en el tratamiento del paciente, en base a la información que se dispone de este, pudiendo realizar una clasificación (*ranking*) de más interés a menos.

2.2. Objetivos secundarios

Para poder cumplir con el objetivo principal **OP1**, desglosaremos los siguientes objetivos secundarios:

OS1 - Extraer la información de la base de datos y tratarla para quedarnos solo con la que consideramos valida.

OS2 - Hacer un análisis de componentes principales (estudio de que atributos son relevantes para alcanzar el objetivo).

OS3 - Enriquecer de los datos (*data augmentation*) prediciendo los resultados que no están indicados si son relevantes o no. Aproximación por Vecinos más próximos (*K-Nearest-Neighbor*).

OS4 - Predecir los resultados usando el algoritmo de aprendizaje supervisado para clasificación llamado Regresión logística '*Logistic regression*'.

OS5 - Predecir los resultados usando el algoritmo de aprendizaje supervisado para clasificación llamado Bosques Aleatorios '*Random Forests*'.

OS6 - Predecir los resultados usando el algoritmo de aprendizaje supervisado para clasificación llamado Máquinas de vector soporte '*Support Vector Machines*'.

OS7 - Comparar los resultados obtenidos de los 3 modelos.

Capítulo 3

Metodología

3.1. Reuniones con el cliente

Para poder comprender y abordar con éxito el [objetivo principal](#) se ha realizado dos reuniones en donde el cliente expuso el [problema](#) a abordar y el origen de los datos para poder realizar el estudio.

En estas reuniones se pudo observar que los datos facilitados por el usuario requerían de una limpieza y tratamiento para poder cumplir el objetivo principal, ya que muchas observaciones tenían información poco relevante que podía generar ruido.

3.2. Limitaciones detectadas

Realizando un primer análisis visual, se detecto que los datos aportados por el cliente eran insuficientes para completar el [OP1](#), ya que solo se disponía de la información respecto de si una referencia bibliográfica había sido útil o no, pero no se disponía de la información suficientemente detallada para saber si había sido muy útil o poco útil para poder llegar a realizar una clasificación (*ranking*). El cliente nos comenta que en el momento actual no dispone de ese nivel de detalle.

Se acuerda con el cliente que intentara conseguir más volumen de información y con más detalle para poder realizar la clasificación (*ranking*) que necesita. Todo y con eso, se comenta que se realizara una aproximación para indicar si una referencia bibliográfica es útil o no dejando para mas adelante la opción de poder realizar la clasificación (*ranking*) si se consigue ese nivel de detalle por parte del cliente.

También se pudo comprobar que el cliente disponía de un volumen de observaciones bajo por lo que se planteo la posibilidad de, o intentar obtener más observaciones facilitadas por el cliente, o intentar enriquecer las observaciones actuales generando nuevos datos por aproximación a los reales.

Finalmente se decidió estudiar si era viable generar nuevos valores por aproximación, debido a que en el momento en que se trato el problema, el cliente no podía facilitar más datos. Si a lo largo del estudio, el cliente conseguía obtener nuevas observaciones, estas serian añadidas al estudio para aproximar mejor la solución final.

3.3. Análisis de datos

Tal y como comentamos en el capítulo de [Metodología](#), al recibir los datos por parte del cliente y realizar el primer análisis superficial, se detecto que el volumen de datos no era suficiente como para que los resultados que pudieran surgir del estudio, fueran concluyentes. Todo y con eso se acordó (tal y como se refleja en el [OS1](#)) hacer una valoración de los datos para ver si se podrían enriquecer de algún modo, mientras el cliente intenta obtener más volumen de datos.

Para ello, transformamos los datos desde el origen a un formato columnar en el que poder aplicar los modelos predictivos. Una vez tenemos los datos transformados, observamos que existen ciertos atributos que contienen listas de opciones como son los atributos *pubmed_keys* (que corresponde a las palabras clave que la API[6] de pubmed nos devuelve para esta observación), *articles* (que corresponde a los ids de los artículos relacionados con esa observación), *articles-RevisedYear* y *articlesRevisedMonth* (que corresponde a los años y meses de los artículos según están ordenados en el atributo *articles*). Para poder recomendar artículos útiles, necesitamos tener una observación por artículo, para poder analizar posteriormente de manera independiente si ese articulo fue útil o no para la observación a la que hace referencia, por lo que aplicaremos las transformaciones necesarias para acabar obteniendo la representación de un articulo con su correspondiente *feedback* por parte de los profesionales sanitarios. Todos los pasos para realizar las transformaciones se pueden consultar en el anexo ([Procesado de datos](#)).

Aplicados las transformaciones anteriores, nos quedan los siguientes atributos que describimos a continuación:

- **pedido.data.attributes.age:** Nos indica la edad del paciente.

- **pedido.data.attributes.diagnostic_main:** Nos indica el diagnostico principal que dio el profesional sanitario para la enfermedad que tiene el paciente
- **pedido.data.attributes.gender:** Nos indica el sexo del paciente
- **articulo:** Nos indica el identificador del articulo relacionado con el diagnostico principal.
- **respuesta.articlesRevisedYear:** Nos indica el año de revisión del articulo referenciado por el identificador del campo articulo.
- **respuesta.articlesRevisedMonth:** Nos indica el mes de revisión del articulo referenciado por el identificador del campo articulo.
- **respuesta.pubmed_keys:** Nos indica las palabras clave relacionadas con el articulo del campo articulo
- **utilidad:** Nos indica si un profesional sanitario ha considerado o no que el articulo referenciado en el campo articulo es útil para el diagnostico principal relacionado con la enfermedad del paciente. Este campo puede tener 3 valores: 1 para saber que el articulo es útil, 0 para indicar que no es útil y null para indicar que aun no se ha valorado.

3.4. Análisis de componentes principales

El análisis de componentes principales (o como se le conoce en ingles por '*Principal Component Analysis*' o *PCA*), es un componente fundamental en el análisis de los datos, ya que permite reducir el número de atributos de un conjunto de datos para eliminar el ruido que los posteriores análisis/modelos predictivos funcionen con mejor precisión[7].

Por poner un ejemplo sencillo, imaginemos que tenemos un conjunto de datos de modelos de coche con muchísimos atributos de estos, como por ejemplo, el nombre del modelo, el color, el número de puertas, la cilindrada y la potencia, entre otros. Probablemente si intentáramos analizar los datos en busca de cual es el modelos que menos consume y quisiéramos crear un modelo predictivo con este objetivo, podríamos observar a simple vista que tenemos atributos que no nos son necesarios y que generaría distracción (ruido) a la hora de conseguir nuestro objetivo, como es el caso del atributo color, o el nombre del modelo.

El *PCA* nos ayudara a encontrar cuales son los atributos más significativos del conjunto de datos para conseguir predecir el atributo que queremos[7], que en el caso que nos toca, es el atributo '*utilidad*'.

Es cierto que, en nuestro conjunto de datos, no tenemos un gran volumen de atributos, pero este proceso nos puede ayudar a eliminar atributos con poca relevancia, para así, poder simplificar el modelo predictivo final.

Para realizar el *PCA* nos ayudaremos de la librería *sklearn.decomposition* que ya nos ofrece implementada la lógica para ejecutarlo. Además realizaremos la transformación de todos los atributos de Categóricos (texto) a Continuos (números continuos). Este paso se realiza para que el *PCA* pueda realizar operaciones matemáticas sobre los valores de las observaciones. A este proceso se le conoce como factorización (*factorize*).

También estandarizaremos los valores a un rango de entre 1 y -1. Esto se realiza para igualar la importancia de todos los atributos, ya que en el paso anterior, al realizar la factorización, se nos puede dar el caso de tener valores muy altos (por ejemplo, al factorizar un atributo con 100 valores diferentes, se nos dará casos de observaciones que en un atributo tienen el valor 100, que puede ser más alto que otros valores que no se han transformado). Al realizar el *PCA*, si no se hace esta estandarización de los datos, los valores más altos se les dará más peso, pero no por eso pueden ser relevantes. Por lo que es imperativo el realizar esta estandarización.

Para realizar el *PCA* se han seguido 3 estrategias, para valorar que impacto tiene las observaciones según el conjunto de datos a analizar. Hemos realizado 3 iteraciones con el siguiente conjunto de datos:

- **1:** El dataset completo con las transformaciones mencionadas en el apartado [Análisis de datos](#)
- **2:** El mismo dataset de la anterior iteración pero cogiendo solo las observaciones que se ha informado el atributo utilidad.
- **3:** El mismo dataset realizado en la iteración 2 pero añadiendo el mes y año del artículo, eliminando los atributos *gender* (ya que todas las observaciones tienen el mismo valor) y artículo (que solo representa un identificador de un artículo). Además se expande el atributo *respuesta.pubmed.keys* para que exista una observación por cada keyword.

Una vez realizado esa transformación ya podemos ejecutar el PCA. Todos los pasos para realizar el PCA se pueden consultar en el anexo ([Iteración 1](#)) para la iteración 1, el anexo ([Iteración 2](#)) para la iteración 2 y el anexo ([Iteración 3](#)) para la iteración 3.

3.5. Enriquecimiento de los datos.

Debido a que tenemos un conjunto de datos con poco volumen, se acuerda con el cliente el intentar enriquecer los datos ([OS3](#)) con los conjuntos de datos que tenemos que no tienen informado el atributo a predecir.

Para realizar esta operación utilizaremos el algoritmo de Enriquecimiento por Aproximación de Vecinos más próximos (conocido por K-Nearest-Neighbor o K-NN). Este Algoritmo pretende asociar si un registro es de un conjunto de datos o del otro en base la a aproximación de otros resultados que si se conoce su resultado [\[8\]](#).

Para poder ejecutar el algoritmo de K-NN necesitamos realizar la transformación de todos los atributos Categóricos (texto) a Continuos (números continuos) igual que se realizo para el caso del [PCA](#). Esto es debido a que el algoritmo necesita tener valores numéricos para poder trabajar con los datos[\[9\]](#). Con esta transformación ya podemos realizar operaciones matemáticas con el conjunto de datos.

Para realizar el K-NN se han seguido 2 estrategias, para valorar que impacto tiene las observaciones según el conjunto de datos a analizar. Hemos realizado 2 iteraciones con el siguiente conjunto de datos:

- **1:** El dataset completo con las transformaciones mencionadas en el apartado [Análisis de datos](#) eliminando las observaciones que no tenían identificado el atributo edad.
- **2:** El mismo dataset realizado en la iteración 1 pero eliminando los atributos *gender* (ya que todas las observaciones tienen el mismo valor) y *articulo* (que solo representa un identificador de un articulo). Ademas se expande el atributo *respuesta.pubmed.keys* para que exista una observación por cada keyword.

Una vez realizado esa transformación ya podemos ejecutar el K-NN. Todos los pasos para realizar el K-NN se pueden consultar en el anexo ([Iteración 1](#)) para la iteración 1 y el anexo ([Iteración 2](#)) para la iteración 2.

Capítulo 4

Planificación

4.1. Fechas importantes

A continuación se detalla en la tabla 4.1 las fechas clave del proyecto.

Nombre	Fecha de inicio	Fecha de fin	Duración
Reunión inicial con el cliente	22/09/20	22/09/20	1
Redacción de objetivos	23/09/20	27/09/20	5
Análisis de mercado	28/09/20	18/10/20	21
Enriquecer el DataSet	19/10/20	01/11/20	14
Diseño del Modelo 1	02/11/20	15/11/20	14
Diseño del Modelo 2	16/11/20	29/11/20	14
Diseño del Modelo 3	30/11/20	13/12/20	14
Redacción de conclusiones	14/12/20	27/12/20	14
Preparación de la Defensa	28/12/20	01/01/21	5

Cuadro 4.1: *Tabla de fechas clave del proyecto.*

4.2. Diagrama Gantt

A continuación se muestra en la figura 4.1 el diagrama Gantt del proyecto.

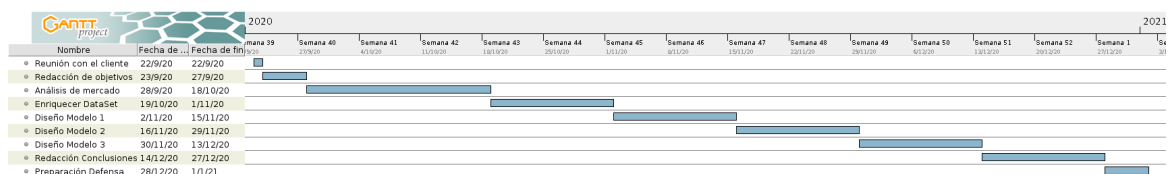


Figura 4.1: Gantt del proyecto.

Capítulo 5

Estado del arte

5.1. Recomendadores actuales

A día de hoy, existen multitud de clasificadores y recomendadores[10] que se usan a diario para realizar todo tipo de recomendaciones, que van desde cual es el mejor producto en base a unas necesidades, pasando por clasificar piezas o productos en una fabrica para organizar la disponibilidad de estas, hasta incluso a recomendadores de canciones, series o películas en base a unas características de quien solicita estas.

Esto se ha conseguido gracias a la investigación y al IOT (Internet de las cosas 'Internet of Things') que facilita la vida a los investigadores y sobre todo a las personas. Antiguamente eran los usuarios los que, si querían comparar algo, tenían que buscar y buscar por tiendas para encontrar, o bien el mismo producto a distintos precios o bien productos alternativos al que quería adquirir o que directamente desconocía de la existencia de esa alternativa. Esto le consumía muchas horas (o incluso días) y ademas no tenia la garantía de haber abarcado todas las posibilidades[11].

Hoy en día, a través de Internet y gracias a los Recomendadores/Clasificadores actuales, el usuario puede consultar todas las alternativas y precios y estar mas seguro de que hace una buena elección antes de hacer la inversión[12].

Los sistemas de recomendación/clasificación existen desde hace muchos años[13]. Han estado ahí para ayudarnos sin darnos cuenta, aunque al principio fueron muy rudimentarios. A medida que las características y el volumen de opciones crecía, se hacia mas complicado y lento realizar esa clasificación/recomendación. Y en ese punto es donde entro los algoritmos de aprendizaje

automático (*Machine Learning*) para ayudar a que esas clasificaciones y recomendaciones fueran posibles en un tiempo razonable.

Dos empresas que quizás merezcan mención especial (gracias a la fama que consiguieron al implementar estos algoritmos) son *Spotify*[14] con su recomendador de canciones y *Netflix*[15] con su recomendador de series. Ambos utilizan el poder de los algoritmos de recomendadores de *Machine Learning* para analizar el comportamiento que tienen sus usuarios en las correspondientes plataformas y sugerir recomendaciones basándose en estos comportamientos.

5.1.1. Recomendadores en el ámbito de la medicina

Pese a que en el ámbito de la medicina se ha utilizado mucho la estadística, el uso de estos recomendadores ha sido mas discreto, o al menos no se ha dado a conocer su uso al publico potencial, llegando a quedarse solo en la fase de caso de estudio sin llegar a convertirse en un producto final. Un ejemplo de esto es el recomendador de medicamentos personalizado[16] que se basa en el historial clínico del paciente para recomendar la medicación que se ha de tomar este. Este es un caso de uso muy similar a nuestro [OP1](#), aunque en nuestro caso, buscamos recomendar referencias bibliográficas clínicas.

También se conocen casos donde se han usado estos Clasificadores/Recomendadores para clasificar muestras o grupos de individuos para realizar un posterior análisis del grupo, como por ejemplo este caso de uso[17] que utiliza un clasificador para clasificar plaquetas para un posterior ensayo medico.

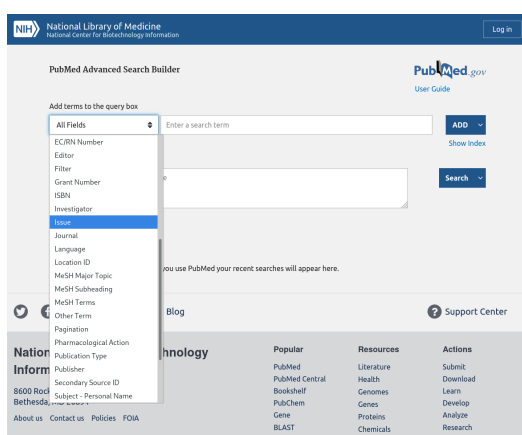
Todo y con eso, como comentamos anteriormente, el uso de Clasificadores/Recomendadores en este ámbito suele pasar a estar en segundo plano por lo que se ha de intuir su uso para darse cuenta que sin ellos, muchos estudios no serian posibles a día de hoy.

5.1.2. Recomendadores en el ámbito de la salud

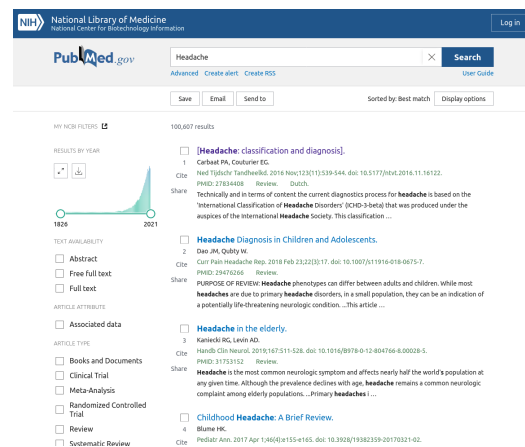
Al igual que en el caso anterior, el uso de recomendadores en el ámbito de la salud pasa desapercibido y es raro la vez que se destaque como solución final para la ayuda a los profesionales sanitarios.

A continuación destacamos los más conocidos en este ámbito:

● **Pubmed[18]:** Esta web Creada por el Centro Nacional de Biotecnología *National Center for Biotechnology Information (NCBI)* Permite hacer búsquedas por texto libre y permite realizar filtros por año, tipo de artículo y disponibilidad del texto. En el listado también se puede observar el estado de los artículos, si la publicación es definitiva o aun se esta revisando su contenido. El orden de los resultados es por aproximación al texto que se busque (mejor acierto), aunque permite modificar la ordenación por fecha del artículo, nombre de autor u origen del artículo. En las figuras 5.1a y 5.1b puede ver el aspecto de la pagina y los resultados de una búsqueda.



(a) Vista principal de Pubmed.



(b) Resultados de búsqueda de Pubmed.

Figura 5.1: Web de Pubmed.

● **MedlinePlus**[20]: Esta web creada por Biblioteca Nacional de Medicina de EE. UU. (NLM, por sus siglas en inglés), permite buscar artículos médicos basándose en texto relacionado. Por ejemplo si ponemos un tipo de dolencia en el buscador, este nos muestra artículos relacionados con ese texto como si de un buscador de Internet común se tratara, aunque a diferencia de estos, MedlinePlus promete brindar información de calidad y relevante de salud y bienestar que sea confiable y fácil de entender. Aparentemente este buscador muestra los resultados en orden de relevancia aunque, igual que en los casos anteriores, no permite valorar si esos artículos son realmente interesantes para la búsqueda relacionada. En las figuras 5.4a y 5.4b puede ver el aspecto de la pagina y los resultados de una búsqueda.



(a) Vista principal de MedlinePlus.



(b) Resultados de búsqueda de MedlinePlus.

Figura 5.4: Web de MedlinePlus.

5.2. Recomendadores en el futuro

Los recomendadores/clasificadores han ayudado mucho en la medicina, sobre todo en el ámbito de la investigación, en donde se requiere un trabajo previo para poder tener un conjunto de datos valido para el posterior estudio que se quiera realizar. En el caso que nos toca, ayudara a los profesionales sanitarios a poder encontrar otras referencias bibliográficas clínicas de otros casos similares al que tienen, para así poder valorar alternativas, cosa que sin estos recomendadores/clasificadores, seria un trabajo eterno de realizar, por no decir casi imposible, debido al gran volumen de información que cada día se genera. Estos algoritmos de recomendación/clasificación, no solo seguirán siendo un pilar básico en el ámbito medico, sino que seguramente aparecerán nuevos algoritmos de clasificación/recomendación que facilitaran las futuras investigaciones, como ya se esta observando en los nuevos recomendadores para el control de la glucosa en la sangre[\[21\]](#).

Capítulo 6

Resultados

6.1. Análisis de componentes principales

6.1.1. Conjunto de datos completo

Para este caso, al ejecutar el *PCA* con el conjunto de datos completo, este nos muestra tal y como podemos ver en la figura 6.1, que con solo 3 atributos, el modelo es capaz de explicar (predecir) el 95 % de las observaciones.

```
: print('explained variance ratio (first three components): %s' %  
    str(pca.explained_variance_ratio_))  
print('sum of explained variance (first three components): %s' %  
    str(sum(pca.explained_variance_ratio_)))  
  
explained variance ratio (first three components): [0.44726263 0.25669502 0.25009763]  
sum of explained variance (first three components): 0.9540552760689917
```

Figura 6.1: Resultado del *PCA*.

Dado que el *PCA* solo nos muestra 3 atributos como relevantes, podemos representar los resultados de las observaciones a un gráfico que mostramos en la figura 6.2 dibujando los 3 componentes principales como si fueran 3 dimensiones, y pintamos los valores a predecir sobre estos (aplicando el color rojo a las observaciones que no consideran útiles, y azul a los que si se consideran útiles), podemos observar que la gran mayoría de los rojo se encuentran en el plano del componente principal 1 y dos, mientras que el azul esta mas repartidos en los dos planos.

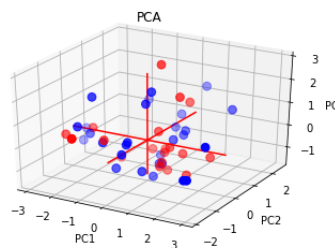


Figura 6.2: Proyección de los resultados del *PCA* en una gráfica.

Si le pedimos al modelo que nos indique cuáles son los 3 componentes que ha detectado que son los principales, este nos muestra los resultados que se pueden ver en la figura 6.3. Para saber cuáles son los componentes, tendremos que quedarnos para cada componente principal, el valor del atributo que más se acerque a 1.

```
pd.DataFrame(pca.components_, columns=features, index = ['PC1', 'PC2', 'PC3'])
```

	pedido.data.attributes.age	pedido.data.attributes.diagnostic_main	pedido.data.attributes.gender	respuesta.pubmed_keys	articulo
PC1	-0.050066	0.705007	-1.110223e-16	0.705462	0.052744
PC2	0.760475	0.071805	-1.318390e-16	0.030411	-0.644668
PC3	-0.630542	-0.130930	-1.110223e-16	0.142297	-0.751682

Figura 6.3: Proyección de los resultados del *PCA* sobre los atributos del conjunto de datos.

Para este caso, el *PCA* nos muestra que los 3 componentes principales son, el atributo *pubmed_keys*, el atributo *age* y el atributo *diagnostic_main*.

6.1.2. Conjunto de datos solo con atributo utilidad definido

En este caso, el *PCA* nos muestra tal y como podemos ver en la figura 6.4, que necesitamos 5 atributos para que el modelo sea capaz de explicar (predecir) el 97 % de las observaciones.

```
print('explained variance ratio (first three components): %s' %
      str(pca.explained_variance_ratio_))
print('sum of explained variance (first three components): %s' %
      str(sum(pca.explained_variance_ratio_)))
```

```
explained variance ratio (first three components): [0.32833981 0.23094825 0.17265137 0.14192254 0.09557452]
sum of explained variance (first three components): 0.9694364891147347
```

Figura 6.4: Resultado del *PCA* de la segunda ejecución.

Para este caso, el *PCA* nos muestra tal y como podemos ver en la figura 6.5, que los 5 componentes principales son, el atributo *diagnostic_main*, el atributo *Year*, el atributo *pubmed_keys*, el atributo *age* y el atributo *articulo*.

```
pd.DataFrame(pca.components_, columns=features, index = ['PC1', 'PC2', 'PC3', 'PC4', 'PC5'])
```

	pedido.data.attributes.age	pedido.data.attributes.diagnostic_main	pedido.data.attributes.gender	respuesta.articilesRevisedYear	respuesta.articilesRevisedMont
PC1	-0.094054	0.603912	-1.110223e-16	-0.365651	0.33837
PC2	0.135318	0.327845	-5.273559e-16	0.362710	-0.55606
PC3	-0.891259	-0.174252	-5.828671e-16	-0.236812	-0.13360
PC4	0.382620	-0.115174	-1.929013e-15	-0.621075	0.16245
PC5	0.098479	0.051453	3.774758e-15	-0.540677	-0.72925

PC1 => diagnostic_main PC2 => pubmed_keys/Year PC3 => pubmed_keys PC4 => age PC5 => articulo

Figura 6.5: Proyección de los resultados del *PCA* sobre los atributos del conjunto de datos para la segunda ejecución.

6.1.3. Conjunto de datos solo con atributo utilidad definido, añadiendo el mes y año del articulo, eliminando los atributos *gender* y *articulo* y expandiendo el atributo *respuesta.pubmed_keys*

En este caso, el *PCA* nos muestra tal y como podemos ver en la figura 6.6, que necesitamos 4 atributos para que el modelo sea capaz de explicar (predecir) el 90 % de las observaciones.

```
print('explained variance ratio (first three components): %s' %
      str(pca.explained_variance_ratio_))
print('sum of explained variance (first three components): %s' %
      str(sum(pca.explained_variance_ratio_)))
```

explained variance ratio (first three components): [0.31205207 0.24763151 0.19472659 0.14857077]
 sum of explained variance (first three components): 0.9029809474497172

Figura 6.6: Resultado del *PCA* de la tercera ejecución.

Para este caso, el *PCA* nos muestra como podemos ver en la figura 6.7, que los 4 componentes principales son, el atributo *diagnostic_main*, el atributo *Month*, el atributo *pubmed_keys* y el atributo *age*.

```
pd.DataFrame(pca.components_, columns=features, index = ['PC1', 'PC2', 'PC3', 'PC4'])
```

	pedido.data.attributes.age	pedido.data.attributes.diagnostic_main	respuesta.articlesRevisedYear	respuesta.articlesRevisedMonth	respuesta.pubmed_keys
PC1	0.202186	0.688327	-0.236152	-0.080358	0.650462
PC2	-0.299028	-0.000864	-0.648047	0.698266	-0.055149
PC3	-0.908866	-0.006936	0.153411	-0.221720	0.318152
PC4	-0.015928	-0.145721	-0.702451	-0.673071	-0.179024

PC1 => diagnostic_main / pubmed_keys PC2 => articlesRevisedMonth PC3 => pubmed_keys PC4 => age

Figura 6.7: Proyección de los resultados del *PCA* sobre los atributos del conjunto de datos para la tercera ejecución.

6.2. Enriquecimiento de los datos. Aproximación por Vecinos más próximos (K-NN)

6.2.1. Características de los conjuntos de datos a analizar

Una vez realizado todas las transformaciones necesarias para poder ejecutar el algoritmo de *K-NN* podemos observar que tanto el '*Conjunto de datos solo con atributo utilidad definido*' (Caso de estudio 1) como el '*Conjunto de datos solo con atributo utilidad definido, añadiendo el mes y año del artículo, eliminando los atributos gender y artículo y expandiendo el atributo respuesta.pubmed_keys*' (Caso de estudio 2) **se observa el atributo a predecir (utilidad) sesgado** (no tienen un balanceo de datos equilibrado en el atributo a predecir), teniendo más resultados con valor 0 que con valor a uno. Esto puede suponer que el modelo resultante tienda a estar mas sobre ajustado hacia los resultados con valor 1 que hacia los resultados con valor 0 (debido a que tiene más datos en los que basarse para dar el resultado 1 frente a los otros datos como podemos ver en la figura 6.8).

También, para el **caso de estudio 2**, se da el caso que la distribución del **atributo pubmed_keys**, una vez expandido, también **se observa sesgado** mostrando algunos resultados con pocas *keywords* y otros como '*Diuresis*' o '*Abdomen*' con muchas observaciones como podemos ver en la figura 6.9. Esto puede condicionar el comportamiento del modelo final igual que en el caso anterior, dando sospechas a que el modelo podría estar sobre ajustado[22].

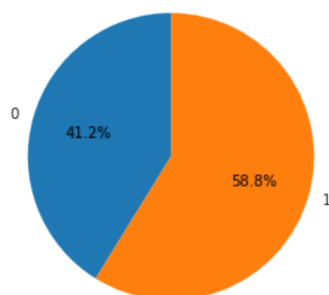


Figura 6.8: Distribución del atributo a predecir (utilidad) en el conjunto de entrenamiento.

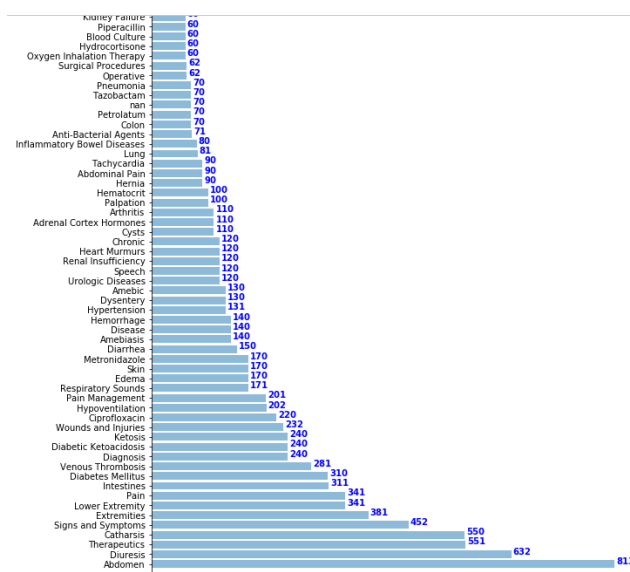


Figura 6.9: Distribución del atributo *pubmed_keys* en el conjunto de entrenamiento.

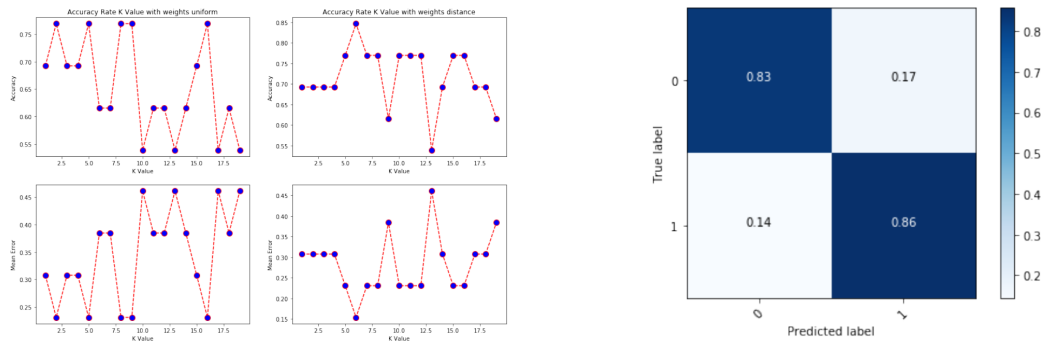
6.2.2. (Caso de estudio 1) Resultados del entrenamiento con el Conjunto de datos solo con atributo utilidad definido

Antes de poder entrenar el algoritmo *K-NN* es necesario realizar un estudio para detectar a partir de que número de vecinos mas próximos (valor *K*), podemos dar por supuesto que el resultado esta relacionado con ellos. El algoritmo también nos permite poder decidir como calcular la distancia entre los vecinos (bien dando más relevancia a los vecinos con menos distancia '*distance*'[23] o bien dando la misma relevancia a todos los vecinos que están próximos '*uniform*'[23]).

Para realizar este estudio ejecutaremos el modelo con diferentes valores de *K* y contrastaremos que porcentaje de acierto tiene el modelo con un subconjunto del propio conjunto de datos que conocemos su valor a predecir correcto. Como podemos ver en la figura 6.10a, el valor más optimo de *K* es 6 utilizando el calculo de la distancia '*distance*', con un **porcentaje de acierto**

del 85 %.

Si mostramos la matriz de confusión (figura 6.10b) para poder ver que contrastar que porcentaje de aciertos tiene el modelo con el conjunto de test, vemos que tiene un porcentaje de aciertos bastante aceptable pese a tener mas acierto para los casos con resultado 1 vrs. los resultados con valor 0 como comentamos anteriormente.



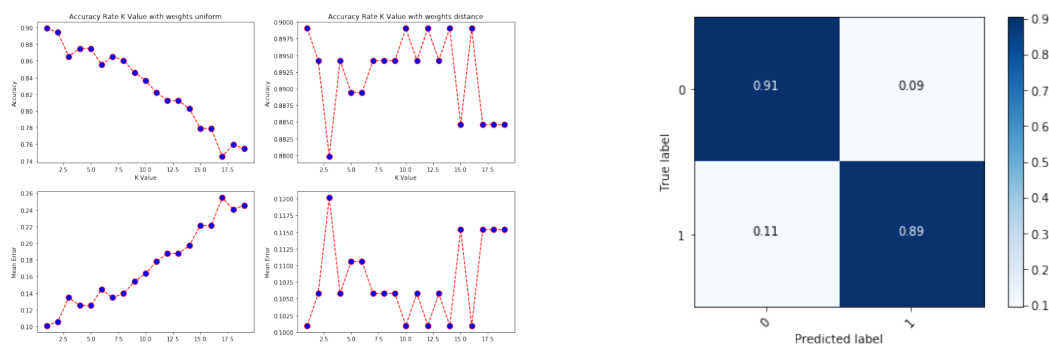
(a) Calculo del valor K para el caso de estudio 1. (b) Matriz de confusión para el modelo K-NN en el caso de estudio 1.

Figura 6.10: Modelo K-NN en el caso de estudio 1.

6.2.3. (Caso de estudio 2) Conjunto de datos solo con atributo utilidad definido, añadiendo el mes y año del articulo, eliminando los atributos *gender* y articulo y expandiendo el atributo *respuesta.pubmed_keys*

Igual que en el caso anterior, realizaremos el estudio para detectar a partir de que número de vecinos mas próximos (valor K), podemos dar por supuesto que el resultado esta relacionado con ellos. A diferencia del caso anterior, podemos ver en la figura 6.11a, que el valor más optimo de K es 1 utilizando el calculo de la distancia '*uniform*', con un **porcentaje de acierto del 90 %**.

Si mostramos la matriz de confusión (figura 6.11b) para poder ver que contrastar que porcentaje de aciertos tiene el modelo con el conjunto de test, vemos que tiene un porcentaje de aciertos bastante elevado por lo que podemos afirmar que **el modelo es sospechoso de estar sobre ajustado**[22].



(a) Calculo del valor K para el caso de estudio 2. (b) Matriz de confusión para el modelo K-NN en el caso de estudio 2.

Figura 6.11: Modelo K-NN en el caso de estudio 2.

6.2.4. Conclusiones del algoritmo K-NN

Pese a que aparentemente, el caso de estudio 1 tiene unos resultados aceptables, **aconsejamos** al cliente **no utilizar el algoritmo de K-NN** para enriquecer los datos del dataset, prediciendo las observaciones que no tienen informado el atributo a predecir. Esto es **debido** a que detectamos que **el conjunto de datos se encuentra sesgado** en varios atributos, lo que hará que el modelo tienda a dar falsos positivos (como hemos podido ver en el caso de estudio 2), perjudicando el entrenamiento de posteriores modelos predictivos que estudiaremos en este documento.

Capítulo 7

Conclusiones

7.1. Objetivos Secundarios

En el apartado [Análisis de datos](#) se ha cumplido el [OS1](#), aplicado las transformaciones necesarias a los datos que tenemos para poder trabajar posteriormente con ellos con los modelos predictivos.

En el apartado [Análisis de componentes principales](#) se ha cumplido el [OS2](#), pudiendo observar en el apartado de [resultados](#), que hay mucha diferencia según si se usa todas las observaciones o solo las que tienen informado el atributo a predecir (*utilidad*), por lo que se decide junto con el cliente, intentar enriquecer los datos que no tienen informado el atributo (*utilidad*) definiendo este paso como objetivo secundario [OS3](#). Se analizara en el siguiente paso si el resultado de intentar enriquecer estos datos aporta algo al conjunto de datos.

En el apartado [Enriquecimiento de los datos](#) no se ha podido cumplir el [OS3](#), pudiendo observar en el apartado de [resultados](#) que el dataset se encuentra [sesgado](#), por lo que intentar enriquecer los datos lo único que generaría son falsos positivos, perjudicando a los posteriores entrenamientos de modelos predictivos que se estudian en este documento.

Capítulo 8

Anexos

8.1. Procesado de datos

[S0 - HOPE extract data.pdf](#)

8.2. Análisis de componentes principales

8.2.1. Iteración 1

[S1 - HOPE PCA.pdf](#)

8.2.2. Iteración 2

[S1 - HOPE PCA_v2.pdf](#)

8.2.3. Iteración 3

[S1 - HOPE PCA_v3.pdf](#)

8.3. Enriquecimiento de los datos

8.3.1. Iteración 1

[S2 - HOPE kNN.pdf](#)

8.3.2. Iteración 2

[S2 - HOPE kNN_v2.pdf](#)

Bibliografía

- [1] R. J. W. Cline and K. M. Haynes, “Consumer health information seeking on the internet: the state of the art.” <https://academic.oup.com/her/article/16/6/671/571640>, 12 2001.
- [2] C. AR, “Los buscadores en la recuperación de información en salud.” <https://www.medigraphic.com/cgi-bin/new/resumen.cgi?IDARTICULO=35274>, 2011.
- [3] “Hope project website.” <http://proyectohope.centromedicodespierta.es/>.
- [4] C. L. S. Bocanegra, J. L. S. Ramos, C. Rizo, A. Civit, and L. Fernandez-Luque, “Healthrecsys: A semantic content-based recommender system to complement health videos.” <https://bmcmmedinformdecismak.biomedcentral.com/articles/10.1186/s12911-017-0431-7>, 05 2017.
- [5] C. Vargas and K. Azucena, “Hope desarrollo de una metodología piloto para vinculación a clinical trials como herramienta de medicina basada en evidencias, para proporcionar medicina centrada en el paciente.” http://discovery.uoc.edu/iii/encore/record/C__Rx1033890__Orightresult__U__X2?lang=spi&suite=def, 07 2017.
- [6] “Documentación api pubmed.” <https://www.ncbi.nlm.nih.gov/CBBresearch/Lu/Demo/tmTools/RESTfulAPIs.html>.
- [7] “Definición de pca (análisis de componentes principales.” https://es.wikipedia.org/wiki/An%C3%A1lisis_de_componentes_principales.
- [8] “Definición de k-nn (k vecinos más próximos).” https://es.wikipedia.org/wiki/K_vecinos_m%C3%A1s_pr%C3%B3ximos.
- [9] “Why is scaling required in knn and k-means?.” <https://medium.com/analytics-vidhya/why-is-scaling-required-in-knn-and-k-means-8129e4d88ed7>.

- [10] G. Adomavicius and A. Tuzhilin, “Toward the next generation of recommender systems: a survey of the state-of-the-art and possible extensions.” <https://ieeexplore.ieee.org/abstract/document/1423975>, 04 2005.
- [11] S.-J. Yoon and J.-H. Kim, “Is the internet more effective than traditional media? factors affecting the choice of media.” <http://www.journalofadvertisingresearch.com/content/41/6/53>, 11 2001.
- [12] C. MaféS and S. Blas, “Implicaciones del uso de buscadores en el comportamiento de compra online.” <https://www.sciencedirect.com/science/article/pii/S1135252312601014>, 02 2009.
- [13] V. M. and T. K., *History and Overview of the Recommender Systems*. Vishal Bhatnagar, IGI Global, 2017.
- [14] A. aron van den Oord, S. Dieleman, and B. Schrauwen, “Deep content-based music recommendation.” <http://papers.nips.cc/paper/5004-deep-content-based-music-recommendation.pdf>, 2013.
- [15] D. Chong, “Deep dive into netflix’s recommender system.” <https://towardsdatascience.com/deep-dive-into-netflixs-recommender-system-341806ae3b48>, 04 2020.
- [16] S. Bhoi, L. M. Li, and W. Hsu, “Premier: Personalized recommendation for medical prescriptions from electronic records.” <https://arxiv.org/abs/2008.13569>, 08 2020.
- [17] J. C. Clauser, J. Maas, J. Arens, T. Schmitz-Rode, U. Steinseifer, and B. Berkels, “Automation of hemocompatibility analysis using image segmentation and a random forest.” <https://arxiv.org/abs/2010.06245>, 10 2020.
- [18] “Pubmed website.” <https://pubmed.ncbi.nlm.nih.gov/advanced/>.
- [19] “Clinicaltrials website.” <https://clinicaltrials.gov>.
- [20] “Medlineplus website.” <https://medlineplus.gov>.
- [21] T. Yamagata, A. O’Kane, A. Ayobi, D. Katz, K. Stawarz, P. Marshall, P. Flach, and R. Santos-Rodríguez, “Model-based reinforcement learning for type 1diabetes blood glucose control.” <https://arxiv.org/abs/2010.06266>, 10 2020.
- [22] “Que es el sobreajuste (en inglés overfitting).” <https://es.wikipedia.org/wiki/Sobreajuste>.

-
- [23] “Documentación del algoritmo kneighborsclassifier.” <https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html>.