



UTPL

UNIVERSIDAD TÉCNICA PARTICULAR DE LOJA
La Universidad Católica de Loja

ÁREA TÉCNICA

**INGENIERO EN SISTEMAS INFORMÁTICOS Y
COMPUTACIÓN**

TRABAJO DE TITULACIÓN

Método de enriquecimiento de Grafos de Conocimiento de indicadores de ODS usando un enfoque de inferencia de entidades semánticas equivalentes desde fuentes de datos abiertos

Autor: Pinto Orellana, Diego Fernando

Director: Piedra Pullaguari, Nelson Oswaldo

LOJA - ECUADOR
2021



Esta versión digital, ha sido acreditada bajo la licencia Creative Commons 4.0, CC BY-NC-SA: Reconocimiento-No comercial-Compartir igual; la cual permite copiar, distribuir y comunicar públicamente la obra, mientras se reconozca la autoría original, no se utilice con fines comerciales y se permiten obras derivadas, siempre que mantenga la misma licencia al ser divulgada. <http://creativecommons.org/licenses/by-nc-sa/4.0/deed.es>

2021

Aprobación del director del Trabajo de Titulación

Loja, 01, de febrero, 2021

Magíster

Fernanda Maricela Soto Guerrero

Coordinadora de carrera de Sistemas Informáticos y Computación

Ciudad. -

De mi consideración:

El presente trabajo de titulación denominado: Método de enriquecimiento de Grafos de Conocimiento de indicadores de ODS usando un enfoque de inferencia de entidades semánticas equivalentes desde fuentes de datos abiertos, realizado por Diego Fernando Pinto Orellana, ha sido orientado y revisado durante su ejecución, por cuanto se aprueba la presentación de este. Así mismo, doy fe que dicho trabajo de titulación ha sido revisado por la herramienta antiplagio institucional.

Particular que comunico para los fines pertinentes.

Atentamente,

Firma:.....

Nelson Oswaldo Piedra Pullaguari.

C.I: 1102809462

Declaración de autoría y cesión de derechos

“Yo, Diego Fernando Pinto Orellana, declaro y acepto en forma expresa lo siguiente:

- Ser autor(a) del Trabajo de Titulación denominado: Nombre del trabajo, de la Titulación Método de enriquecimiento de Grafos de Conocimiento de indicadores de ODS usando un enfoque de inferencia de entidades semánticas equivalentes desde fuentes de datos abiertos de la Titulación Sistemas Informáticos y Computación, específicamente de los contenidos comprendidos en: Capítulo 1. Estado del Arte, Capítulo 2. Propuesta, Capítulo 3. Desarrollo de la propuesta, Capítulo 4. Pruebas y Resultado, siendo Nelson Oswaldo Piedra Pullaguari; y, en tal virtud, eximo expresamente a la Universidad Técnica Particular de Loja y a sus representantes legales de posibles reclamos o acciones judiciales o administrativas, en relación a la propiedad intelectual. Además, ratifico que las ideas, conceptos, procedimientos y resultados vertidos en el presente trabajo investigativo son de mi exclusiva responsabilidad.
- Que mi obra, producto de mis actividades académicas y de investigación, forma parte del patrimonio de la Universidad Técnica Particular de Loja, de conformidad con el artículo 20, literal j), de la Ley Orgánica de Educación Superior; y, artículo 91 del Estatuto Orgánico de la UTPL, que establece: “Forman parte del patrimonio de la Universidad la propiedad intelectual de investigaciones, trabajos científicos o técnicos y tesis de grado que se realicen a través, o con el apoyo financiero, académico o institucional (operativo) de la Universidad”.
- Autorizo a la Universidad Técnica Particular de Loja para que pueda hacer uso de mi obra con fines netamente académicos, ya sea de forma impresa, digital y/o electrónica o por cualquier medio conocido o por conocerse, sirviendo el presente instrumento como la fe de mi completo consentimiento; y, para que sea ingresada al Sistema Nacional de Información de la Educación Superior del Ecuador para su difusión pública, en cumplimiento del artículo 144 de la Ley Orgánica de Educación Superior.

Firma:

Autor: Diego Fernando Pinto Orellana

C.I.: 1105642332

Dedicatoria

Este trabajo lo dedico:

A Dios por haberme dado la fortaleza en los momentos difíciles y haberme permitido llegar hasta este momento tan importante de mi formación profesional.

A mis padres, quienes me han acompañado durante todo este transcurso de mi vida, por sus consejos y apoyo incondicional.

A mi hermana, quien siempre me ha apoyado y motivado para seguir adelante para lograr este gran sueño.

A mis abuelitas, tíos, primos y a toda mi familia, por haberme brindado su apoyo incondicional y por compartir conmigo buenos y malos momentos.

Diego Fernando Pinto Orellana

Agradecimiento

A Dios por guiar mi camino y haberme dado fuerzas para saber llevar cada obstáculo y las dificultades a lo largo de mi formación.

A Nelson Piedra, por brindarme las oportunidades que he tenido en varias etapas de mi formación, por la confianza en mi capacidad para desarrollar cada una de las ideas que tiene en mente, por dedicar su tiempo para revisar y corregir mis errores y enseñarme como tener una mejor versión de mi trabajo.

Agradezco la confianza y el apoyo de mis padres para alcanzar mis metas. A mi hermana, abuelitas, tíos y a toda mi familia por haberme brindado muchos momentos de alegría.

A todos mis docentes quienes fueron parte de mi formación académica, por compartir su conocimiento y experiencias para mi formación profesional.

A todos mis amigos/as, que siempre me animaron a ser mejor cada día y se alegran de cada uno de mis logros.

Índice de Contenido

Carátula.....	I
Aprobación del director del Trabajo de Titulación	II
Declaración de autoría y cesión de derechos	III
Dedicatoria	V
Agradecimiento	VI
Índice de Contenido	VII
Resumen	1
Abstract	2
Introducción.....	3
Capítulo uno.....	5
Estado del Arte.....	5
1.1 Revisión Sistemática	5
1.1.1 Metodología.....	6
1.1.2 Definiendo las preguntas de investigación	6
1.1.3 Búsqueda y selección	7
1.1.4 Extracción de artículos y documentos	8
1.1.5 Evaluación de los datos extraídos	9
1.1.6 Año de publicación y distribución geográfica.....	11
1.1.7 Tipos de publicaciones.....	12
1.1.8 Áreas de aplicación o investigación.....	13
1.1.9 Discusión y Resultados.....	14
1.2 Conceptos	18
1.2.1 Datos abiertos	18
1.2.2 Objetivos de desarrollo sostenibles.	23
1.2.3 Web Semántica	24
1.2.4 Procesamiento de lenguaje natural	29
1.2.5 Base de datos semántica con GraphDB.....	30
1.2.6 Discusión.....	33
Capítulo dos.....	34
Propuesta.....	34
2.1 Contexto	34
2.2 Problema	34
2.3 Propuesta de solución	35
2.4 Componentes arquitectónicos	36
Capítulo tres.....	38
Desarrollo de la propuesta.....	38
3.1 Configuración y extracción	38

3.1.1	Creación de repositorios y carga de información	38
3.1.2	Anotaciones semánticas de los objetivos de desarrollo sostenibles	44
3.1.3	Transformación de información de ODS	45
3.2	Análisis - Creación de enlaces entre ods y dataset semánticos	48
3.2.1	Traducción.....	48
3.2.2	NER y NLP	49
3.2.3	Creación de nuevas relaciones.....	57
3.3	Visualización	59
Capítulo Cuatro		60
Pruebas y Resultados		60
4.1	Casos de aplicación	60
4.2	Visualización del resultados obtenidos del enriquecimiento datos	65
4.2.1	Sitio web.....	65
4.2.2	Otros Resultados	73
4.3	Discusión	75
Conclusiones.....		77
Recomendaciones		79
Referencias		80
Apéndice		83

Índice de Tablas

Tabla 1	8
Tabla 2	9
Tabla 3	21

Índice de Figuras

Figura 1	7
Figura 2	10
Figura 3	11
Figura 4	12
Figura 5	13
Figura 6	14
Figura 7	21
Figura 8	25
Figura 9	31
Figura 10	32
Figura 11	36
Figura 12	39
Figura 13	40
Figura 14	40
Figura 15	41
Figura 16	42
Figura 17	42
Figura 18	43
Figura 19	45
Figura 20	46
Figura 21	46
Figura 22	49

Figura 23	51
Figura 24	52
Figura 25	56
Figura 26	56
Figura 27	57
Figura 28	58
Figura 29	61
Figura 30	62
Figura 31	63
Figura 32	64
Figura 33	66
Figura 34	67
Figura 35	68
Figura 36	68
Figura 37	70
Figura 38	70
Figura 39	71
Figura 40	72
Figura 41	73
Figura 42	74
Figura 43	74
Figura 44	83
Figura 45	84
Figura 46	85
Figura 47	86
Figura 48	87
Figura 49	88
Figura 50	89
Figura 51	90
Figura 52	91
Figura 53	92
Figura 54	93

Resumen

El objetivo de este trabajo es crear un método para enriquecer un grafo de conocimiento relacionado con los Objetivos de Desarrollo Sostenible (ODS). Los ODS buscan garantizar que las personas gocen de paz y prosperidad. Este trabajo describe una propuesta para el enriquecimiento de grafos de conocimiento relacionados con datos ODS. El contexto de fuentes de información y ODS se aborda a través de la gestión de datos heterogéneos y distribuidos a través del modelo de datos orientados a grafos. Los grafos aportan mayor nivel de abstracción que facilita la descripción semántica de información, en términos de nodos y relaciones, tanto del modelo de datos, como de las interacciones que se producen entre componentes de conjuntos de datos locales como distribuidos.

El propósito central del trabajo consiste en incorporar anotaciones a un grafo de conocimiento semántico con información relacionada con ODS, de manera que enriquezca la descripción semántica de los nodos con componentes de grafos semánticos externos. Finalmente, la propuesta desarrollada para este trabajo implementa una aplicación web piloto que permite interactuar y acceder a los recursos enriquecidos.

Palabras claves: Objetivos de Desarrollo Sostenible, Enriquecimiento de grafos, Web Semántica.

Abstract

The objective of this work is to create a method to enrich a knowledge graph related to the Sustainable Development Goals (SDGs). The SDGs aim to ensure that people enjoy peace and prosperity. This paper describes a proposal for the enrichment of knowledge networks related to SDGs data. The context of information sources and SDGs is addressed through the management of heterogeneous and distributed data via the graph-oriented data model. Networks provide a higher level of abstraction that facilitates the semantic description of information, in terms of nodes and relationships, both of the data model and of the interactions that occur between components of local and distributed datasets.

The central purpose of the work consists of incorporating annotations to a semantic knowledge graph with information related to SDGs, in a way that enriches the semantic description of the nodes with components of external semantic graphs. Finally, the proposal developed for this work implements a pilot web application that allows interacting and accessing the enriched resources.

Keywords: Sustainable Development Goals, Graphics Enrichment, Semantic Web.

Introducción

Los portales de datos de acceso público tienen como misión recopilar y compartir datos de múltiples fuentes y de una variedad de dominios de información. Los portales de datos abiertos son claves en acciones que contribuyan a la transparencia, el mejoramiento de la gestión participativa, la innovación, la promoción de la colaboración ciudadana y académica, entre otras características importantes. Existen avances importantes relacionados con la apertura y el reuso: licenciamiento abierto, uso de estándares abiertos, mejoras en las capacidades de integración e interoperabilidad de datos, buenas prácticas para mejorar la calidad de datos, la accesibilidad, mejorar la colaboración, el uso y reutilización de información abierta.

Los ODS comprenden 17 objetivos y 169 metas, y un conjunto preliminar de 330 indicadores, propuestos en marzo de 2015 y con un plazo para que se cumpla la resolución hasta el año 2030. Los portales de datos relacionados con los Objetivos de Desarrollo Sostenible (ODS) plantea retos en torno al acceso e interoperabilidad de datos a través de los diferentes mecanismos e interfaces de acceso a la información.

Actualmente no existe un proceso único que permita conectar conjuntos de datos abiertos relacionados con ODS con otras fuentes externas. Este trabajo busca crear un método de enriquecimiento de grafos de tal manera poder observar y determinar la calidad del grafo de conocimiento construido a partir de conjuntos de datos abiertos relacionado con los ODS, y que se publican a través de la plataforma CKAN. El trabajo se enfoca en desarrollar un método para enriquecer un grafo de conocimiento semántico, e implementa una interfaz de navegación Web sobre el grafo que forma parte del alcance de este trabajo. El sistema de enriquecimiento de grafos de conocimiento proporciona una forma de generar candidatos para un nodo y que puede ser desplegado en el grafo existente, y también una nueva forma de conectar nodos a través de equivalencia semántica inferida a través del enfoque propuesto.

El resultado del enriquecimiento de datos de este trabajo se mostrará a través de una aplicación Web. Se ha validado que la aplicación sea de fácil uso y acceso.

Este trabajo inició con una revisión sistemática relacionada con datos, grafos de conocimiento semántico, interoperabilidad y enriquecimiento. En la revisión sistemática de literatura se obtuvo información que permitió construir una forma para convertir recursos abiertos de múltiples fuentes en información interoperable y esta pueda ser enriquecida con más información. Entonces mediante un correcto análisis y procesamiento de datos se puede mejorar la accesibilidad, uso y reutilización.

Capítulo uno

Estado del Arte

1.1 Revisión Sistemática

El crecimiento intensivo de la información heterogénea, que se encuentra en la Web, plantea la necesidad de contar con maneras efectivas para mejorar la capacidad de descubrimiento, re-uso, analítica, personalización, y en general de obtención de valor.

A través de la Web, se acceden a diferentes plataformas de información abierta. Algunas de estas plataformas están destinadas a la publicación de datos abiertos (en inglés, Open Data - OD). Estas plataformas están relacionadas con la recolección y la distribución de información de diferente tipo, recolectada por diferentes instituciones. Las plataformas de datos abiertos publican información en forma de conjuntos de datos (denominados en inglés como *datasets*). Los *datasets* pueden ser de mucha ayuda en tareas de re-uso de información. RDF es un método para describir semánticamente piezas de conocimiento, a partir de fuentes de información estructurada o no estructurada. RDF se ha convertido en el estándar de la descripción semántica de recursos de la Web Semántica. Con RDF se puede tener un conjunto de *datasets* y estos conformarían repositorios semánticos como un gran conocimiento.

La gestión de la información contenida actualmente en la Web avanza hacia un paradigma de Web Semántica (WS). La WS busca que cada dato se pueda extraer, analizar, publicar y enlazar entre varios datos. De hecho, el interés en esta área sigue siendo alta debido a la abundancia de aplicaciones prácticas que ayudan a los usuarios a lidiar con diferentes tipos de información.

A lo largo de este tiempo en la web ha ido generando grandes volúmenes de información y a su vez esta siendo tratada con el fin de obtener conocimiento de ellos de diferente índole. Lo que se busca es con este trabajo de titulación es establecer un método para enriquecer un grafo de conocimiento con diferentes fuentes externas.

1.1.1 **Metodología**

La revisión sistemática es realizó con las siguientes pautas proporcionadas por (Okoli & Schabram, 2010, p. 7) y (Kitchenham & Charters, 2007) con la finalidad de buscar documentos relevantes de investigación. Los resultados de las búsquedas permitirán identificar desafíos o direcciones futuras en la investigación, así proporcionando un método verificable para resumir la búsqueda de información en la orientación de métodos para enriquecer grafos de conocimiento. La **Figura 1** presenta la metodología a seguir para la revisión sistemática.

Las siguientes subsecciones describen los pasos que se llevaron a cabo en la revisión sistemática del trabajo de titulación.

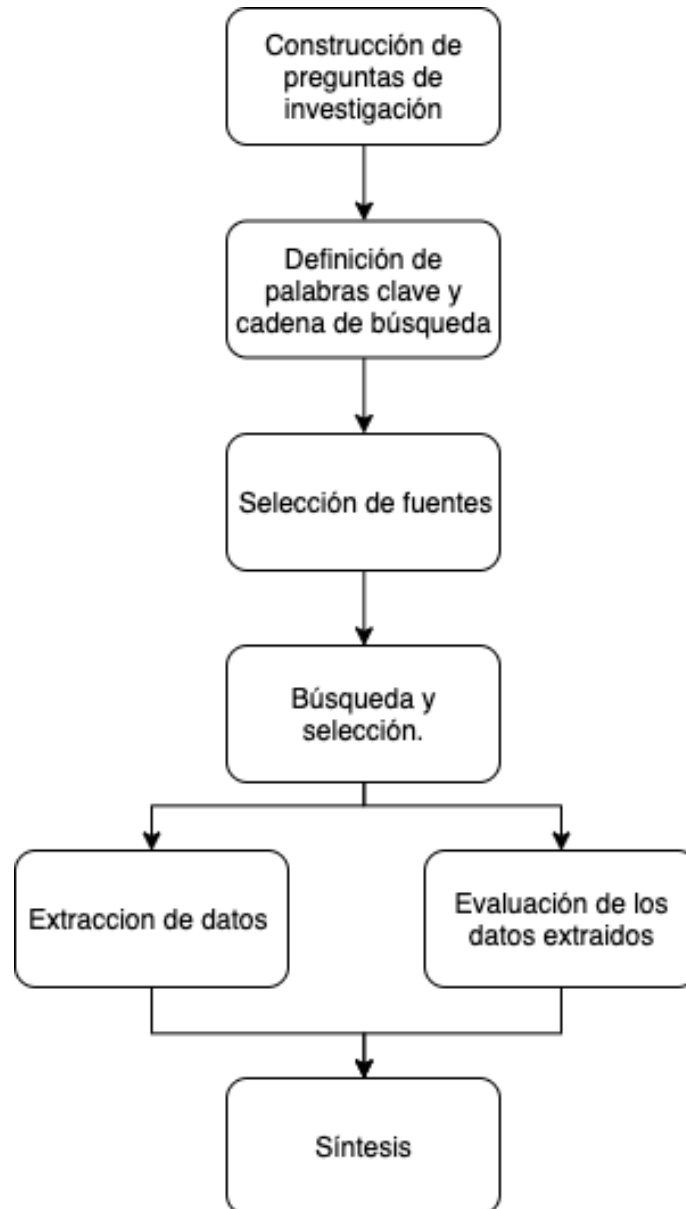
1.1.2 **Definiendo las preguntas de investigación**

El objetivo de la revisión sistemática es comprender como un conjunto de datos o *datasets* pueden ser almacenados y esto a su vez con un método de enriquecimiento que pueden ser explotados para crear un conocimiento más grande. Por tanto, hemos definido las siguientes preguntas de investigación:

- a) ¿Qué repositorios semánticos han desarrollado un método para enriquecer su grafo de conocimiento?
- b) ¿Qué desafíos y problemas se han enfrentado para enriquecer los datos?
- c) ¿Qué criterios y técnicas se utilizan para la evaluación el conjunto de datos?
- d) ¿Qué direcciones son las más prometedoras para futuras investigaciones?

Figura 1

Revisión sistemática de manera general.



Nota. Elaboración propia

1.1.3 **Búsqueda y selección**

La obtención de conjuntos de documentos relevantes para esta revisión sistemática requiere de búsquedas en las bases de datos científicas: Scopus¹, IEEE Xplore² y Google Scholar³. Estas bases de datos fueron elegidas en base a las exigencias que requiere la

investigación y porque estas son mundialmente reconocidas por sus investigaciones y publicaciones. La **Tabla 1**, se define un conjunto de cadenas de búsquedas para determinar los tipos de documentos serán óptimos.

Tabla 1

Cadenas de consultas

Base de datos científicas	Cadena de búsquedas
- SCOPUS ¹	(data AND enrichment AND in AND semantic AND web)
-IEEE Xplore ²	(enrichment AND of AND open AND data AND graphs)
-Google Scholar ³	(adding AND value AND to AND linked AND open AND data)
	(data AND enrichment AND in AND semantic AND web) AND PUBYEAR > 2015
	(data AND enrichment AND in AND semantic AND web) AND DOCTYPE (ar OR re) AND PUBYEAR > 2015

Nota: Representación de las cadenas de consultas en las bases de datos científicas.

El objetivo es encontrar artículos y documentos de calidad con la finalidad de obtener un refinamiento en los resultados de búsquedas y a su vez los resultados sean relacionados al tema de investigación. Sin embargo, debido al gran número de publicaciones de todo tipo se optimizarán las búsquedas, especificando los años y los tipos de documentos requeridos. Adicionalmente, que la etapa de búsqueda y selección de artículos fue de gran ayuda debido a que con esta forma se pudo definir de mejor manera que información será de gran relevancia y de que bases de datos científicas serán parte para el TT.

1.1.4 **Extracción de artículos y documentos**

Una vez determinado e identificado cada uno de los documentos científicos los cuales serán la base necesario para el desarrollo de la investigación. Y después de finalizar la búsqueda, se procede a extraer los metadatos de cada uno de los artículos consultados. Al extraer esta información de los diferentes sitios científicos estos contienen

¹ <https://www.scopus.com/>

² <https://ieeexplore.ieee.org/>

³ <https://scholar.google.com/>

diferentes metadatos. La **Tabla 2** describe las características que se toma en cuenta para la extracción por cada búsqueda, los cuales son:

Tabla 2

Elementos de los resultados de la búsqueda

Elementos de datos	Descripción
Título del documento	Título del trabajo
Autor(es)	Los nombres del autor
Año	Año de publicación
Fuente y tipo de documento	El tipo de la investigación o área de interés
DOI	identificar un objeto digital
Tipo de acceso	Si el documento puede ser accedido libremente por cualquier persona sin pago.

Nota: Representación de resultados de las búsquedas en las bases de datos científicas.

Una vez extraída la información de cada una de las búsquedas en las diferentes bases de datos científicas mencionadas. Se procede a excluir ciertos resultados con el fin de evitar la duplicidad de datos, limpieza de datos de los resultados obtenidos en las tres plataformas que fueron de ayuda. Por último, la estrategia de análisis fue leer el abstract, introducción y las conclusiones de los diferentes artículos o documentos para cerciorarse de que contiene todos los datos necesarios para este TT. Estos artículos o documentos se encuentran en el idioma inglés.

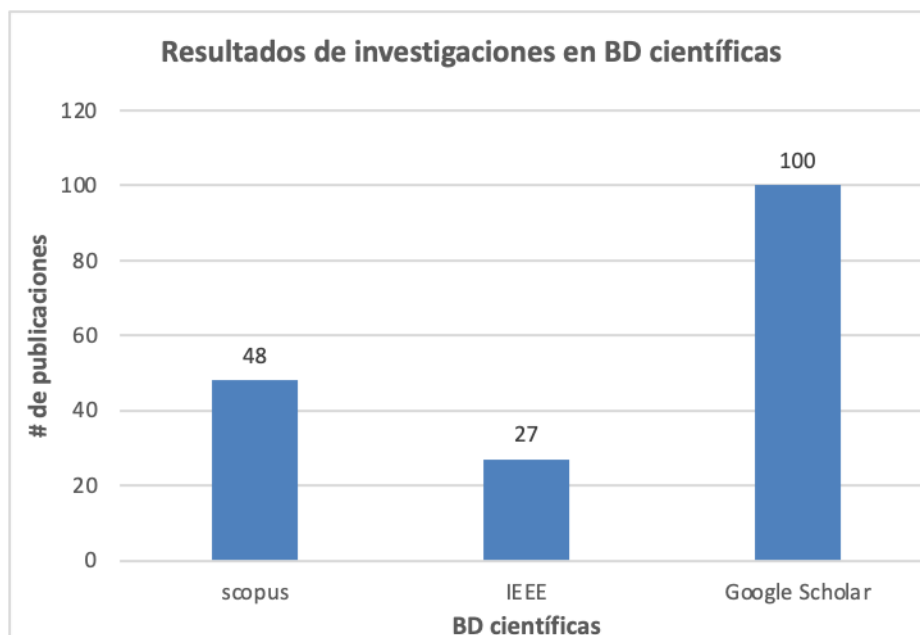
1.1.5 Evaluación de los datos extraídos

Después de extraer y limpiar los datos de los artículos y documentos se procedió al análisis los datos artículos extraídos con el fin de tener un criterio más claro y determinar que información es más relevante. Además para el análisis se utilizan criterios que surgieron de las preguntas planteadas al inicio de la investigación. Los criterios usados para el análisis se los puede definir como temáticas o palabras claves de tal manera que se pueda obtener un resultado conciso y más acercado a lo que se pretende revisar. Las temáticas contienen las siguientes palabras claves: semantic web, linked data, knowledge graphs, enriching a knowledge graph.

Con base a los datos extraídos, las búsquedas y selección de los resultados se presentan en la **Figura 2**.

Figura 2

Resultado de búsquedas en BD científicas



Como resultado de las búsquedas se obtuvo un total de 175 artículos de los cuales se usarán 12 documentos que serán los base para la TT y la definición del método de enriquecimiento de grafos de conocimiento. Para la selección de estos 12 documentos se tomó en cuenta al año de publicación y el enfoque de la investigación que estos tienen en relación al trabajo de TT.

La **Figura 2** muestra el resultado final de las búsquedas obtenidas en cada una de las bases de datos científicas. Hay que tener en cuenta que la búsqueda en Google Scholar arroja un resultado total de búsqueda que inicialmente se obtuvo un valor aproximado de 1710 documentos, entonces se asumió que las diez primeras hojas del resultado de búsqueda son las óptimas, esto es los 100 primeros documentos.

Por otro lado, los criterios que fueron parte esencial para determinar los tipos de documentos que son claves para la investigación teniendo en cuenta los años de publicación, la participación de los estudios primarios como se muestra en la **Figura 3**.

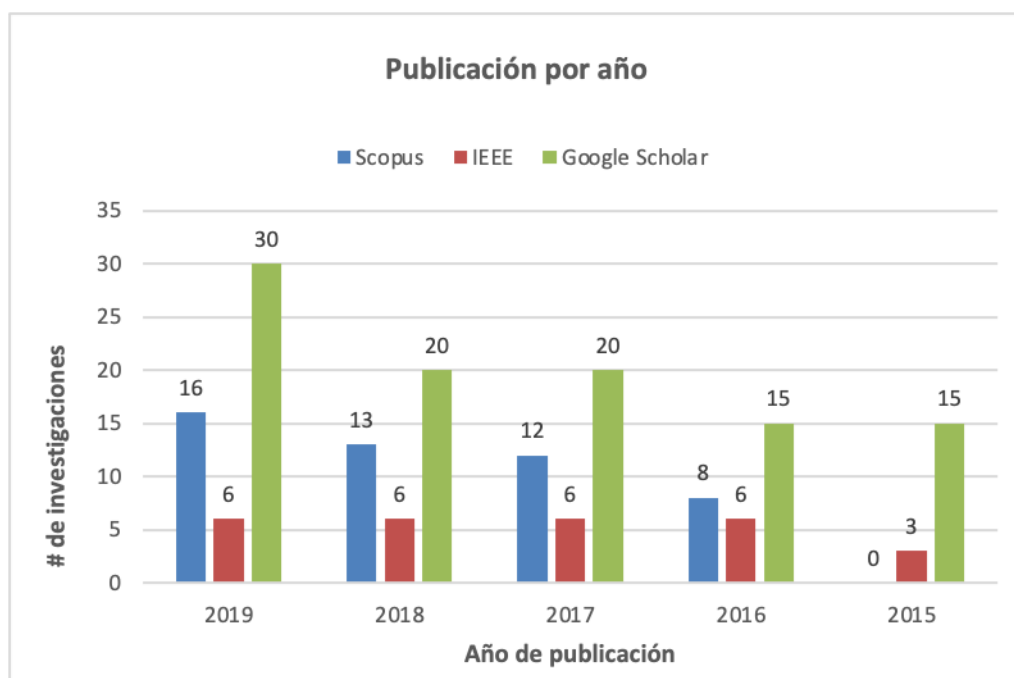
Cabe destacar que los documentos científicos seleccionados se publicaron desde el 2015 hasta la actualidad 2019 o 2020 en caso de que existan algunos.

1.1.6 Año de publicación y distribución geográfica

La **Figura 3** demuestra que esta área de investigación es moderna, novedosa y lleva varios años en la búsqueda de aportar conocimiento. En la **Figura 3** se puede observar el resultado de la distribución del año de publicación de todos los trabajos seleccionados tiene un incremento relativo con relación al año anterior donde se publicaron un total de 18 trabajos (10,2%) en el 2015, 29 trabajos (16,4%) en el 2016, 38 trabajos (21,5%) en 2017, 39 trabajos (22%) 2018 y 52 trabajos (29,4%) hasta noviembre del 2019. Por lo tanto, el número de publicaciones aumenta cada año, lo que refleja el aumento de interés en esta área de investigación.

Figura 3

Año de publicación de los documentos primarios seleccionados.



La **Figura 4** muestra la distribución geográfica de los países que están en la búsqueda sobre este tema. Existen más 20 países que tienen trabajos relacionados, pero se toma en cuenta los 10 primeros países que tienen gran número de trabajos para la

representación gráfica. En esta **Figura 4** se muestra la distribución geográfica donde la atención de la investigación tiene más impacto. El mayor porcentaje de trabajos o documentos es de 21,74% en Estados Unidos, y fueron publicados por universidades o empresas en función del proyecto de investigación. Después de esto, tenemos España Alemania y Grecia que comparten un porcentaje del 10,87%. En el resto de los países comparten valores similares entre ellos.

Figura 4

Distribución geográfica de los países en relación con el tema.



1.1.7 Tipos de publicaciones

La **Figura 5** muestra el total de publicaciones de las cuales fueron seleccionadas para formar parte de esta investigación. Con un 62,29% del total de todos los documentos seleccionados fueron publicados en artículos de conferencia, seguido del 27,43% en

artículos científicos y 6,86% se publicaron en artículos de revisión. Una de las bases de datos científica o bibliografía con más número de publicaciones es Scopus, en este se indexan gran cantidad de documentos de investigación de diferentes categorías de las cuales se basan para determinar que método de enriquecimiento será usado para el TT.

Figura 5

Tipos de publicaciones

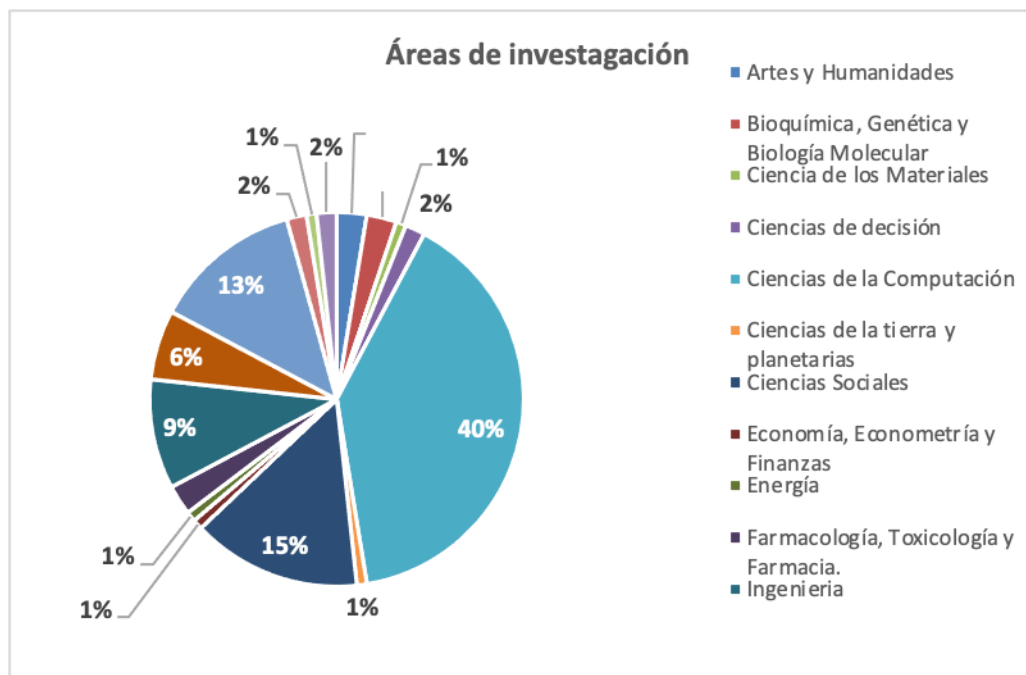


1.1.8 Áreas de aplicación o investigación.

La **Figura 6** se aprecia mediante un diagrama de pastel como esta temática se ha ido aplicando en cada una de las áreas de investigación.

Figura 6

Áreas de investigación



Nota. Elaboración propia

La gran mayoría de los temas de investigación de las universidades se centran en las áreas técnicas y socio humanísticas, donde se puede apreciar gran cantidad de investigaciones relacionadas a este tema. Las áreas con mayor aporte es ciencias de computación que tiene un 43% y representa las áreas que están interesadas en resolver o buscar métodos para implementar y adaptar las diferentes tecnologías para este tema del TT, le sigue ciencias sociales con un 15% y luego el área medica que comparte un 13% de todas las investigaciones en este apartado.

1.1.9 **Discusión y Resultados.**

Finalmente, después de todas las fases anteriores se llegó a la selección de los diferentes métodos mencionados en cada artículo analizado, basándose desde la fase de recolección de los datos, transformación y hasta el enriquecimiento de grafos de conocimiento, siendo esto la parte fundamental para el desarrollo del TT.

Son 12 artículos científicos que son base para este TT. Para la selección de esto se basa en el año de publicación y el enfoque del contenido que estos tienen. De estos 12, en

5 artículos mencionan diferentes procesos o métodos de enriquecimiento de grafos y los restantes son los que permiten entender sobre la web semántica, grafos de conocimiento y más temas relacionados al TT. A continuación se describen algunas formas para realizar el proceso de enriquecimiento.

Según (Quattrini et al., 2017) propone un modelo de enriquecimiento de datos basándose en usar un etiquetado semántico, una organización semántica a todo el conjunto de datos y para relacionar o enlazar los *datasets* con otros usaron *Revit* el usar parámetros compartidos para enriquecer datos. Estos parámetros usados por *Revit* son de crear, agrupar por áreas temáticas, agrupar conjuntos de parámetros de una ontología de dominio y diferentes formas de agrupar como texto, números, referencias de enlaces externos, entre otros.

Otras formas de enriquecer grafos son planteadas por (Jabbar & Bulbul, 2019) y (Sánchez et al., 2018), los cuales proponen dos formas, la primera definida por *Jabbar* y *Bulbul* mencionan tener un almacén de datos semánticos donde tiene un motor de transformación a *linked data*, y el enriquecimiento se realiza usando el motor de enriquecimiento semántico de *Jena* y así tener un servicio de *linked data* para luego su explotación.

Por otro lado, (Sánchez et al., 2018) en su trabajo proponen una vinculación híbrida entre vinculación automática y manual. La vinculación automática se realiza para todos los recursos mediante el servicio de reconciliación ofrecido por *RDFRefine*, donde se establece *Sparql Endpoint* y la propiedad con la que se buscan los recursos relacionados y por otro lado la vinculación manual trata de enlazar datos, pero con los criterios que requiere para la vinculación de datos extendidos.

En (Escobar et al., 2020) se plantea un método de enriquecimiento de grafos, partiendo desde la extracción de datos hasta el enriquecimiento. Primero se debe especificar las fuentes de datos (utilizan un proceso ETL para normalización de datos obtenidos de fuente de datos heterogéneos), establecen un modelo RDF de datos (transformación de datos usando vocabulario *Data Cube*), generación de datos (*OpenRefine*)

esta es la parte que se encarga de enriquecer los datos y finalmente la publicación y explotación de los datos.

Otro método propuesto por (Song & Park, 2018) menciona en usar un modelo basado en traducción, esto consiste en integrar entidades y relaciones en un solo conjunto de datos de la misma área o del mismo vector. Con la traducción de entidades busca mapear vectores de entidad en diferentes vectores en espacios de relación según la entidad y la relación.

Todos estos métodos planteados por los investigadores mencionados serán de gran ayuda para definir un método que se adapte a la necesidad para resolver en problema de este trabajo de titulación.

Por otro lado, al concluir este análisis e interpretación de los métodos usados en las investigaciones mencionadas se tiene como resultado las respuestas a la incógnitas planteadas para esta revisión sistemática. La primera pregunta planteada menciona lo siguiente, ¿Qué repositorios semánticos han desarrollado un método para enriquecer su grafo de conocimiento? En la web existen múltiples repositorios semánticos que se usan para tareas de enlazado y enriquecimiento, por ejemplo en (Escobar et al., 2020) se presenta el proceso o método para la recolección, tratamiento, formulación de estructura semántica y el enriquecimiento con otras fuentes como GeoNames, RDF Data Cube con la finalidad de ampliar el repositorio Barcelona Open Data (BOD)..

Después de haber revisado todos los posibles desafíos y problemas que se puede encontrar al momento de enriquecer los datos se resumen en lo siguiente. El problema es que no existe un proceso estándar que permita que todos los datos existen en la web puedan ser estructurados semánticamente, llevando esto a un nuevo desafío de poder generar un método que permita estructurar y esto a su vez permita añadir más información a los datos existentes. De la misma manera los criterios y técnicas se utilizan para la evaluación el conjunto de datos se basa en los criterios que los investigadores planteaban en sus artículos.

Finalmente las direcciones más prometedoras para plantear las investigaciones acorde al tema es que se pueda crear modelos estándares que permitan que cualquier tipo de información semántica se pueda enriquecer. Por consiguiente el objetivo de esta revisión sistemática fue identificar investigaciones de temas predominantes dentro del campo de las tecnologías de web semántica y posibles aplicabilidades de métodos que permitan enriquecer un grafo de conocimiento semántico.

En resumen se menciona lo siguiente, que mediante una fase de 5 pasos se puede resolver un poco el problema de los datos no enlazados o enriquecidos, estas 5 fases son las siguientes:

1. Extractor: permite la extracción rutinaria y automatizada de nuevos datos de fuentes específicas.
2. NER (Reconocimiento de Entidades Nombradas): permite la extracción de información que se busca localizar y clasificar en categorías predefinidas, tales como: personas, organizaciones, lugares y las entidades nombradas encontradas en un texto.
3. NLP (Procesamiento del Lenguaje Natural): reconocimiento de entidades en el cual se pueden idear búsquedas lingüísticas específicas y detectar similitudes entre entidades.
4. Flujo de trabajo: los usuarios podrán hacer uso del API de los datos.
5. Publicación: la publicación de los datos en el piloto y el resultado de visualización serán recursos fáciles de procesar por parte del usuario final.

Es importante tener en cuenta que varios de estos los temas son parte fundamental para el desarrollo de la propuesta de solución del TT.

1.2 Conceptos

En esta sección se detallan conceptos claves para el desarrollo del Trabajo de Titulación (TT) que permitieron la mejor comprensión del contexto del TT y la presentación una propuesta relevante que permita resolver o satisfaga a los objetivos del tema. En las secciones 1.2.1 y 1.2.2 se presentan partes fundamentales para el entendimiento del desarrollo del TT. La sección 1.2.3 analiza la importancia de tener claro los conceptos que sean útiles en el campo de la web semántica con finalidad de crear una propuesta de solución para el trabajo. La sección 1.2.4. muestra la técnica para la integración, extracción de información y herramientas que se usaran para resolver el problema. Finalmente, en la sección 1.2.5. realizará una discusión del contenido obtenido.

1.2.1 *Datos abiertos*

Se debe tener en cuenta que los datos abiertos, también conocido como Open Data (OD). Según (Open Data Handbook, 2019) define que los datos abiertos son “*datos que pueden ser utilizados, reutilizados y redistribuidos libremente por cualquier persona*”, nos referimos a datos que están disponibles de forma libre por la Web. Estos datos tienen que pasar por varios niveles para poder ser completamente libres y tienen que cumplir ciertas características como la disponibilidad y acceso, pues tienen que estar disponibles por completo y así estos pueden ser accedidos por cualquier interesado. Otra característica es la reutilización que se tiene que hacer bajo unos permisos o términos que permitan su uso y explotación de forma libre y no pongan ninguna restricción a la hora de trabajar con otros tipos de datos. Todos los interesados que deben poder acceder a los datos y trabajarlos como ellos lo requieran, pero sin ninguna restricción alguna.

1.2.1.1 Importancia

La importancia de OD radica en la posibilidad de acceder a los datos de forma libre y así mejorar el conocimiento⁴. También los datos deben contribuir a mejorar eficientemente los servicios de cualquier entidad o grupo que usen estos datos (Pfenninger et al., 2017). Principalmente permite el acceso, la responsabilidad y la transparencia con la finalidad de que estos recursos abiertos permitan beneficiar a varios individuos o también permite ayudar con búsquedas de soluciones.

Hay muchas áreas donde podemos esperar que los datos abiertos sean valiosos, por lo tanto, (Duus & Cooray, 2016) mencionan en su portal que la disponibilidad de datos abiertos crean oportunidades para todo tipo de organizaciones, gobiernos y organizaciones con o sin fines de lucro para encontrar nuevas formas de abordar los problemas, como esto se promueve que exista transparencia en las acciones que estos grupos realicen.

Los datos abiertos juegan un papel fundamental para poder de ampliar la visión de cada persona en el uso y manejo de los datos. Según (Duus & Cooray, 2016) menciona que los datos describen los patrones detrás de cómo vivimos puede ayudarnos a resolver problemas de formas que no hubiéramos previsto. El autor concluye que la apertura de los datos abiertos es de gran ayuda con la finalidad de tener una perspectiva diferente a lo que cada uno tiene actualmente de la información porque los OD buscan resolver problemas.

1.2.1.2 Calidad de los datos abiertos

Dentro de la calidad de OD vienen también los retos que estos pueden tener como poder realizar la captura de datos y de estos tener obtener una información bien estructura.

Las herramientas de análisis de datos cada vez más sofisticadas, así permitiendo analizar los datos de nuevas maneras para descubrir tendencias y hallazgos que han dado forma a decisiones cruciales. Estas herramientas permiten ver soluciones entre datos que nunca han sido explorados. Para esto la información debe estar abierta y disponible.

⁴ Tecnologías de la información y la comunicación

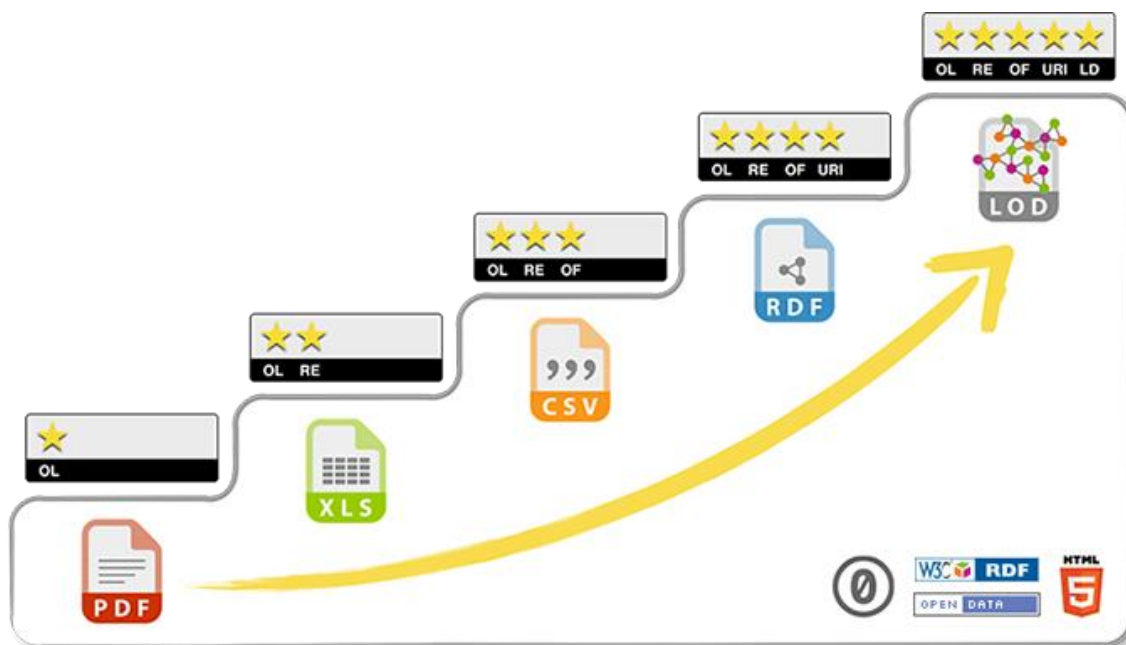
La disponibilidad de información y el acceso público a los datos cada vez mayor y el crecimiento exponencial de la información día a día es muy grande desde el nacimiento de la World Wide Web (WWW) que fue planteado por Tim Berners-Lee. Uno de los mayores retos que menciona (Stauffacher et al., 2012) es que el acceso a datos o información se traduce en el empoderamiento con el fin de poder tomar decisiones informadas, resolver problemas y para mejorar el nivel de vida de las personas.

Para que cada uno de estos datos puedan ser usado debe de pasar por diferentes procesos de análisis con la finalidad de determinar a que estrella o nivel se encuentra y de esta manera saber si esta cumpliendo con cada una de las características que tiene la Open Data. Tim Berners-Lee, el creador de la WWW, sugirió un esquema de desarrollo de 5 estrellas de OD⁵ el cual describe en niveles los estados de los datos en los que se encuentra desde el primer nivel hasta llegar a la ultimo nivel en la cual representa que los datos son abiertos, de acceso libre y con una estructura fácil de manejar y mantener. La **Figura 7** muestra de manera resumida los niveles de estrellas en los que se encuentra los datos.

⁵ <https://5stardata.info/en/>

Figura 7

Calidad de OD en 5 niveles



Nota. Adaptado de 5 stard data [Fotografía], por (Berners-Lee, 2015)

Como se aprecia en la **Figura 7** cada nivel está representado por una estrella a través de la que se evalúa que tan accesible es la información. Según (Berners-Lee, 2015) describe los niveles de la siguiente manera en la **Tabla 3**:

Tabla 3

Niveles de las 5 estrellas de la Open Data

★☆☆☆☆	Los datos deben estar disponibles en la Web (cualquier formato) bajo una licencia abierta
★★☆☆☆	Deben ser datos estructurados, por ejemplo, un archivo usado en Excel.
★★★☆☆	Deben estar disponible en un formato abierto no propietario, por ejemplo, CSV en lugar de Excel.
★★★★☆	El uso de URI para denotar cosas, de modo que las personas puedan señalar sus cosas.
★★★★★	Los datos deben estar vinculados entre datos para proporcionar contexto.

Nota: Presenta los 5 niveles de los datos abiertos planteados por (Berners-Lee, 2015)

En este trabajo se ha considerado una versión ampliada del modelo de cinco estrellas de calidad de los datos abiertos, reportada en (Alesso & Smith, 2006).

1.2.1.3 Plataformas de datos abiertos

Existen diversas plataformas de datos abiertos o Open Data Platform (ODP), las cuales permiten obtener datos o a su vez publicar todo tipo de información relevante de varias temáticas, son una parte fundamental en el proceso de apertura de datos según (Techopedia, n.d.).

Las ODP sirven para almacenar, compartir, conectar y visualizar todo tipo de datos en una base de datos que estas posean. También son las puertas de entrada a la búsqueda de estrategias de una organización o grupos con el beneficio de abrir sus datos y así los ciudadanos pueden ver esta información (Techopedia, n.d.).

Dentro de las plataformas no existe un único modelo para los datos, sino los portales poseen varias características con algunas secciones concretas, pues todo esto no permite tener una misma estructura de la información de una plataforma con otra.

Hay información y herramientas que permiten trabajar con datos que son abiertos o públicos. Algunas aplicaciones de varias funciones en la sociedad se han desarrollado gracias a estos datos, ha permitido también que esta información ayude a predecir eventos que puedan suceder. A continuación, enumeraremos las partes fundamentales que tiene un portal de datos abiertos (Danneels et al., 2017):

a. Catálogo de datos:

Los catálogos de datos son la parte esencial de las plataformas de datos abiertos, estos poseen una lista en una tabla los cuales describen el posible contenido que estas puedan ofrecer, el nombre de los responsables de la información, el formato en la que se encuentran (pdf, xlsx, csv, json entro otros), la frecuencia que estos son actualizados, número de visitas entre más características que cambian según los portales.

b. Colección de datos:

La colección de datos es un listado de información. La colección puede tener un enfoque con temas relacionados (salud, economía, clima, educación) o puede tener alguna organización que esta a cargo de estos datos.

c. Visualización de datos:

Es el sitio donde se da a conocer de otra manera los datos, es decir, por medio de gráficas o vista de datos que fueron desarrollados por los usuarios o administrador del sitio, todo esto depende si la plataforma nos permite.

d. Historias de datos:

Las historias basadas en datos son consideradas una prioridad para muchas organizaciones. Estas historias de datos se conforman por versiones previas que se han llevado durante un tiempo. Estas historias de datos pueden ayudar a tomar decisiones o también a determinar como ha ido cambiando la organización a lo largo de este tiempo ya que nos permite dar una retroalimentación de como estas han ido avanzando a lo largo del tiempo y ver si las cosas que han hecho han sido la mejor opción o no.

1.2.2 Objetivos de desarrollo sostenibles.

En el 2015 se reunieron en las Naciones Unidas⁶ (UN) todos los líderes de las diferentes naciones para adoptar un conjunto de objetivos con una visión de mejorar la vida en un futuro para todos, estas metas están planteadas en la Agenda 2030 de la siguiente manera, 17 Objetivos de Desarrollo sostenibles⁷ (ODS) con 169 objetivos y 330 indicadores (Griggs et al., 2013). Las Naciones Unidas ha estimado un incremento de la

⁶ <https://www.un.org/es/>

⁷ <https://www.un.org/sustainabledevelopment/sustainable-development-goals/>

población de 66% para el 2050 lo que esto implica más retos y tratar de solucionar de manera rápida y oportuna a estos problemas.

1.2.2.1 Indicadores de los ODS alineados con los datos abiertos

A través de los ODS, los países cuentan con indicadores específicos, observables y medible que permiten mostrar el cambio y progreso en las actividades relacionadas a los ODS (*Indicadores de ODS*, 2010).

En la búsqueda de ayudar a mejorar la calidad de vida de todos con los ODS, varios grupos o personas han desarrollado herramientas que permite tener de cierta manera características de cada uno de los ODS y así estos puedan ser accedidos y poder relacionar grupos de datos abiertos a cada criterio con los indicadores y así saber de manera mas rápida el crecimiento que tiene el objetivo de desarrollo sostenible. Las Naciones Unidas con el objetivo de ayudar a determinar esto desarrollaron una API⁸ que permite saber tantos detalles de los objetivos como de los indicadores. Estos recursos pueden ayudar de gran manera a integrar varios recursos abiertos a cada uno de los ODS e indicadores.

1.2.3 Web Semántica

A lo largo de estos años la web ha ido variando desde su primera versión la 1.0 basándose en que las personas puedan intercambiar información en la web. Después pasó a la web 2.0 que se describe como la segunda generación de World Wide Web (WWW)⁹ donde las páginas estáticas son reemplazadas por páginas interactivas para una experiencia dinámica para los usuarios. La Web 2.0 es creada con el objetivo de permitir a las personas compartir su información en la red a través de las plataformas de redes sociales, blogs, chats y poder contribuir en la web (Cheaney, 2012).

La Web 3.0 es el tercer nivel de desarrollo de la web a lo que se denomina la Web Semántica (WS) la que demuestra la naturaleza semántica de los contenidos. Tim

⁸ <https://unstats.un.org/SDGAPI/swagger/>

⁹ Es un término que es utilizado para referirse a la World Wide Web, también conocida como Internet.

Berbers-Lee fue el mentor del término de web semántica lo que nos dice que los datos que se encuentran en la Web es información que se puede tratar y buscar relaciones entre ellas con situaciones reales (Sharma, 2015).

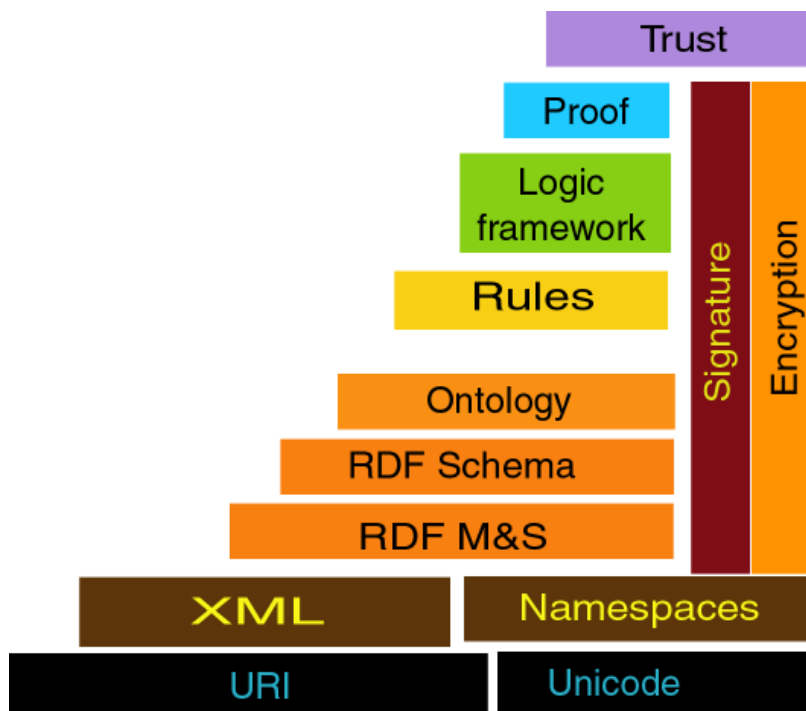
Tim Berners-Lee, el padre la web, señala que *“La Web Semántica es una extensión de la web actual en la que la información tiene un significado bien definido, permitiendo que las computadoras y las personas trabajen mejor en cooperación”*, (Berners-Lee et al, 2000) esta visión, evoluciona la Web actual en una red de datos y significados global, con un enorme potencial para apoyar a la resolución de problemas que requieren gestionar conocimiento global, de manera rápida y oportuna.

1.2.3.1 Tecnologías de la Web Semántica

La Web Semántica desde sus inicios se ha planteado varios criterios que se deben considerar el desarrollo de la WS y sus posibles aplicaciones. La **Figura 8** presenta la organización por capas de las tecnologías de la WS cual será el criterio clave para el desarrollo del TT.

Figura 8

Las capas de las tecnologías web semánticas



Nota. Adaptado de w3: capas de las tecnologías [Fotografía], por (Hazaël-Massieux & Berners-Lee, 2003)

Las capas fundamentales para el desarrollo del TT se las describe a continuación:

a. URI

Uniform *Resource Identifier* (URI), su significado en español es Identificador de Recursos Uniforme, su función es nombrar recursos en la web, lo que un navegador recibe es información o le proporciona un método/ubicación de acceso al recurso (Miessler, 2020).

Esté termino surge de la unión de Uniform Resource Locator (URL) y Uniform Resource Name (URN).

b. XML

XML es conocido como eXtensible Markup Language, traducido como Lenguaje de Marcado Extensible o Lenguaje de Marcas Extensible, el cual es un lenguaje de marcado similar a HTML y es una especificación de W3C, es decir que XML no está predefinido, por lo cual debe ser definido las propias etiquetas (Gavin, 2018).

c. RDF

Resource Description Framework (RDF) es un estándar que permite la interoperabilidad entre datos en la Web. Según (Lapiente & Lamarca, 2018) “RDF es un lenguaje para representar información sobre recursos, ... particularmente sobre metadatos de recursos web, tales como el título, autor, modificaciones de los datos de la página web, copyright y otras licencias de información sobre documentos web, así como la disponibilidad para algunos recursos compartidos”. Todo esto se trata de un modelo de datos para objetos o recursos.

1.2.3.2 Ontologías y vocabularios

Las ontologías tienen componentes que sirven para representar el conocimiento de algún dominio. Son descritas como una estructura para la comunicación entre humano y máquinas. Estas permiten mejor las búsquedas, la comunicación entre conocimientos y entre otros aspectos.

Según (W3C, 2015) los vocabularios definen conceptos y relaciones que son usados para describir y representar cosas. Estos vocabularios pueden ser muy complejos a la hora de usar una gran cantidad de datos.

La W3C dice que no existe la división exacta entre vocabulario y ontologías, ya que la tendencia es usar la palabra *ontología* para una colección de términos más complejos. Mientras que *vocabulario* es usado para la construcción de bloques con técnicas básica a seguir.

1.2.3.3 Linked data

Linked data es un término que hace referencia a un método de publicación de datos, pero estructurados. Esta tecnología tiene dos ideas, las cuales engloban en ser más grandes y poderosas.

Este método permite interrelacionar conjuntos de datos, por lo que nos permite que sea mas fácil y útil para la búsqueda de información. Según (Bizer, 2009) menciona que la linked data se refiere “*al conjunto datos que buscar mejorar la información para publicar y conectar datos estructurados en la web*”, siendo este el siguiente paso de la open data con

la finalidad de que estos datos tengan una estructura pero que estos se puedan relacionar con más conjuntos de datos abiertos.

Linked data (LD) tiene algunas capacidades o características que según (Bizer, 2009) pueden ser las siguientes:

- Cualquiera puede publicar datos en la Web de datos vinculados.
- Los enlaces conectan entidades creando un grafo que abarca fuentes de datos y permite el descubrimiento de nuevas fuentes de datos.
- Los datos se auto describen. Si una aplicación encuentra datos representados usando un vocabulario desconocido, la aplicación puede resolver los URI que identifican los términos del vocabulario para encontrar sus definiciones.
- La Web de datos vinculados está abierta, lo que significa que las aplicaciones pueden descubrir nuevas fuentes de datos en tiempo de ejecución siguiendo los enlaces.

El enfoque de datos enlazados puede ser aplicado en cualquier tipo de información. A través de LD la información se estructura como piezas de conocimiento, que se pueden enlazar con otros datos enlazados, y así expandir el conocimiento (Piedra, N. & Suárez, J.P, 2017)..

1.2.3.4 Aplicaciones de Web Semántica

A la medida que pasa el tiempo, la Web Semántica ha ido teniendo mayor auge y se ha visto que está siendo aplicada en diferentes áreas de investigación. Al igual que las necesidades han incrementado, el volumen de los datos también, en los cuales su estructura se ha vuelto más heterogénea y difícil de procesar, pues gracias a esto se ha visto la necesidad de tratar y dar un valor agregado a esta información con la finalidad de encontrar nuevas oportunidades y tratar de crear nuevas aplicaciones con esta información como: Buscadores semánticos, servicios de web semántica, e-Learning, entre otras.

Finalmente, las aplicaciones web semánticas se las pueden definir como aplicaciones basadas en la web que aprovechan el contenido semántico y estructurado

que estas poseen incluyendo no solo información, sino también metadatos que describen cada uno de los recursos obtenidos según (Alesso & Smith, 2006).

1.2.4 *Procesamiento de lenguaje natural*

El procesamiento de lenguaje natural (NLP¹⁰) es la capacidad de interpretar o analizar el lenguaje humano con la finalidad de saber de lo que están hablando en el conjunto de datos que son analizados.

Otra definición que nos entrega (Wonderflow, 2018) menciona que es la “interacción entre el lenguaje humano con las computadoras”. Una de las características es que el NLP es superior a las actividades de un ser humano en la capacidad de interpretar gran cantidad lenguajes y datos. El NLP es una rama de la inteligencia artificial (AI¹¹).

El desarrollo del NLP es desafiante debido a los algoritmos que son usados para realizar los procesos que necesitan para una comprensión más fuerte para poder interpretar la información por un computador. Para que esto se realice debe ser mediante comandos de voz o texto, los cuales son interpretados por los computadores. Para el desarrollo del NLP existen varias librerías o framework en los diferentes lenguajes los cuales permiten realizar el análisis de la información que se quiere interpretar.

En la web existen varios métodos que son fáciles de aplicar, por ejemplo, tenemos a **fasttext**¹² desarrollado para Python y el cual es una librería de aprendizaje de incrustaciones de palabras y de clasificación de texto, este modelo permite un aprendizaje supervisado y no supervisado dando un resultado de una representación de vectores para las palabras y soporta 294 idiomas.

Con esta librería se busca entender el contexto de cada palabra y así formar la oración que serán usadas para la búsqueda, de tal manera refinar la búsqueda y obtener un resultado más acorde a lo que se quiere.

¹⁰ Natural Language Processing

¹¹ Artificial Intelligence

¹² <https://fasttext.cc/>

Finalmente, el Procesamiento de Lenguaje Natural permite realizar varias actividades de manera más rápida y oportuna. El ser humano no puede realizar análisis o interpretar grandes volúmenes de datos en pocos segundos.

1.2.5 Base de datos semántica con GraphDB

GraphDB es una base de datos grafos que sirve a las organizaciones para almacenar, organizar y gestionar contenido en forma de datos inteligentes semánticamente enriquecidos y es altamente eficiente, escalable y robusta con soporte para RDF y SPARQL. GraphDB incluye integración de datos e interconexión, cumplimiento de estándares W3C, modelo de datos expresivo, rico y flexible, espacio único de información interconectado formado por datos estructurados y documentos de texto, razonamiento, compatibilidad con datos abiertos vinculados y procedencia de datos. Según (Ontotext, 2020), empresa que creó GraphDB menciona que tiene tres tipos de licencias, una gratis, la estándar y versión empresarial; en la que se trabajó es con la licencia gratis.

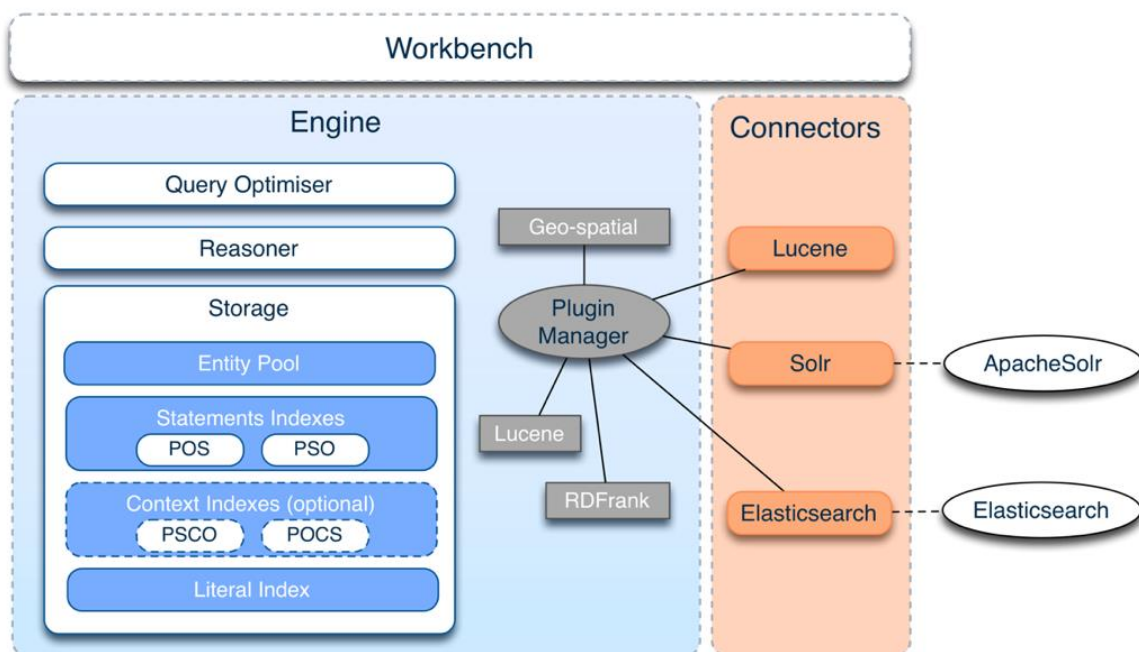
1.2.5.1 Arquitectura y componentes

GraphDB se estructura con una capa de almacenamiento e inferencia (SAIL¹³) para esto usa las propiedades del framework RDF4J y hace un amplio uso de las características y la infraestructura de RDF4J, especialmente el RDF modelo, analizadores RDF y motores de consulta.

¹³ Storage And Inference Layer

Figura 9

Arquitectura de alto nivel GraphDB



Nota. Tomado de Ontotext architecture GraphDB [fotografía], por (Ontotext, 2020)

1.2.5.1.1 RDF4J

RDF4J es un framework para almacenar, consultar y razonar con datos RDF. Es compatible con el lenguaje de consulta W3C SPARQL. También es compatible con los formatos de archivo RDF más populares y los formatos de resultados de consultas incrustado en una aplicación como una biblioteca Java. Para comunicarse con RDF4J se puede utilizar una API con las funciones de JDBC y de RESTful HTTP.

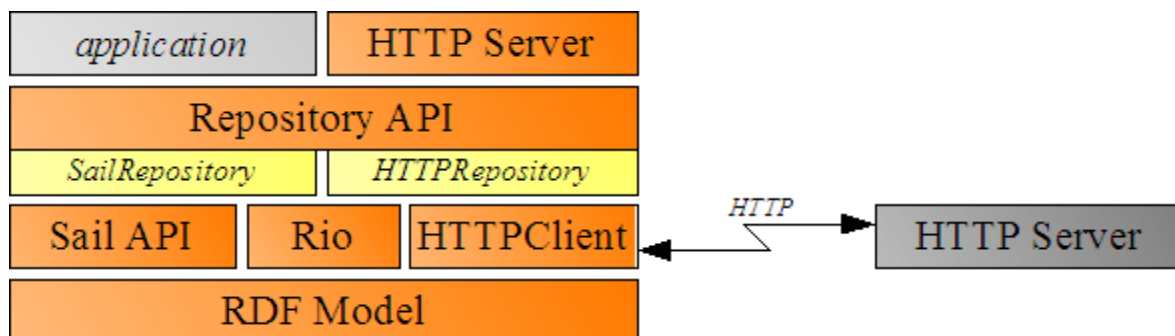
La **Figura 10** nos muestra una representación esquemática de la arquitectura de RDF4J y una breve descripción de los componentes principales.

Las aplicaciones normalmente se comunicarán con RDF4J a través de la API del repositorio y este proporciona un nivel de abstracción lo suficientemente alto como para que los detalles de los componentes subyacentes particulares permanezcan ocultos, es decir, se pueden intercambiar diferentes componentes sin necesidad de modificar la aplicación.

Este TT está basado en Python y la comunicación que se usa para el desarrollo se el framework SPARQLWrapper. El formato de comunicación del framework sigue manteniendo la misma estructura arquitectónica que plantea RDF4J.

Figura 10

Arquitectura RDF4J



Nota. Tomado de Ontotext architecture *RDF4J* [Fotografía], por (Ontotext, 2020)

1.2.5.1.2 Workbench

Workbench es una herramienta de GraphDB que permite la administración de su *triplestore*, Según (Ontotext, 2020) proponen los siguientes puntos en lo que se puede usar:

- Gestionar repositorios GraphDB
- Carga y exportación de datos
- Ejecutar consultas y actualizaciones SPARQL
- Gestión de espacios de nombres
- Manejo de contextos
- Visualización / edición de recursos RDF
- Consultas de seguimiento
- Monitoreo de recursos
- Gestión de usuarios y permisos
- Conectores de gestión
- Proporciona API REST para automatizar varias tareas para administrar y administrar repositorios

1.2.6 *Discusión*

Este trabajo propone como obtener las ventajas que entrega la Web Semántica con esto mejorar todas las formas de relacionar su contenido con más fuentes de datos semánticos. Se busca realizar el proceso de enriquecimiento de datos para esto se debe tener en cuenta que los datos deben de estar en un formato abierto y semántico con la finalidad de poder realizar consultas SPARQL sobre el contenido que se quiera enriquecer.

Para esto se debe tener en cuenta todos los criterios que se plantear en la sección 1.2, siendo estos la base fundamental para tener la idea principal del trabajo de titulación.

Al definir cada uno de los conceptos que se mencionan en esta sección (1.2), se busca identificar que aspectos son los necesarios para estar al tanto de que va el TT. Por otra parte, se define ciertas herramientas que serán usadas para el desarrollo de modelo de solución del TT.

Las herramientas mencionadas en esta sección cumplen el rol fundamental para el manejo y uso de los datos CKAN semánticos y así sea posible poner la solución de TT.

Capítulo dos

Propuesta

La propuesta presentada en este capítulo está fundamentada en la revisión de la literatura realizada en el capítulo anterior, y enfocada en alcanzar los objetivos planteados para el Trabajo de Titulación.

En la propuesta se define un método de enriquecimiento de un grafo de conocimiento relacionado con indicadores de ODS, se define los conjuntos de metadatos que serán usados para el análisis e integración con otros metadatos y así estos se puedan relacionar con los ODS, se crea una base de conocimiento, se aplica un método de enriquecimiento, y se desarrolla una herramienta Web para visualización de datos.

En la sección 2.1 trata del contexto del trabajo, la sección 2.2 se menciona cuál es el problema general del Trabajo de Titulación, la sección 2.3 se trabaja en la propuesta del método de solución y en la sección 2.4 se presenta el modelo de la arquitectura que se emplea para el método de enriquecimiento de grafos.

2.1 *Contexto*

Los Objetivos de Desarrollo Sostenible son indicadores que permiten medir el progreso de la información en la web en relación de los objetivos, metas e indicadores que estos presentan. Desde el punto de vista de los datos se busca saber si la información semántica está relacionada con los indicadores y determinar si está surgiendo un cambio o no, debemos tener en cuenta que todos los datos relacionados están en formatos abiertos en la web. Sin embargo, existen más fuentes de datos abiertos que podrían ser de gran ayuda para aplicar la visión de los datos semánticos en la web.

2.2 *Problema*

La web ha ido creciendo rápidamente y la información en ella también, y eso ha generado datos heterogéneos; por consiguiente, un enorme volumen y variedad de datos inadecuados para los ordenadores tradicionales, lo cual a dado origen a la Big Data. Los datos abiertos son aquellos que pueden ser accedidos y utilizados cuando están disponibles en un formato común legible por máquina, lo que permite que una

computadora los lea y procese automáticamente. Tim Berners-Lee propuso un nuevo modelo llamado Linked Data para publicar información legible por máquina como datos estructurados, basado en RDF. Los datos abiertos vinculados (LOD) son datos que se publican bajo una licencia abierta.

El problema que al existir grandes volúmenes de datos y de diferentes fuentes es que estos no pueden incrementar el conocimiento que estas poseen (Eguiguren Palacios, 2019). Por otro lado, se tiene a los ODS que permiten analizar que datasets pueden estar vinculados a sus objetivos, metas e indicadores. Por consiguiente, lo que se busca es este Trabajo de Titulación es enriquecer un dataset con otros datos vinculados, para lo cual se crea un método que permite enriquecer los datos y se enlacen con otros datasets y a su vez con los Objetivos de Desarrollo Sostenibles.

2.3 ***Propuesta de solución***

Se requiere desarrollar un método de enriquecimiento de Grafos de Conocimiento de indicadores de ODS usando un enfoque de inferencia de entidades semánticas equivalentes desde fuentes de datos abiertos. Para lograr esto se realizó previamente una revisión sistemática en bases de datos científicas para encontrar investigaciones relacionadas que describe diferentes métodos usados para el enriquecimiento de grafos de conocimiento ya que estas investigaciones son la base para plantear la solución.

Para enriquecer los datos del datasets CKAN semántico, se usará el dataset de la Dbpedia debido a que este posee una gran información de varios temas o recursos semánticos que tienen relación al dataset de estudio que es el CKAN semántico.

Este trabajo se basa en los conjuntos de datos reportados en (Eguiguren Palacios, 2019), y en (Eguiguren Palacios, 2019). Estos conjuntos de datos están descritos en tripletas y en diferentes idiomas según el conjunto de datos. Sin embargo, con el fin de ampliar el alcance de búsqueda a otros lenguajes, se incorporó un componente de traducción al inglés como parte de la solución propuesta, también una aproximación basada en NLP (procesamiento de lenguaje natural) y en NER (Reconocimiento de entidades nombradas).

Una vez hecho los pasos anteriores se procede a determinar que valores o entidades tienen mayor relación y con las que se puedan crear una relación entre un dataset A y un dataset B que describen o se relación con el contexto. De esta manera se puede tener la inferencia semántica entre entidades y las nuevas relaciones entre diferentes conjuntos de datos y de esta forma cumplir con los objetivos del trabajo de titulación.

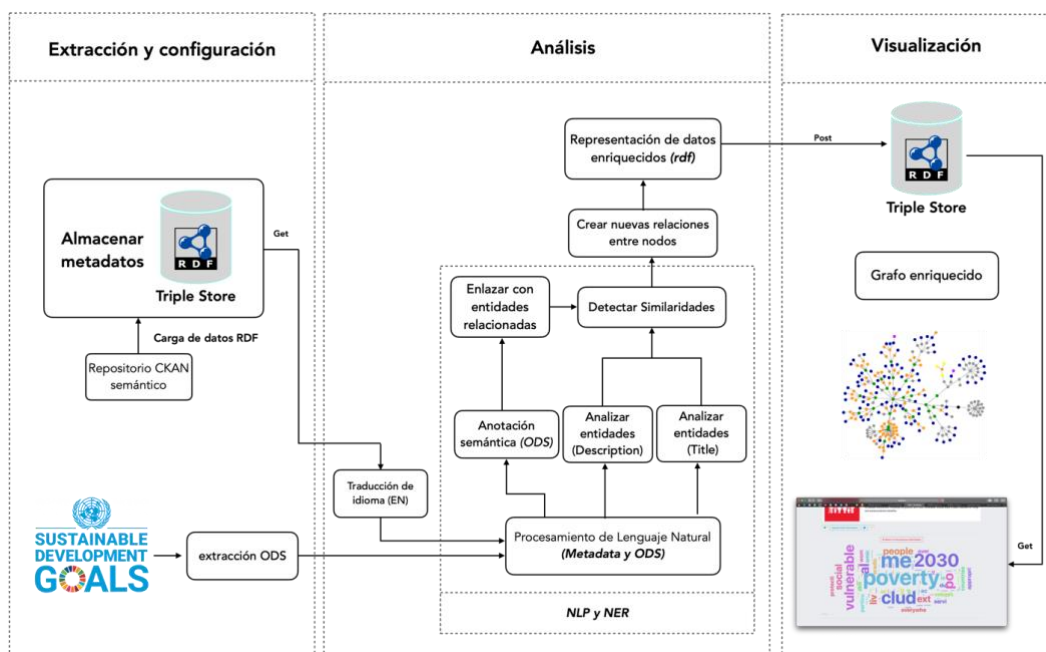
2.4 Componentes arquitectónicos

Para cumplir con los objetivos planteados para este Trabajo de Titulación es necesario crear una arquitectura que cumpla con cada uno de los elementos o características mencionadas en la propuesta de solución. Con esta arquitectura se busca establecer el proceso para el enriquecimiento de grafos semánticos a partir de otros conjuntos de datos abiertos y ver sus posibles relaciones con los Objetivos de Desarrollo Sostenible.

La **Figura 11** representa cada uno de los elementos esenciales para determinar el proceso del método de enriquecimiento de grafos y la asociación con los ODS. Esta representación arquitectónica muestra de manera general el funcionamiento del prototipo para el método de enriquecimiento de grafos con indicadores de ODS e inferencias semánticas con fuentes de datos abiertos.

Figura 11

Arquitectura de prototipo de método de enriquecimiento de grafos con ODS



Primero se realiza la representación de los ODS, lo que se busca es la extracción de cada uno de los objetivos, metas e indicadores. Una vez obtenida esta información se somete a un análisis de Procesamiento de Lenguaje Natural (NLP) con la finalidad de identificar palabras o partes de oraciones que permitan relacionar cada ODS a un recurso o entidad de otro conjunto de datos semánticos y a su vez crear un enlace entre ambos y estos sean almacenados en un triplestore.

De la misma manera se emplea el NLP para el conjunto de repositorios que es la base de estudio y así buscar que entidades semánticas que se relacionan entre si y de esta forma crear una nueva relación semántica.

Finalmente, lo que se busca es tener una nueva base de conocimiento que contenga no solo un conocimiento, sino varios y esta permita tener una mejor precisión en la toma de decisiones o mejor conocimiento en sus relaciones entre datos. También permite identificar las relaciones que cada dataset o cada recurso tiene con los Objetivos de Desarrollo Sostenible y de esta manera se pueda medir que el crecimiento la información enriquecida.

Capítulo tres

Desarrollo de la propuesta

Este capítulo especifica la implementación de la propuesta realizada en el capítulo anterior. Tomando en consideración cada una de las fases desarrolladas para el TT en base a la arquitectura planteada en la **Figura 11**. La sección 3.1 describe la configuración que se tiene que realizar con la herramienta de almacenamiento orientado a grafos GraphDB y la extracción de los ODS. En la sección 3.2 se describe el proceso de creación de enlaces entre los ODS y el dataset semántico. Finalmente, la sección 3.3 describe la visualización de los datos enriquecidos.

3.1 *Configuración y extracción*

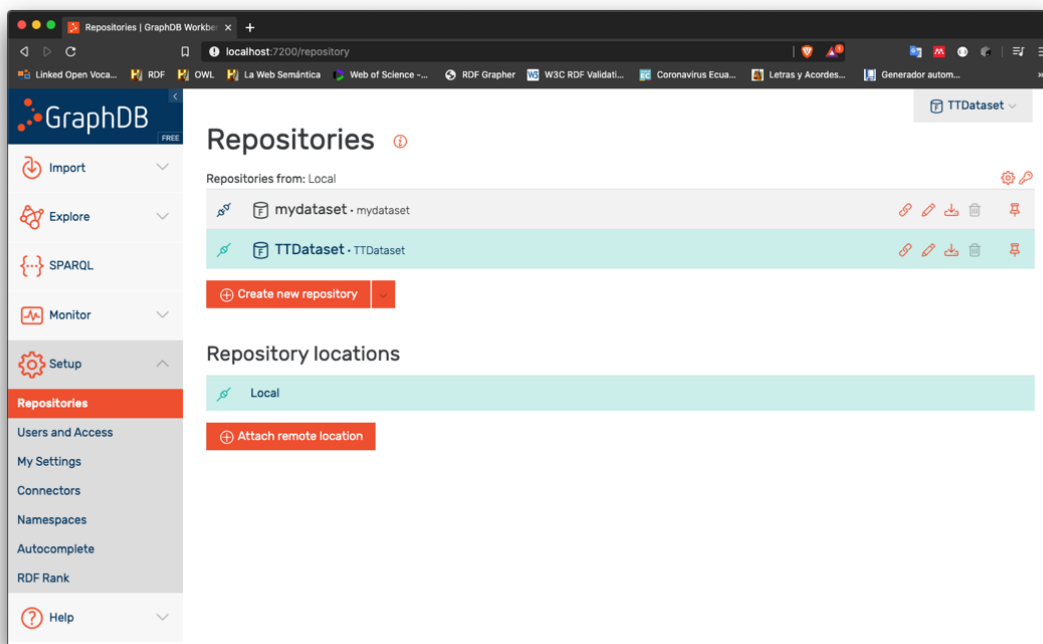
Previamente al almacenamiento de la información en GraphDB se tiene que crear el repositorio en el que se importará la información para obtener el tripleStore y de esta manera trabajar para el consumo de los datos en tiempo real como se plantea en la arquitectura de la **Figura 11** y la extracción de la información acerca de los objetivos, metas e indicadores de los ODS.

3.1.1 *Creación de repositorios y carga de información*

Dentro del entorno de GraphDB está Workbench, este permite de manera fácil y rápida la carga de la información. La **Figura 12** presenta el administrador para la creación del repositorio, y para la creación de este se da clic en la opción “create new repository” y luego se le asigna un ID o identificador único para el repositorio y continuar con los demás ítems que se necesita para la creación.

Figura 12

Administrador para la creación del repositorio



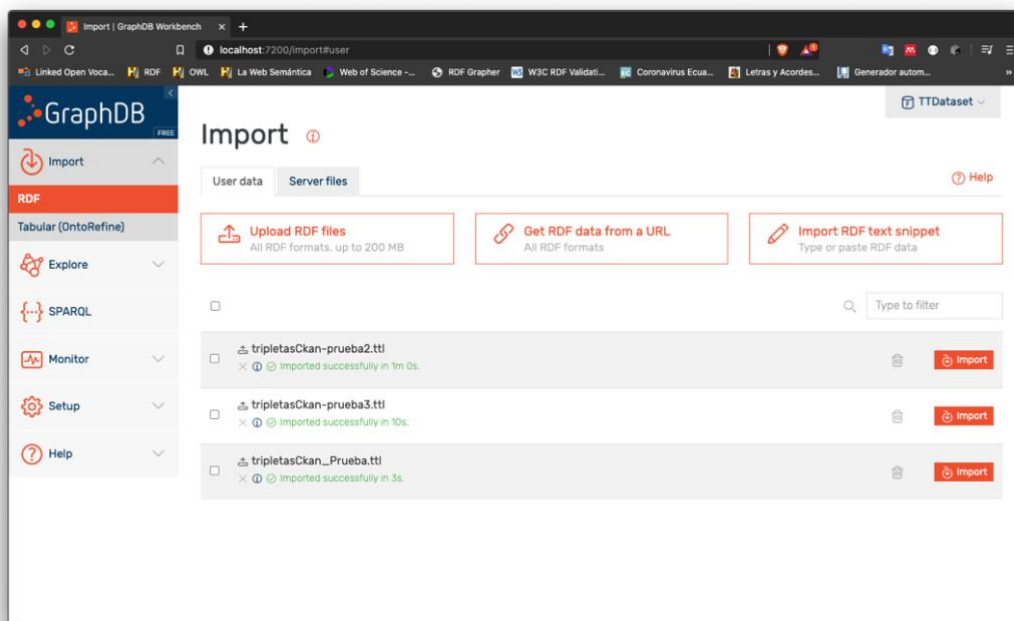
Una vez creado el repositorio se procede a la carga de la información. Esta sección se puede hacer mediante la carga directa desde la aplicación dando clic en “Upload RDF files” como se muestra en la **Figura 13**, o la otra forma mediante el uso del Endpoint de cada repositorio al que se quiera añadir la información, primero se debe obtener la URL del repositorio como se aprecia en la **Figura 13** y la configuración para la carga es la siguiente:

```
url -X POST --header "Content-Type:multipart/form-data" -F "config=@./config.ttl"  
"http://localhost:7200/<ID>/repositories"
```

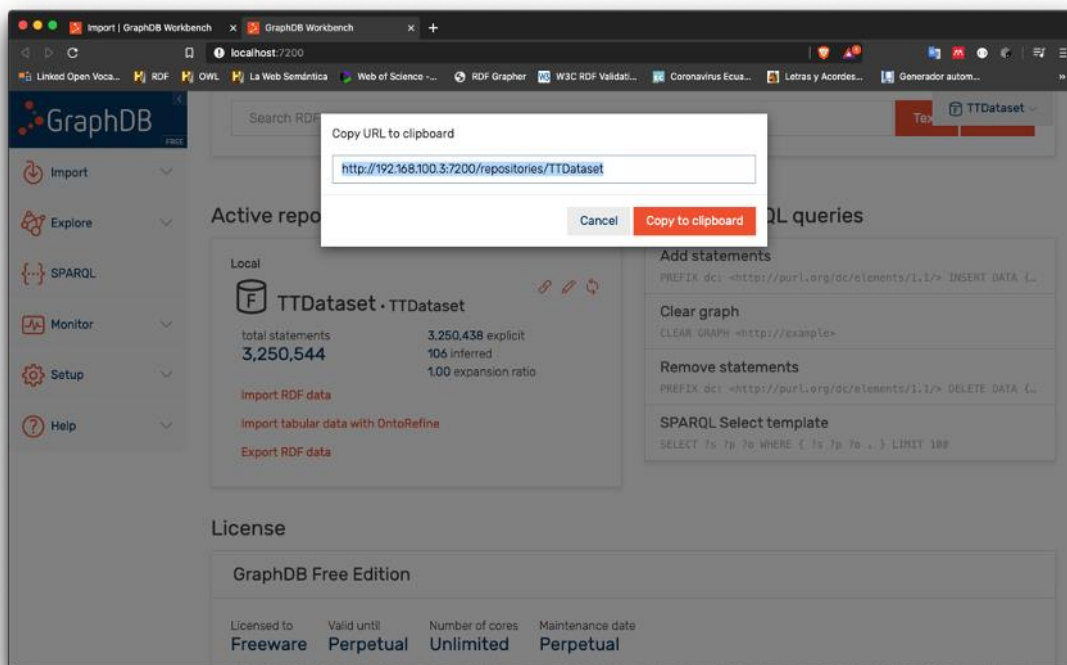
Con este comando se debe reemplazar el ID de la URL por el nombre del repositorio que usará y también se debe especificar el formato serializable RDF del archivo.

Figura 13

Administrador para la carga de la información en el repositorio

**Figura 14**

Administrador para obtener la URL del repositorio



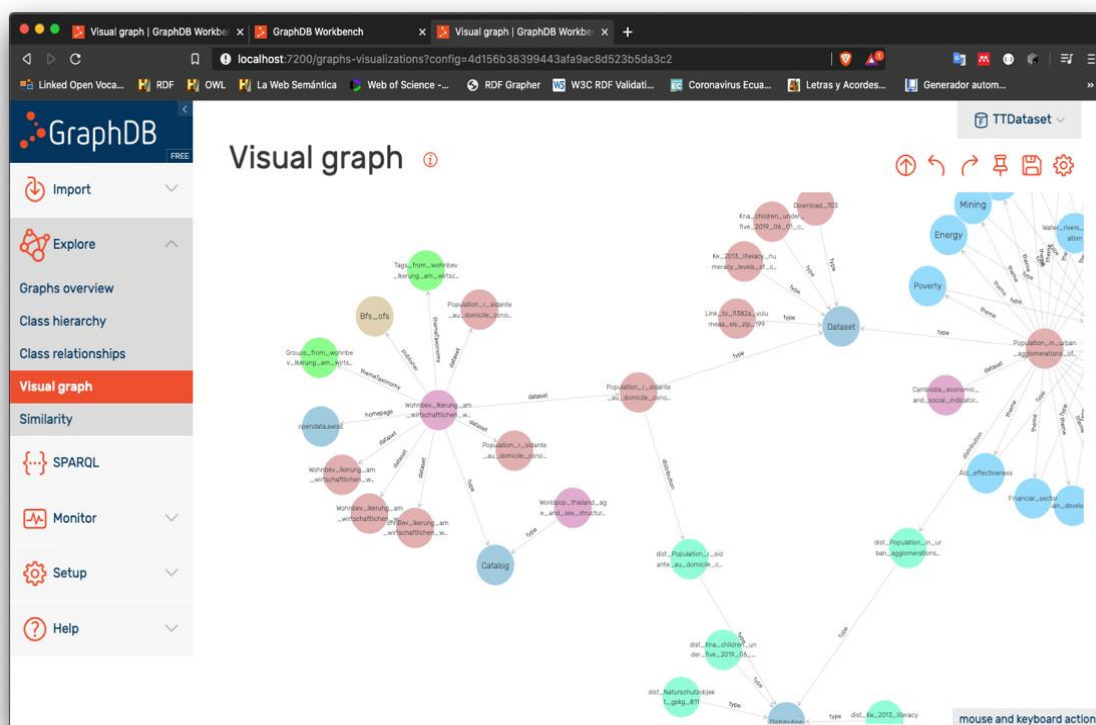
3.1.1.1 Exploración y visualización de datos

Esta plataforma GraphDB tiene muchas funcionalidades que son muy potentes, siendo así fácil de explorar los datos y sus herramientas que nos permiten visualizar instancias, clases, propiedades y sus relaciones a través de una jerarquía de clases del conjunto de datos que dan lugar al grafo de conocimiento.

Dentro de la aplicación Workbench también permite realizar la visualización del grafo con los datos existen del repositorio que se creó. Para la construcción del grafo se busca los recursos que están dentro de los datos almacenados o también se puede crear uno en base a una configuración previa. También se busca a través de una instancia poder explorar todas las relaciones y dando clic se puede ver más información de la instancia principal y de esta manera generar el grafo que se muestra en la **Figura 15**.

Figura 15

Visualizador de grafos



3.1.1.2 Jerarquía y relaciones de clases

Como se viene apreciando en las gráficas anteriores que GraphDB es una plataforma muy robusta que permite realizar múltiples cosas y en este caso se explica la jerarquía de las clases en base al número de instancias almacenadas como se lo presenta en la **Figura 16**. Dado este grafo se observa las clases hijas y padres.

La **Figura 17** presenta un diagrama que muestra las relaciones y la cantidad de enlaces que están asociadas a las clases. La **Figura 17** muestra el número total de las relaciones existentes a las clases como ejemplo en número de relaciones asociadas a `dcat:Dataset`.

Figura 16

Visualizador de grafos

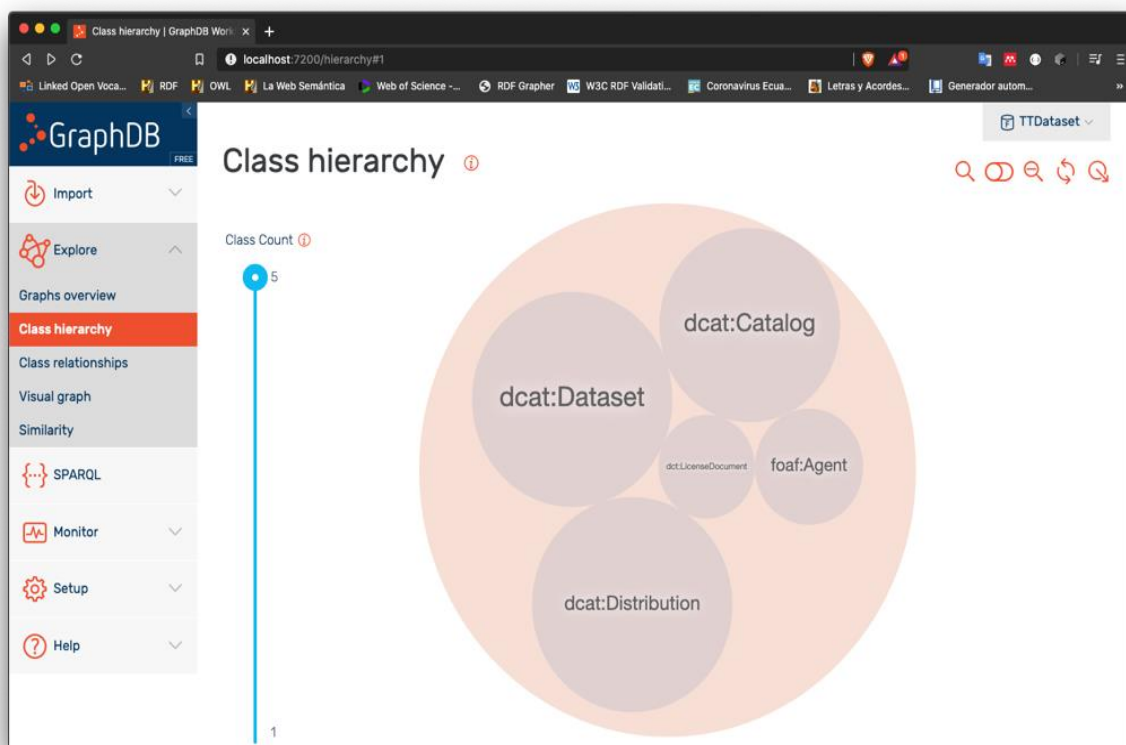
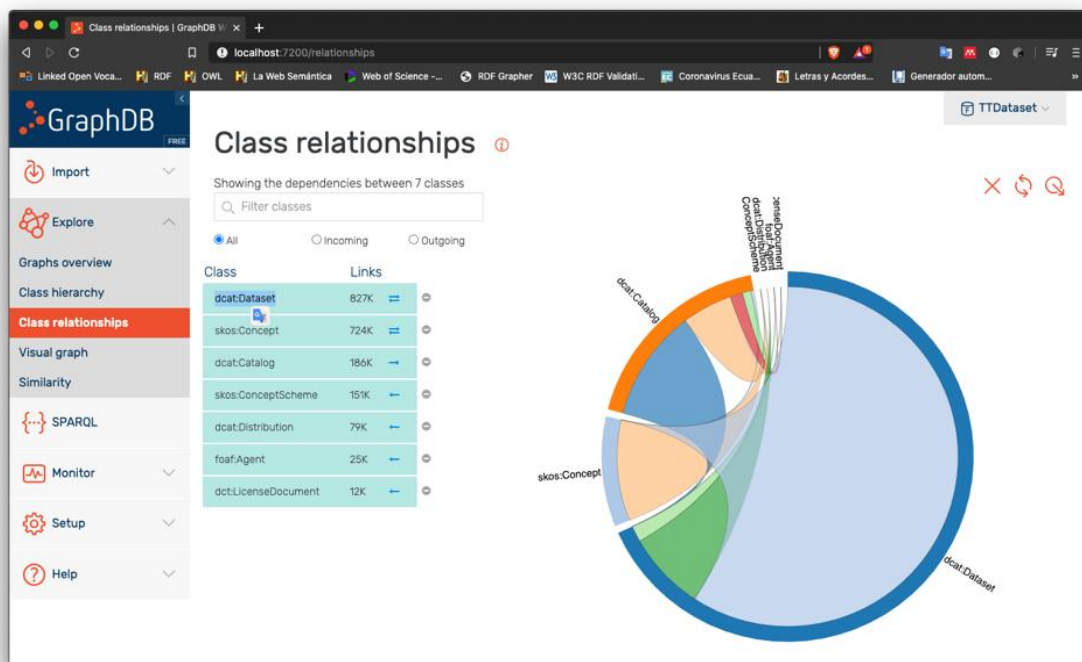


Figura 17

Diagrama de relaciones de clases

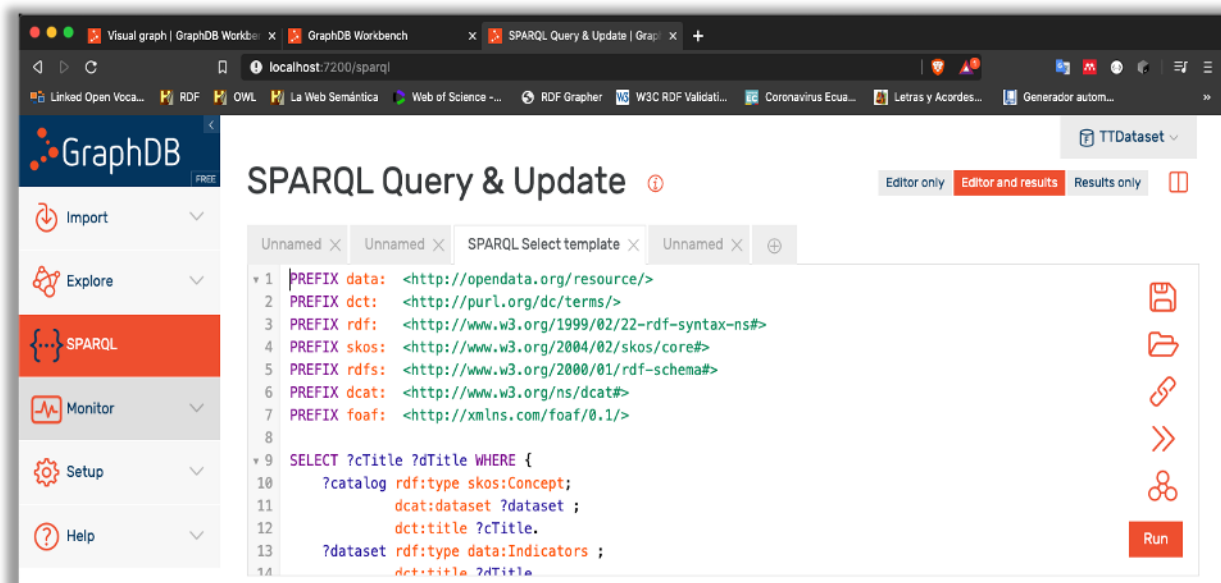


3.1.1.3 Consulta de datos

Finalmente se tiene las consultas de los datos que se encuentran en el repositorio de GraphDB, estos pueden ser usados para realizar las consultas mediante el editor SPARQL en Workbench. La **Figura 18** muestra un ejemplo de como se ven las consultas en el entorno. Las consultas SPARQL son súper sencillas y fácil de hacer con todas las funciones que da Workbench.

Figura 18

Diagrama de consultas Sparql



Nota: (Ontotext, 2020)

3.1.2 Anotaciones semánticas de los objetivos de desarrollo sostenibles

Los datos son usables si están en un formato abierto y con una estructura. La estructura que poseen los recursos de los Objetivos de Desarrollo Sostenibles (ODS) no son lo suficientemente abiertos para su análisis y reutilización. De esta manera se busca que estos datos se los pueda procesar y luego crear un enlace entre ellos poder obtener las anotaciones semánticas de los ODS.

3.1.2.1 Extracción de información de ODS

Para obtener la información de los ODS acerca de los objetivos, metas e indicadores se realizó una extracción de datos desde el sitio oficial en inglés en el cual se tomó en cuenta todas las características mencionadas para su extracción o también existe otra forma de extracción que es mediante el API que posee Naciones Unidas¹⁴ y para esto se usa los Endpoints¹⁵ que ofrece la ONU.

Con cualquiera de estas dos formas de extracción se obtuvo los 17 objetivos y cada uno de estos con sus características título, objetivos, metas e indicadores y así saber la

¹⁴ <https://unstats.un.org/SDGAPI/swagger/>

¹⁵ <https://unstats.un.org/SDGAPI/v1/sdg/Goal/List?includechildren=true>

finalidad de cada uno de estos objetivos. Si la extracción es mediante la API, el cual permite distribuir las solicitudes en diferentes parámetros para obtener una mejor información y características de los ODS y el resultado obtenido es cada objetivo con sus metas e indicadores como se aprecia en la **Figura 19**.

Figura 19

Resultado de consulta mediante la API de la ONU

```

▼ 1:
  code: "2"
  ▶ title: "End hunger, achieve food...sustainable agriculture"
  ▶ description: "Goal 2 seeks to end hung...nctioning food markets."
  uri: "/v1/sdg/Goal/2"
  ▼ targets:
    ▼ 0:
      goal: null
      code: "2.1"
      ▶ title: "By 2030, end hunger and ...ent food all year round"
      ▶ description: "By 2030, end hunger and ...ent food all year round"
      uri: "/v1/Target/2.1"
      ▼ indicators:

```

Nota. Tomado de ODS ONU [Fotografía], por (ONU, 2019)

3.1.3 Transformación de información de ODS

La transformación de estos datos se realiza para obtener una representación de conocimiento de cada uno de los objetivos, metas e indicadores. En esta fase se realizó dos procesos, uno para la creación de base de conocimiento esta se encuentra en inglés (**Figura 20**) para obtener esta información se realizó un proceso semiautomático el cual consiste en obtener la información manualmente y guardar un archivo; después mediante un algoritmo eliminar las Stop Word existentes en el texto obtenido. Después de esto se tiene correctamente la base de conocimiento. Se procede a realizar el procesamiento de lenguaje natural (NLP¹⁶) usando el motor de fasttext el cual permite interpretar esta información y así mejorar el conocimiento obtenido y buscar posibles relaciones con otra información.

¹⁶ Natural Language Processing

Figura 20

Base de conocimiento de los ODS

```

__label__SDG_1 __label__SDG_1aen End poverty forms
__label__SDG_2 __label__SDG_2aen End hunger achieve food security improved nutrition promote sustainable agriculture
__label__SDG_3 __label__SDG_3aen Ensure healthy lives promote well-being ages
__label__SDG_4 __label__SDG_4aen Ensure inclusive equitable quality education promote lifelong learning opportunities
__label__SDG_5 __label__SDG_5aen Achieve gender equality empower women girls
__label__SDG_6 __label__SDG_6aen Ensure availability sustainable management water sanitation
__label__SDG_7 __label__SDG_7aen Ensure access affordable reliable sustainable modern energy
__label__SDG_8 __label__SDG_8aen Promote sustained inclusive sustainable economic growth productive employment decent work
__label__SDG_9 __label__SDG_9aen Build resilient infrastructure promote inclusive sustainable industrialization foster innovation
__label__SDG_10 __label__SDG_10aen Reduce inequality countries
__label__SDG_11 __label__SDG_11aen Make cities human settlements inclusive safe resilient sustainable
__label__SDG_12 __label__SDG_12aen Ensure sustainable consumption production patterns
__label__SDG_13 __label__SDG_13aen Take urgent action combat climate change impacts[b]
__label__SDG_14 __label__SDG_14aen Conserve sustainably use oceans seas marine resources sustainable development
__label__SDG_15 __label__SDG_15aen Protect restore promote sustainable use terrestrial ecosystems sustainably manage forests comb
__label__SDG_16 __label__SDG_16aen Promote peaceful inclusive societies sustainable development provide access justice build effe
__label__SDG_17 __label__SDG_17aen Strengthen means implementation revitalize Global Partnership Sustainable Development
__label__SDG_1 __label__SDG_1aen __label__target_1-1 __label__target_1-1aen By 2030 eradicate extreme poverty people everywhere c
__label__SDG_1 __label__SDG_1aen __label__target_1-2 __label__target_1-2aen By 2030 reduce half proportion men women children age
__label__SDG_1 __label__SDG_1aen __label__target_1-3 __label__target_1-3aen Implement nationally appropriate social protection sy
__label__SDG_1 __label__SDG_1aen __label__target_1-4 __label__target_1-4aen By 2030 ensure men women particular poor vulnerable e
__label__SDG_1 __label__SDG_1aen __label__target_1-5 __label__target_1-5aen By 2030 build resilience poor vulnerable situations r
__label__SDG_1 __label__SDG_1aen __label__target_1-a __label__target_1-aen Ensure significant mobilization resources variety sou
__label__SDG_1 __label__SDG_1aen __label__target_1-b __label__target_1-baen Create sound policy frameworks national regional inte
__label__SDG_2 __label__SDG_2aen __label__target_2-1 __label__target_2-1aen By 2030 end hunger ensure access people particular po
__label__SDG_2 __label__SDG_2aen __label__target_2-2 __label__target_2-2aen By 2030 end forms malnutrition including achieving 20
__label__SDG_2 __label__SDG_2aen __label__target_2-3 __label__target_2-3aen By 2030 double agricultural productivity incomes sma
__label__SDG_2 __label__SDG_2aen __label__target_2-4 __label__target_2-4aen By 2030 ensure sustainable food production systems in
__label__SDG_2 __label__SDG_2aen __label__target_2-5 __label__target_2-5aen By 2020 maintain genetic diversity seeds cultivated p
__label__SDG_2 __label__SDG_2aen __label__target_2-a __label__target_2-aen Increase investment including enhanced international
__label__SDG_2 __label__SDG_2aen __label__target_2-b __label__target_2-baen Correct prevent trade restrictions distortions world
__label__SDG_2 __label__SDG_2aen __label__target_2-c __label__target_2-caen Adopt measures ensure proper functioning food commod

```

Lo que se busca el NLP es poder identificar del dataset principal y analizar que recursos tienen una probabilidad de ser semejantes y de esta manera se puedan relacionar los ODS a los demás recursos semánticos.

Y en la segunda parte una vez identificadas las entidades más importantes que se está relacionando a cada uno de los 17 objetivos, 33 metas y 21 indicadores. Para esto se usó la librería `rdflib` la cual permite analizar las entidades resultantes del NLP y buscar las posibles relaciones a cada ODS el cual contiene un parser/serializador a RDF/XML. Este resultado del uso de esta librería es un conjunto de tripletas serializadas en RDF/XML, el resultado está almacenado en un archivo denominado `SDG.rdf` (Figura 21).

Luego de obtener el análisis por cada ODS, se transformó y consolidó con la taxonomía de SKOS y para relacionar sus clases se usó DCAT y Schema.

Figura 21

Resultado de análisis de los ODS

```

1 <?xml version="1.0" encoding="utf-8" ?>
2 <rdf:RDF xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#"
3     xmlns:skos="http://www.w3.org/2004/02/skos/core#"
4     xmlns:ns0="http://dbpedia.org/ontology/"
5     xmlns:schema="http://schema.org/">
6
7     <skos:Collection rdf:about="http://ld.utpl.edu.ec/resource/ods">
8         <ns0:abbreviation>ODS</ns0:abbreviation>
9         <skos:perflabel>Objetivos de Desarrollo Sostenible</skos:perflabel>
10        <rdf:type type>
11            <rdf:Description rdf:about="http://ld.utpl.edu.ec/resource/nameES">
12                <schema:description>Objetivos de Desarrollo Sostenible</schema:description>
13                <skos:topConceptOf>
14                    <rdf:Description rdf:about="http://ld.utpl.edu.ec/resource/odsES1">
15                        <rdf:type>fin de la pobreza</rdf:type>
16                    </rdf:Description>
17                </skos:topConceptOf>
18
19                <skos:topConceptOf>
20                    <rdf:Description rdf:about="http://ld.utpl.edu.ec/resource/odsES2">
21                        <rdf:type>hambre cero</rdf:type>
22                    </rdf:Description>
23                </skos:topConceptOf>
24
25                <skos:topConceptOf>
26                    <rdf:Description rdf:about="http://ld.utpl.edu.ec/resource/odsES3">
27                        <rdf:type>salud y bienestar</rdf:type>
28                        <rdf:type>salud</rdf:type>
29                        <rdf:type>bienestar</rdf:type>
30                    </rdf:Description>
31                </skos:topConceptOf>
32
33                <skos:topConceptOf>
34                    <rdf:Description rdf:about="http://ld.utpl.edu.ec/resource/odsES4">
35                        <rdf:type>educación de calidad</rdf:type>
36                    </rdf:Description>

```

Con estas tripletas en formato RDF/XML se puede identificar fácilmente los objetivos, metas e indicadores que han sido representados por el nombre y características que se agruparon con todos los conceptos. De tal manera se puede usar el análisis y buscar si existe alguna relación con otros textos. Finalmente, estos datos pueden ser usados para buscar si existe la probabilidad de tener relaciones entre los ODS y un recurso, y estos permitirán identificar como están incrementando el conocimiento de los ODS con relación a los dataset.

3.2 Análisis - creación de enlaces entre ods y dataset semánticos

La creación de enlaces entre el Dataset ckan Semántico (dataset base) con otros datasets externos y los ODS, de esta manera obtener una nueva representación semántica del conjunto de datos ckan.

Inicialmente para la creación de estos enlaces se busca identificar cuales son las entidades principales y cuales son literales, con estos se pueden interpretar que cada uno de los recursos que serán procesados en la parte del procesamiento de lenguaje natural (NLP). Con esto se busca interpretar el contenido de las entidades, de los literales y determinar la similitud entre los conjuntos de datos de tal manera que permitan obtener nuevos enlaces entre ellos. En los siguientes ítems se describe todos los pasos para la creación de los enlaces.

3.2.1 Traducción

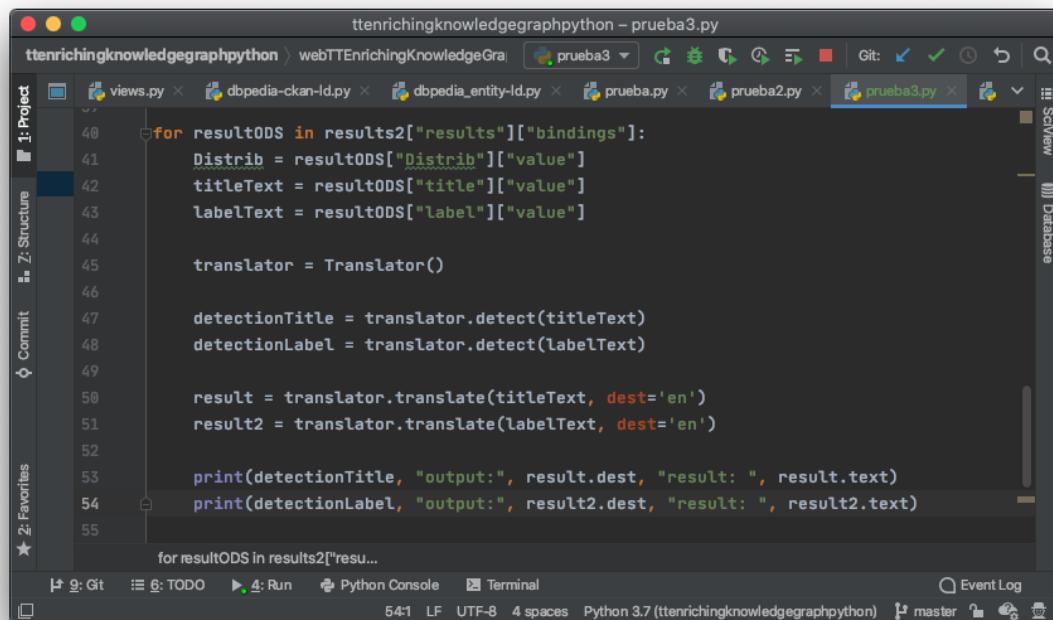
Conceptualmente, la traducción es la comunicación del significado de un texto en el idioma de origen, por medio de un texto equivalente en el idioma de destino. En este trabajo debido a que los recursos provienen desde diferentes fuentes de información y en varios idiomas, a cada recurso se aplicó un proceso de traducción. Se ha realizado un proceso de análisis y interpretación de idioma con el framework de google traductor (googletrans¹⁷). Esta herramienta permite realizar una interpretación del texto que entra y así determinar el idioma en que se encuentra, dando como resultado la similitud de idioma y traducción. Como parte del alcance de este trabajo se usó la traducción al idioma inglés, debido a que este idioma es el más usado en los datos utilizados.

En la **Figura 22** se aprecia un fragmento del código de la parte de la traducción que se realiza previo al análisis de traducción.

¹⁷ <https://pypi.org/project/googletrans/>

Figura 22

Traducción de texto con google traductor api



```

ttenrichingknowledgegraphpython - prueba3.py
ttenrichingknowledgegraphpython webTTErichingKnowledgeGra prueba3
views.py dbpedia-ckan-ld.py dbpedia_entity-ld.py prueba.py prueba2.py prueba3.py
40 for result0DS in results2["results"]["bindings"]:
41     Distrib = result0DS["Distrib"]["value"]
42     titleText = result0DS["title"]["value"]
43     labelText = result0DS["label"]["value"]
44
45     translator = Translator()
46
47     detectionTitle = translator.detect(titleText)
48     detectionLabel = translator.detect(labelText)
49
50     result = translator.translate(titleText, dest='en')
51     result2 = translator.translate(labelText, dest='en')
52
53     print(detectionTitle, "output:", result.dest, "result: ", result.text)
54     print(detectionLabel, "output:", result2.dest, "result: ", result2.text)
55
for result0DS in results2["resu...
54:1 LF UTF-8 4 spaces Python 3.7 (ttenrichingknowledgegraphpython) master

```

3.2.2 NER y NLP

El Reconocimiento de Entidades Nombradas (NER¹⁸) permite la extracción de información que busca localizar y clasificar en categorías predefinidas, tales como: personas, organizaciones, lugares y las entidades nombradas encontradas en un texto. Lo que se busca es identificar diversas características que tiene la información que se está enriqueciendo de tal manera que en el NLP sea más rápido y el porcentaje de similitud sea mayor. Y por otra parte lo que busca el procesamiento de lenguaje natural es el reconocimiento de entidades en el cual quiere idear búsquedas lingüísticas específicas y detectar similitudes entre entidades.

¹⁸ Reconocimiento de Entidades Nombradas.

3.2.2.1 FastText

Teniendo estas premisas, para el NLP en primera instancia se usó fasttext como herramienta para este proceso, el cual busca analizar los datos y determinar cuál es el rango de similitud entre un texto con base de conocimiento que fue entrenada previamente y aplicado el NER, en la **Figura 20** se presenta un fragmento de la base de conocimiento de los ODS¹⁹.

Con la base de conocimiento previamente estructurada y revisada se procede a realizar el entrenamiento con la finalidad de obtener un rango de similitud para esto en la documentación de fasttext muestra que se debe usar para wordNgrams=1 y se tiene que trabajar con un rango de aprendizaje entre 0.5 a 0.6, cuando se usa wordNgrams=2 el rango debe ir entre 0.7 a 1.0 como mejor modelo de entrenamiento. En este trabajo se usó el wordNgrams=1 con el fin de obtener un acercamiento más real a la información obtenida y el entrenamiento con 0.5 en el rango de similitud siendo esto la media para determinar si ese valor es similar o no. La **Figura 23** muestra un fragmento del código en Python del proceso que se usó para el entrenamiento supervisado, para esto se realizó la herramienta de fasttext.

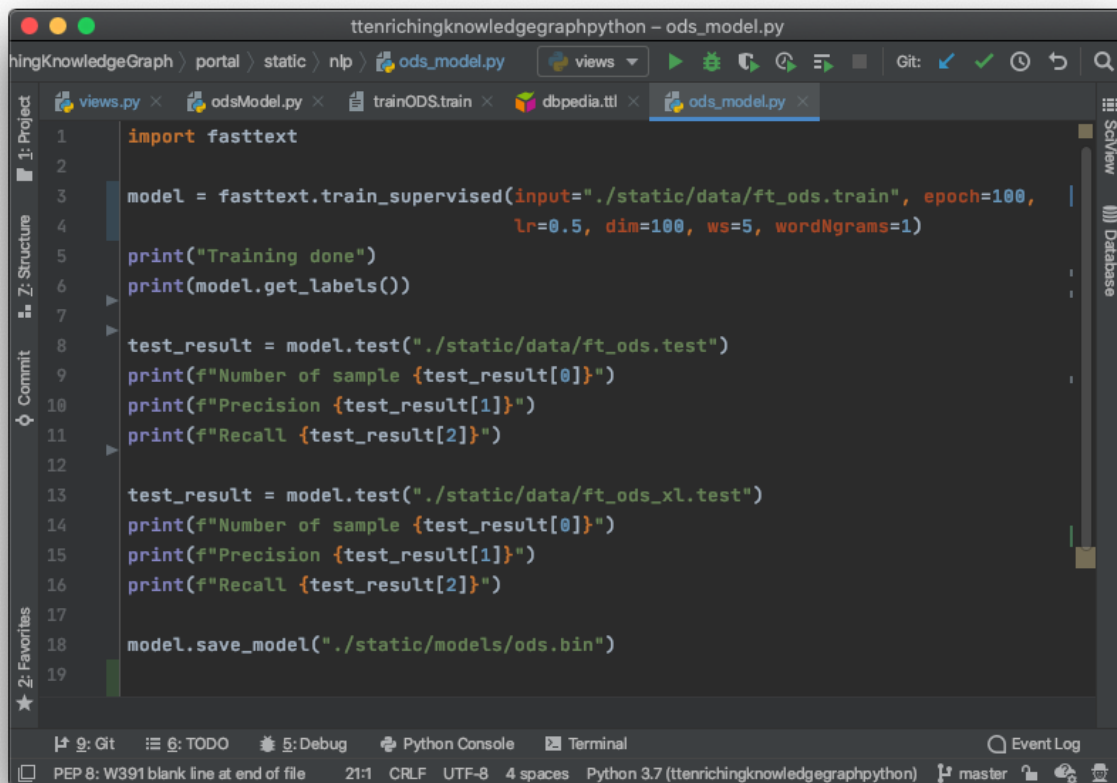
En el **Apéndice 1** se presenta un fragmento del modelo resultante del entrenamiento.

Una vez realizado este proceso nos da como resultado un modelo, el cual será la base en la que dando los parámetros en texto los analizar y tomará de referencia para determinar a que se puede relacionar o que pueda presentar alguna similitud que sea igual o mayor a 0.5 hasta 1, con esto se busca que el resultado de la interpretación sea más preciso.

¹⁹ Objetivos de Desarrollo Sostenibles

Figura 23

Proceso de entrenamiento de los ODS



```

1  import fasttext
2
3  model = fasttext.train_supervised(input="./static/data/ft_ods.train", epoch=100,
4                                  lr=0.5, dim=100, ws=5, wordNgrams=1)
5  print("Training done")
6  print(model.get_labels())
7
8  test_result = model.test("./static/data/ft_ods.test")
9  print(f"Number of sample {test_result[0]}")
10 print(f"Precision {test_result[1]}")
11 print(f"Recall {test_result[2]}")
12
13 test_result = model.test("./static/data/ft_ods_xl.test")
14 print(f"Number of sample {test_result[0]}")
15 print(f"Precision {test_result[1]}")
16 print(f"Recall {test_result[2]}")
17
18 model.save_model("./static/models/ods.bin")
19

```

Para la detección de las similitudes entre los recursos y los ODS, se realiza un análisis previo con fasttext el cual me permite interpretar el contexto y por consiguiente obtener una interpretación mejor. Antes de empezar la detección de similitudes se hizo la traducción del contenido con la finalidad de que todos los recursos sean analizados y de esta manera obtener una mayor calidad de los resultados.

En el análisis se tiene que cargar primero el modelo entrenado el cual servirá para el proceso de NLP supervisado, lo que se realiza en esta parte es adecuar la información que le llega al analizador, luego de esto se crea un variable donde se almacena el resultado que es un arreglo, el cual tiene como salida la predicción, el label (url) y el porcentaje de similitud siendo este mayor a 0.5.

La **Figura 24** presenta un fragmento del código que permite la detección y obtener de la similitud del recurso con un ODS.

Figura 24

Proceso de predicción con fastText

```

ttenrichingknowledgegraphpython - views.py
ttenrichingknowledgegraphpython > webTTErichingKnowle prueba3
views.py dbpedia-ckan-ld.py dbpedia_entity-ld.py graphdb_conect.py models.py
368 text_filter = str(kdata4)
369 text_filterlabel = str(kdataSujeto4)
370
371 search = model.predict(text=text_filter, k=-1, threshold=0.5)
372 for i in search[0]:
373     try:
374         time.sleep(1)
375
376         iNew = i.replace('__label__', 'http://opendata.org/resource/')
377
378         # reemplazar valor que añadie "empowermentwomen"
379         iAll = iNew.replace('empowermentwomen', '')
380
381         executor.submit(saveData, str(Distrib), str(iAll))
382
383     except Exception as x:
384         print("[ERROR] Oops! Se perdió la conexión... ", x)
385
fasttext_nlp()
9: Git 6: TODO 4: Run Python Console Terminal Event Log
401:25 LF UTF-8 4 spaces Python 3.7 (ttenrichingknowledgegraphpython) master

```

Una vez obtenido el resultado del análisis se procede a interpretar el resultado. Lo que se busca con esto es la interpretación del contenido y saber que información es la que me sirve o no, y así enviar a un método que se lo denomina *saveData* el cual realiza el proceso de almacenar la información directamente en la base de datos semántica en este caso en GraphDB en un formato de tripletas.

Luego de haber sido guardada la información esta ya puede ser consultada y consumida por una API para luego su explotación. Lo que se busca con esta herramienta de fasttext es que se pueda analizar cada recurso sea con el label, title o description y como resultado nos devuelva a que Objetivos de Desarrollo Sostenible se relaciona y de esta manera enriquecer el grafo principal con una mayor información.

En el **Apéndice 2** se presenta un fragmento del resultado final del enriquecimiento con fasttext.

3.2.2.2 DBpedia spotlight

Primero la DBpedia es un proyecto para la extracción de datos desde la Wikipedia y ellos proponen una versión de la web semántica en el contenido. En el proceso de estructuración del sitio hacia la web semántica, Dbpedia desarrollo una API que la denominan dbpedia-spotlight el cual es una herramienta para anotar automáticamente las menciones de los recursos de DBpedia en un texto. Y así proporcionar una solución para vincular fuentes de información no estructuradas a la nube Linked Open Data a través de DBpedia.

Esta herramienta funciona con un enfoque de cuatro pasos establecidos por la DBpedia. DBpedia Spotlight realiza la extracción de entidades nombradas, incluyendo la detección de entidades y la resolución de nombres. También se puede utilizar para el reconocimiento de entidades nombradas, entre otras tareas de extracción de información.

Sabiendo los beneficios de este framework se procede a realizar la primera actividad que es el NLP y NER, en este caso usando spacy y de esta existe una extensión en la que también esta trabajando dbpedia-spotlight la cual se denomina *spacy_dbpedia_spotlight*, con esta biblioteca de software de código abierto de procesamiento avanzado de lenguaje natural, esta busca la interpretación sintáctica y semántica de los parámetros que se le pase desde los recursos que se analizó.

Lo que se busca al usar spacy es que el procesamiento de texto en nuestro caso el texto del label, title o Description. Estos textos se encuentran sin procesar, la mayoría de las palabras son raros y es común para las palabras que se ven completamente diferente a significar casi lo mismo. Las mismas palabras en un orden diferente pueden significar algo completamente diferente. Incluso dividir el texto en unidades útiles similares, estos pueden resultar difícil en muchos idiomas.

Si bien es posible resolver algunos problemas comenzando solo con los caracteres en bruto, generalmente es mejor usar el conocimiento lingüístico para agregar información útil.

Para poder identificar las palabras que están relacionadas entre si, spacy realiza una tokenización, con el fin de que esta pueda analizar y etiquetar un documento determinado. Aquí es donde entra el modelo estadístico o en otras palabras el modelo de etiquetas previamente entrenado, que permite a spacy hacer una predicción y el resultado que es mayor a 0.5 siendo esta la media que permita etiquetar que palabra se relaciona con la probabilidad del contexto. Un modelo consta de datos binarios y en diferentes idiomas.

Las anotaciones lingüísticas están disponibles como atributos del token en la herramienta. Entonces, para obtener la representación de cadena legible de un atributo necesitamos agregar un guión bajo “_” a su nombre como por ejemplo:

- *token.lemma_*
- *token.pos_*
- *token.tag_*
- *token.dep_*
- *token.shape_*
- *token.kb_id_*
- *token.is_alpha*
- *token.is_stop*

Cada uno de estos traen un resultado como:

- Texto: el texto de la palabra original.
- Lemma: La forma básica de la palabra.
- POS: La etiqueta simple de parte del discurso de Universal POS tags.
- Tag: la etiqueta detallada de la parte del discurso.
- Dep: Dependencia sintáctica, es decir, la relación entre tokens.

- Shape: la forma de la palabra: mayúsculas, puntuación, dígitos.
- Kb id: URL resultado de las palabras similares.
- is alpha: ¿el token es un carácter alfa?
- is stop: ¿El token forma parte de una lista de *stop*, es decir, las palabras más comunes del idioma?

Entonces, la mayoría de las etiquetas se ven bastantes abstractas y varían entre los idiomas. De la misma manera trabaja la librería de *dbpedia-spotlight* pero en este caso usando la herramienta de *spacy* ayudará a la interpretación del contexto del contenido entrante lo que hace *spacy_dbpedia_spotlight* trabajo con las anotaciones lingüísticas y con reconocimiento de entidad nombrada.

Lo que *spacy* hace con el NER es reconocer una entidad nombrada, el cual es un "objeto del mundo real" al que se le asigna un nombre, por ejemplo, una persona, un país, un producto, el título de un libro y en nuestro caso la URI del recurso que se encuentre estrechamente relacionada al recurso analizado para esto se usa la etiqueta *ent.kb_id_*. Lo que se busca es poder identificar las entidades que se encuentran en el texto y una vez reconocidas, estas pueden ser interpretadas y ver que información contiene y con esto crear el enriquecimiento del recurso con el recurso de la DBpedia. A continuación se describe un ejemplo para tener una mayor claridad del proceso.

De todos los recursos que se tiene se tomó un recurso de ejemplo, el cual dice "*Distribution of: Costs of proceedings - Australia*", este habla sobre costos de distribución en Austria, de este recurso se toma la **Description**, la **Figura 25** muestra el texto de entrada sin procesar, en este caso será el texto de entrada que recibe *spacy* para su análisis e interpretación.

Dentro de la documentación de *spacy* especifica las variables que serán usadas para la interpretación de texto de entrada entre estas se encuentra *ent.text* (texto resultado), *ent.label_* (Valor de análisis sea verbo o algún resultado que nos de *spacy*) y *ent.kb_id_* (token resultado) los cuales son los principales recursos que voy a tomar en

cuenta del resultado. Con estas variables de análisis se activa el método para que comience a trabajar e interpretar el texto.

La **Figura 26** presenta el valor resultante de manera visual, para obtener esto se realiza un mapeo con las anotaciones lingüísticas y el reconocimiento de entidades nombradas que se encuentre en el texto. Con estas nuevas etiquetas ya se puede saber que información me sirve para la creación de la nueva relación o el enriquecimiento de los datos y este nuevo recurso para su presentación se lo puede identificar con la URL de la DBpedia en el resultado y también lo presenta un como *dbr* que se hace referencia la URI “<http://dbpedia.org/resource/NameResource>” o “<http://dbpedia.org/page/NameResource>”, y de esta manera se obtenga el resultado *dbr* y enviado para su almacenamiento. La **Figura 27** muestra un fragmento del código que nos permite realizar este proceso.

En el **Apéndice 3** se presenta un fragmento del resultado final del enriquecimiento usando spacy Dbpedia-spotlight.

Figura 25

Fragmento de texto de entrada ejemplo

Authoritative descriptive metadata for: Farm
share and price spread in Australia's beef
supply chain

Figura 26

Ejemplo de interpretación del NER y el NLP con spacy Dbpedia-spotlight

Authoritative **descriptive verb** metadata **dbr** for: Farm
share and **price \$** spread in **Australia's dbr** **beef food**
supply **chain dbr**

Figura 27

Proceso de predicción con *dbpedia-spotlight*

```

ttenrichingknowledgegraphpython - prueba3.py
ttenrichingknowledgegraphpython webTTErichingKnowledgeGra| prueba3
views.py x prueba3.py x prueba2.py x prueba.py x sqarqL_graphdb.py x dbpedia-ckan-ld.py x dbpedia
45 for result in results2["results"]["bindings"]:
46     Distrib = result["Distrib"]["value"]
47     titleText = result["title"]["value"]
48     labelText = result["label"]["value"]
49
50     # if para filter de la data Objeto
51     if titleText == titleText:
52         try:
53             # carga del modelo en el idioma
54             nlp = spacy_dbpedia_spotlight.load('en')
55
56             # pasar el texto a las funcion de parametrizar (parameter) con la dbpedia
57             text = text_filter
58
59             # cargar el texto al modelo
60             doc = nlp(text)
61
62             for ent in doc.ents:
63                 # resultado de analisis
64                 entitiesList = (ent.text, ent.start_char, ent.end_char, ent.kb_id_, ent.label_)
65                 # valor de URI
66                 entitiesId = ent.kb_id_
67                 # Valor de analisis
68                 has_vectord = ent.label_
69
70                 executor.submit(save_dbpedia, str(Distrib), str(entitiesId))
71
72             except Exception as x:
73                 print("[ERROR] Oops! Se perdió la conexión...: ", x)
74
for result in results2["results..."] if titleText == titleText > try
Python Console Terminal Event Log
56:35 LF UTF-8 4 spaces Python 3.7 (ttenrichingknowledgegraphpython) master

```

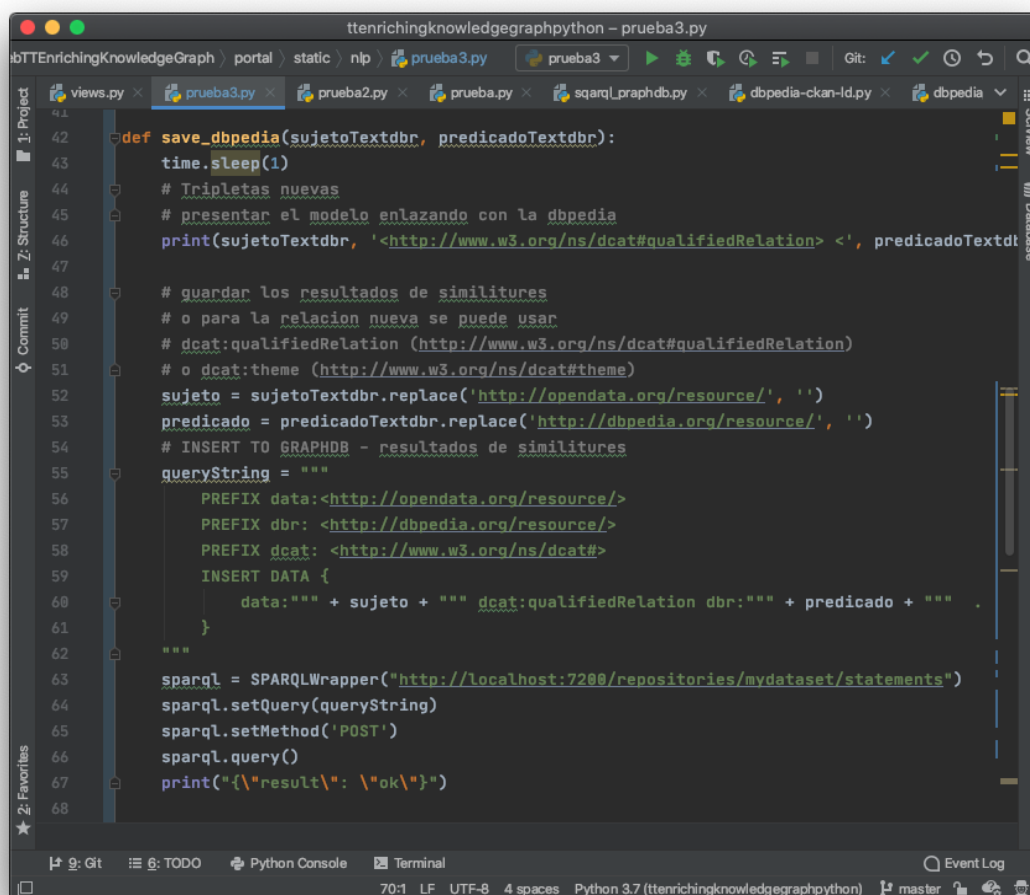
3.2.3 Creación de nuevas relaciones

Para la creación o enlazamiento de nuevas entidades, se creó un método en el que se envía los parámetros resultantes de los análisis realizados con las dos herramientas. En este método se crea un proceso de inserción de datos directamente a la base de datos semántica, en este caso GraphDB, para esto se crea una consulta SPARQL de inserción para que se pueda almacenar el resultado final. Los parámetros que recibe son la URI del recurso interpretado más en resultado obtenido del análisis es decir la URI del nuevo

recurso que se asocia a este. La **Figura 28** presenta el proceso para almacenar la información en el triple store.

Figura 28

Método de almacenamiento en GraphDB



```

42 def save_dbpedia(sujetoTextdbr, predicadoTextdbr):
43     time.sleep(1)
44     # Tripletas nuevas
45     # presentar el modelo enlazando con la dbpedia
46     print(sujetoTextdbr, '<http://www.w3.org/ns/dcat#qualifiedRelation> <', predicadoTextdbr)
47
48     # guardar los resultados de similitudes
49     # o para la relacion nueva se puede usar
50     # dcat:qualifiedRelation (http://www.w3.org/ns/dcat#qualifiedRelation)
51     # o dcat:theme (http://www.w3.org/ns/dcat#theme)
52     sujeto = sujetoTextdbr.replace('http://opendata.org/resource/', '')
53     predicado = predicadoTextdbr.replace('http://dbpedia.org/resource/', '')
54     # INSERT TO GRAPHDB - resultados de similitudes
55     queryString = """
56     PREFIX data:<http://opendata.org/resource/>
57     PREFIX dbr: <http://dbpedia.org/resource/>
58     PREFIX dcat: <http://www.w3.org/ns/dcat#>
59     INSERT DATA {
60         data:"" + sujeto + "" dcat:qualifiedRelation dbr:"" + predicado + "" .
61     }
62     """
63     sparql = SPARQLWrapper("http://localhost:7200/repositories/mydataset/statements")
64     sparql.setQuery(queryString)
65     sparql.setMethod('POST')
66     sparql.query()
67     print("{\"result\": \"ok\"}")
68
  
```

Esto se aplica para los dos procesos de análisis sea para los ODS o la Dbpedia y adicionalmente se lo almacena en una base de datos SQL para revisar el contenido que esta siendo almacenado permitiendo analizar si lo obtenido está acorde al contenido analizado o no.

En el **Apéndice 5** se presenta de manera visual un fragmento de todos los datos que fueron insertados en GraphDB.

3.3 Visualización

La visualización de los datos enriquecidos se presenta en la sección 4.2 con los casos de aplicación en la detección de similitudes del dataset y los ODS. El resultado de la visualización se presenta con el sitio web piloto desarrollando.

Capítulo Cuatro

Pruebas y Resultados

En este capítulo se especifica un caso de aplicación para la detección de posibles relaciones con diferentes umbrales de entrenamiento en el proceso de NLP. Además, se evalúa el sitio piloto web. La sección 4.1 describe el caso de aplicación para el proceso de creación de enlaces entre los ODS y el dataset. Y en la sección 4.2 describe el resultado final de TT que es la visualización de los datos enriquecidos.

4.1 Casos de aplicación

Esta sección pretende mostrar varios ejemplos de algunos datasets que se tiene y presentar el proceso de enriquecimiento de los datos CKAN con la Dbpedia y con los Objetivos de Desarrollo Sostenibles.

Para el proceso de enriquecimiento con fasttext y la Dbpedia se trabajar con los 20 primeros dataset. A estos se les aplicará todos los procesos mencionados en las secciones anteriores como el NER y el NLP. Como primeros pasos se realiza una consulta básica con SPARQL para obtener los valores que van a ser analizados, en este caso serian el título y la descripción de recurso como se muestra en la **Figura 29**.

Figura 29

Consulta SPARQL para enriquecimiento

The screenshot displays the GraphDB SPARQL Query & Update interface. The query editor shows the following SPARQL query:

```

1 PREFIX data: <http://opendata.org/resource/>
2 PREFIX dcat: <http://www.w3.org/ns/dcat#>
3 PREFIX skos: <http://www.w3.org/2004/02/skos/core#>
4 PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
5 PREFIX schema: <http://schema.org/>
6 PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
7 PREFIX dct: <http://purl.org/dc/terms/>
8 PREFIX foaf: <http://xmlns.com/foaf/0.1/>
9 SELECT DISTINCT * WHERE {
10   ?Distrib rdf:type dcat:Catalog;
11   dct:title ?title ;
12   dct:description ?description .
13 }

```

The results table shows the following data:

	Distrib	title	description
1	data:Employees_devices_to_acce	"Employees - devices to access the internet (NACE Rev. 1.1 activity)"	"Employees - devices to access the internet (NACE Rev. 1.1 activity)"
2	data:Dg_near_pre_accession_ass	"DG NEAR - Pre-accession assistance to Iceland"	"Information on projects financed through the Instrument for Pre-accession Assistance"

Usando el método de fasttext para la predicción de similitud que se presenta en la **Figura 24**, en este caso se trabaja con un rango de 0,5 y con ese rango obtendremos que algunas relaciones sean semejantes o no. La **Figura 30** muestra la ejecución del enriquecimiento de los datasets.

Figura 30

Captura del proceso de enriquecimiento con un rango de 0.5 de similitud

```

[INFO] ...
[INFO] ....
[INFO] Traducción y enriquecimiento
(('__label__odsInEN3',), array([0.81227142]))
(('__label__odsInEN3',), array([0.81227142]))
http://opendata.org/resource/Employees_devices_to_access_the_internet_nace_rev_1_1_activity_catalog <http://www.w3.org/ns/dcat#theme> < http://opendata.org/z
esource/odsInEN3 > .
{"result": "ok"}
http://opendata.org/resource/Employees_devices_to_access_the_internet_nace_rev_1_1_activity_catalog <http://www.w3.org/ns/dcat#theme> < http://opendata.org/z
esource/odsInEN3 > .
{"result": "ok"}
(('__label__odsInEN11',), array([0.74275362]))
(('__label__odsInEN4',), array([0.6884319]))
http://opendata.org/resource/Dg_near_pre_accession_assistance_to_iceland_catalog <http://www.w3.org/ns/dcat#theme> < http://opendata.org/resource/odsInEN11 >
.
{"result": "ok"}
http://opendata.org/resource/Dg_near_pre_accession_assistance_to_iceland_catalog <http://www.w3.org/ns/dcat#theme> < http://opendata.org/resource/odsInEN4 >
.
{"result": "ok"}
(('__label__odsEN5empowermentwomen',), array([0.99121553]))
((), array([], dtype=float64))
http://opendata.org/resource/Stress_test_for_bank_banco_comercial_portugues_sa_bcp_or_millennium_bcp_catalog <http://www.w3.org/ns/dcat#theme> < http://opend
ata.org/resource/odsEN5 > .
{"result": "ok"}
((), array([], dtype=float64))
(('__label__odsEN11',), array([0.96681297]))
http://opendata.org/resource/Rem_data_bank_year_1997_catalog <http://www.w3.org/ns/dcat#theme> < http://opendata.org/resource/odsEN11 > .
{"result": "ok"}
((), array([], dtype=float64))
((), array([], dtype=float64))
(('__label__odsEN5',), array([0.57148886]))
((), array([], dtype=float64))
http://opendata.org/resource/Wms_reb_kataster_ffentlich_rechtliche_eigentumsbeschr_nkungen_ogd_catalog <http://www.w3.org/ns/dcat#theme> < http://opendata.or
g/resource/odsEN5 > .
{"result": "ok"}
(('__label__odsEN5empowermentwomen',), array([0.71693534]))
((), array([], dtype=float64))
http://opendata.org/resource/Notariatskreis_catalog <http://www.w3.org/ns/dcat#theme> < http://opendata.org/resource/odsEN5 > .
{"result": "ok"}
((), array([], dtype=float64))
((), array([], dtype=float64))
(('__label__odsEN5empowermentwomen',), array([0.99121553]))
(('__label__odsInEN17',), array([0.80294442]))

```

Con el rango de similitud de 0.5 se obtiene 28 nuevas tripletas de relación de los 20 recursos analizados. Se observa que un recurso puede tener más de dos relaciones alineadas a los ODS. Para esto se requiere otro análisis o interpretación humana con el fin de obtener una información más precisa o que tenga relación semánticamente.

A continuación, se ejecuta el mismo proceso de enriquecimiento con fasttext pero en este caso el rango de similitud se incrementa a 0.8 para obtener una mayor precisión.

Figura 31

Captura del proceso de enriquecimiento con un rango de 0.8 de similitud

```

[INFO] ...
[INFO] ...
[INFO] Traducción y enriquecimiento
({'_label__odsInEN3',), array([0.81227142])}
({'_label__odsInEN3',), array([0.81227142])}
http://opendata.org/resource/Employees_devices_to_access_the_internet_nace_rev_1_1_activity_catalog <http://www.w3.org/ns/dcat#theme> < http://opendata.org/r
esource/odsInEN3 > .
{"result": "ok"}
http://opendata.org/resource/Employees_devices_to_access_the_internet_nace_rev_1_1_activity_catalog <http://www.w3.org/ns/dcat#theme> < http://opendata.org/r
esource/odsInEN3 > .
{"result": "ok"}
([], array([], dtype=float64))
([], array([], dtype=float64))
({'_label__odsENSEmpowermentwomen',), array([0.99121553])}
([], array([], dtype=float64))
http://opendata.org/resource/Stress_test_for_bank_banco_comercial_portugues_sa_bcp_or_millennium_bcp_catalog <http://www.w3.org/ns/dcat#theme> < http://opend
ata.org/resource/odsEN5 > .
{"result": "ok"}
([], array([], dtype=float64))
({'_label__odsEN11',), array([0.96681297])}
http://opendata.org/resource/Rom_data_bank_year_1997_catalog <http://www.w3.org/ns/dcat#theme> < http://opendata.org/resource/odsEN11 > .
{"result": "ok"}
([], array([], dtype=float64))
([], array([], dtype=float64))
([], array([], dtype=float64))
([], array([], dtype=float64))
([], array([], dtype=float64))
([], array([], dtype=float64))
([], array([], dtype=float64))
([], array([], dtype=float64))
([], array([], dtype=float64))
({'_label__odsENSEmpowermentwomen',), array([0.99121553])}
({'_label__odsInEN17',), array([0.88294442])}
http://opendata.org/resource/Wcs_digitales_terrain_dtm_und_oberfl_chenmodell_dom_catalog <http://www.w3.org/ns/dcat#theme> < http://opendata.org/resource/ods
InEN17 > .
{"result": "ok"}
([], array([], dtype=float64))
({'_label__odsENSEmpowermentwomen',), array([0.99121553])}
http://opendata.org/resource/Wah_sozialhilfebeziehende_von_18_bis_64_jahren_nach_erwerbssituation_und_besch_ftigungsgrad_catalog <http://www.w3.org/ns/dcat#t
heme> < http://opendata.org/resource/odsEN5 > .
{"result": "ok"}

```

Luego de análisis con el grado de similitud mayor, en este caso mayor a 0.8 de los 20 recursos analizados se implementan nuevas tripletas de relación. Esto nos asegura que las nuevas relaciones son más precisas para el conjunto de datos. Para asegurar esta similitud se requiere un análisis de un humano para verificar si las relaciones nuevas son las correctas en las que el valor de similitud se encuentre entre 0.8 hasta 0.89, y si esta es mayor o igual a 0.9 las relaciones son 95% más confiables y no requieren un análisis humano.

Una vez obtenido el resultado del análisis se envía al método save-Data, siendo este el que permita almacenar las nuevas relaciones obtenidas.

Con estos 20 datasets analizados con el rango de 0.5 se obtuvo un gran número de relaciones nuevas, en cambio con el rango de 0.8 las nuevas relaciones fueron menores debido a que el análisis tiene una mejor interpretación y por ende los resultados serán más

exactos en las relaciones de los ODS. Con el valor de similitud que sea mayor los datos serán más precisas y confiables, en cambio con el rango menor se requiere un análisis de una persona para validar cierta información.

El enriquecimiento de datos con *spacy_dbpedia_spotlight* se aplicará el en mismo método que se muestra en la **Figura 27**. En este modelo spacy define el rango de similitud acorde a la expresión de las entidades nombradas y las relaciones lingüísticas que obtiene. La **Figura 32** muestra el proceso de enriquecimiento que esta realizando *Dbpedia-spotlight*.

Figura 32

Captura del proceso de enriquecimiento de la Dbpedia-spotlight

```

webTTEnrichingKnowledgeGraph — Python · Python manage.py runserver — 157x44
[INFO] ...
[INFO] ...
[INFO] Traducción y enriquecimiento
/Users/diepinto30/OneDrive - Universidad Técnica Particular de Loja - UTPL/TESIS-DIEGO/BitBucket/ttenrichingknowledgegraphpython/envTT/lib/python3.7/site-packages/spacy/language.py:639: UserWarning: [W033] Training a new parser or NER using a model with no lexeme normalization table. This may degrade the performance of the model to some degree. If this is intentional or the language you're using doesn't have a normalization table, please ignore this warning. If this is surprising, make sure you have the spacy-lookups-data package installed. The languages with lexeme normalization tables are currently: da, de, el, en, id, lb, pt, ru, sr, ta, th.
**kwargs
http://opendata.org/resource/Dg_near_pre_accession_assistance_to_iceland_catalog <http://www.w3.org/ns/dcat#qualifiedRelation> < http://dbpedia.org/resource/International_Phonetic_Alphabet > .
{"result": "ok"}
http://opendata.org/resource/Stress_test_for_bank_banco_comercial_portugues_sa_bcp_or_millennium_bcp_catalog <http://www.w3.org/ns/dcat#qualifiedRelation> < http://dbpedia.org/resource/Psychological_stress > .
http://opendata.org/resource/Stress_test_for_bank_banco_comercial_portugues_sa_bcp_or_millennium_bcp_catalog <http://www.w3.org/ns/dcat#qualifiedRelation> < http://dbpedia.org/resource/Portuguese_Commercial_Bank > .
{"result": "ok"}
{"result": "ok"}
http://opendata.org/resource/Stress_test_for_bank_banco_comercial_portugues_sa_bcp_or_millennium_bcp_catalog <http://www.w3.org/ns/dcat#qualifiedRelation> < http://dbpedia.org/resource/European_Banking_Authority > .
http://opendata.org/resource/Stress_test_for_bank_banco_comercial_portugues_sa_bcp_or_millennium_bcp_catalog <http://www.w3.org/ns/dcat#qualifiedRelation> < http://dbpedia.org/resource/Eba > .
http://opendata.org/resource/Stress_test_for_bank_banco_comercial_portugues_sa_bcp_or_millennium_bcp_catalog <http://www.w3.org/ns/dcat#qualifiedRelation> < http://dbpedia.org/resource/Psychological_stress > .
http://opendata.org/resource/Stress_test_for_bank_banco_comercial_portugues_sa_bcp_or_millennium_bcp_catalog <http://www.w3.org/ns/dcat#qualifiedRelation> < http://dbpedia.org/resource/Psychological_stress > .
{"result": "ok"}
{"result": "ok"}
{"result": "ok"}
{"result": "ok"}
http://opendata.org/resource/Stress_test_for_bank_banco_comercial_portugues_sa_bcp_or_millennium_bcp_catalog <http://www.w3.org/ns/dcat#qualifiedRelation> < http://dbpedia.org/resource/Systemic_risk > .
{"result": "ok"}
http://opendata.org/resource/Rem_data_bank_year_1997_catalog <http://www.w3.org/ns/dcat#qualifiedRelation> < http://dbpedia.org/resource/Roentgen_equivalent_man > .
{"result": "ok"}
http://opendata.org/resource/Rem_data_bank_year_1997_catalog <http://www.w3.org/ns/dcat#qualifiedRelation> < http://dbpedia.org/resource/Environmental_radioactivity > .
http://opendata.org/resource/Rem_data_bank_year_1997_catalog <http://www.w3.org/ns/dcat#qualifiedRelation> < http://dbpedia.org/resource/Radioactive_decay > .
http://opendata.org/resource/Rem_data_bank_year_1997_catalog <http://www.w3.org/ns/dcat#qualifiedRelation> < http://dbpedia.org/resource/Environmental_monitoring > .
http://opendata.org/resource/Rem_data_bank_year_1997_catalog <http://www.w3.org/ns/dcat#qualifiedRelation> < http://dbpedia.org/resource/Roentgen_equivalent_man > .

```

4.2 Visualización de los resultados obtenidos del enriquecimiento datos

La visualización del grafo enriquecido se puede presentar de varias maneras como el mismo grafo general, pero esto sería muy difícil de entender a simple vista. Se ha pensado en crear un sitio web piloto que permitirá visualizar la información y la interpretación de los datos para el usuario de mejor manera. A continuación, se detalla el uso de las interfaces o vistas de usuario creadas el sitio web piloto y los resultados obtenidos del proceso de implementación de la fase de desarrollo.

Se tiene que tomar en cuenta que para la visualización se realizó consultas SPARQL en tiempo real para recuperar la información directamente desde el triple Store y así esta sea consumida por un api/json, siendo este formato uno de los mas fáciles para usar, mantener y consumir la información que devuelve la consulta.

Dentro de esta API se toman en cuenta los valores que se quiere recibir para la presentación en el JSON. Mediante este formato se hace el consumo con JavaScript (js) para su presentación en el sitio web.

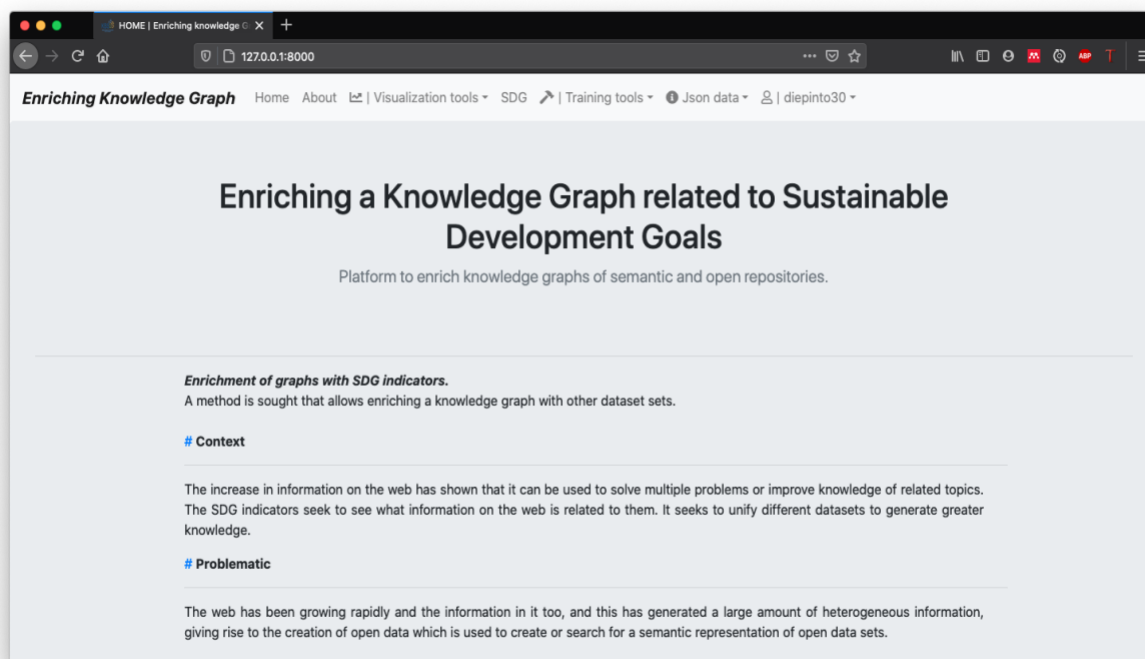
En el **Apéndice 6** se presenta un fragmento de una de las consultas SPARQL que se hizo para obtener los recursos que se esta usando. Y en el **Apéndice 7** se muestra una de las api/json, esta api se la proceso para obtener la información precisa y estructurar el contenido que se muestra en la misma, siendo así usado para el consumo con js.

4.2.1 *Sitio web*

En las siguientes figuras se describen las características y las acciones que cada una de las interfaces que se han creado para el sitio web piloto. En la **Figura 33** muestra el diseño del home o página principal de la aplicación web. En esta se presenta los objetivos de la aplicación, la estructura general de arquitectura en la que se basa para el desarrollo y las opciones del menú de navegación que tiene el sitio.

Figura 33

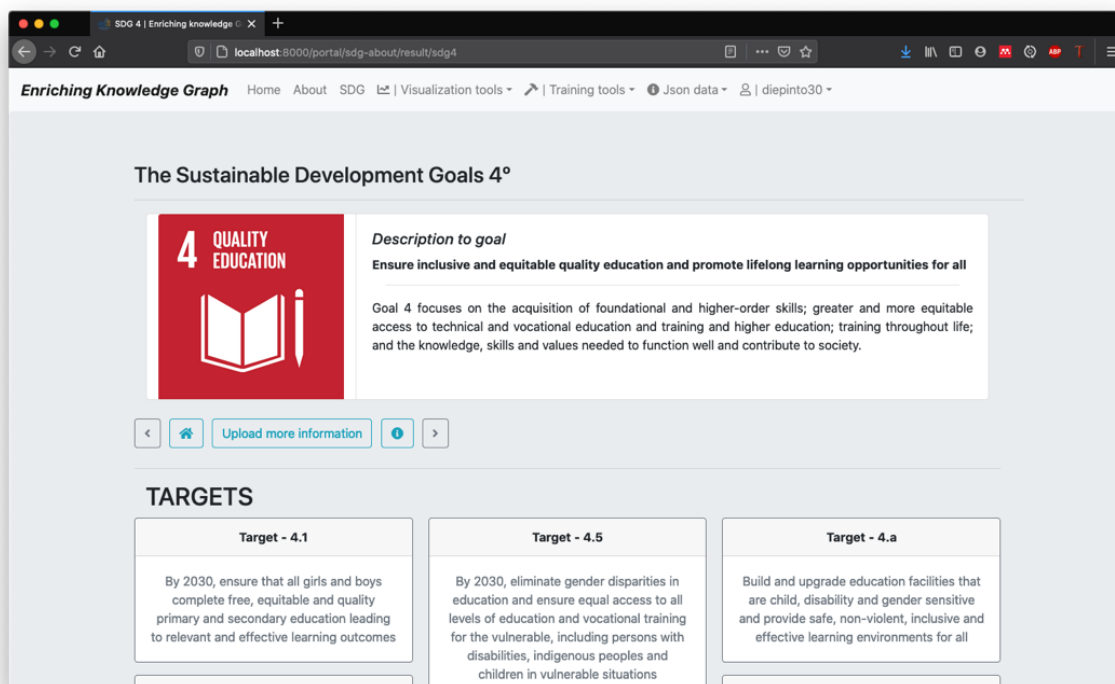
Página principal de la aplicación Web



Dentro del menú de navegación se tiene una opción que dice “SDG”, este punto se presenta cada uno los de los objetivos de desarrollo de sostenible. Al hacer clic en el icono del objetivo nos lleva a una página en la cual nos permite ver las características de sus indicadores y metas que estos objetivos quieren realizar. La **Figura 34** presenta este apartado.

Figura 34

Página de los ODS, en la aplicación Web



En este sitio se presenta una opción de “Upload more information” y presenta gráficas de información que los recursos que estén relacionados a cada objetivo en una Word Cloud y más gráficas que permitan ver los datos relacionados y enriquecidos. La **Figura 35** se muestra los resultados obtenidos y estos se los muestra en una nube de palabras que presenta los conceptos de los ODS, en esta gráfica se busca ver la mayor incidencia de palabras mencionadas en este ODS.

También en la misma sección tiene una opción que presenta una tabla con buscador (**Figura 36**), el cual permite ver todos los recursos que se encuentren relacionados a cada Objetivo de Desarrollo Sostenible y de esta manera poder apreciar más información enriquecida.

Esta tabla permite redirigir a otra página a visualizar toda la información del recurso asociado. La **Figura 36** se muestra un fragmento esta tabla.

Otra visualización que se tiene es con un buscador y diferentes gráficas de resumen de los datos. Si nos dirigimos al menú de navegación existe un apartado de “Visualization tools” y en esta existe dos herramientas de visualización de los datos enriquecidos, a continuación de las describe a cada una de ellas.

4.2.1.1 Search

En este apartado se muestra un buscador (**Figura 37**), el cual permite realizar consultas sobre la información de los recursos en este caso sobre el título, descripción y sobre los datos enriquecidos del recurso, en esta búsqueda se aplicará la traducción en caso de que el texto de la búsqueda sea diferente del idioma base que es el inglés y adicionalmente la corrección del texto si esta mal escrito dando lugar a la corrección de texto de búsqueda. El resultado obtenido de la búsqueda nos llevará a una nueva pestaña denominada browser.

El browser busca presentar todos los recursos que han sido ya enriquecidos; para esto se usa un API del resultado de la consulta en tiempo real hacia el tripleStore, de esta manera el resultado es consumido con Ajax para crear la tabla resultante del contenido.

La **Figura 38** presenta el resultado obtenido de la búsqueda, también cabe mencionar que en la parte derecha se encuentran algunos filtros que se aplican a los resultados existen de la búsqueda, los filtros implementados en el resultado son selección por página de destino, por ODS, por idioma del recurso y un buscador que ayuda filtrar la información resultado según lo que busco.

Figura 37

Buscador de palabras relacionadas

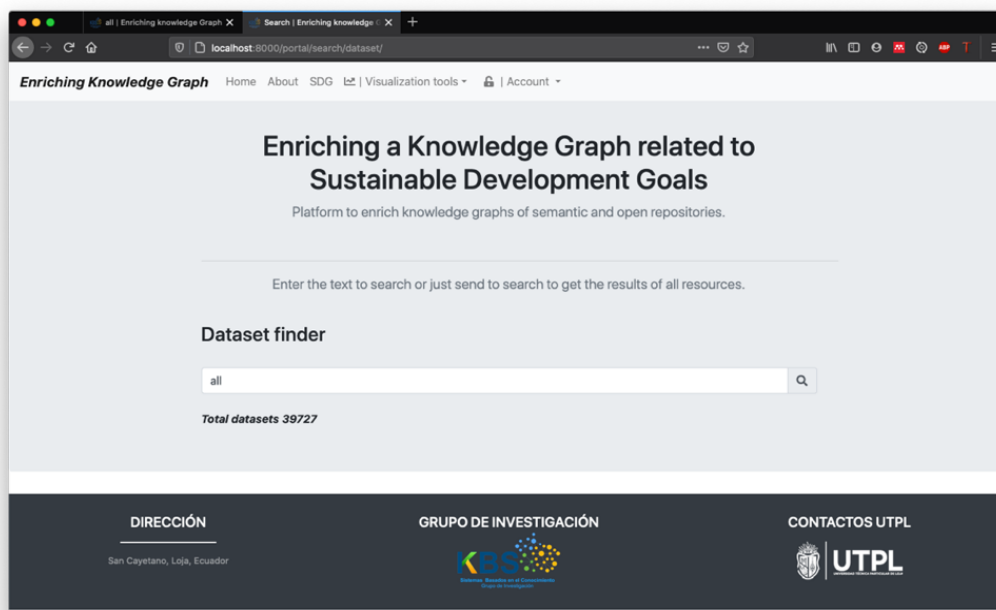
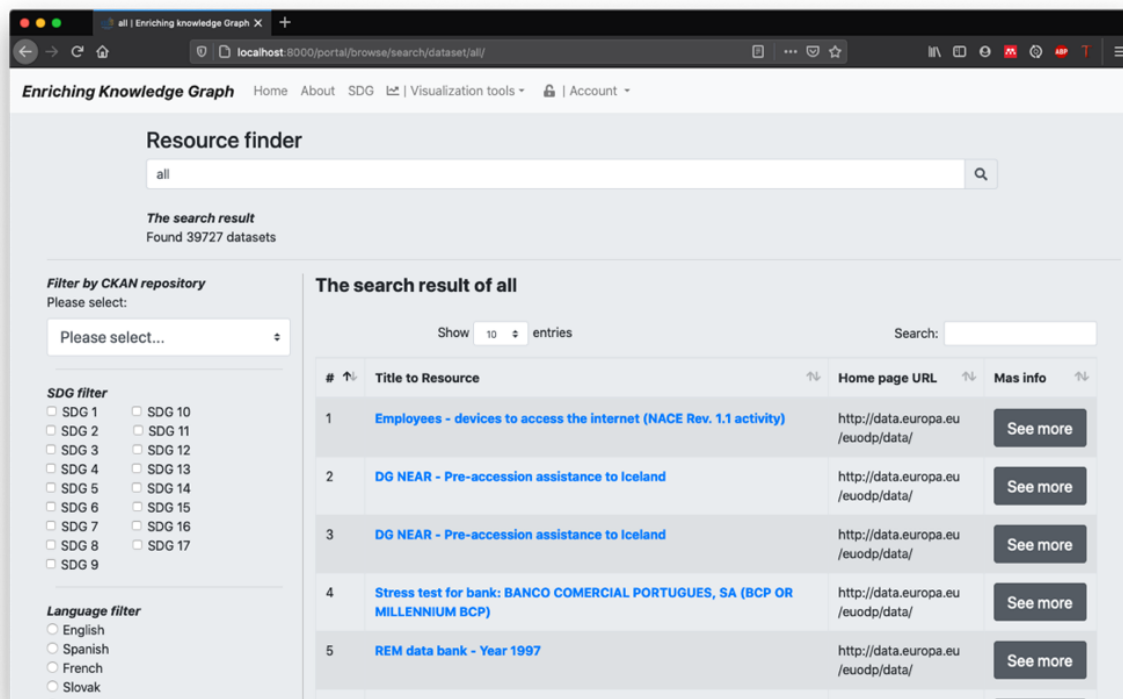


Figura 38

Buscador de recursos enriquecidos

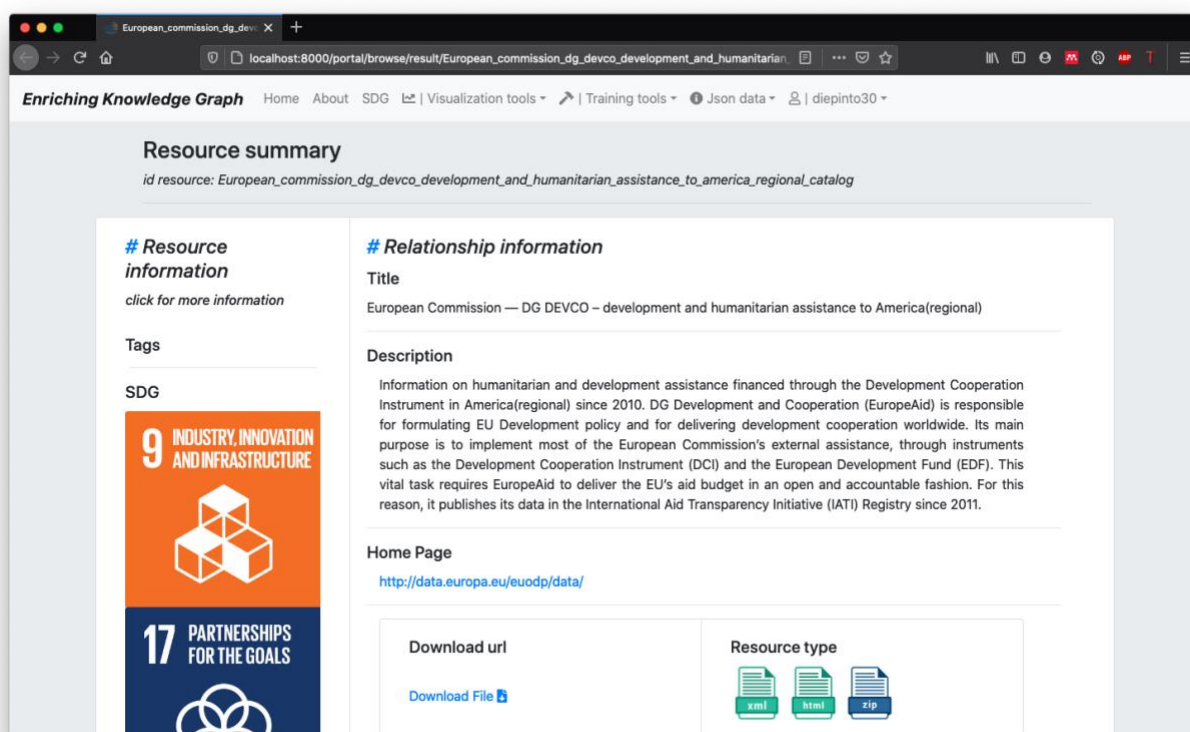


En los apéndices del 8 al 11 se presenta la evolución del buscador desde su primera versión hasta la versión 4 que incluye un buscador y más flitros de búsqueda.

Una vez que se tiene visible todos los datos, en el botón de “see more” o en el mismo título nos lleva a una nueva pestaña en que se cargan los datos del recurso al que se hizo clic, en la **Figura 39** se muestra el contenido resultante.

Figura 39

Resultado de búsqueda en nueva pestaña



4.2.1.2 Graphics Dataset New

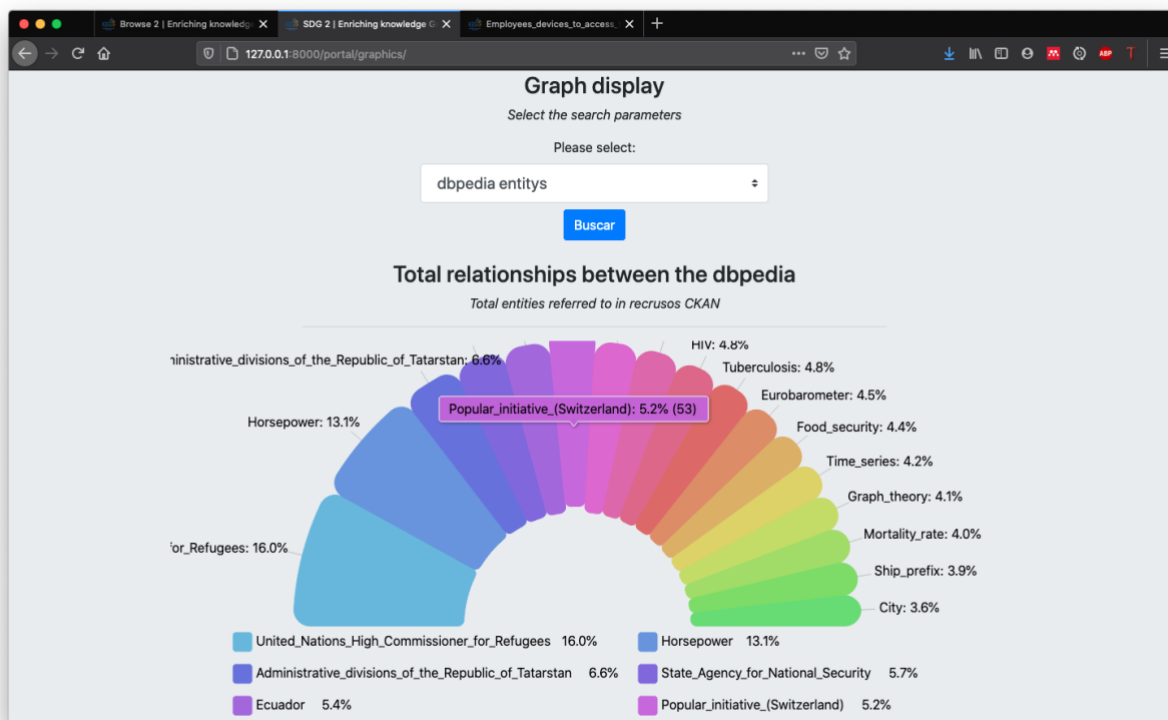
En esta opción se representa con unas gráficas varios resultados generales relacionados con las nuevas características del proceso de enriquecimiento de datos.

La **Figura 40** presenta un conteo total de las 25 entidades más nombras de la DBpedia que se encuentra en cada uno de los recursos. Esta gráfica denominada semicírculo de sectores, busca representar las etiquetas que sean más usadas en los

recursos y determinar de que temáticas están abordando estos recursos que fueron analizados y mapeados.

Figura 40

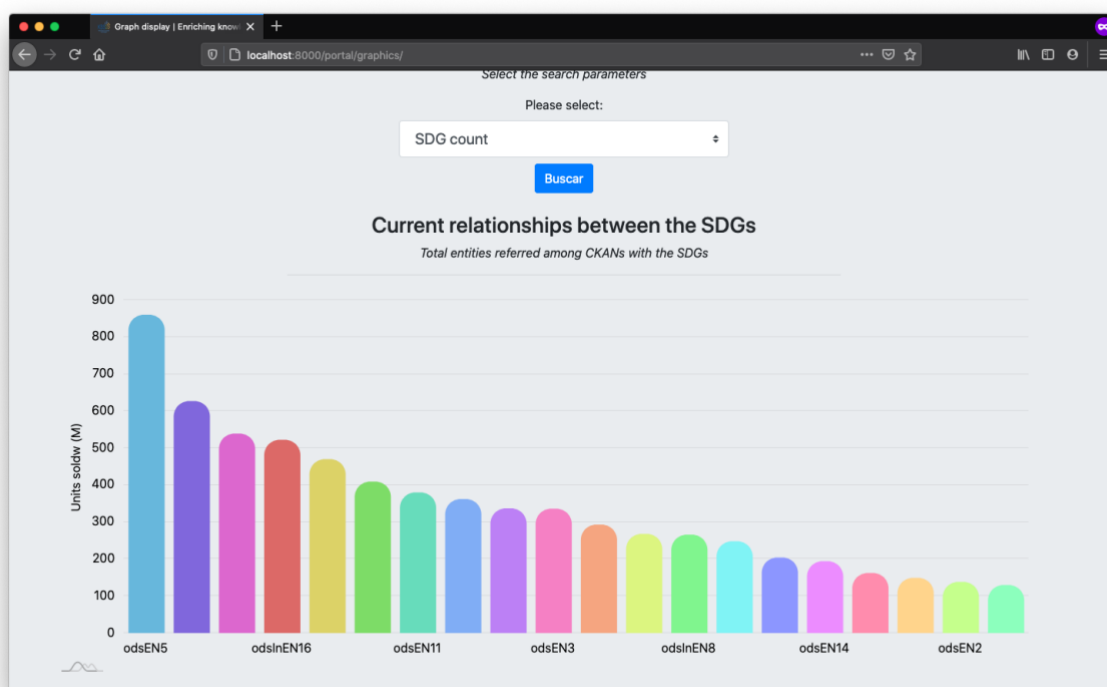
Semi-círculo de sectores de mayores entidades mencionadas de la dbpedia



La **Figura 41** presenta el total de recursos analizados y el total de enlaces que se tiene con cada ODS.

Figura 41

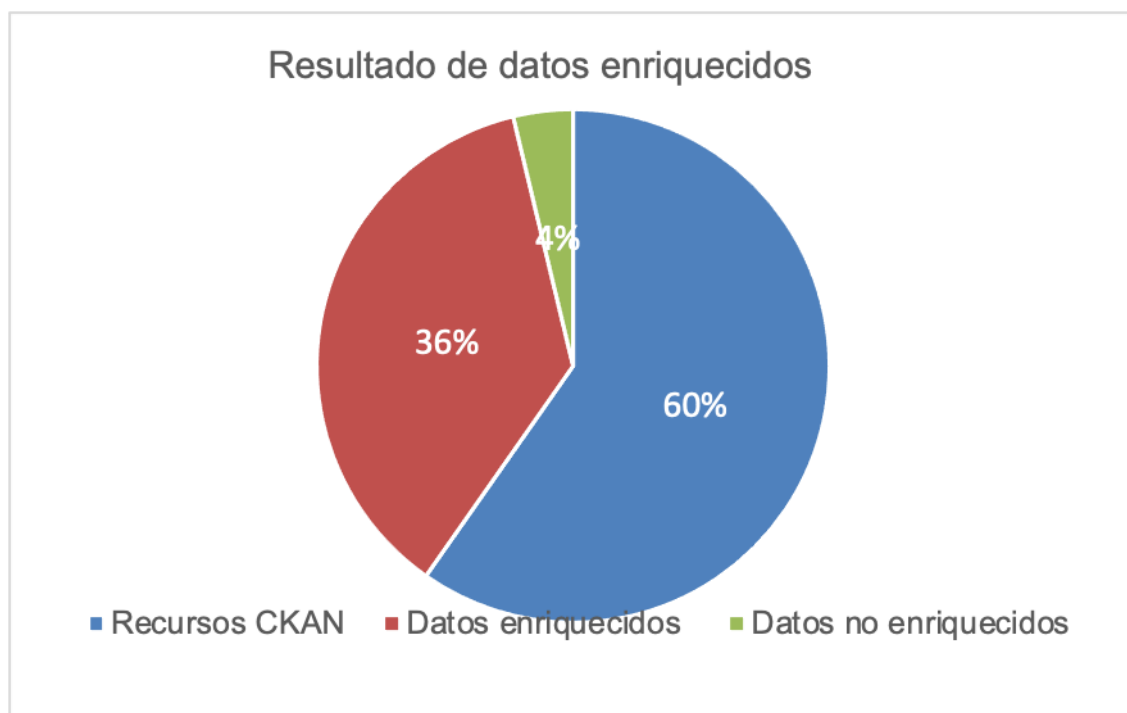
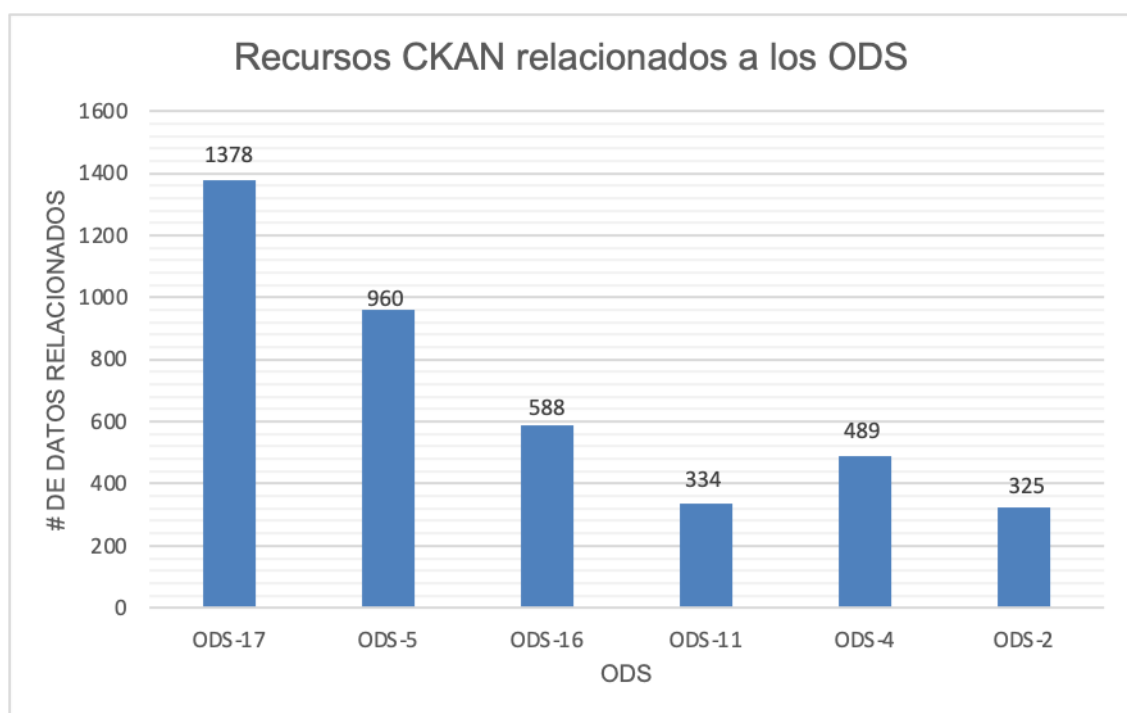
Total de entidades referenciadas entre los recursos CKAN con los ODS



4.2.2 Otros Resultados

Una vez concluido el sitio web piloto y el proceso de enriquecimiento de los datos se puede resumir lo siguiente. De los 39727 recursos CKAN semánticos se ha podido enriquecer más de un 36% de todo el contenido. También un 4% del dataset principal no se ha encontrado alguna coincidencia en los ODS y la Dbpedia. Y el 60% representa el resto de los datos que faltan por enriquecer. Este resumen se lo puede apreciar en la **Figura 42**.

Después de pasar por el proceso de análisis mediante los ODS se puede mencionar que los ODS más usados en el enriquecimiento fueron los siguientes. El ODS 17 que representa a *Alianza para lograr los objetivos* en el cual se tiene más de 1300 recursos asociados al CKAN. Otro más usado es el ODS 5 que representa *la igualdad de género* con más 900 recursos enlazados y otro que es el más mencionado es el ODS 16 que es *Promover sociedades justas, pacíficas e inclusivas* con más de 500 recursos enlazados. En la figura se puede ver más datos de uso de los ODS.

Figura 42*Resultado de datos enriquecidos***Figura 43***Recursos CKAN relacionados a los ODS*

Finalmente, a través de la arquitectura o del método planteado es capaz de acoplarse a diversos datos que requieran ser enriquecidos, tomando en cuenta la estructura que estos posean y de esta manera poder incrementar el conocimiento actual en los datos. El algoritmo es autónomo, esto quiere decir que es capaz de trabajar por sí solo durante un tiempo prolongado.

4.3 Discusión

Después de terminar la aplicación se analiza lo siguiente, la solución propuesta para este trabajo cumple con todos los objetivos del trabajo de titulación. Para esto se logró proponer y crear un método de enriquecimiento de Grafos de Conocimiento mediante la arquitectura propuesta en la **Figura 11**. Este busca unir cualquier recurso abierto a uno y varios indicadores de ODS y también con otras fuentes de datos abiertos en este caso con los recursos de la DBpedia.

En primera instancia se almacenó los datos en GraphDB la primera versión y una vez analizados los datos y procesados se generaron nuevas tripletas siendo estas el resultado del enriquecimiento de los nuevos datos y luego fueron almacenadas en la base de datos semántica GraphDB (segunda versión o datos enriquecidos). Esto permite el consumo en tiempo real de la información que se requiere, pero para esto se necesita generar consultas basadas en SPARQL de esta forma permitiendo visualizar el nuevo contenido añadido a los recursos CKAN.

El proceso de enriquecimiento de datos se lo realizó usando la herramienta de fasttext. La parte fundamental para el proceso es tener una base de conocimiento bien estructurada que permita un análisis lógico y semántico. Y también se usó la herramienta de la Dbpedia que se llama dbpedia-spotlight. Este Framework permitirá dar un valor agregado a la información que es enriquecida con el fin de determinar que palabras o frases tienen incidencia en el texto con la Dbpedia.

Por otra parte, con las herramientas usadas se lograron establecer nuevos enlaces entre diferentes datasets y enlaces con el dataset externo con la dbpedia siendo estos de acceso abierto y determinado su relación con los ODS. Toda esta información enriquecida

puede ser visualizada en el sitio web piloto mediante las diferentes APIs creadas. Este sitio permite que cualquier persona que pueda hacer uso de estos resultados, ayudar al análisis de las mismas y usarla fácilmente.

En estebuye con la capacidad de crear un método que permita en 5 fases (la extracción, NER, NLP, Flujo de trabajo y publicación) extraer información, analizar, crear nuevas relaciones con diferentes conjuntos de datos relacionados con la objetivos, metas e indicadores de los ODS, y publicar estos datos en la web.

Finalmente, este trabajo buscar ser la base para investigaciones futuras, y aportar a nuevas iniciativas para enlazar a múltiples conjunto de datos relacionados con la Agenda 2030 y así mejorar la información que se encuentra en la web. Adicionalmente, los procesos realizados en la aplicación pueden ser adaptados para cualquier tipo de conjunto de datos semánticos y así mejorar el conocimiento actual e ir incrementado la información de forma rápida y segura.

Conclusiones

La motivación principal del trabajo de investigación es crear un método para enriquecer conjuntos de datos abiertos relacionados con la agenda 2030, extraídos desde plataformas CKAN, y que se han descrito como Grafos de Conocimiento Semántico anotados a una base de conocimiento de ODS y a recursos de la Dbpedia.

A partir de los ODS, indicadores y metas se ha creado una base de conocimiento (BC) descrita en RDF, que se usó para detección de similitudes, la recomendación y etiquetado semántico a través de *fasttext*. Esta base de conocimiento se creó a través de la combinación de procesos manuales (extracción) y automatizados (transformación). El idioma utilizado en la BC es el inglés, y contiene 852 entradas relacionadas con objetivos, metas e indicadores.

Se usó Dbpedia-spotlight para tareas de análisis semántico: Reconocimiento de Entidades Nombradas, Procesamiento de Lenguaje Natural y Mapeo Semántico. El análisis de texto ha resultado en la creación de más de 8 mil nuevas relaciones a partir de los datasets descritos como parte de este trabajo.

La creación del proceso el enriquecimiento ha sido evolutivo. La versión final de este proceso implica realizar un proceso de traducción para identificar el idioma del texto de entrada, en caso de que sea diferente del idioma inglés el texto es traducido al idioma base. A continuación, se realiza una limpieza de caracteres especiales. El siguiente paso consiste en aplicar procesamiento de lenguaje natural para reconocer entidades nombradas.

Para la creación de las nuevas relaciones se usó vocabularios semánticos, *dcat:theme* dando lugar a 6680 nuevas relaciones resultantes del uso de *fasttext*; *dcat:qualifiedRelation* para anotar 8054 nuevas relaciones resultantes del uso de *dbpedia-spotlight* y finalmente *dct:language* para la anotación de 7593 nuevas relaciones que determinan el idioma del grafo de conocimiento inicial (Eguiguren Palacios, 2019).

El enriquecimiento y la visualización del conjunto de datos abiertos enriquecidos dio lugar al desarrollo de una aplicación web. Dicha aplicación web permite navegar conjuntos de datos abiertos anotados con los ODS. El sitio es interactivo y accesible para todo tipo de usuario.

Recomendaciones

A continuación, se describen algunas recomendaciones al término del presente trabajo y tomando en cuenta trabajos futuros.

Se recomienda realizar un análisis previo de los datos que se van a procesar con la finalidad de determinar los valores serán los usados para el enriquecimiento del grafo. Es importante tener en cuenta como adaptar sus estrategias de información con los datos con las nuevas relaciones y esas estén adaptadas al vocabulario del grafo principal.

El re-uso de recursos ontológicos es clave en el proceso de enriquecimiento semántico. Sin embargo, para ampliar las posibilidades de enriquecimiento de conjuntos de datos abiertos relacionados con los ODS, sería conveniente crear una ontología o vocabularios semánticos específicos para la anotación de recursos con elementos de información propios de la Agenda 2030.

Además de las tareas NER y NLP reportadas en este trabajo, consideramos importante combinar estos con SPACY con el fin de mejorar la calidad del enriquecimiento de los grafos de conocimiento.

En cuanto a la base de conocimiento usado para fasttext, como trabajo futuro se tiene previsto ampliarla a otros idiomas, esto aumentará la capacidad de anotar semánticamente conjuntos de datos que usan lenguas diferentes al inglés.

En el análisis que se realiza para identificar las entidades nombradas de los conjuntos de datos e interpretar el texto de entrada se recomienda que el índice de similitud sea superior a 0,8.

La propuesta de enriquecimiento de grafos semánticos presentada en este trabajo al usar servicios y recursos de información en la nube y en tiempo real (dbpedia-spotlight, googleTranslate, fasttext) demanda calidad de conexión a Internet. Este es uno de los desafíos que debe ser gestionado si se busca contar con un entorno basado en grafos de conocimiento semántico a gran escala.

Referencias

- Alesso, H. P., & Smith, C. F. (2006). *THINKING ON THE WEB Berners-Lee, Gödel, and Turing*.
- Berners-Lee, T. (2015). *5-star Open Data*. <https://5stardata.info/en/>
- Berners-Lee, T., Connolly, D., Stein, L.A., Swick R. (2000). *The Semantic Web*. <https://www.w3.org/2000/Talks/0906-xmlweb-tbl/text.htm>
- Bizer, C. (2009). The emerging web of linked data. *IEEE Intelligent Systems*, 24(5). <https://doi.org/10.1109/MIS.2009.102>
- Cheaney, A. (2012). *What is Web 2.0?* Quora. <https://www.quora.com/What-is-Web-2-0>
- Danneels, L., Viaene, S., & Van den Bergh, J. (2017). Open data platforms: Discussing alternative knowledge epistemologies. In *Government Information Quarterly* (Vol. 34, Issue 3, pp. 365–378). Elsevier Ltd. <https://doi.org/10.1016/j.giq.2017.08.007>
- Duus, R., & Cooray, M. (2016). *The importance of open data | World Economic Forum*. <https://www.weforum.org/agenda/2016/02/the-importance-of-open-data/>
- Eguiguren Palacios, J. E. (2019). *Marco de trabajo para la integración semántica de portales de datos abiertos basados en CKAN Piloto de Datos Abiertos alineados a los indicadores ODS. (Trabajo de Titulación de Ingeniero en Sistemas Informáticos y Computación)*. UTPL, Loja. <http://dspace.utpl.edu.ec/handle/20.500.11962/24629>
- Escobar, P., Candela, G., Trujillo, J., Marco-Such, M., & Peral, J. (2020). Adding value to Linked Open Data using a multidimensional model approach based on the RDF Data Cube vocabulary. *Computer Standards and Interfaces*, 68(September 2019), 103378. <https://doi.org/10.1016/j.csi.2019.103378>
- Gavin, B. (2018). *What Is An XML File?* <https://www.howtogeek.com/357092/what-is-an-xml-file-and-how-do-i-open-one/>
- Griggs, D., Stafford-Smith, M., Gaffney, O., & Rockström, J. (2013). *Sustainable development goals for people and planet*. 2013. <https://www.nature.com/articles/495305a>

- Hazaël-Massieux, D., & Berners-Lee, T. (2003). *The Semantic Web and its applications at W3C*. <https://www.w3.org/2003/Talks/simo-semwebapp/all.htm>
- Indicadores de ODS*. (2010). <http://www.endvawnow.org/es/articles/336-indicadores.html>
- Jabbar, J. A., & Bulbul, R. (2019). SEMANTIC ENRICHMENT of ROUTING ENGINES USING LINKED DATA: A CASE STUDY USING GRAPHHOPPER. *International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences - ISPRS Archives*, 42(4/W14), 111–117. <https://doi.org/10.5194/isprs-archives-XLII-4-W14-111-2019>
- Kitchenham, B., & Charters, S. (2007). *Guidelines for performing Systematic Literature Reviews in Software Engineering*.
- Lapuate, M. J., & Lamarca, C. (2018). *RDF*.
- Miessler, D. (2020). *Difference Between a URI and a URL*. <https://danielmiessler.com/study/difference-between-uri-url/>
- Okoli, C., & Schabram, K. (2010). *Working Papers on Information Systems A Guide to Conducting a Systematic Literature Review of Information Systems Research A Guide to Conducting a Systematic Literature Review of Information Systems Research*.
- Ontotext. (2020). *GraphDB Documentation*. <http://graphdb.ontotext.com/documentation/standard/>
- ONU. (2019). *UNSD SDGs API*. <https://unstats.un.org/SDGAPI/swagger/#!/Indicator/V1SdgIndicatorByIndicatorCodeSeriesListGet>
- Open Data Handbook. (2019). *La importancia del Open Data*. <https://www.unblogenred.es/la-importancia-del-open-data/>
- Pfenninger, S., DeCarolis, J., Hirth, L., Quoilin, S., & Staffell, I. (2017). The importance of open data and software: Is energy research lagging behind? *Energy Policy*, 101, 211–215. <https://doi.org/10.1016/j.enpol.2016.11.046>
- Quattrini, R., Pierdicca, R., & Morbidoni, C. (2017). Knowledge-based data enrichment for HBIM: Exploring high-quality models using the semantic-web. *Journal of Cultural*

Heritage, 28, 129–139. <https://doi.org/10.1016/j.culher.2017.05.004>

Sánchez, A., Piedra, N., & Morocho, J. C. (2018). Using linked data to ensure that digital information about historical figures of Loja remains accessible and usable. *CEUR Workshop Proceedings, 2096*, 21–34.

Sharma, V. (2015). *What are WEB 2.0 & WEB 3.0?* <https://www.quora.com/What-are-WEB-2-0-WEB-3-0>

Song, H. J., & Park, S. B. (2018). Enriching translation-based knowledge graph embeddings through continual learning. *IEEE Access*, 6, 60489–60497. <https://doi.org/10.1109/ACCESS.2018.2874656>

Stauffacher, D., Hattotuwa, S., & Weekes, B. (2012). *The potential and challenges of open data for crisis information management and aid efficiency A preliminary assessment*. <http://wbi.worldbank.org/wbi/open-aid-partnership>

Techopedia. (n.d.). *What is the Open Data Platform (ODP)*. Retrieved January 8, 2020, from <https://www.techopedia.com/definition/31438/open-data-platform-odp>

W3C. (2015). *Ontologies*. <https://www.w3.org/standards/semanticweb/ontology>

Wonderflow. (2018). *12 NLP Examples: How Natural Language Processing is Used*. Wonderflow. <https://www.wonderflow.co/blog/natural-language-processing-examples>

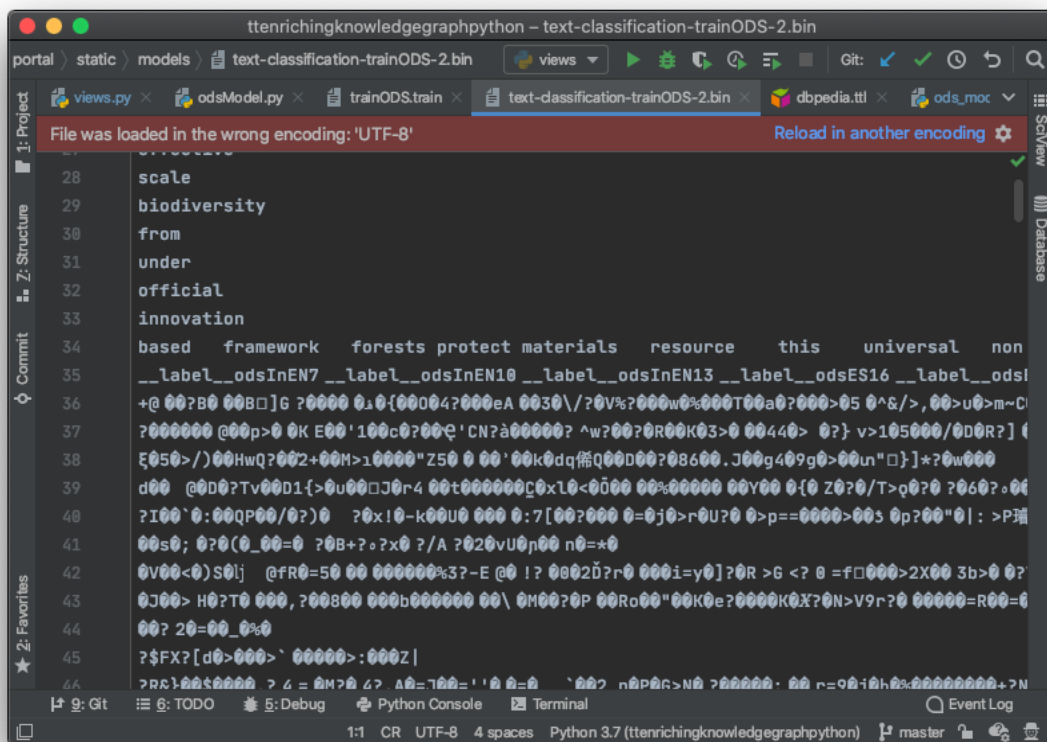
Apéndice

Se incluye de acuerdo al orden citado en el cuerpo del Trabajo de Titulación.

Apéndice 1:

Figura 44

Modelo entrenado con Fasttext.



The image shows a screenshot of a code editor window titled "ttenrichingknowledgegraphpython - text-classification-trainODS-2.bin". The editor displays a Python file with a red error message at the top: "File was loaded in the wrong encoding: 'UTF-8'". The code visible includes:

```

28 scale
29 biodiversity
30 from
31 under
32 official
33 innovation
34 based framework forests protect materials resource this universal non
35 __label__odsInEN7 __label__odsInEN10 __label__odsInEN13 __label__odsES16 __label__odsI
36 +( @ 00?B0 00B[]G ?0000 0_0{00004?000eA 0030\/?0V%?000w0%000T00a0?000>05 0^&/>,00>u0>m~C
37 ?000000 @00p>0 0K E00'100c0?00c'CN?à00000? ^w?00?0R00K03>0 00440> 0?} v>105000/0D0R?] (
38 ξ050>/)00HwQ?002+00M>10000"Z50 0 00'00k0dq000D00?08600.J00g409g0>00un"}]*?0w000
39 d00 @0D0?Tv00D1{>0u00□J00r4 00t000000□0x10<0000 00%00000 00Y00 0{0 Z0?0/T>q0?0 ?060?000
40 ?I00`0:00QP00/0?)0 ?0x!0-k00U0 000 0:7[00?000 0=0j0>r0U?0 0>p==0000>003 0p?00`0| : >P環
41 00s0; 0?0(0_00=0 ?0B+?0?x0 ?/A ?020vU0n00 n0=*0
42 0V00<0)S0lj @fR0=50 00 000000%3?-E @0 !? 0002D?r0 000i=y0]?0R >G <? 0 =f□000>2X00 3b>0 0?
43 0J00> H0?T0 000,?00800 000b000000 00\ 0M00?0P 00R000"00K0e?0000K0X?0N>V9r?0 00000=R00=0
44 00? 20=00_0%0
45 ?$FX?[d0>000>` 00000>:000Z|
46 }R& }0000000 ? 4 = 0M?0 4? A0=7000=' '0 0=0 `00? n0P0R>M0 ?0000000· 00 r=900i0h0:000000000+?M

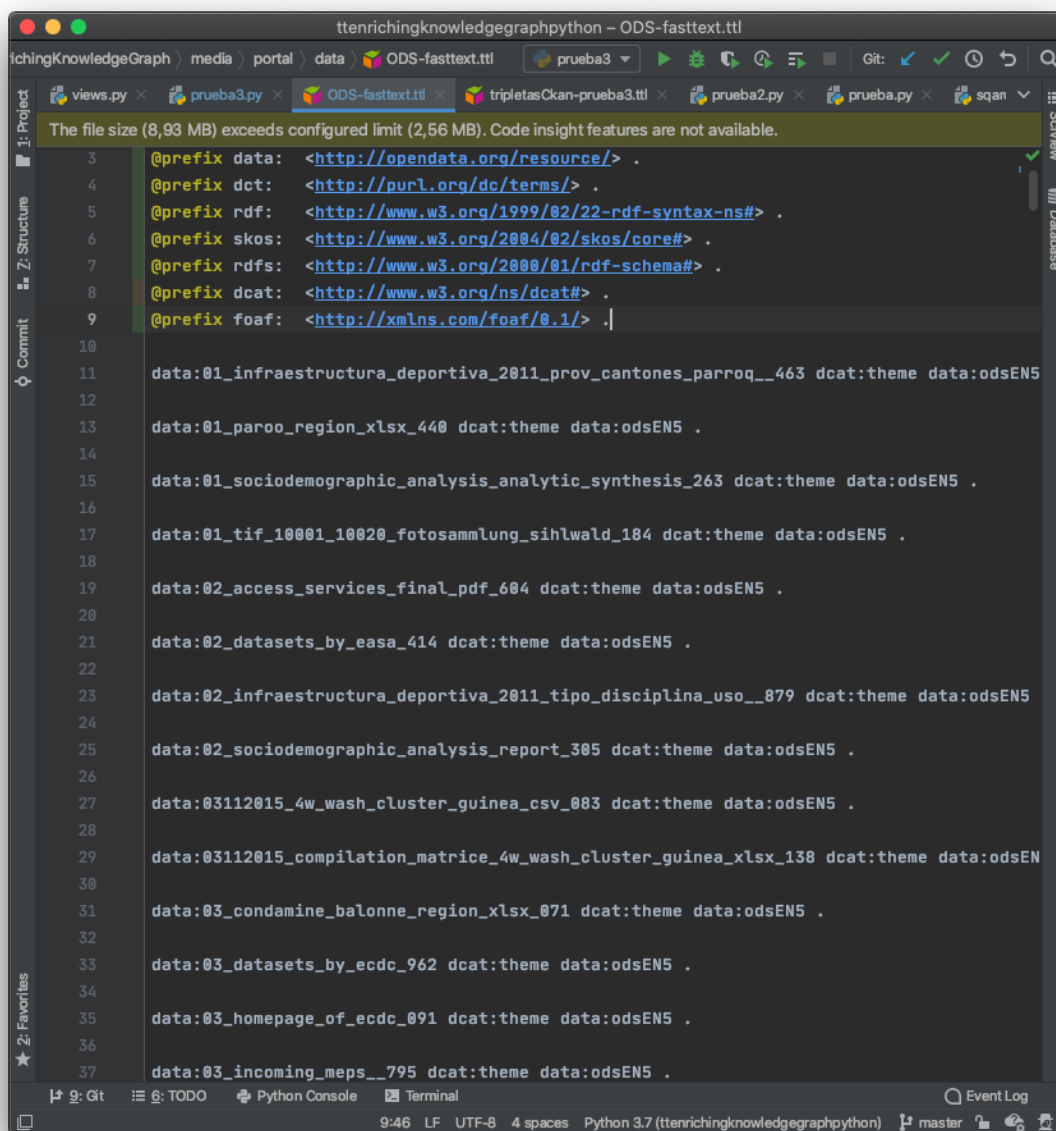
```

The editor interface includes a sidebar with "Project" and "Structure" views, a "Commit" button, and a "Favorites" section. The bottom status bar shows "1:1 CR UTF-8 4 spaces Python 3.7 (ttenrichingknowledgegraphpython) master".

Apéndice 2:

Figura 45

Fragmento de las nuevas relaciones creadas con fasttext



```
ttenrichingknowledgegraphpython - ODS-fasttext.ttl
odingKnowledgeGraph > media > portal > data > ODS-fasttext.ttl
views.py x prueba3.py x ODS-fasttext.ttl x tripletasCkan-prueba3.ttl x prueba2.py x prueba.py x sqan
The file size (8,93 MB) exceeds configured limit (2,56 MB). Code insight features are not available.
3 @prefix data: <http://opendata.org/resource/> .
4 @prefix dct: <http://purl.org/dc/terms/> .
5 @prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#> .
6 @prefix skos: <http://www.w3.org/2004/02/skos/core#> .
7 @prefix rdfs: <http://www.w3.org/2000/01/rdf-schema#> .
8 @prefix dcat: <http://www.w3.org/ns/dcat#> .
9 @prefix foaf: <http://xmlns.com/foaf/0.1/> .
10
11 data:01_infraestructura_deportiva_2011_prov_cantones_parroq__463 dcat:theme data:odsEN5
12
13 data:01_paroo_region_xlsx_440 dcat:theme data:odsEN5 .
14
15 data:01_sociodemographic_analysis_analytic_synthesis_263 dcat:theme data:odsEN5 .
16
17 data:01_tif_10001_10020_fotosammlung_sihlwald_184 dcat:theme data:odsEN5 .
18
19 data:02_access_services_final_pdf_604 dcat:theme data:odsEN5 .
20
21 data:02_datasets_by_easa_414 dcat:theme data:odsEN5 .
22
23 data:02_infraestructura_deportiva_2011_tipo_disciplina_uso__879 dcat:theme data:odsEN5
24
25 data:02_sociodemographic_analysis_report_305 dcat:theme data:odsEN5 .
26
27 data:03112015_4w_wash_cluster_guinea_csv_083 dcat:theme data:odsEN5 .
28
29 data:03112015_compilation_matrice_4w_wash_cluster_guinea_xlsx_138 dcat:theme data:odsEN
30
31 data:03_condamine_balonne_region_xlsx_071 dcat:theme data:odsEN5 .
32
33 data:03_datasets_by_ecdc_962 dcat:theme data:odsEN5 .
34
35 data:03_homepage_of_ecdc_091 dcat:theme data:odsEN5 .
36
37 data:03_incoming_meps__795 dcat:theme data:odsEN5 .
9:46 LF UTF-8 4 spaces Python 3.7 (ttenrichingknowledgegraphpython) master
```

En esta imagen se aprecia en un formato turtle (ttl) en resultado de las nuevas relaciones entre los recursos y los ODS.

Apéndice 3:

Figura 46

Fragmento de las nuevas relaciones creadas con spacy.

```

ttenrichingknowledgegraphpython - dbpedia.ttl
TEnrichingKnowledgeGraph > media > portal > data > dbpedia.ttl
views.py x dbpedia.ttl x tripletasCkan-prueba3.ttl x prueba2.py x prueba.py x sqarj_lpraphdb.py x gr:
5 | }prefix dct: <http://purl.org/dc/terms/> .
6 | }prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#> .
7 | }prefix skos: <http://www.w3.org/2004/02/skos/core#> .
8 | }prefix rdfs: <http://www.w3.org/2000/01/rdf-schema#> .
9 | }prefix dcat: <http://www.w3.org/ns/dcat#> .
10 | }prefix foaf: <http://xmlns.com/foaf/0.1/> .
11 |
12 | |lata:02_full_list_of_meps_elected_in_2014_583 dcat:qualifiedRelation dbr:Member_of_the_Eu
13 |
14 | |lata:03112015_4w_wash_cluster_guinea_csv_083 dcat:qualifiedRelation <[dbr:]Cooper_(profes
15 | |dbr:Eromanga_Basin .
16 |
17 | |lata:03112015_compilation_matrice_4w_wash_cluster_guinea_xlsx_138 dcat:qualifiedRelation
18 |
19 | |lata:03_condamine_balonne_region_xlsx_071 dcat:qualifiedRelation dbr:Eurobarometer,
20 | |dbr:Tobacco .
21 |
22 | |lata:03_homepage_of_ecdc_091 dcat:qualifiedRelation dbr:European_Centre_for_Disease_Preve
23 |
24 | |lata:04_datasets_by_echa_022 dcat:qualifiedRelation dbr:European_Chemicals_Agency .
25 |
26 | |lata:05_homepage_of_eea_615 dcat:qualifiedRelation dbr:Outerwear .
27 |
28 | |lata:05_o_meps_with_alphabet_letter_o_846 dcat:qualifiedRelation dbr:Currency,
29 | |dbr:Current_account .
30 |
31 | |lata:05_rehabilitation_and_return_to_work_analysis_report_on_eu_and_member_states_policie
32 | |<http://dbpedia.org/resource/Swimming_(sport)> .
33 |
34 | |lata:05_s_meps_with_alphabet_letter_s_622 dcat:qualifiedRelation dbr:Alphabet,
35 | |dbr:Member_of_the_European_Parliament .
36 |
37 | |lata:05_z_meps_with_alphabet_letter_z_927 dcat:qualifiedRelation dbr:Alphabet,
38 | |dbr:Member_of_the_European_Parliament .
39 |
40 | |lata:12_20150115_dtm_dataset_round_12_xlsx_314 dcat:qualifiedRelation dbr:Asteroid_family)

```

Reference should use prefix: dbr. Unresolved refs. 10:46 LF UTF-8 4 spaces Python 3.7 (ttenrichingknowledgegraphpython) master

En esta imagen se aprecia en un formato turtle (ttl) en resultado de las nuevas relaciones entre los recursos y los recursos con la dbpedia.

Apéndice 4:

Figura 47

Fragmento de los datos antes de ser enriquecidos.

```

@prefix data: <http://opendata.org/resource/> .
@prefix dct: <http://purl.org/dc/terms/> .
@prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#> .
@prefix skos: <http://www.w3.org/2004/02/skos/core#> .
@prefix rdfs: <http://www.w3.org/2000/01/rdf-schema#> .
@prefix dcat: <http://www.w3.org/ns/dcat#> .
@prefix foaf: <http://xmlns.com/foaf/0.1/> .

data:
  Groups_from_volksinitiative_abtreibungsfinanzierung_ist_privatsache_entlastung_der_krankenversicherung_durch_streichung_der_kosten_des_schwangerschaftsabbruchs_aus_der_obligatorischen_grundversicherung_nach_bezirken
    a
      skos:ConceptScheme ;
    rdfs:label "Groups from Volksinitiative «Abtreibungsfinanzierung ist Privatsache - Entlastung der Krankenversicherung durch Streichung der Kosten des Schwangerschaftsabbruchs aus der obligatorischen Grundversicherung», nach Bezirken" ;
    dct:title "Groups from Volksinitiative «Abtreibungsfinanzierung ist Privatsache - Entlastung der Krankenversicherung durch Streichung der Kosten des Schwangerschaftsabbruchs aus der obligatorischen Grundversicherung», nach Bezirken" .

data:Worldpop_thailand_age_and_sex_structures_summary_page_catalog
  a
    dcat:Catalog ;
  rdfs:label "WorldPop Thailand Age and sex structures Summary Page" ;
  dct:identifier "65885" ;
  dct:title "WorldPop Thailand Age and sex structures Summary Page" ;
  dcat:themeTaxonomy data:Groups_from_worldpop_thailand_age_and_sex_structures_summary_page ,
data:Tags_from_worldpop_thailand_age_and_sex_structures_summary_page ;
  foaf:homepage <https://data.humdata.org/> .

data:Tags_from_wfs_ausnahmetransportrouten
  a
    skos:ConceptScheme ;
  rdfs:label "Tags from WFS Ausnahmetransportrouten" ;
  dct:title "Tags from WFS Ausnahmetransportrouten" .

data:dist_Naturschutzobjekt_gpkg_011
  a
    dcat:Distribution ;
  rdfs:label "Distribution of: naturschutzobjekt.gpkg" ;

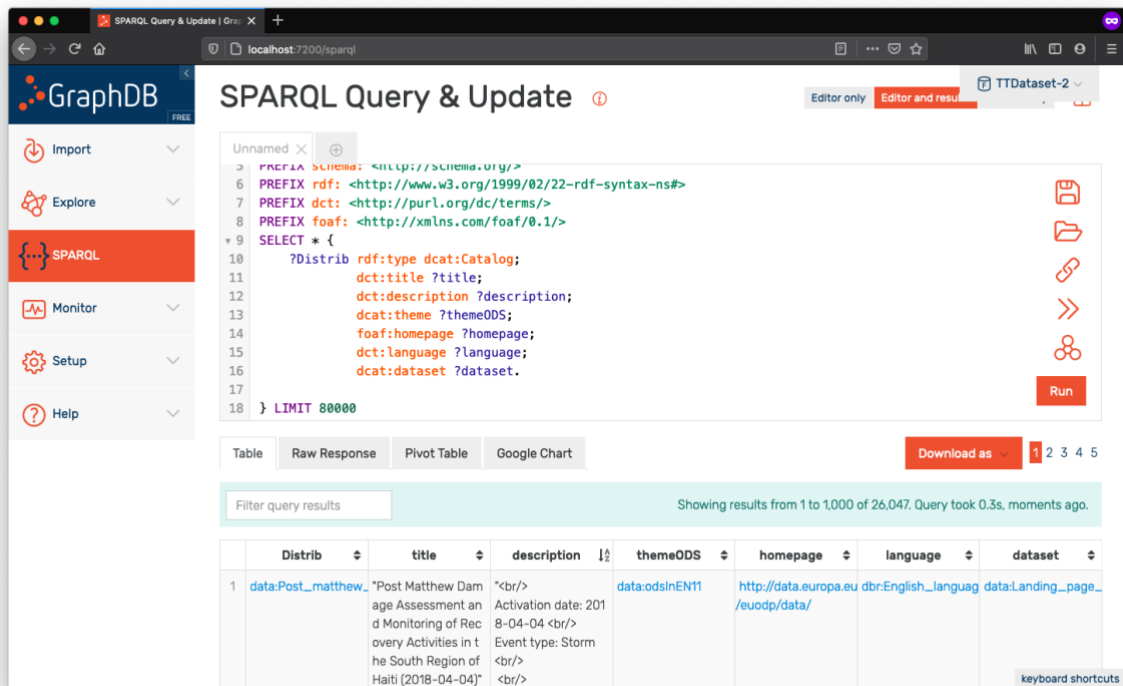
```

En esta imagen se aprecia en un formato turtle (ttl) los recursos CKAN base.

Apéndice 6:

Figura 49

Consulta SPARQL para obtener los resultados de análisis.



The screenshot shows the GraphDB SPARQL Query & Update interface. The query is as follows:

```

PREFIX schema: <http://schema.org/>
PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
PREFIX dct: <http://purl.org/dc/terms/>
PREFIX foaf: <http://xmlns.com/foaf/0.1/>
SELECT * {
  ?Distrib rdf:type dcat:Catalog;
           dct:title ?title;
           dct:description ?description;
           dcat:theme ?themeODS;
           foaf:homepage ?homepage;
           dct:language ?language;
           dcat:dataset ?dataset.
} LIMIT 80000
  
```

The results are displayed in a table format, showing the first result:

	Distrib	title	description	themeODS	homepage	language	dataset
1	data:Post_matthew_	"Post Matthew Damage Assessment and Monitoring of Recovery Activities in the South Region of Haiti (2018-04-04)"	 Activation date: 2018-04-04 Event type: Storm 	data:odsinEN11	http://data.europa.eu/dbr:English_language/euodp/data/		data:Landing_page_

The interface also shows a sidebar with navigation options (Import, Explore, SPARQL, Monitor, Setup, Help) and a top navigation bar with 'Editor only' and 'Editor and results' tabs. The results are shown from 1 to 1,000 of 26,047 total results.

En esa consulta se busca obtener los resultados que nos permita saber sobre los recursos que están enriquecidos y saber más de esos recursos analizados.

Apéndice 7:

Figura 50

Api/Json de los resultados de la consulta Sparql en tiempo real.

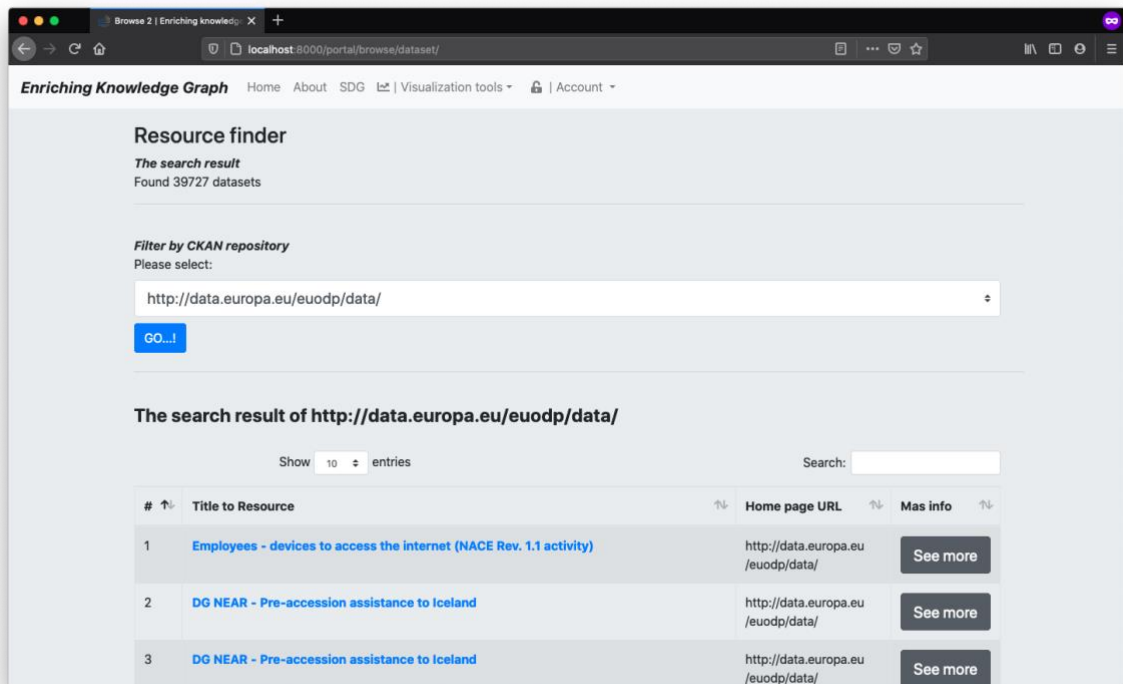
```
localhost:8000/portal/api/apiBrowseNew/
JSON
Datos sin procesar
Cabeceras
Guardar Copiar Contraer todo Expandir todo (lento) Filtrar JSON
9:
  Distrib: "European_commission_dg_devco_development_and_humanitarian_assistance_to_america_regional_catalog"
  title: "European Commission - DG DEVCO - development and humanitarian assistance to America(regional)"
  themeODS: "http://opendata.org/resource/odsInEN9"
  homepage: "http://data.europa.eu/euodp/data/"
  language: "http://dbpedia.org/resource/English_language"
10: {...}
11: {...}
12: {...}
13:
  Distrib: "Enterprises_with_broadband_access_catalog"
  title: "Enterprises with broadband access"
  themeODS: "http://opendata.org/resource/odsEN9"
  homepage: "http://data.europa.eu/euodp/data/"
  language: "http://dbpedia.org/resource/English_language"
14: {...}
15: {...}
16: {...}
17: {...}
18: {...}
19:
  Distrib: "Preliminary_results_on_employer_enterprise_deaths_presented_by_legal_form_until_2007_nace_rev_1_1_catalog"
  title: "Preliminary results on employer enterprise deaths presented by legal form (until 2007, NACE Rev. 1.1)"
  themeODS: "http://opendata.org/resource/odsEN3"
  homepage: "http://data.europa.eu/euodp/data/"
  language: "http://dbpedia.org/resource/English_language"
20: {...}
21: {...}
22: {...}
23: {...}
24: {...}
25: {...}
26: {...}
```

Este json busca estructurar y analizar los resultados de manera mas fácil para el consumo del mismo en el sitio.

Apéndice 8:

Figura 51

Buscador en la primera versión



Resource finder
The search result
Found 39727 datasets

Filter by CKAN repository
Please select:

GO...

The search result of http://data.europa.eu/euodp/data/

Show 10 entries Search:

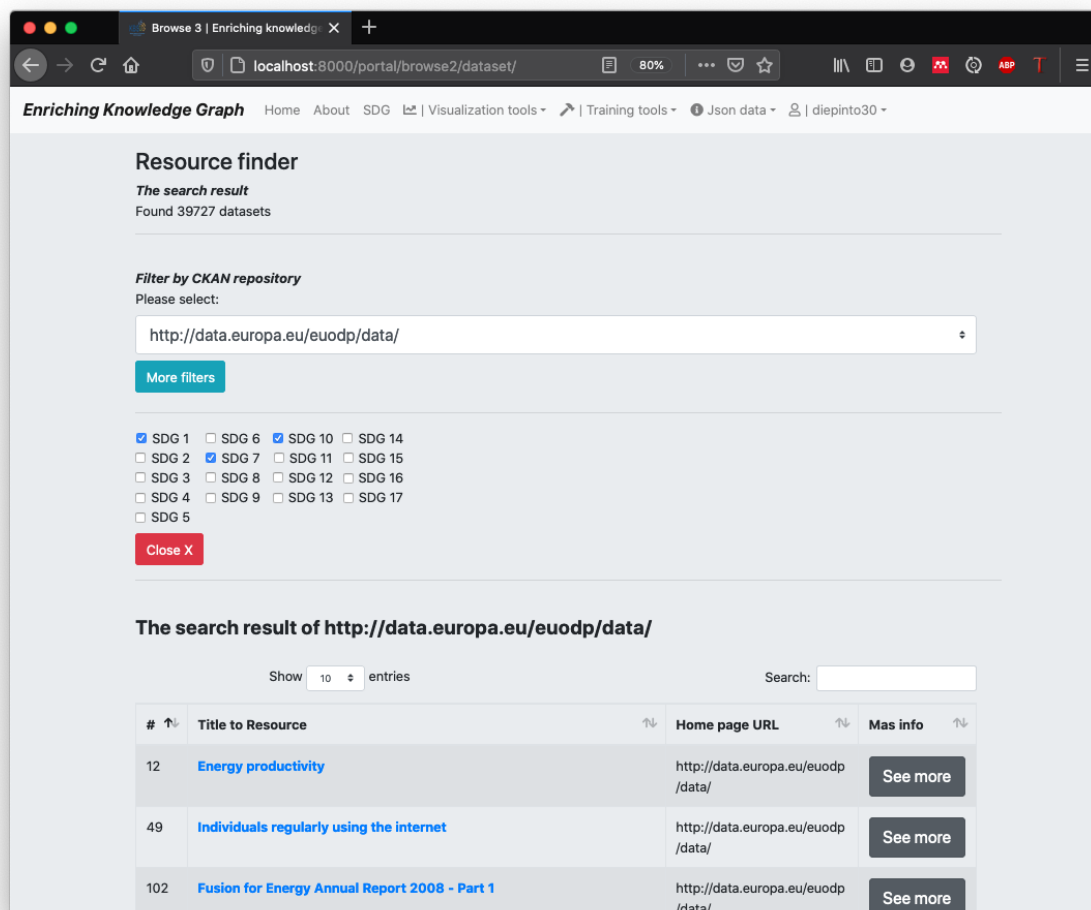
#	Title to Resource	Home page URL	Mas info
1	Employees - devices to access the internet (NACE Rev. 1.1 activity)	http://data.europa.eu/euodp/data/	See more
2	DG NEAR - Pre-accession assistance to Iceland	http://data.europa.eu/euodp/data/	See more
3	DG NEAR - Pre-accession assistance to Iceland	http://data.europa.eu/euodp/data/	See more

Este buscador tiene filtro por los homePage de los recursos y un buscador por entrada de texto.

Apéndice 9:

Figura 52

Buscador en la segunda versión



Resource finder

The search result
Found 39727 datasets

Filter by CKAN repository
Please select:

[More filters](#)

SDG 1 SDG 6 SDG 10 SDG 14
 SDG 2 SDG 7 SDG 11 SDG 15
 SDG 3 SDG 8 SDG 12 SDG 16
 SDG 4 SDG 9 SDG 13 SDG 17
 SDG 5

[Close X](#)

The search result of <http://data.europa.eu/euodp/data/>

Show entries Search:

#	Title to Resource	Home page URL	Mas info
12	Energy productivity	http://data.europa.eu/euodp/data/	See more
49	Individuals regularly using the internet	http://data.europa.eu/euodp/data/	See more
102	Fusion for Energy Annual Report 2008 - Part 1	http://data.europa.eu/euodp/data/	See more

Este buscador tiene filtro por los homePage de los recursos, un checkbox para filtrar por los ODS relacionados con uno o varios recursos y un buscador por entrada de texto.

Apéndice 10:

Figura 53

Buscador en la tercera versión

The screenshot displays the 'Resource finder' interface of the 'Enriching Knowledge Graph' application. The page title is 'Resource finder' and it indicates 'The search result' with 'Found 39727 datasets'. On the left side, there are three filter sections: 'Filter by CKAN repository' with a dropdown menu, 'SDG filter' with checkboxes for SDG 1 through SDG 17, and 'Language filter' with radio buttons for English, Spanish, French, Slovak, Hungarian, and All. The main content area, titled 'The search result of all', shows a table of search results. The table has columns for '#', 'Title to Resource', 'Home page URL', and 'Mas info'. The search results are as follows:

#	Title to Resource	Home page URL	Mas info
1	Employees - devices to access the internet (NACE Rev. 1.1 activity)	http://data.europa.eu/euodp/data/	See more
2	DG NEAR - Pre-accession assistance to Iceland	http://data.europa.eu/euodp/data/	See more
3	DG NEAR - Pre-accession assistance to Iceland	http://data.europa.eu/euodp/data/	See more
4	Stress test for bank: BANCO COMERCIAL PORTUGUES, SA (BCP OR MILLENNIUM BCP)	http://data.europa.eu/euodp/data/	See more
5	REM data bank - Year 1997	http://data.europa.eu/euodp/data/	See more
6	Electricity grid mix; AC; consumption mix, at consumer; 230V (Location: SE)	http://data.europa.eu/euodp/data/	See more

Este buscador tiene filtro por los homePage de los recursos, un checkbox para filtrar por los ODS relacionados con uno o varios recursos, un radiobox para filtrar por idiomas y un buscador por entrada de texto.

Apéndice 11:

Figura 54

Buscador en la cuarta versión

The screenshot shows a web browser window with the URL `localhost:8000/portal/browse/search/dataset/all/`. The page title is "Enriching Knowledge Graph" and the navigation menu includes "Home", "About", "SDG", "Visualization tools", and "Account".

The main content area is titled "Resource finder" and contains a search bar with the text "all". Below the search bar, it says "The search result" and "Found 39727 datasets".

On the left side, there are two filter sections:

- Filter by CKAN repository:** A dropdown menu with the text "Please select..." and a downward arrow.
- SDG filter:** A list of checkboxes for SDG 1 through SDG 17.
- Language filter:** A list of radio buttons for English, Spanish, French, Slovak, and Hungarian.

The main content area is titled "The search result of all" and contains a table with the following columns: "#", "Title to Resource", "Home page URL", and "Mas info". The table shows 5 results:

#	Title to Resource	Home page URL	Mas info
1	Employees - devices to access the internet (NACE Rev. 1.1 activity)	http://data.europa.eu/euodp/data/	See more
2	DG NEAR - Pre-accession assistance to Iceland	http://data.europa.eu/euodp/data/	See more
3	DG NEAR - Pre-accession assistance to Iceland	http://data.europa.eu/euodp/data/	See more
4	Stress test for bank: BANCO COMERCIAL PORTUGUES, SA (BCP OR MILLENNIUM BCP)	http://data.europa.eu/euodp/data/	See more
5	REM data bank - Year 1997	http://data.europa.eu/euodp/data/	See more

Este buscador tiene un buscador de coincidencia de palabras y filtros mejorados por el homePage de los recursos, un checkbox para filtrar por los ODS relacionados con uno o varios recursos, también otro con checkbox para filtrar por idiomas y un buscador por entrada de texto del resultado obtenido.