



UNIVERSIDAD TÉCNICA PARTICULAR DE LOJA

La Universidad Católica de Loja

ÁREA TÉCNICA

**TITULACIONES DE INGENIERÍA EN CIENCIAS DE LA
COMPUTACIÓN**

**Desarrollo de Servicios Web para el proceso de Enlace y Enriquecimiento
de Datos Enlazados**

TRABAJO DE FIN DE TITULACIÓN

AUTOR: Montaña Sozoranga, Wilmer Fabricio

DIRECTOR: Piedra Pullaguari, Nelson Oswaldo, Ing.

LOJA - ECUADOR

2014

APROBACIÓN DEL DIRECTOR DEL TRABAJO DE FIN DE TITULACIÓN

Ingeniero.

Nelson Oswaldo Piedra Pullaguari.

DOCENTE DE LA TITULACIÓN

De mi consideración:

El presente trabajo de fin de titulación: Desarrollo de Servicios Web para el proceso de Enlace y Enriquecimiento de Datos Enlazados. Piloto: dominio de datos Universitarios, realizado por Montaña Sozoranga Wilmer Fabricio , ha sido orientado y revisado durante su ejecución, por se aprueba la presentación del mismo.

Loja, abril de 2015

f)

DECLARACIÓN DE AUTORÍA Y CESIÓN DE DERECHOS

“Yo Montaña Sozoranga Wilmer Fabricio declaro ser autor del presente trabajo de fin de titulación: Desarrollo de Servicios Web para el proceso de Enlace y Enriquecimiento de Datos Enlazados. Piloto: dominio de datos Universitarios, de la Titulación de Ingeniería en Sistemas Informáticos y Computación, siendo el Ing. Nelson Oswaldo Piedra Pullaguari director del presente trabajo; y eximo expresamente a la Universidad Técnica Particular de Loja y a sus representantes legales de posibles reclamos o acciones legales. Además certifico que las ideas, conceptos, procedimientos y resultados vertidos en el presente trabajo investigativo, son de mi exclusiva responsabilidad.

Adicionalmente declaro conocer y aceptar la disposición del Art. 88 del Estatuto Orgánico de la Universidad Técnica Particular de Loja que en su parte pertinente textualmente dice: “Forman parte del patrimonio de la Universidad la propiedad intelectual de investigaciones, trabajos científicos o técnicos y tesis de grado o trabajos de titulación que se realicen con el apoyo financiero, académico o institucional (operativo) de la Universidad”

f.

Autor: Montaña Sozoranga Wilmer Fabricio

Cédula: 11104634421

INDICE DE CONTENIDOS

APROBACIÓN DEL DIRECTOR DEL TRABAJO DE FIN DE TITULACIÓN	ii
DECLARACIÓN DE AUTORÍA Y CESIÓN DE DERECHOS	iii
INDICE DE CONTENIDOS.....	iv
Indice de Figuras.....	vii
Índice de tablas.....	ix
CAPITULO I: MARCO TEÓRICO	1
1. Datos Enlazados	2
1.1. Introducción.	2
1.2. Principios de Datos Enlazados.	2
1.3. Tecnologías.....	3
1.3.1. URI.....	3
1.3.2. RDF	5
1.3.3. Especificación de formatos de serialización de RDF	8
1.3.4. SPARQL Query Language for RDF	14
1.4. Acerca de DBpedia.....	15
1.4.1. Framework extracción	15
1.4.2. DBpedia Dataset	16
1.4.3. Acceso a DBpedia Dataset	16
2. Procesamiento de Lenguaje Natural (PLN).....	17
2.1. Introducción	17
2.2. Part of Speech Tagger	18
2.3. Chunking	18
2.4. Desambiguación	19
2.4.1. Métodos basados en el conociendo.	19
3. Servicios Web.....	21
3.1. Introducción	21
3.2. Tipos de servicios web	21
3.2.1. SOAP AND THE WS-* STACK	22
3.2.2. REST.....	22
3.3. Recursos y representaciones.....	23
CAPITULO 2: PROBLEMÁTICA	25
1. Estado actual.....	26
2. Justificación.....	26

3.	Objetivo General.....	28
4.	Objetivos Específicos.....	28
CAPITULO 3: Solución.....		29
5.	30	
1.	Propuesta.....	30
2.	Metodología.....	31
2.1.	Fases de desarrollo	31
3.	Desarrollo	32
3.1.	Análisis de requerimientos.....	32
3.1.1.	Requerimientos	32
3.1.2.	Modelo de Dominio	33
3.1.3.	Modelo de caso de Uso	34
3.2.	Análisis y diseño preliminar.....	35
3.2.1.	Especificación de casos de uso	35
3.3.	Diseño	41
3.3.1.	Arquitectura	41
3.3.2.	Componentes	42
3.3.3.	Diagrama de secuencia	45
3.4.	Implementación.....	48
3.4.1.	Servidor	48
3.4.2.	Servidor Dataset DBpedia Local.....	51
3.4.3.	Cliente web.....	52
3.4.4.	Resumen de prototipos.....	59
CAPITULO 4: VALIDACIÓN Y PRUEBAS		63
4.	Validación de resultados	64
4.1.	Objetivo.....	64
4.2.	Contexto.....	64
4.3.	Pruebas sobre el abstract de publicación.	64
4.3.1.	Descripción de publicaciones.	64
4.3.2.	Resultados de los servicios web	65
4.3.3.	Comparación y enriquecimientos de datos.....	66
4.4.	Pruebas en base a los contenido del proyecto SMARTLAND	67
4.4.1.	Datos del proyecto	67
4.4.2.	Contabilización de Resultados.....	68

4.4.3.	Comparación y enriquecimiento de datos.	69
5.	Pruebas funcionales	72
5.1.	Objetivo.....	72
5.2.	Escenario	73
5.3.	Pruebas sobre el servicio web de Segmentación en Sentencias	73
5.4.	Prueba sobre el servicio web de Tokenización.....	74
5.5.	Servicio web de Etiquetado.....	74
5.6.	Servicio web de Extracción.....	75
5.7.	Servicio web de Desambiguación y Enlace.....	76
	DISCUSIÓN.....	77
	CONCLUSIONES.....	79
	RECOMENDACIONES.....	80
	Bibliografía	81
	Anexos.....	84
1.	Anexo 1: Especificación de Requerimientos de Software (ERS)	85
2.	Anexo 2: Especificación de Caso de Uso (ECS) - Tokenización en Sentencias	93
3.	Anexo 3: Especificación de Caso de Uso (ECS) - Tokenización en Palabras.....	96
4.	Anexo 4: Especificación de Caso de Uso (ECS) - Etiquetado.....	99
5.	Anexo 5: Especificación de Caso de Uso (ECS) - Extracción de Entidades	102
6.	Anexo 6: Especificación de Caso de Uso (ECS) - Desambiguación y Enlace.....	105
7.	Anexo 7: Pseudocódigo del algoritmo de Lesk encontrado en (Satanjeev, 2002)	109
8.	Anexo 8: Pseudocódigo del algoritmo de Lesk encontrado en (Satanjeev, 2002)	110
9.	Anexo 9: Resultado completo de la prueba sobre el artículo: Consuming and producing linked open data: The case of OpenCourseWare	111
9.1.	Sentencias	111
9.2.	Tokens y Etiquetado	111
9.3.	Extracción	116
9.4.	Desambiguación y Enlace.....	117
10.	Anexo 10: Descripción del Proyecto SMARTLAND	118
10.1.	Sentencias	118
10.2.	Tokens y etiquetas	118
10.3.	Extracción	124
10.4.	Desambiguación y Enlace.....	126

INDICE DE FIGURAS

Figura 1: Relación entre URI, URL y URN.	4
Figura 2: Estructura de la URI (Scheme URI).	5
Figura 3: RDF 1.0 y 1.1 formatos de serialización.....	6
Figura 4: RDF con dos nodos (sujeto y objetos) y una arista (predicado).	7
Figura 5: Representación gráfica RDF	16
Figura 6: Arquitectura de provisión de Datos de Dbpedia.	17
Figura 7: Ejemplo POS Tag y Chunking.....	19
Figura 8: Lógica propuesta para la Aplicación.....	31
Figura 9: Modelo de dominio.....	34
Figura 10: Modelo de casos de uso	35
Figura 11. Arquitectura.	42
Figura 12. Dependencia de servicios web.....	44
Figura 13. Diagrama de secuencias de tokenización de sentencias	45
Figura 14. Diagrama de secuencias de tokenización en palabras.....	46
Figura 15. Diagrama de secuencias de etiquetado de palabra.....	46
Figura 16. Diagrama de secuencias de extracción.....	47
Figura 17. Diagrama de secuencia de desambiguación y enlace.....	47
Figura 18. Captura de consulta de recursos de DBpedia	48
Figura 19. Consulta de abstracts de recursos de DBpedia.....	48
Figura 20. Consulta para extraer el tipo de recurso.....	49
Figura 21. Resultado de servicio de tokenización en sentencias.....	50
Figura 22. Resultado del servicio web de desambiguación y enlace	50
Figura 23. Captura de la interfaz es su estado inicial.	53
Figura 24. Momento previo a la selección de la funcionalidad de etiquetado	53
Figura 25. Función etiquetado seleccionado junto a las funcionalidades dependientes ...	54
Figura 26. Momento previo a la deselección de la funcionalidad de tokenización	54
Figura 27. Función de tokenización deseleccionada junto las funciones que depende de esta.....	54
Figura 28. resultado del procesamiento del texto	55
Figura 29. Captura de la tabla con datos cuantitativos de los servicios invocados.	55
Figura 30. Menú construido con todos los servicios	55
Figura 31. Resultado de la función de tokenización	56
Figura 32. Resultado de la funcionalidad de tokenización.....	56
Figura 33. Resultado de la funcionalidad de etiquetado.	57

Figura 34. Resultado del servicio de extracción.	57
Figura 35. Resultado del servicio de enlace.....	58
Figura 36. Captura de la visualización del JSON.	58
Figura 37. Captura del resultado del servicio de etiquetado de palabra	59
Figura 38. Tabla y menú generado de la llamada al servicio de etiquetado.	59
Figura 39: Comparación grafica entre elementos extraídos y enlazados	66
Figura 40: Comparación de datos extraídos y enlazados proyecto SMARTLAND.....	69

ÍNDICE DE TABLAS

Tabla 1. Ejemplos de Prefijos de Espacios de Nombres e IRIs	8
Tabla 2. Ejemplo N-triple.....	8
Tabla 3. Ejemplo Turtle	9
Tabla 4. Ejemplo TriG	10
Tabla 5. Ejemplo N-Quads.....	11
Tabla 6. Ejemplo JSON-LD.....	11
Tabla 7. Ejemplo RDFa.....	12
Tabla 8. Ejemplo RDF/XML	13
Tabla 9. Resultado consulta SPARQL	14
Tabla 10: Fases de desarrollo del proyecto	31
Tabla 11: Resumen de requerimientos funcionales.....	32
Tabla 12: Requerimiento de tokenización de sentencias	36
Tabla 13. Requerimiento de tokenización en palabras.....	37
Tabla 14: Requerimiento de etiquetado	38
Tabla 15. Requerimiento de extracción de entidades.....	39
Tabla 16. Especificación del requerimiento de desambiguación y enlace	40
Tabla 17. Propiedades del JSON resultado de los servicios web.....	50
Tabla 18. Tabla resumen del prototipo 1	60
Tabla 19. Tabla resumen del prototipo 2.....	60
Tabla 20. Tabla resumen del prototipo 3.....	61
Tabla 21. Tabla resumen del prototipo 4.....	61
Tabla 22. Tabla resumen del prototipo 5.....	62
Tabla 23. Tabla resumen del prototipo 6.....	62
Tabla 24: Resumen de resultados del procesamiento de la publicación	65
Tabla 25: comparación de resultados y keywords dados por los autores	66
Tabla 26: Elementos enlazados con recursos de DBpedia	67
Tabla 27: Resultados de servicios del procesamiento de la descripción del proyecto SMARTLAND.....	69
Tabla 28: Comparación de elementos extraídos y los paquetes de trabajo del proyecto SMARTLAND.....	69
Tabla 29: Elementos enlazado hacia recursos de DBpedia del análisis del proyecto SMARTLAND.....	70
Tabla 30: Pruebas de funcionalidad sobre servicio web de segmentación en sentencias.	73

Tabla 31: Pruebas del servicio web de Tokenización.....	74
Tabla 32: Pruebas del servicio web de Etiquetado.....	74
Tabla 33: Pruebas del servicio web de Extracción	75
Tabla 34: Pruebas del servicio web de Desambiguación y Enlace.....	76

CAPITULO I: MARCO TEÓRICO

1. Datos Enlazados

1.1. Introducción.

En sus inicios la web en su primera versión 1.0, donde web era rígida en cuanto a la entrega de información, además de poco actualizada, convertía al visitante de un sitio web un simple lector, restringido de cualquier interacción. Se puede decir que la web no era más que paginas enlazadas mediante hipervínculos, limitada inclusive por la tecnología existente.

La web que está en constante crecimiento y evolución, alcanza su siguiente versión conocida como la web 2.0, en donde usuario juega el papel más importante, se convierte en gestor del contenido al que accede, puede calificar, compartir, rectificar, y aún más importante retroalimentar y subir su propia información a la web. Esto producto de la aparición de nuevas tecnologías y estandarización¹.

Los datos enlazados (Linked Data, por su nombre en inglés) llegan para dar forma a una nueva versión de la web, la web semántica. La W3C² los define así : “Linked Data se refiere a la utilización de las mejores prácticas para publicación, estructuración de los datos en la web, de tal forma que puedan ser enlazados entre sí, utilizando tecnología propias de web semántica como RDF, OCW, SPARQL, etc.” Se refiere a una estructuración de la web hasta convertirse en una web semántica, que en si busca que la información que se en publica en internet pueda no solo ser entendida por seres humanos sino también por la maquinas que navegan en la web, en donde a partir de un dato se puede descubrir otros datos a través sus relaciones, y el poder resolver inquietudes, encontrar información y compartirla sea fácil cosa que con una web carente de semántica se pude volver una tarea complicada y frustrante³.

1.2. Principios de Datos Enlazados.

La implementación de este concepto de datos relacionados entre sí, se basan

(Berners-Lee, Linked Data - Design Issues, 2006) Tim Berners Lee en su publicación Linked Data - Design Issues describe cuatro reglas base para la publicación de datos enlazados:

¹ <http://www.evolutionoftheweb.com>

² <http://www.w3.org/>

³ <http://www.w3c.es/Divulgacion/GuiasBreves/WebSemantica>

1. Usar URIs como nombre de las cosas
2. Usar URIs HTTP par que esas cosas puedan ser referenciadas
3. Representar los datos en RDF y SPARQL como lenguaje de consulta
4. Incluir enlaces hacia otra cosas, para descubrir más cosas

La utilizar de estas reglas para la publicación de datos, permite que estos por las características propias de las tecnologías sobre las cuales se construyen como:

- Las cosas nombradas por URIs son inequívocas y estos serán recursos.
- Los detalles o atributos y las relaciones de los datos van a estar descritos y estructurados en formato RDF
- Se puede acceder o realizar consultas sobre estos mediante SPARQL
- Los recursos publicados estar relacionados entre sí.

1.3. Tecnologías.

1.3.1. *URI.*

URI (Uniform Resource Identifier) permite mediante una secuencia de caracteres la identificación de recursos dentro de la web, como: imágenes, videos, programa, páginas web, correos electrónicos, servicios, etc., para cumplir con esta función utiliza URN (por sus siglas en ingles de Uniform Resource Name, en español “Nombre de Recurso Uniforme”) que cumple con la función de nombrar a los recursos y URL (Uniform Resource Locator, en español “Localizador Uniforme de Recursos”) que permite “apuntar” hacia la ubicación de un recurso. Sobre al acceso a recursos por medio de URIs en (Lapiente, 2013) aporta con las siguiente síntesis “(...). Los URIs hacen posible encontrar los recursos bajo una gran variedad de esquemas definidos y métodos de acceso tales como HTTP, FTP, Gopher, news, telnet o correos electrónicos localizables siempre de la misma manera, ya que a un mismo documento se puede acceder desde distintos protocolos.”

En la Figura 1 se representa la relación existente entre estas tecnologías, podemos observar que URI abarca a URL y URN para nombrar a los recursos.

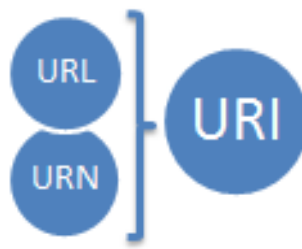


Figura 1: Relación entre URI, URL y URN.

Fuente: propio

Para entender mejor la función que cumple en la web se explican los términos que lo conforman tomado del RFC⁴⁵ (Request for Comments) que lo describe como: “Un identificador uniforme de recursos (URI) proporciona un medio simple y extensible para la identificación de un recurso.”

Uniforme: Uniformidad ofrece varios beneficios. Permite diferentes tipos de identificadores de recursos que se utilizarán en el mismo contexto, aun cuando los mecanismos utilizados para acceder a esos recursos pueden ser diferentes. Permite la interpretación semántica uniforme de convenciones sintácticas comunes a través de diferentes tipos de identificadores de recursos.

Recurso: Esta especificación no limita el alcance de lo que podría ser un recurso; más bien, el término "recurso" se utiliza en un sentido general de lo que pudiera ser identificado por un URI. Ejemplos conocidos incluyen un documento electrónico, una imagen, una fuente de información con un propósito consistente (por ejemplo, "parte meteorológico de hoy para Los Ángeles"), un servicio (por ejemplo, una puerta de enlace HTTP a SMS), y una colección de otros recursos. Un recurso no es necesariamente accesible a través de Internet; por ejemplo, los seres humanos, las empresas y los libros encuadrados en una biblioteca también pueden ser recurso (...).

Identificador: Un identificador encarna la información necesaria para distinguir lo que se identificó a partir de todas las otras cosas dentro de su ámbito de aplicación de la identificación. Nuestro uso de los términos "identificar" y "identificación" se refieren a este fin de distinguir un recurso de todos los demás recursos, independientemente de cómo se logra ese propósito (ejemplos, nombre, dirección, o el contexto). Estos términos no deben

⁴ <http://www.ietf.org/rfc/rfc3986.txt>

⁵ http://es.wikipedia.org/wiki/Request_for_Comments

confundirse con la presunción de que un identificador define o encarna la identidad de lo que se hace referencia, aunque esto puede ser el caso de algunos identificadores. Tampoco debe asumirse que un sistema que utiliza los URI tendrá acceso al recurso identificado: en muchos casos, los URI se utilizan para referirse a los recursos sin ninguna intención de que se puede acceder. Del mismo modo, "un" recurso identificado podría no ser singular en la naturaleza (ejemplo, un recurso puede ser un conjunto con nombre o una asignación que varía con el tiempo)" (Berners-Lee, Fielding, & Masinter, 2005).

La estructura de una URI es explicada en el Figura 2 de (Albahari & Albahari, 2012), todos los componente aquí identificados son explicados ampliamente en el RFC3986⁶ en el cual se definen las URIs.

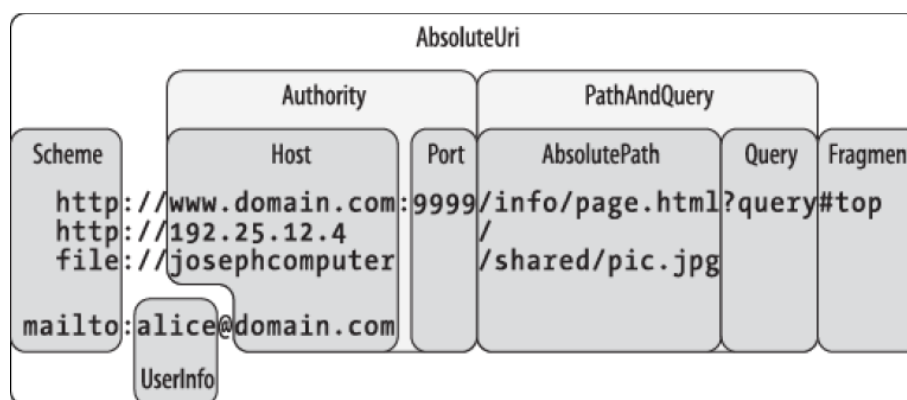


Figura 2: Estructura de la URI (Scheme URI).
Fuente: (Albahari & Albahari, 2012)

1.3.2. RDF

En (McBride, 2004) se encuentran las especificaciones sobre *Resource Description Framework* (RDF), en su traducción al español *Marco de Descripción de Recursos* que es un lenguaje para representar la información sobre recursos en la World Wide Web. Busca guardar metadatos de recursos que se encuentran disponibles en web, como título, autor, fechas relacionadas, derechos del autor y la información de la licencia del recurso.

Existen recursos que no necesariamente son digitales y que por lo tanto no pueden ser recuperados a través de la red, en la web encontramos representaciones de los mismos, como los que existen en un sitio de compras, de estos recursos se puede guardar

⁶ <https://tools.ietf.org/html/rfc3986#section-3>

información relevante para la operaciones de adquisición, por ejemplo, disponibilidad, características específicas de cada producto. Otro recurso que se rescata de este ejemplo son los usuarios de quienes se puede rescatar información relevantes de para las diferentes cuentas que pueden poseer.

Desde la primera publicación de RDF en su versión 1.0 en febrero del 2004 pasaron 10 años para que en la primera mitad del 2014 se conoció la versión 1.1, la cual se encuentra redactada en *What's New in RDF 1.1*⁷. Una de las principales novedades es la incorporación de IRIs⁸ para referenciar los recursos, esto en lugar de URIs además de nuevos formatos de serialización, este último se visualizan en la Figura 3.

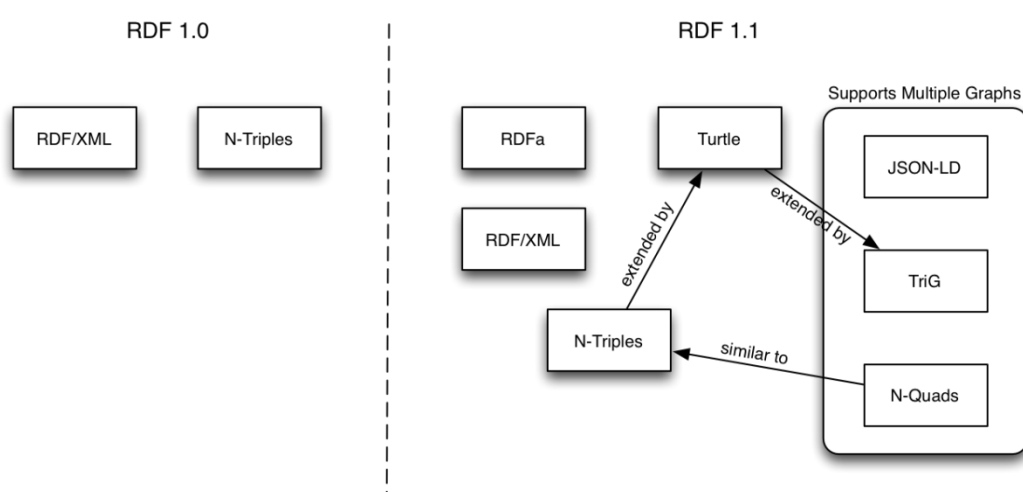


Figura 3. RDF 1.0 y 1.1 formatos de serialización.

Fuente: (Wood, 2014)

1.3.2.1. Modelo de datos RDF

La estructura esencial que define la sintaxis abstracta de un RDF es el conjunto formado por un sujeto, un predicado y un objeto, constituyendo una tripleta, a la vez un conjunto de tripletas se denominan grafo RDF, cada tripleta puede ser visualizada como se representa como un enlace nodo-arco-nodo. A una tripleta la conforman tres elementos denominados Sujeto, Predicado y Objeto, que por convencionalismo se listan en ese orden, en la Figura 4 se visualiza de forma gráfica. Los nodos pueden ser de tres tipos diferentes: IRIS, Literales, Nodos en blanco. (Cyganiak, Wood, & Lanthaler, 2014).

⁷ <http://www.w3.org/TR/2014/REC-rdf11-concepts-20140225/>

⁸ <http://www.ietf.org/rfc/rfc3987.txt>

IRI (International Resource Identifier) son cadenas de caracteres UNICODE⁹ que nos permiten referenciar los recursos web. Como parte de una tripleta estos pueden ser sujeto, predicado u objeto.

Literales son cadena de caracteres UNICODE a los cuales se determina su tipo de dato, coda fechas, valores numéricos o cadenas de caracteres y en caso de que este último adicionalmente se puede determinar en que idioma se encuentra redactado. Por su propósito en si estos, como parte de una tripleta solo ocupar el lugar de un objeto.

Nodos en blanco que permiten el poder referirnos a recursos son la necesidad utilizar identificadores globales. Debido a sus características propias los nodos en blanco solo pueden ser utilizados ya sea como sujetos o predicados.

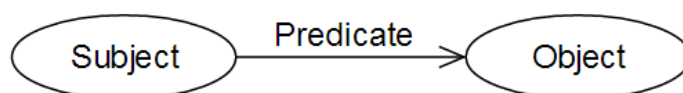


Figura 4: RDF con dos nodos (sujeto y objetos) y una arista (predicado).

Fuente: (Cyganiak, Wood, & Lanthaler, 2014)

1.3.2.2. Vocabularios RDF y Espacio de Nombres IRIs

Un RDF vocabulario es una colección de IRIs destinados para ser usados en grafos RDF. Por ejemplo, los IRIs documentados en [RDF11-SCHEMA]¹⁰ son el Vocabulario RDF Esquema (RDF Schema Vocabulary). Son una colección de "términos" que definen conceptos y relaciones que sirven para representar un área de conocimiento, para un propósito en particular. Ayudan a la integración de datos y a aumentar el conociendo al descubrir una nueva relación. (Cyganiak, Wood, & Lanthaler, 2014)

Los vocabularios pueden ir desde simples a complejos, ampliamente usados como Schema RDF utilizado, FOAF¹¹ y Dublin Core Metadata Element Set¹² para vocabularios con miles de términos, tales como los utilizados en la asistencia sanitaria para describir síntomas, enfermedades y tratamientos. Vocabularios juegan un papel muy importante en Linked Data, específicamente para ayudar con la integración de datos. El uso de este término se superpone con la Ontología.

⁹ <http://www.unicode.org/versions/Unicode7.0.0/>

¹⁰ <http://www.w3.org/TR/2014/REC-rdf-schema-20140225/>

¹¹ <http://www.foaf-project.org/>

¹² <http://dublincore.org/documents/dces/>

Los IRIs dentro de un vocabulario RDF a menudo comienzan con una subcadena común conocido como un espacio de nombres IRI (Namespaces IRI). Algunos IRIs de espacio de nombres se asocian por convención con un nombre corto conocido como un prefijo de espacio de nombres (namespace prefix). (W3C, 2013)

El término "espacio de nombres" por sí no tiene un significado bien definido en el contexto de la RDF, pero a veces se utiliza de manera informal en el sentido de "espacio de nombres IRI" o "vocabulario RDF". (W3C, 2013)

Tabla 1. Ejemplos de Prefijos de Espacios de Nombres e IRIs

Namespace prefix	Namespace IRI	RDF vocabulary
rdf	http://www.w3.org/1999/02/22-rdf-syntax-ns#	The RDF built-in vocabulary [RDF11-SCHEMA] ¹³
rdfs	http://www.w3.org/2000/01/rdf-schema#	The RDF Schema vocabulary [RDF11-SCHEMA] ¹⁴
xsd	http://www.w3.org/2001/XMLSchema#	The RDF-compatible XSD types ¹⁵

Fuente: (Cyganiak, Wood, & Lanthaler, 2014)

1.3.3. Especificación de formatos de serialización de RDF

1.3.3.1. N-Triples

Tripletas N-Triples son una secuencia de términos RDF que representan al sujeto, predicado y objeto de una Tripletta RDF. Estos pueden estar separados por espacios en blanco (espacios U +0020 o tabulaciones U +0009). Esta secuencia es terminada por un '.' y una nueva línea (opcional al final de un documento). (Carothers & Seaborne, 2014)

N-triple en formato en texto plano para grafos RDF, en la Tabla 2 se ejemplifica su estructura.

Tabla 2. Ejemplo N-triple

01	< http://example.org/bob#me > < http://www.w3.org/1999/02/22-rdf-syntax-ns#type > < http://xmlns.com/foaf/0.1/Person > .
----	---

¹³ <http://www.w3.org/TR/2014/REC-rdf11-concepts-20140225/#bib-RDF11-SCHEMA>

¹⁴ <http://www.w3.org/TR/2014/REC-rdf11-concepts-20140225/#bib-RDF11-SCHEMA>

¹⁵ <http://www.w3.org/TR/2014/REC-rdf11-concepts-20140225/#dfn-rdf-compatible-xsd-types>

02	<http://example.org/bob#me> <http://xmlns.com/foaf/0.1/knows> <http://example.org/alice#me> .
03	<http://example.org/bob#me> <http://schema.org/birthDate> "1990-07-04"^^<http://www.w3.org/2001/XMLSchema#date> .
04	<http://example.org/bob#me> <http://xmlns.com/foaf/0.1/topic_interest> <http://www.wikidata.org/entity/Q12418> .
05	<http://www.wikidata.org/entity/Q12418> <http://purl.org/dc/terms/title> "Mona Lisa" .
06	<http://www.wikidata.org/entity/Q12418> <http://purl.org/dc/terms/creator> <http://DBpedia.org/resource/Leonardo_da_Vinci> .
07	<http://data.europeana.eu/item/04802/243FA8618938F4117025F17A8B813C5F9AA4D61 9> <http://purl.org/dc/terms/subject> <http://www.wikidata.org/entity/Q12418> .

Fuente: (Schreiber & Raimond, 2014)

1.3.3.2. Turtle

Permite a un grafo RDF a ser completamente escrito en un formulario de texto compacto y natural, con las abreviaturas para los patrones y tipos de datos de uso común. Turtle ofrece niveles de compatibilidad con el formato N-Triples, así como con la sintaxis de patrón de triplas de la Recomendación de la W3C de SPARQL¹⁶. (Beckett, Berners-Lee, Prud'hommeaux, Carothers, & Machina., 2014)

Turtle es una extensión del N-Triples. Además de la sintaxis básica N-Triples, Turtle introduce una serie de atajos sintácticos, como el soporte para prefijos de espacio de nombres, listas y abreviaturas para datos tipo literales. Turtle ofrece una compensación entre la facilidad de la escritura, la facilidad de análisis y facilidad de lectura. (Schreiber & Raimond, 2014)

Tabla 3. Ejemplo Turtle

01	<http://example.org/>
02	PREFIX foaf: <http://xmlns.com/foaf/0.1/>
03	PREFIX xsd: <http://www.w3.org/2001/XMLSchema#>
04	PREFIX schema: <http://schema.org/>
05	PREFIX dcterms: <http://purl.org/dc/terms/>
06	PREFIX wd: <http://www.wikidata.org/entity/>
07	
08	<bob#me>
09	a foaf:Person ;
10	foaf:knows <alice#me> ;
11	schema:birthDate "1990-07-04"^^xsd:date ;
12	foaf:topic_interest wd:Q12418 .
13	
14	wd:Q12418
15	dcterms:title "Mona Lisa" ;
16	dcterms:creator <http://DBpedia.org/resource/Leonardo_da_Vinci> .
17	

¹⁶ <http://www.w3.org/TR/sparql11-query/>

18	<http://data.europeana.eu/item/04802/243FA8618938F4117025F
19	17A8B813C5F9AA4D619> dcterms:subject wd:Q12418 .

Fuente: (Schreiber & Raimond, 2014)

1.3.3.3. TriG

La sintaxis de la Turtle sólo soporta la especificación de grafos simples sin un medio para "nombrarlos". TriG es una extensión de la Turtle que permite la especificación de múltiples grafos en forma de un conjunto de datos RDF. (Schreiber & Raimond, 2014)

Un documento TriG permite escribir un conjunto de datos RDF en una forma textual compacta. Se consiste de una sucesión de directivas, declaraciones triples, declaraciones de grafos que contienen declaraciones triple-generación y líneas en blanco opcionales. (Bizer & Cyganiak, 2014)

Tabla 4. Ejemplo TriG

01	BASE <http://example.org/>
02	PREFIX foaf: <http://xmlns.com/foaf/0.1/>
03	PREFIX xsd: <http://www.w3.org/2001/XMLSchema#>
04	PREFIX schema: <http://schema.org/>
05	PREFIX dcterms: <http://purl.org/dc/terms/>
06	PREFIX wd: <http://www.wikidata.org/entity/>
07	
08	GRAPH <http://example.org/bob>
09	{
10	<bob#me>
11	a foaf:Person ;
12	foaf:knows <alice#me> ;
13	schema:birthDate "1990-07-04"^^xsd:date ;
14	foaf:topic_interest wd:Q12418 .
15	}
16	
17	GRAPH <https://www.wikidata.org/wiki/Special:EntityData/Q12418>
18	{
19	wd:Q12418
20	dcterms:title "Mona Lisa" ;
21	dcterms:creator <http://DBpedia.org/resource/Leonardo_da_Vinci> .
22	
23	<http://data.europeana.eu/item/04802/243FA8618938F4117025F17A8B813C5F9AA4D619>
24	dcterms:subject wd:Q12418 .
25	}
26	
27	<http://example.org/bob>
28	dcterms:publisher <http://example.org> ;
29	dcterms:rights <http://creativecommons.org/licenses/by/3.0/> .

Fuente: (Schreiber & Raimond, 2014)

1.3.3.4. N-Quads

N-Quads es una simple extensión de N-Triples para permitir el intercambio de RDF Datasets. N-Quads le permite a uno agregar un cuarto elemento a una línea, capturando en la gráfica IRI la tripleta descrito en esa línea. (Schreiber & Raimond, 2014)

Tabla 5. Ejemplo N-Quads

01	<http://example.org/bob#me> <http://www.w3.org/1999/02/22-rdf-syntax-ns#type> <http://xmlns.com/foaf/0.1/Person> <http://example.org/bob> .
02	<http://example.org/bob#me> <http://xmlns.com/foaf/0.1/knows> <http://example.org/alice#me> <http://example.org/bob> .
03	<http://example.org/bob#me> <http://schema.org/birthDate> "1990-07-04"^^<http://www.w3.org/2001/XMLSchema#date> <http://example.org/bob> .
04	<http://example.org/bob#me> <http://xmlns.com/foaf/0.1/topic_interest> <http://www.wikidata.org/entity/Q12418> <http://example.org/bob> .
05	<http://www.wikidata.org/entity/Q12418> <http://purl.org/dc/terms/title> "Mona Lisa" <https://www.wikidata.org/wiki/Special:EntityData/Q12418> .
06	<http://www.wikidata.org/entity/Q12418> <http://purl.org/dc/terms/creator> <http://DBpedia.org/resource/Leonardo_da_Vinci> <https://www.wikidata.org/wiki/Special:EntityData/Q12418> .
07	<http://data.europeana.eu/item/04802/243FA8618938F4117025F17A8B813C5F9AA4D619> <http://purl.org/dc/terms/subject> <http://www.wikidata.org/entity/Q12418> <https://www.wikidata.org/wiki/Special:EntityData/Q12418> .
08	<http://example.org/bob> <http://purl.org/dc/terms/publisher> <http://example.org> .
09	<http://example.org/bob> <http://purl.org/dc/terms/rights> <http://creativecommons.org/licenses/by/3.0/> .

Fuente: (Schreiber & Raimond, 2014)

JSON-LD (JSON-based RDF syntax);

Proporciona una sintaxis JSON para gráficos y conjuntos de datos RDF. JSON-LD puede ser utilizado para transformar documentos JSON a RDF con cambios mínimos. JSON-LD ofrece identificadores universales de objetos JSON, un mecanismo en el que un documento JSON se puede referir a un objeto descrito en otro documento JSON en otros lugares en la Web, así como el tipo de datos y el lenguaje de manipulación.

Tabla 6. Ejemplo JSON-LD

01	{
02	"@context": "example-context.json",

```

03  "@id": "http://example.org/bob#me",
04  "@type": "Person",
05  "birthdate": "1990-07-04",
06  "knows": "http://example.org/alice#me",
07  "interest": {
08    "@id": "http://www.wikidata.org/entity/Q12418",
09    "title": "Mona Lisa",
10    "subject_of": "http://data.europeana.eu/item/04802/
243FA8618938F4117025F17A8B813C5F9AA4D619",
11    "creator": "http://DBpedia.org/resource/Leonardo_da_Vinci"
12  }
13  }

```

Fuente: (Schreiber & Raimond, 2014)

1.3.3.5. RDFa

Es una sintaxis de RDF que se puede utilizar para insertar datos RDF dentro de los documentos HTML y XML. Esto permite, por ejemplo, a los motores de búsqueda agregar estos datos al rastrear la Web y enriquecer los resultados de búsqueda. (Schreiber & Raimond, 2014)

Tabla 7. Ejemplo RDFa

```

01  <body prefix="foaf: http://xmlns.com/foaf/0.1/
02          schema: http://schema.org/
03          dcterms: http://purl.org/dc/terms/">
04  <div resource="http://example.org/bob#me" typeof="foaf:Person">
05    <p>
06      Bob knows <a property="foaf:knows"
07      href="http://example.org/alice#me">Alice</a>
08      and was born on the <time property="schema:birthDate"
09      datatype="xsd:date">1990-07-04</time>.
10    </p>
11    <p>
12      Bob is interested in <span property="foaf:topic_interest"
13      resource="http://www.wikidata.org/entity/Q12418">the Mona Lisa</span>.
14    </p>
15  </div>
16  <div resource="http://www.wikidata.org/entity/Q12418">
17    <p>
18      The <span property="dcterms:title">Mona Lisa</span> was painted by
19      <a property="dcterms:creator"
20      href="http://DBpedia.org/resource/Leonardo_da_Vinci">Leonardo da Vinci</a>
21      and is the subject of the video
22      <a
23      href="http://data.europeana.eu/item/04802/243FA8618938F4117025F17A8B813
24      C5F9AA4D619">'La Joconde à Washington'</a>.
25    </p>
26  </div>

```

	resource="http://data.europeana.eu/item/04802/243FA8618938F4117025F17A8B813C5F9AA4D619">
23	<link property="dcterms:subject"
24	href="http://www.wikidata.org/entity/Q12418"/>
25	</div>
	</body>

Fuente: (Schreiber & Raimond, 2014)

1.3.3.6. RDF/XML

Proporciona una sintaxis XML para grafos RDF. Cuando RDF fue desarrollado originalmente en la década de 1990, esta fue su única sintaxis, y algunas personas siguen llamando esta sintaxis "RDF". En 2001, se propuso un precursor de la Tortuga llamado "N3", y poco a poco los otros idiomas que aparecen aquí se han adoptado y normalizado. (Schreiber & Raimond, 2014)

Tabla 8. Ejemplo RDF/XML

01	<?xml version="1.0" encoding="utf-8"?>
02	<rdf:RDF
03	xmlns:dcterms="http://purl.org/dc/terms/"
04	xmlns:foaf="http://xmlns.com/foaf/0.1/"
05	xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#"
06	xmlns:schema="http://schema.org/">
07	<rdf:Description rdf:about="http://example.org/bob#me">
08	<rdf:type rdf:resource="http://xmlns.com/foaf/0.1/Person"/>
09	<schema:birthDate
	rdf:datatype="http://www.w3.org/2001/XMLSchema#date">1990-07-
10	04</schema:birthDate>
	<foaf:knows rdf:resource="http://example.org/alice#me"/>
11	<foaf:topic_interest
	rdf:resource="http://www.wikidata.org/entity/Q12418"/>
12	</rdf:Description>
13	<rdf:Description rdf:about="http://www.wikidata.org/entity/Q12418">
14	<dcterms:title>Mona Lisa</dcterms:title>
15	<dcterms:creator
	rdf:resource="http://DBpedia.org/resource/Leonardo_da_Vinci"/>
16	</rdf:Description>
17	<rdf:Description
	rdf:about="http://data.europeana.eu/item/04802/243FA8618938F4117025F17A8B813C5F9AA4D619">
18	<dcterms:subject
	rdf:resource="http://www.wikidata.org/entity/Q12418"/>
19	</rdf:Description>
20	</rdf:RDF>

Fuente: (Schreiber & Raimond, 2014)

1.3.4. SPARQL Query Language for RDF

SPARQL se puede utilizar para expresar consultas que permiten interrogar diversas fuentes de datos, si los datos se almacenan de forma nativa como RDF o son definidos mediante vistas RDF a través de algún sistema middleware. SPARQL contiene las capacidades para la consulta de los patrones obligatorios y opcionales de grafo, junto con sus conjunciones y disyunciones. SPARQL también soporta la ampliación o restricciones del ámbito de las consultas indicando los grafos sobre los que se opera. Los resultados de las consultas SPARQL pueden ser conjuntos de resultados o grafos RDF. (Prud'hommeaux & Seaborne, 2008)

La mayoría de las formas de consulta en SPARQL contienen un conjunto de patrones de tripleta (triple patterns) denominadas patrón de grafo básico. Los patrones de tripleta son similares a las tripletas RDF, excepto que cada sujeto, predicado y objeto puede ser una variable. Un patrón de grafo básico concuerda con un subgrafo de datos RDF cuando los términos RDF (RDF terms) de dicho subgrafo pueden ser sustituidos por las variables y el resultado es un grafo RDF equivalente al subgrafo en cuestión. (Prud'hommeaux & Seaborne, 2008)

A continuación se redacta un ejemplo de una consulta SPARQL, tomado de (Prud'hommeaux & Seaborne, 2008).

Datos:

```
<http://example.org/book/book1> <http://purl.org/dc/elements/1.1/title> "SPARQL Tutorial" .
```

Consulta:

```
SELECT ?title
WHERE
{
  <http://example.org/book/book1> <http://purl.org/dc/elements/1.1/title> ?title.
}
```

Resultado de la consulta:

Tabla 9. Resultado consulta SPARQL

title
"SPARQL Tutorial"

1.4. Acerca de DBpedia

DBpedia¹⁷ da la siguiente definición sobre si misma: “Es un esfuerzo de la comunidad crowd-sourced¹⁸ para extraer información estructurada de Wikipedia¹⁹ y hacer esta información disponible en la web. DBpedia permite que hacer consultas sofisticadas contra Wikipedia.” El conocimiento extraído de Wikipedia es publicado cumpliendo los estándares de la Web Semántica y las mejores prácticas de Linkend Data.

1.4.1. *Framework extracción*

Los artículos de Wikipedia consisten sobre todo en texto libre, pero también comprenden diversos tipos de información estructurada en forma de wiki markup²⁰. Dicha información incluye plantillas infobox, información de categorización, imágenes geo-coordenadas, enlaces a páginas web externas, páginas de desambiguación, redirecciones entre páginas y vínculos a través de las diferentes ediciones lingüísticas de Wikipedia. El marco de la extracción DBpedia extrae esta información estructurada de Wikipedia y lo convierte en una rica base de conocimientos (Lehmann, y otros, 2012)

En la figura 5 se observa el marco de trabajo necesario para lograr que todo el proceso partiendo de la extracción de información desde Wikipedia hasta poder disponer de ella como datos enlazados.

¹⁷ <http://dbpedia.org/About>

¹⁸ <http://es.wikipedia.org/wiki/Crowdsourcing>

¹⁹ <http://www.wikipedia.org/>

²⁰ http://en.wikipedia.org/wiki/Help:Wiki_markup

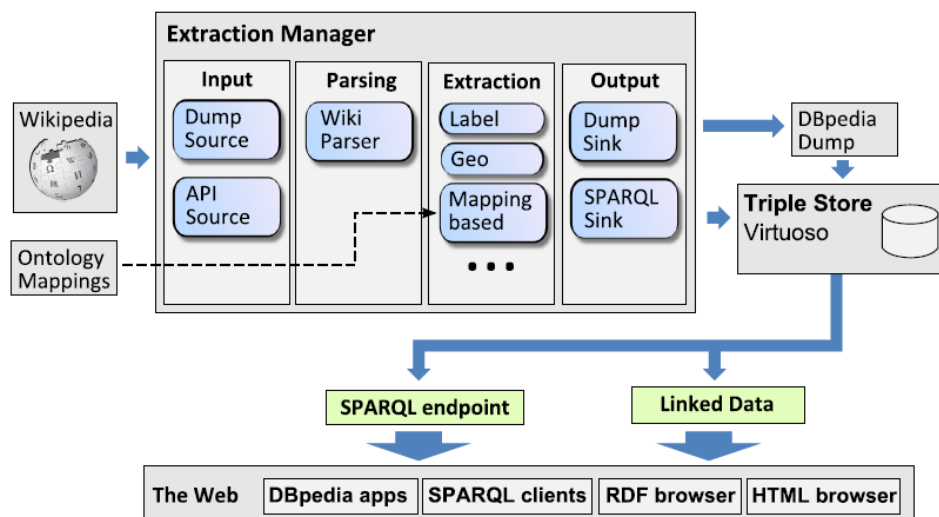


Figura 5: Representación gráfica RDF
Fuente: (Lehmann, y otros, 2012)

1.4.2. *DBpedia Dataset*

DBpedia se trata de una base de conocimiento (en inglés knowledge base) que se encuentra distribuida en 119 idiomas que en total describen 12.6 millones de cosas únicas, 24.6 millones de enlaces a imágenes, 27.6 millones de enlaces a fuentes externas, 45 millones de enlaces a fuentes externas de datos RDF y 67 millones de enlaces a categorías Wikipedia, 42.1 millones de enlaces a categorías YAGO²¹.

Lo cual la establece como una fuente muy buena de información sobre cualquier ámbito de conocimiento, esto gracias al continuo crecimiento de la Wikipedia, su fuente de información. Pero no esto no quiere decir que la única base de conocimiento disponible en la web, se encuentran disponibles otras como YAGO.

1.4.3. *Acceso a DBpedia Dataset*

El Dataset de DBpedia se almacena y publica mediante OpenLink Virtuoso. La infraestructura de Virtuoso permite el acceso a los datos RDF de DBpedia a través de un SPARQL endpoint, junto al soporte HTTP para cualquier GET estándar de cliente Web para HTML o representación RDF de un recurso DBpedia. (Bizer, DBpedia, 2009).

²¹ <http://www.mpi-inf.mpg.de/departments/databases-and-information-systems/research/yago-naga/yago/>

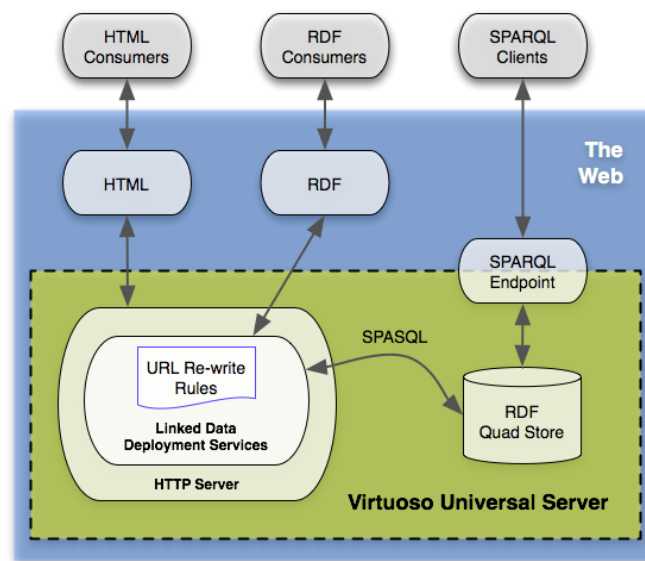


Figura 6. Arquitectura de provisión de Datos de Dbpedia.

Fuente: (Bizer, Dbpedia, 2009)

Se exponen dos formas para acceder a la Dataset de DBpedia:

1. Querying DBpedia: se puede acceder a través del SPARQL endpoint público <http://DBpedia.org/sparql> proporcionado por OpenLink Virtuoso. Por este método se puede acceder enviando query Sparql para hacer consultar sobre Dataset.
2. Linked Data: se refiere a la aplicación de los principios de datos enlazados revisados en 1.1.2. para nombrar y referenciar los recursos dentro de Dataset de DBpedia como por ejemplo: http://DBpedia.org/resource/The_Lord_of_the_Rings

Además de estas opciones se puede descargar el Dataset de DBpedia en diferentes idiomas teniendo en cuenta de que el número de recursos puede cambiar de idioma a idioma puesto que no se trata de una traducción sino de una recopilación de información de Wikipedia la cual se encuentra más extendida en inglés que otros lenguas,

2. Procesamiento de Lenguaje Natural (PLN)

2.1. Introducción

El procesamiento de lenguaje natural se preocupa por entender el lenguaje humano, la comunicación una tarea que para las personas e inclusive animales es tan natural y cotidiana, se vuelve un reto al tratar de interpretarlo mediante procesos computacionales a fin de comprenderlo y poder replicarlo.

La dificultad de la construcción de una aplicación de la ingeniería lingüística variara de acuerdo a objetivo que se persiga, esto explicado por (García, 2005) en donde ejemplifica: “un sistema de generación de cartas personalizado no precisa ningún tratamiento de comprensión, o un sistema de identificación de la lengua (o un detector de errores ortográficos) no necesitan generar lenguaje humano. La mayoría de las aplicaciones incluyen, sin embargo, alguna forma más o menos precisa de comprensión. Así, un sistema de consulta en lenguaje humano a una base de datos precisa un nivel muy alto de comprensión de las expresiones del interlocutor humano para que la respuesta del sistema sea de utilidad. En cambio, en un sistema de traducción o de resumen automáticos se pueden lograr niveles de corrección muy notables con niveles de comprensión bajos. Es decir, no es preciso comprender totalmente una oración para ser capaz de traducirla correctamente.”

2.2. Part of Speech Tagger

Permite distinguir la función de una palabra en un determinado contexto mediante la asignación de una etiqueta predefinida. “Una part of speech tagger es un sistema que usa el contexto para asignar parte de un discurso a una palabra”, (Cutting, D., Kupiec, J., Pedersen, J., & Sibun, P., 1992).

El etiquetado de palabra ya permite una primera desambiguación en cuanto a la función de la palabra en un contexto. Así se puede por ejemplo²² ver que la palabra “*dado*” que si bien es nombre en singular, también puede ser una forma del verbo dar.

Pero antes de poder etiquetar una palabra por su función es necesaria una Tokenización del texto que va a analizar, que consiste en separarlo en palabras individuales reconociendo un token para palabra o carácter extraído.

2.3. Chunking

Text Chunking consiste en dividir un texto en frases de tal manera que palabras sintácticamente relacionadas sean miembros de la misma clase. Estas frases no se superponen es decir que una sola palabra puede ser miembro de un chunk. (Tjong Kim Sang, E. F., & Buchholz, S., 2000)

Este proceso es básico al momento de detectar entidades dentro de un texto, este proceso lo se puede observar en la figura 1 en donde la sentencia, *We saw the yellow*

²² <http://es.wikipedia.org/wiki/Ambig%C3%BCedad>

dog, está separada en cuadros en los más pequeños se observa etiquetas de POS Tag y las más grades al nivel de chunking. Una vez la frase ha pasado por el proceso de chunking se puede rescatar dos entidades dentro de la sentencia como *We* y *the yellow dog*.

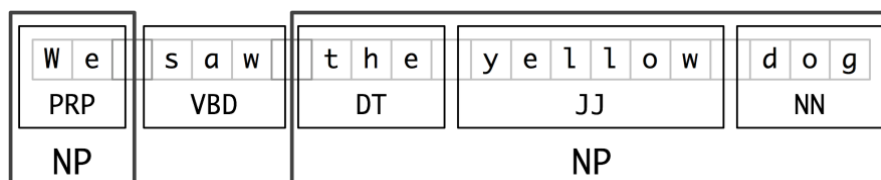


Figura 7: Ejemplo POS Tag y Chunking

Fuente: <http://www.nltk.org/book/ch07.html>

Las etiquetas utilizadas el idioma inglés se encuentran descritas en el anexo 10.

2.4. Desambiguación

La polisemia es un fenómeno muy común, que se refiere a cuando una palabra tiene varios significados, la desambiguación busca descifrar que significado una palabra está siendo utilizado de acuerdo a un contexto en específico, se denomina Desambguacion del sentido de la palabra (Word Sense Disambiguation (WSD), en inglés), este es un problema propio del procesamiento de lenguaje natural (PLN). El descifrar estos distintos significados para los seres humanos es muy común, lo resolvemos de forma cotidiana y pasa casi desapercibida.

La desambiguación de recursos en contexto con Linked Open Data segun (Peláez, Morocho, & Malla, 2012) encontramos “Al realizar la desambiguación estamos consiguiendo un paso clave en la recuperación de la información ya que numerosas palabras cambian de sentido según el contexto en el que nos encontremos trabajando y es sumamente necesario aclarar con el usuario cuál es el sentido que le otorga al término que está utilizando.”

2.4.1. Métodos basados en el conociendo.

Utilizan fuentes de léxicas estructuras existentes para resolver el significado de las palabras, (Tello Leal, 2009) lo define de la siguiente forma: “Estos métodos utilizan un conocimiento lingüístico previamente adquirido. La idea básica consiste en utilizar recursos externos para desambiguar las palabras, tales como diccionarios, tesauros

(vocabularios controlados que representan las relaciones semánticas con otras palabras y sus significados), textos sin ningún tipo de etiquetado e incluso recursos de la Web”

Algunos recursos lingüísticos que se utilizan para la desambiguación de sentido de la palabra en PNL se describen a continuación.

2.4.1.1. *WordNet*:

“WordNet es una base de datos léxico-conceptual del inglés estructurada en forma de red semántica y construida manualmente. Es el lexicón relacional en formato electrónico más completo y extenso existente, comparable sólo con el diccionario bilingüe para el japonés y el inglés EDR Electronic Dictionary²³. La unidad básica en que se estructura WordNet es el synset²⁴, un conjunto de sinónimos representando un concepto.” (Mihaela, 2004)

2.4.1.2. *EuroWordNet*

“El propósito de EuroWordNet ha sido construir una base de datos léxica multilingüe para diferentes lenguas europeas, siguiendo la metodología de WordNet. EuroWordNet es una base de datos multilingüe, con wordnets para varias lenguas (holandés, italiano, español, inglés, alemán, francés, estoniano y checo), compatibles entre sí en cobertura e interpretación de las relaciones.” (Mihaela, 2004)

2.4.1.3. *Extended WordNet*

“Debido a que WordNet ha sido construida como una base de datos léxica, hay limitaciones en su uso para ciertas aplicaciones del procesamiento del conocimiento; por ejemplo, no es posible extraer palabras relacionadas temáticamente. El propósito del proyecto Extended WordNet es transformar WordNet en un formato que permita la derivación de relaciones semánticas y lógicas adicionales.” (Mihaela, 2004)

2.4.1.4. *MindNet*

“MindNet. Una alternativa a WordNet y EuroWordNet es MindNet (Dolan et al., 2000). Aunque tiene un núcleo derivado a partir de diccionarios, la red se construye a base de oraciones nuevas de un corpus que, una vez analizadas, se incorporan a la red. En la visión de sus autores, la dinamicidad y la continua ampliación, permitirán a MindNet perfilarse como un sistema de amplia cobertura.” (Mihaela, 2004)

²³ Para detalles, consúltese el sitio: <http://www.ijnet.or.jp/edr/>.

²⁴ De synonym set ‘conjuntos de sinónimos’.

2.4.1.5. Algoritmo de Lesk

Este algoritmo se basa en diccionarios para resolver la ambigüedad, en (Pérez, 2009) una descripción de la Algoritmo de Lesk 1986, “es uno de los primeros algoritmo desarrollado para la desambiguación semántica de toda las palabras en cualquier texto. El único recurso requerido por el algoritmo es un conjunto de entradas en un diccionario, una por cada posible sentido y conocimiento sobre el contexto inmediato donde se desarrolla la desambiguación”

3. Servicios Web

3.1. Introducción

La W3C²⁵ (World Wide Web Consortium) encarda de estandarización de las tecnologías en la web aborda este tema de la siguiente forma: “Los servicios web proporcionan un medio estándar de interoperabilidad entre las distintas aplicaciones de software, que se ejecuta en una variedad de plataformas y/o marcos de trabajo. Los servicios Web se caracterizan por su gran interoperabilidad y extensibilidad. Se pueden combinar en una forma de acoplamiento flexible con el fin de lograr operaciones complejas. Programas que prestan servicios simples pueden interactuar entre sí con el fin de ofrecer servicios de valor añadido sofisticados.” Los servicios web permiten la colaboración entre aplicaciones independientemente de la plataforma en las que están desarrolladas, utiliza protocolos y normas estandarizadas en la web, además esto permite la reutilización de código, además de disminuir el coste de integración.

3.2. Tipos de servicios web

Dos tipos de servicios web se pueden encontrar de acuerdo con la forma en que se puede implementar abarcado diferentes tecnologías: *RESTful Web Services* y “*Big*”²⁶ *Web Services* (o también, The “Big” Web services technology stack, debido a la diversas tecnologías en las que se implementar como: SOAP, WSDL, WS-Addressing, WS-ReliableMessaging, WSSecurity, etc), estos dos tipos son expuesto en (Pautasso, Zimmermann, & Leymann, 2008)

²⁵ <http://www.w3.org/>

²⁶ Nombrado así en (Richardson & Ruby, RESTful Web Services, 2007)

3.2.1. SOAP AND THE WS-* STACK

Proporcionar interoperabilidad sin fisuras entre los heterogéneos pilas de tecnología de middleware y el fomento de la articulación flexible de servicio al consumidor (solicitante, cliente) y proveedor de servicios son los principales objetivos de diseño de arquitectura orientada a servicios (SOA) conceptos y tecnologías de servicios Web. (Pautasso, Zimmermann, & Leymann, 2008)

En el plano conceptual, un servicio es un componente de software que se proporciona a través de un endpoint²⁷ accesible en la red. Consumidores de servicios y proveedores usan mensajes para intercambiar solicitudes e información de respuesta en forma de *documentos self-containing*²⁸ que hacen muy pocas suposiciones sobre las capacidades tecnológicas del receptor. En particular, no hay noción de una referencia de objeto remoto que requeriría un corredor de objeto para gestionar un espacio distribuido dirección de memoria. En el nivel de la tecnología, SOAP es un lenguaje XML que define una arquitectura de mensajes y formatos de mensaje, por lo tanto, proporcionar un protocolo de procesamiento rudimentario. El documento SOAP define un elemento XML de nivel superior llamada sobre, que contiene un encabezado y un cuerpo. El encabezado SOAP es un contenedor de información de infraestructura extensible de capa de mensajes que se puede utilizar para fines de enrutamiento (por ejemplo, hacer frente) y Calidad de Servicio (QoS) de configuración (por ejemplo, las transacciones, la seguridad, la fiabilidad). El cuerpo contiene la carga útil del mensaje. Esquema XML se usa para describir la estructura del mensaje SOAP, por lo que los motores de jabón en los dos puntos finales pueden Marshall y Resolver referencia el contenido del mensaje y la ruta a la aplicación apropiada. (Pautasso, Zimmermann, & Leymann, 2008)

3.2.2. REST

Transferencia de estado representacional (REST) se introdujo originalmente como un estilo de arquitectura para la construcción de sistemas hipermedia distribuidos a gran escala. Este estilo arquitectónico es una entidad más abstracta, cuyos principios se han utilizado para explicar la excelente escalabilidad del protocolo HTTP 1.0 y también han limitado el diseño de su siguiente versión, HTTP 1.1. Por lo tanto, el término REST muy a menudo se utiliza junto con HTTP. (Pautasso, Zimmermann, & Leymann, 2008)

El estilo arquitectónico REST se basa en cuatro principios:

²⁷ <http://www.w3.org/TR/ws-gloss/#endpoint>

²⁸ <http://www.thefreedictionary.com/self-contained>

Identificación de recursos a través de URI. Un servicio web RESTful expone un conjunto de recursos que identifican los objetivos de la interacción con sus clientes. Los recursos son identificados por URI, que proporcionan un espacio de direccionamiento global de los recursos y de descubrimiento de servicios.

Interfaz uniforme. Los recursos son manipulados utilizando un conjunto fijo de cuatro crear, leer, actualizar, eliminar operaciones: PUT, GET, POST y DELETE. PUT crea un nuevo recurso, que puede ser luego borrar con DELETE. GET recupera el estado actual de un recurso en alguna representación. POST transfiere un nuevo estado sobre un recurso.

Mensajes de auto-descriptivo. Recursos están desconectados de su representación para que su contenido se puede acceder en una variedad de formatos (por ejemplo, HTML, XML, texto plano, PDF, JPEG, etc.) Metadatos sobre el recurso está disponible y se utiliza, por ejemplo, para controlar el almacenamiento en caché, detectar errores de transmisión, negociar el formato de representación adecuada, y llevar a cabo la autenticación o controlar el acceso. Interacciones con estado a través de hipervínculos. Cada interacción con un recurso no tiene estado, es decir, los mensajes de solicitud son autónomos.

Interacciones con estado se basan en el concepto de transferencia de estado explícito. Existen varias técnicas para el intercambio de estado, por ejemplo, la reescritura de URI, cookies, y los campos de formulario ocultos. Estado puede ser embebido en los mensajes de respuesta para señalar válidos estados futuros de la interacción.

3.3. Recursos y representaciones

(Richardson & Amundsen, RESTful Web APIs, 2013) Rest denomina recursos a los datos estructurados que son objetos de las interacciones entre métodos de HTTP, y se dice que cualquier cosa que pueda ser almacenado de un computador puede ser un recurso, como documento electrónico, una fila de una base de datos o el resultado de un algoritmo

No solo las cosas almacenadas en un computador pueden ser llamados recursos también pueden ser recursos artículos tangibles como frutas por ejemplos, y es posible representarlo como recursos a través de la web como por ejemplo como un artículo en de venta o una imagen binaria depende de la aplicación así que por eso decimos sobre las

representaciones que puede ser cualquier documento legible que contenga información acerca de un recurso.

CAPITULO 2: PROBLEMÁTICA

1. Estado actual

La documentación dentro del desarrollo de trabajos educativos universitarios es indispensable para la difusión de los avances y resultados de investigaciones y experimentos realizados en pro del desarrollo de las ciencias en sus distintas áreas. El desarrollo de estos documentos es estructurado pero no deja de ser texto plano el cual se encuentra en algún formato digital.

En la web las publicaciones son almacenadas y presentadas de tal forma que son alcanzados por motores de búsquedas (si es lo que se desea), que si bien realizan un procesamiento del texto esto no permite identificar de forma clara los principales elementos que intervienen.

En este estado el contenido de los documentos no es explotado apropiadamente puesto que en la web existen recursos que son representaciones de estos elementos, a los que se hace referencia en las publicaciones, lo cuales en primera instancia permiten ampliar la información sobre sí mismos y abrir conexiones con otros recursos relacionados con el tema que se trata en la publicación, pero no son directamente referenciados desde la misma.

El procesamiento de estos documentos es un aspecto natural para los seres humanos, que son capaces de identificar los elementos que intervienen, pero se debe tener en cuenta que en la actualidad no solo humanos navegan por la web haciendo necesario desarrollo de herramientas que permitan a las “maquinas” procesar estos contenidos.

2. Justificación

Teniendo en cuenta la tendencia actual de web, la web semántica que se basa en los principios de datos enlazados, los datos toma un factor importante, por lo cual que estos se encuentren “ocultos” dentro del texto no hace posible su enlace e impide la apertura hacia otras fuentes de información.

En el contexto educativo existen esfuerzos para la publicación de Datos Enlazados, uno de los problemas es las fuentes y sus estructura, como lo exponen (Moroch, Piedra, & Valverde, 2012), “La información que se posee no se encuentra en un formato estructurado, y la encontramos en documentos de tipo pdf, hojas de excel, word, texto plano y medios digitales como DVD. Adicionalmente, se completó información faltante, a través de búsquedas manuales en páginas web de organismos de educación de cada país. Al contar con información en diversos formatos y no estructurada, se dificulta

realizar una extracción automática de la información sobre el ámbito que nos interesa, la Legislación de Educación Superior.”

Esto hace necesario medios que permitan extraer y relacionar estos datos dentro de las publicaciones, de acuerdo a los principio de la web semántica que se encuentra en construcción y que iniciativas como esta ayudan a su expansión.

La información escrita es de fácil comprensión para las seres humanos, se puede entender palabras por palabra su significado, aun cuando este puede varias de acurdo al contexto en que se encuentre y a la vez modificando en significado de otras palabras.

Dentro de un texto existen palabras que son más representativas que otras al momento de dar sentido a toda una sentencia o frase, esto puede ocurrir debido a que una palabra o varias palabras, más allá de tener un sentido pueden ser representaciones que entidades existentes en mundo real, como: personas, lugares, eventos, organizaciones etc. o representen entidades abstractas como la Web y diferentes tecnologías existentes, en sí, un texto plano como tal pude estar relacionado con diferentes representaciones de entidades del mundo real alojadas en la web por pedio de las palabras.

Pero el sentido de una frase descansa en todas las palabras siendo unas más representativas que otras como ya se ha visto, aunque no necesariamente estas tengas representaciones en la web, su extracción es indispensable para la comprensión del contenido.

Las publicaciones pueden llegar a contener gran cantidad de texto, pero su estructura exige que posea un campo resumen (etiquetado como *abstract* por su escritura en inglés) del trabajo que abarca la publicación, si bien es cierto el mismo no alcanza para redactar todos los esfuerzos realizados en la publicación pero si da una visión general del tema de la misma y nombra los principales “elementos” que han intervenido en trabajo realizado y que se repiten a lo largo del texto de la publicación.

De esta forma diferentes publicaciones, proyectos educativos y documentos en general poseen descripciones o resúmenes que pueden ser objeto de análisis y cuyos resultados aporten datos que permiten en enlace hacia fuentes de datos abiertos y relacionados que pueden ser aprovechados de mejor manera en la nueva tendencia de la web, la web semántica.

Una de las razones y motivaciones de concentrar los esfuerzos en el ámbito académico, a través del procesamiento de su contenido escrito y el enlace con fuentes libres

promoviendo de esa forma Linked Open Data es debido a la gran cantidad de datos que están siendo publicados en las distintas áreas de la ciencias y siguiendo la nueva tendencia de publicación de forma libre y abierta del conocimiento, como son los Open Educational Resources (OER) que son recursos educativos abiertos y los OpenCourseWare (OCW) que es material docente que se publican de forma abierta sin restricciones. Esto abre oportunidad y necesidad de aplicaciones que faciliten el descubrimiento de estos recursos y la asociación entre sí y con fuentes de datos. Ya existen diversos avances e implantaciones que han demostrado que Linked Data permite crear relaciones en te silos OCW, es posible enriquecer los metadatos de estos; un entorno de datos enlazados de OCW y OER permiten en descubrimiento y la reutilización de estos recursos abiertos (Piedra, Chicaiza, & López, 2014)

Ya se han vendido realizando aportes con el afán de extraer datos desde fuentes académicas en la web y estructurarla en base a los principios de Linked Data, aportes como (Piedra, Tovar, Colomo-Palacios, Lopez-Varga, & Chicaiza, 2014) que busca un “marco de trabajo para evolucionar hacia un sistema más interoperable e integrado para compartir, conectar y el descubrimiento de datos y metadatos de las iniciativas OCW .“

3. Objetivo General.

Implementar servicios web para descubrimiento, enlace y enriquecimiento de Datos Enlazados, aplicado a publicaciones universitarias que permitan el enlace con nuevos datos vinculados a LOD Cloud.

4. Objetivos Específicos

- Implementar Base de Conocimiento que sirva como fuente de información para el descubrimiento de datos.
- Desarrollar Servicios Web para Desambiguación, Enlace, Descubrimiento y Enriquecimiento Datos LOD-Cloud
- Implementar Frontal Piloto para Integración de Servicios Web.

CAPITULO 3: SOLUCIÓN

1. Propuesta

Los datos que se encuentran dentro del texto tanto de las publicaciones como de fuentes en general, se encierran relacionados con otros temas y fuentes de datos, a los que por medio este proyecto se tratara de acceder, extraer, relacionar y enlazar con fuentes de información abierta como lo es DBpedia (que se basa en los principios de los Datos Enlazados), esto permitirá el enriquecimiento del contenido.

Esta propuesta utiliza Servicios Web Rest lo que permite una independencia de la fuente de origen de texto a ser analizado, con la lógica de la aplicación propuesta y esta a su vez devuelve un resultado en que en este caso será formato JSON²⁹ (JavaScript Object Notation - Notación de Objetos de JavaScript) que es ampliamente conocido y utilizado para el intercambio de datos, hacia el cliente que consume el servicio.

La lógica que se propone es explicada en la Figura 10, en donde las interacciones inician con el ingreso del texto a ser analizado, en el cual se aplican las diferentes tecnologías revisadas en capítulos 1 de este documento, para obtener como resultado entidades y keywords estructurados en formato JSON, que estarán desambiguadas y enlazadas, de existir un recurso al cual corresponda dentro del Dataset de DBpedia, es decir que no todas las entidades que se encuentren dentro del texto de una publicación es referenciado en DBpedia.

Puesto que los esfuerzos se concentraran en los *abstracts* de las publicaciones y que esto se redactados en idioma ingles a pesar de que se trate de una publicación en español, se limitara el desarrollo de la solución a este idioma.

Esta propuesta se surge como solución para el descubrimiento de datos en el texto de los *abstracts* de las publicaciones universitarias, pero debido a la gran cantidad de recursos de diversos temas que se encuentran actualmente disponibles en DBpedia, se puede aplicar a cualquier texto (en idioma ingles) para esto se implementa un interfaz gráfica web donde el usuario puede insertar su texto y ver el resultado.

²⁹ <http://json.org/json-es.html>

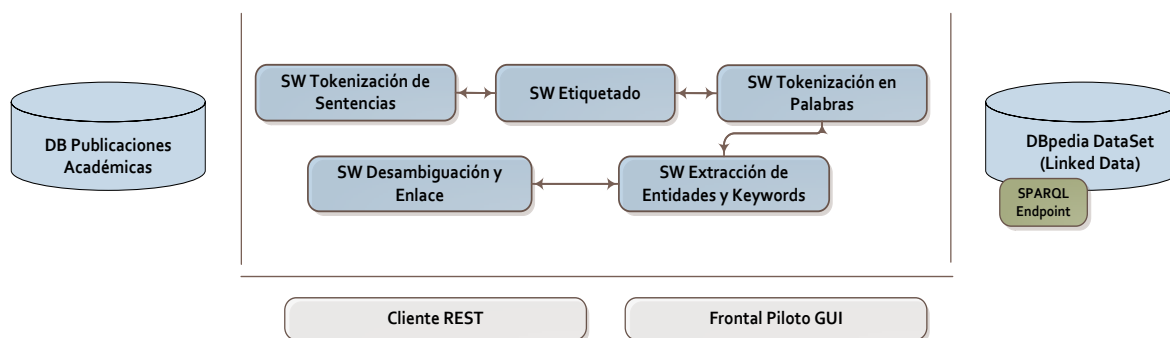


Figura 8. Lógica propuesta para la Aplicación
Fuente: (Propio)

2. Metodología

Para el desarrollo del proyecto de software se propone una metodología de desarrollo en prototipos, con el fin que en el desarrollo del producto de software adicionar funcionalidades al sistema de forma incremental, los cuales se dispondrán en forma de Servicios Web de acuerdo con los objetivos del proyecto.

Por lo cual se decide por el “proceso” de desarrollo ICONIX, que está entre la complejidad de RUP (Rational Unified Processes) y la simplicidad y pragmatismo de XP (Extreme programming), sin eliminar las tareas de análisis y de diseño que XP no completa. (San Martin Oliva)

2.1. Fases de desarrollo

Se establece las siguientes fases de desarrollo de acuerdo con metodología ICONIX, con las respectivas tareas y resultados esperados de las iteraciones de estas fases.

Tabla 10: Fases de desarrollo del proyecto

Fase	Tarea	Resultado
Análisis de requerimientos	<ul style="list-style-type: none"> Identificar objetos del mundo real que interviene en el proceso Identificar actores involucrados Identificar casos de uso del sistema en interacciones con actores identificados 	<ul style="list-style-type: none"> Requerimientos Modelo de dominio Modelo de caso de uso

Diseño preliminar	<ul style="list-style-type: none"> ▪ Describir los casos de uso como un flujo de acciones ▪ Verificar el diseño 	<ul style="list-style-type: none"> ▪ Especificación de casos de uso
Diseño	<ul style="list-style-type: none"> ▪ Especificar comportamiento a través de diagrama de secuencias ▪ Verificar el diseño 	<ul style="list-style-type: none"> ▪ Diagrama de secuencias
Implementación	<ul style="list-style-type: none"> ▪ Escribir código ▪ Realizar pruebas 	<ul style="list-style-type: none"> ▪ Código ▪ Pruebas

Fuente (propio)

3. Desarrollo

3.1. Análisis de requerimientos

A partir de las primeras reuniones se determina los requerimientos funcionales con los que debe contar la propuesta de software a desarrollar y se inicia con la formalización de estos requerimientos y luego con su necesario análisis.

3.1.1. *Requerimientos*

Los requerimientos descritos a continuación definen el comportamiento del sistema para obtener los resultados esperados.³⁰ Estos requerimientos forman la base del desarrollo y se complementa con requerimientos no funcionales que los cuales serán implementados con la finalidad de agregar calidad al producto frente al usuario que lo va a utilizar.

Tabla 11: Resumen de requerimientos funcionales

Código	Requerimiento	Descripción
REQ001	Extraer entidades y palabra relevantes	Descubrir datos relevantes en el texto, a quien se describe y las palabra relevantes que lo acompañan
REQ002	Enlazar entidades y palabra relevantes con LOD Cloud	Se enlazara los términos encontrados en caso de que sea posible con la LOD Cloud

³⁰ Ver en anexos documento de especificación de requerimientos

REQ003	Desambiguar entidades y palabra relevantes	Se determinara el sentido con que las palabras estas siendo usadas en caso de que estas sean ambiguas.
REQ004	Levantar servicios REST separados para los procesos relevantes.	Para que los procesos relevantes dentro del sistema pueda ser consumidos de forma individual y así reutilizados se levantarán servicios separados
REQ005	Frontal UI Web	Construir una interfaz web que permita visualizar el comportamiento del sistema, es decir, la integración de los servicios y su funcionamiento individual

Fuente: (propio)

Los requerimientos redactados establecen las tecnologías necesarias para el desarrollo del producto de software así como las áreas de conocimiento con las que esta ligados siendo el procesamiento de lenguaje natural (PLN) uno de los puntos más fuertes a resolver junto el levantamiento de servicios web y la construcción de un cliente web.

3.1.2. Modelo de Dominio

Uno de los ámbitos más importantes a resolver para el desarrollo de este sistema es el tratamiento de texto que es enviado por un usuario a través de un cliente, que será la entrada y base del procesamiento para descubrir datos relevantes dentro de este.

Se ha procedido a separar en servicios web distintos los procesos de relevancia del producto de software, en vista de que este uno de los requerimientos (REQ004), los servicios web que serán levantados perteneces a fases importantes dentro de las técnicas de procesamiento de lenguaje natural (PLN) aplicadas al texto, así como funciones dadas por los requerimientos, lo cual se observa en la figura 8.

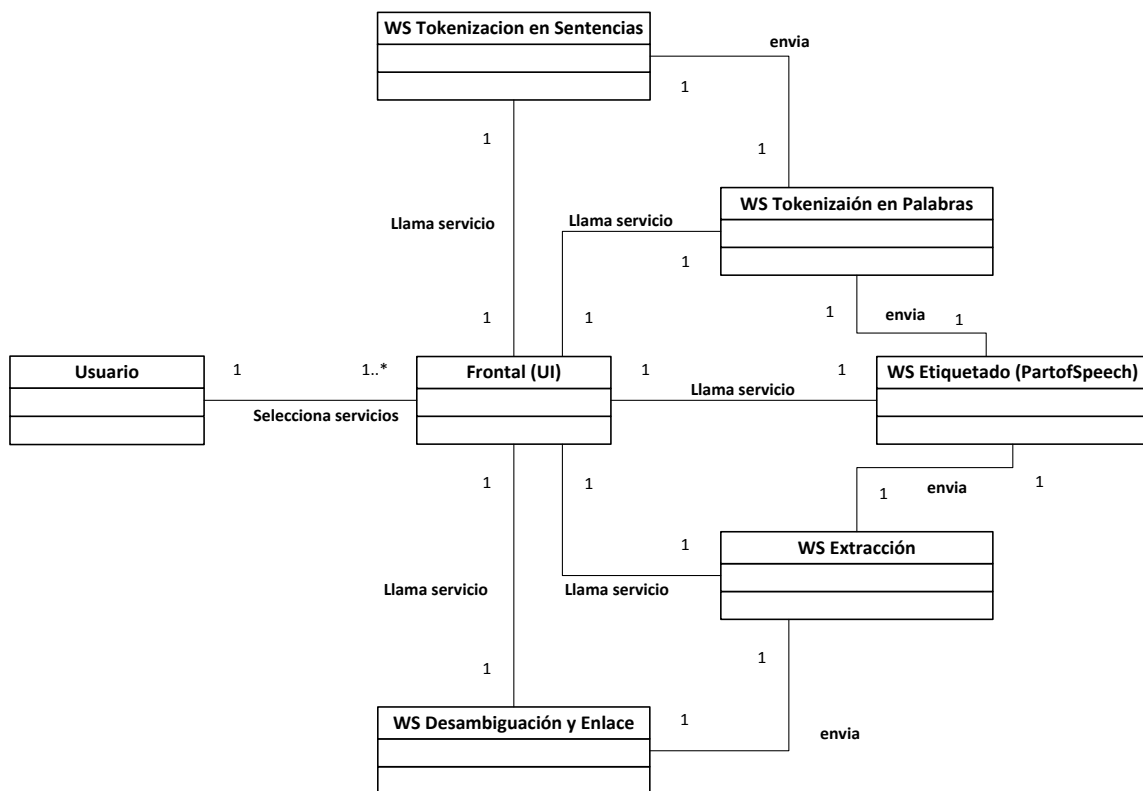


Figura 9: Modelo de dominio
Fuente: (propio)

3.1.3. Modelo de caso de Uso

El comportamiento de los objetos del mundo real, los componentes del sistema a desarrollar y las interacciones que entre estos se pueden realizar son dados por los casos de uso aplicables al sistema dadas las funciones que incorporara basados en el análisis de los requisitos, estas interacciones son visualizados en figura 9.

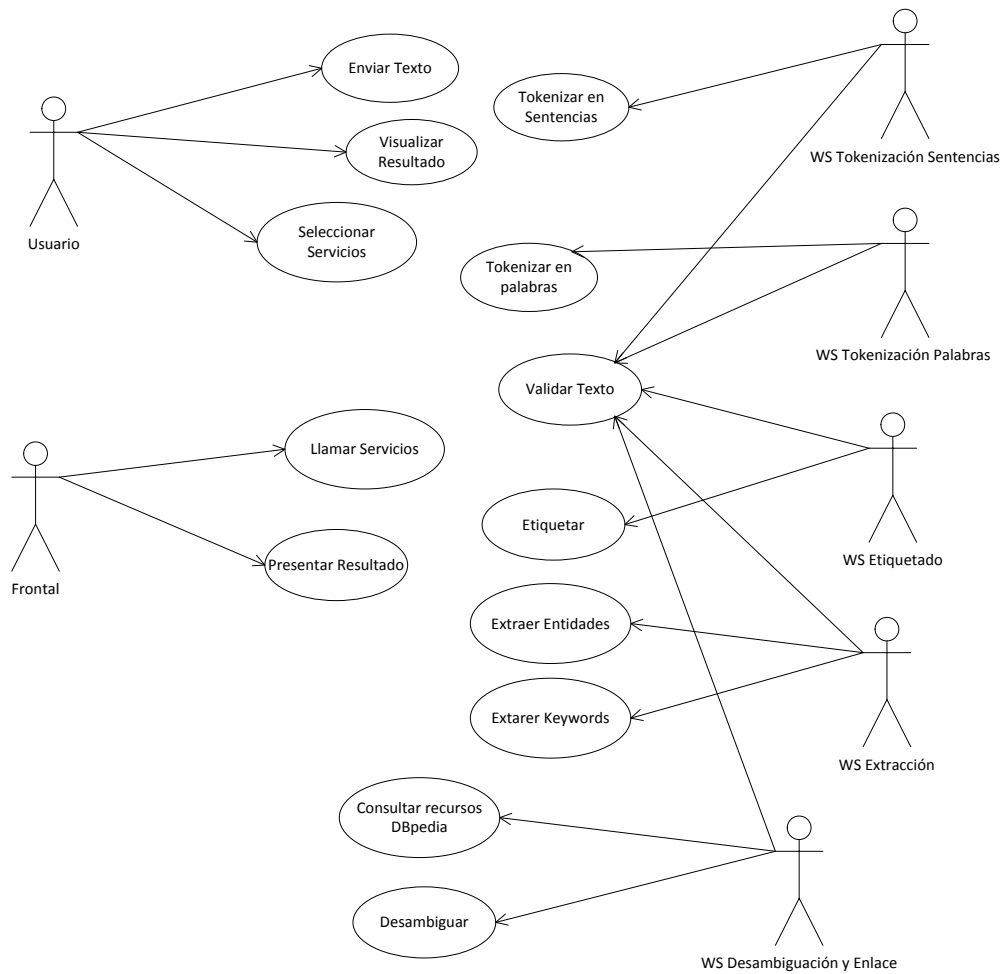


Figura 10: Modelo de casos de uso
Fuente: (propio)

3.2. Análisis y diseño preliminar

3.2.1. Especificación de casos de uso

Una vez determinadas los casos de uso que resolverá el software se presenta a continuación es especificaciones de cada uno, describiendo el comportamiento esperado al realizar una función el software final.

3.2.1.1. Tokenización en Sentencias

Tabla 12: Requerimiento de tokenización de sentencias

Número	ECS-01	
Nombre	Tokenización en Sentencias	
Actores	Usuario, Cliente	
Descripción	Divide el texto de entrada en sentencias cortas separadas por un punto y parte, la salida es una lista de estas sentencias.	
Precondición	<ul style="list-style-type: none"> ▪ Ingreso texto como parámetro de la aplicación ▪ Texto ha sido validado y procesado ▪ Texto segmentado en sentencias 	
Secuencia Normal	Paso	Acción
	1	Entrada del texto validado, como parámetro para el servicio web. SA1
	2	Verifica el número de sentencias que comenten al texto, que estén separadas por un punto seguido (.)
	3	Divide cada una teniendo en cuenta la terminación con punto (.) estructura las sentencias dentro de una lista. SA1
	4	Estructura la lista de elementos formato JSON.
	5	Devuelve el JSON resultante.
Poscondición	<ul style="list-style-type: none"> ▪ El texto dividido en sentencias. 	
Secuencia alternativo	SA1 el número de sentencias es 1 Se estructura una lista de un solo elemento con la sentencia.	
Prioridad	Media	
Requerimientos Especiales	Idioma de texto ingles	
Asunciones y Dependencias		
Notas adicionales		

Fuente: (propio)

3.2.1.2. Tokenización en palabras

Tabla 13. Requerimiento de tokenización en palabras

Número	ECS-02	
Nombre	Tokenización en palabras	
Actores	Cliente, Servicio Web	
Descripción	Divide cada sentencia en palabras validas, tokens.	
Precondición	<ul style="list-style-type: none"> ▪ Ingreso texto como parámetro de la aplicación ▪ Texto ha sido validado y procesado 	
Secuencia Normal	Paso	Acción
	1	Entrada del texto validado, como parámetro para el servicio web.
	2	Segmentación del texto en sentencias. ECS-01
	3	Se recorre la lista de sentencias segmentadas.
	4	Se divide palabra por palabra de la sentencia en una lista, se obtiene una lista de listas.
	5	Se estructura en formato JSON
	6	Retorna el JSON con las sentencias divididas en “tokens”
Poscondición	<ul style="list-style-type: none"> ▪ Texto tokenizado por sentencias y estos a la vez tokenizados en palabras 	
Secuencia alternativo		
Prioridad	Baja	
Requerimientos Especiales	Del funcionamiento del Servicio web de Tokenización en Sentencias	
Asunciones y Dependencias		
Notas adicionales		

Fuente: (propio)

3.2.1.3. Etiquetado de palabra

Tabla 14: Requerimiento de etiquetado

Número	ECS-03	
Nombre	Etiquetado	
Actores	Cliente, Servicio Web	
Descripción	Este servicio permite la tokenización de cada palabra y etiquetación de las mismas de acuerdo a la función que cumplen en el contexto que se encuentra, para hacerlo se apoya en el servicio web de tokenización en sentencias	
Precondición	<ul style="list-style-type: none"> ▪ Ingreso texto como parámetro de la aplicación ▪ Texto ha sido validado y procesado 	
Secuencia Normal	Paso	Acción
	1	Entrada del texto validado, como parámetro para el servicio web.
	2	Tokenización del texto en sentencias. ECS-01
	3	Recorrido de la lista de sentencias
	4	Etiquetado de las palabras que conforman cada sentencia
	5	Estructura y retorna data en JSON
Poscondición	<ul style="list-style-type: none"> ▪ Texto tokenizado a nivel de palabras y etiquetado. 	
Secuencia alternativo		
Prioridad	Alta	
Requerimientos Especiales		
Asunciones y Dependencias		
Notas adicionales	Depende del funcionamiento del servicio web de Etiquetado en Sentencias (ECS-01)	

Fuente:(propio)

3.2.1.4. Extracción de Entidades

Tabla 15. Requerimiento de extracción de entidades.

Número	ECS-04	
Nombre	Extracción de Entidades	
Actores	Cliente, Servicio Web	
Descripción	Permite reconocer y extraer, las entidades y palabras relevantes o claves (keywords) que se encuentran dentro del texto, para lograr se apoya en el servicio web de Etiquetado (y en los que este a su vez , servicio web de tokenización en sentencias)	
Precondición	<ul style="list-style-type: none"> ▪ Ingreso texto como parámetro de la aplicación ▪ Texto ha sido validado y procesado 	
Secuencia Normal	Paso	Acción
	1	Entrada del texto validado, como parámetro para el servicio web.
	2	Tokenización del texto en sentencias. ECS-01
	3	Tokenización y Etiquetado de palabra ECS-03
	4	Reconocimiento de estructuras de Entidades y Keywords
	5	Extracción de Entidades y Keywords
	6	Estructuración de retorno de resultado en formato JSON
Poscondición	<ul style="list-style-type: none"> ▪ Entidades y palabra importantes que las acompañan reconocidos y extraídos 	
Secuencia alternativo		
Prioridad	Alta	
Requerimientos Especiales		
Asunciones y Dependencias		

Notas adicionales	Este servicio depende del funcionamiento del servicio web de Tokenización en Entidades (ECS-01) y Servicio web de Etiquetado (ECS-03)
--------------------------	---

Fuente: (propio)

3.2.1.5. Desambiguación y Enlace

Tabla 16. Especificación del requerimiento de desambiguación y enlace

Número	ECS-05	
Nombre	Desambiguación y Enlace	
Actores	Cliente, Servicio Web	
Descripción	Enlaza las entidades y palabras relevantes (keywords) hacia LOD Cloud, más específicamente DBpedia, esto de existir un recurso al cual vincular, en caso de que una entidad o keyword tuviese más de uno posible recurso al cual enlazar, se realizara un proceso de desambiguación y luego de enlace.	
Precondición	<ul style="list-style-type: none"> ▪ Ingreso texto como parámetro de la aplicación ▪ Texto ha sido validado y procesado 	
Secuencia Normal	Paso	Acción
	1	Entrada del texto validado, como parámetro para el servicio web.
	2	Tokenización del texto en sentencias. ECS-01
	3	Tokenización y Etiquetado de palabra. ECS-03
	4	Extracción de Entidades y keywords. ECS-04
	5	Consulta de recursos a DBpedia.
	6	Consulta de “Abstract” de recurso a DBpedia
	7	Verificar si existen Entidades o keywords ambiguas

	8	Desambiguar Entidades y keywords ambiguos. SA1
	9	Estructurara resultado
	10	Retornar resultado
Poscondición	▪ Entidades y palabra importantes que las acompañan reconocidos y extraídos	
Secuencia alternativo	SA1 Entidades y keywords no ambiguos Se enlaza con los recursos únicos encontrados a las entidades y keywords del texto.	
Prioridad	Alta	
Requerimientos Especiales		
Notas adicionales	Este servicio depende de los servicios web de tokenización en sentencias (ECS-01), etiquetado (ECS-03), extracción de entidades (ECS-04).	

Fuente: (propio)

3.3. Diseño

3.3.1. Arquitectura

Después de haber realizado el análisis y posterior investigación sobre los puntos más relevantes dentro del desarrollo, se propone la siguiente propuesta lógica para el funcionamiento de los componentes a desarrollar y en base a los requerimientos funcionales, casos uso a los que va a responder el software y tecnologías disponibles que se pueden implementar para la resolución de la problemática planteada.

La distribución de los compones implementados se resume en la gráfica 10, haciendo una división de los componentes de la aplicación que interactúan para dar solución a la problemática. Las capas que componen la arquitectura de aplicación se detallan a continuación.

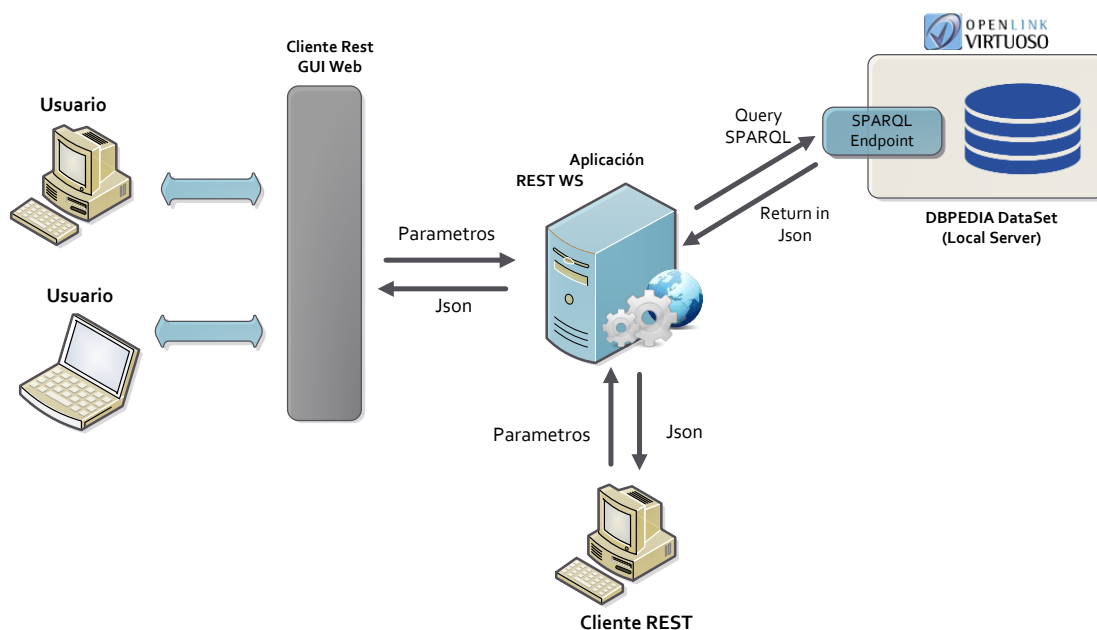


Figura 11. Arquitectura.
Fuente: (Propio)

3.3.2. Componentes

3.3.2.1. Servidor

La lógica de las aplicaciones se la desarrollará en el Lenguaje de alto nivel Python versión 2.7 la elección se ha basado en la familiaridad con la librería, *Natural Language Toolkit* (NLTK), que permite realizar procesamiento de lenguaje natural (PNL) que es una parte fundamental a resolver y de igual forma cuenta con las librerías para satisfacer los requerimientos con los que debe cumplir el software. Las funcionalidades destacadas dentro del desarrollo se resumen a continuación.

Validar texto

Que el texto sea legible para el sistema en funciones siguientes a fin de evitar errores. Comprobar que el texto contenga caracteres a fin de evitar trabajar sobre texto vacío. Es origen obligatorio por el cual los todos los módulos deben pasar para la comprobación de los requisitos del texto que ingresa como parámetro para ser procesado.

Tokenización en sentencias

Dividir el texto ingresado en sentencias u oraciones, generalmente separados por un punto (.) significando el final de esta. Este componente forma parte de los procesos relevantes dentro de procesamiento de lenguaje natural siendo el inicio de estos procesos.

Tokenización en palabras

Una vez divididos el texto en sentencias, se realiza en mismo proceso para las palabras que lo conforman así como los signos de computación, obteniendo un “token” por cada palabra o carácter reconocido. Parte del procesamiento de lenguaje natural.

Etiquetado

Este componente realiza etiquetado de las palabras, de acuerdo al contexto dentro de la sentencia en la que se encuentra, tomando una función específica como verbo, sustantivo, etc. la etiquetación de palabras (Part of Speech, por su nombre en inglés) forma parte del procesamiento de lenguaje natural.

Extracción de entidades y palabras claves

Permite el reconociendo y extracción de entidades y palabra claves dentro del contexto de una oración, para lo cual realiza un reconocimiento de los etiquetas de las palabras y un análisis de su estructura para determinar que palabra es una entidad o una palabra relevante, este análisis puede extraer una palabra o un conjunto de estas formando una entidad o palabras claves. Se basa en componentes posteriores para poder realizar sus operaciones, es decir, necesita de un texto que haya sido tokenizado y etiquetado.

Extracción recursos DBpedia

Una vez que los entidades y palabra claves hayan sido extraídos de las sentencias. Se realizan consultas con estos términos hacia el servidor local de DBpedia, para encontrar los recursos que se denominen igual que estos términos, así como las descripciones rápidas de estos recursos.

Desambiguación de recursos

A través del algoritmo de **Lesk** se analiza el contexto del término y las descripciones de los recursos de DBpedia, para determinar que cual de estos concuerda mejor con el uso que se le está dando al término en la sentencia.

Servicios web

De conformidad con los requerimientos funcionales se levantan servicios diferentes para algunos de los proceso relevantes, permitiendo que puedan ser consumidos y analizados de forma individual. Los servicios levantados serán:

- Tokenización en sentencias
- Tokenización en palabras
- Etiquetado
- Extracción de entidades
- Desambiguación y enlace

Existe una fuerte dependencia entre los servicios, puesto que la salida de unos se convierte en la entrada de otros en forma de secuencia a través de los servicios, el servicio de Tokenización en sentencias, después de proceso de validación, es el primero en trabajar sobre el texto hasta llegar al servicios de Desambiguación y enlace. La gráfica 12, presenta la forma en que los servicios se relacionan e interactúan entre sí.

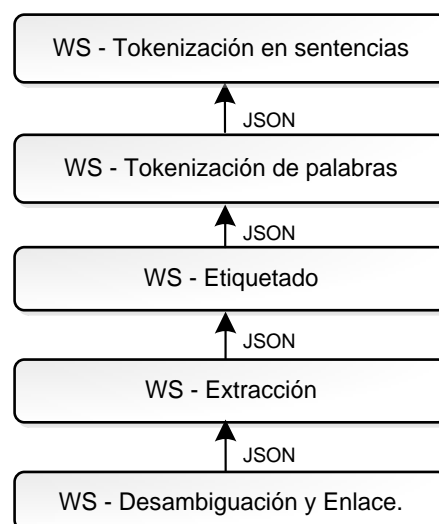


Figura 12. Dependencia de servicios web
Fuente: (propio)

3.3.2.2. Servidor Dataset DBpedia Local (SPARQL endpoint)

Parte de la propuesta consiste en acceder a los recursos de DBpedia por lo cual se implementará un servidor local con los datos, de estos recursos, necesarios para desarrollar la propuesta y así evitar cualquier fallo por problemas de conexión recurrentes al tratar de consultar directamente con su servidor, esto es posible gracias a que DBpedia misma proporciona los medios para descargar sus recursos.

3.3.2.3. Cliente

La vista es un Interfaz Web que permite al usuario la facilidad de la comunicación entre la aplicación y el usuario. Instruye al usuario sobre el uso de la herramienta, identificando con facilidad los parámetros necesarios y en especial la facilidad de la presentación de los resultados. Permite la integración de los componentes del sistema en un entorno amigable para el usuario.

Prototipo de interfaz

3.3.3. Diagrama de secuencia

Al observar el modelado de las interacciones entre los componentes del sistema se evidencia una clara dependencia e interoperabilidad entre los estos, iniciando siempre con el módulo validación del texto a ser procesado.

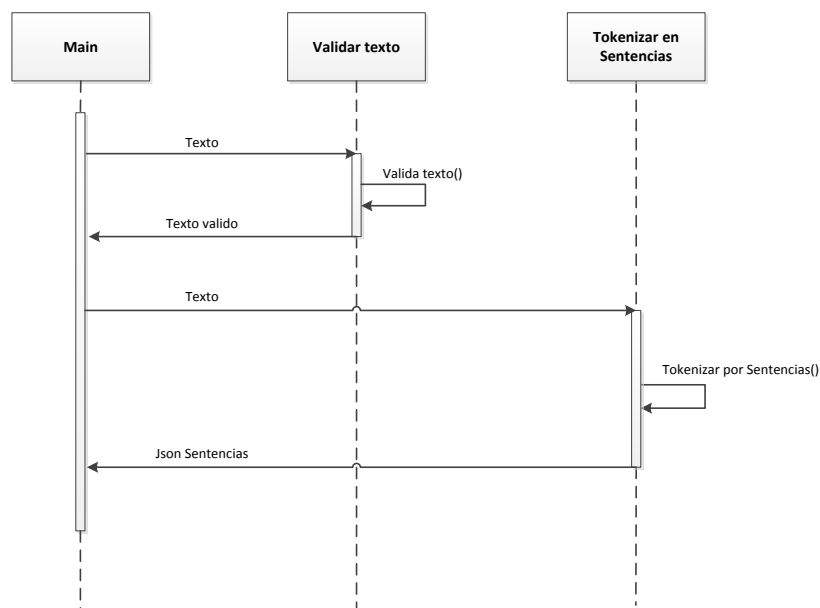


Figura 13. Diagrama de secuencias de tokenización de sentencias
Fuente: (propio)

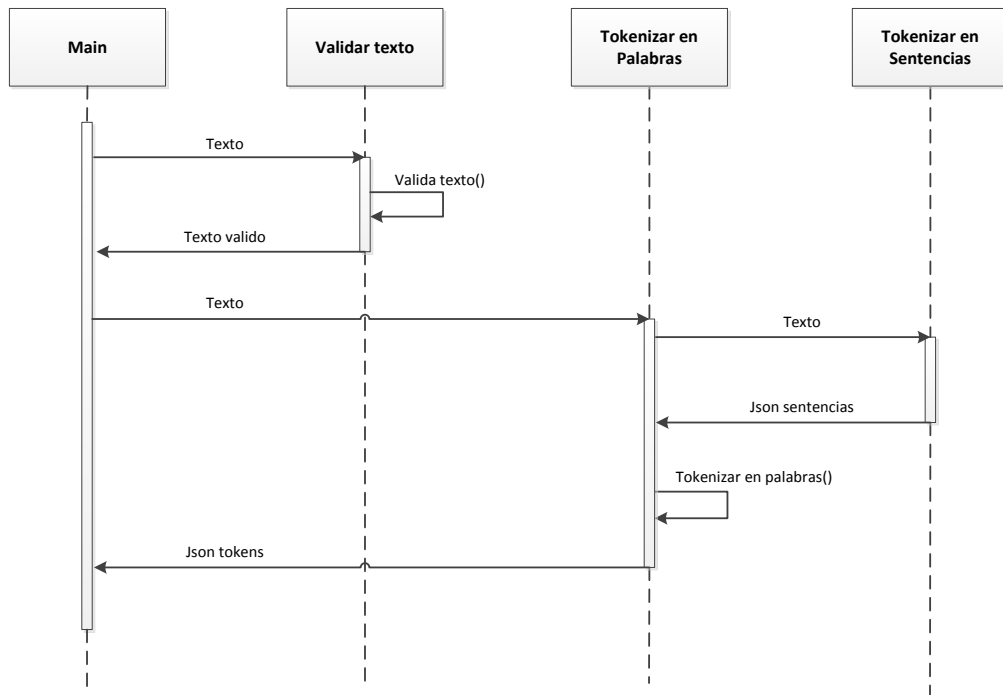


Figura 14. Diagrama de secuencias de tokenización en palabras
Fuente: (propio)

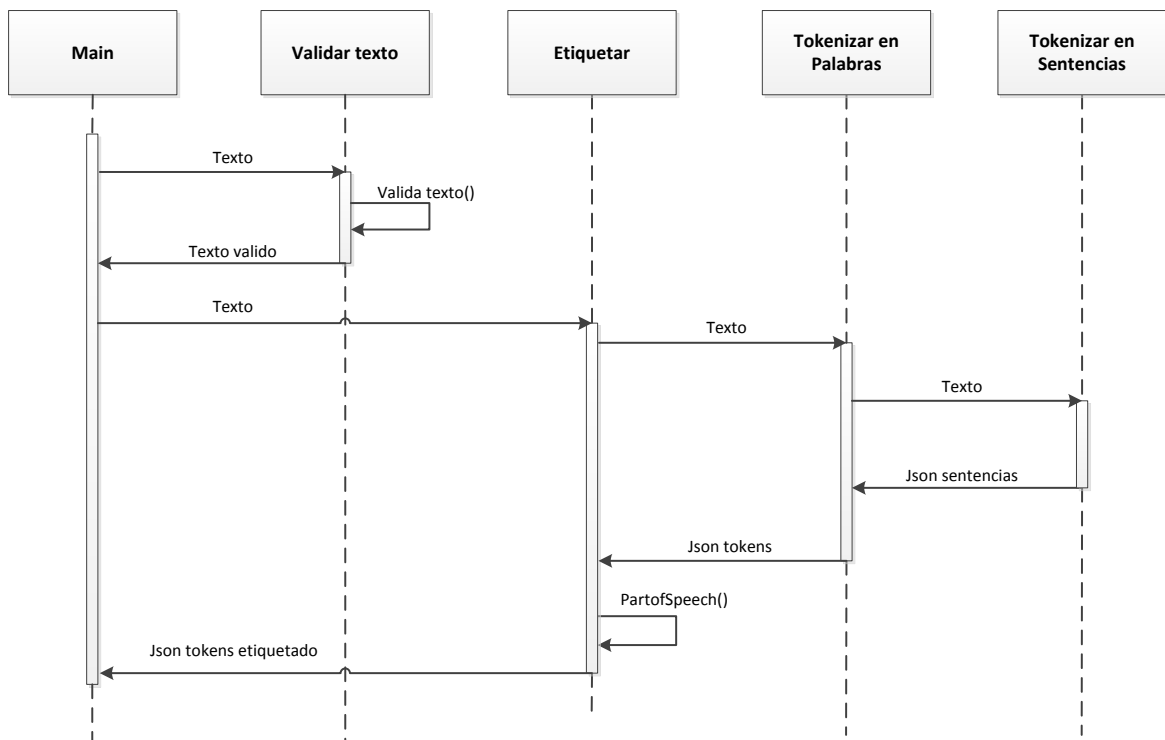


Figura 15. Diagrama de secuencias de etiquetado de palabra
Fuente: (propio)

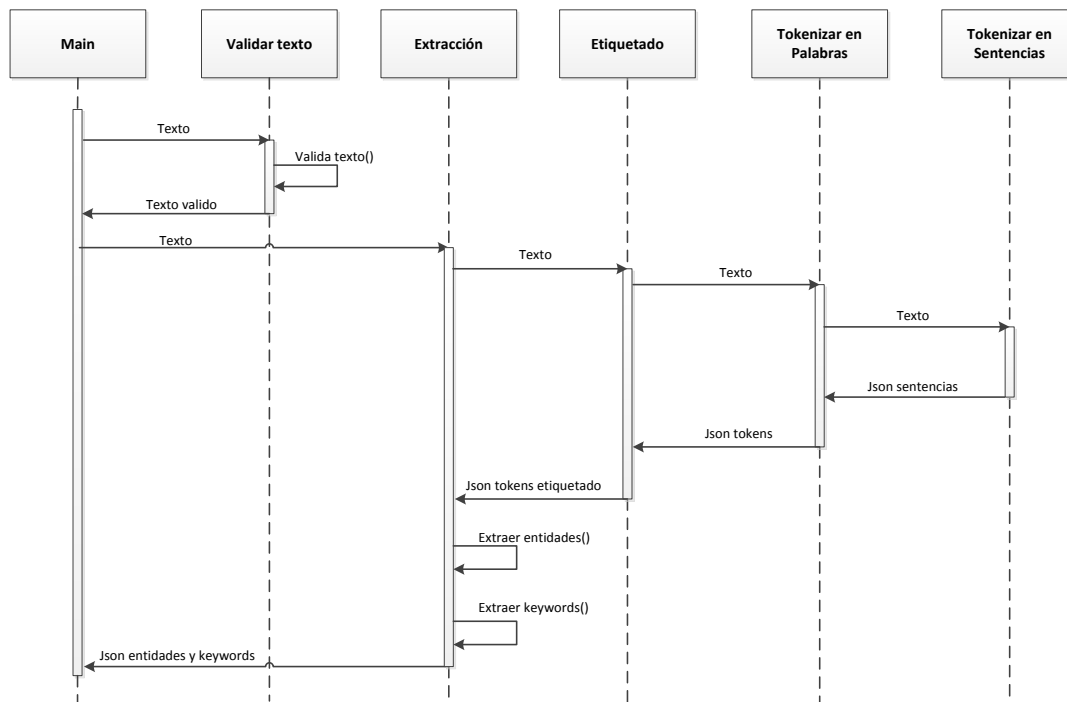


Figura 16. Diagrama de secuencias de extracción
Fuente: (propio)

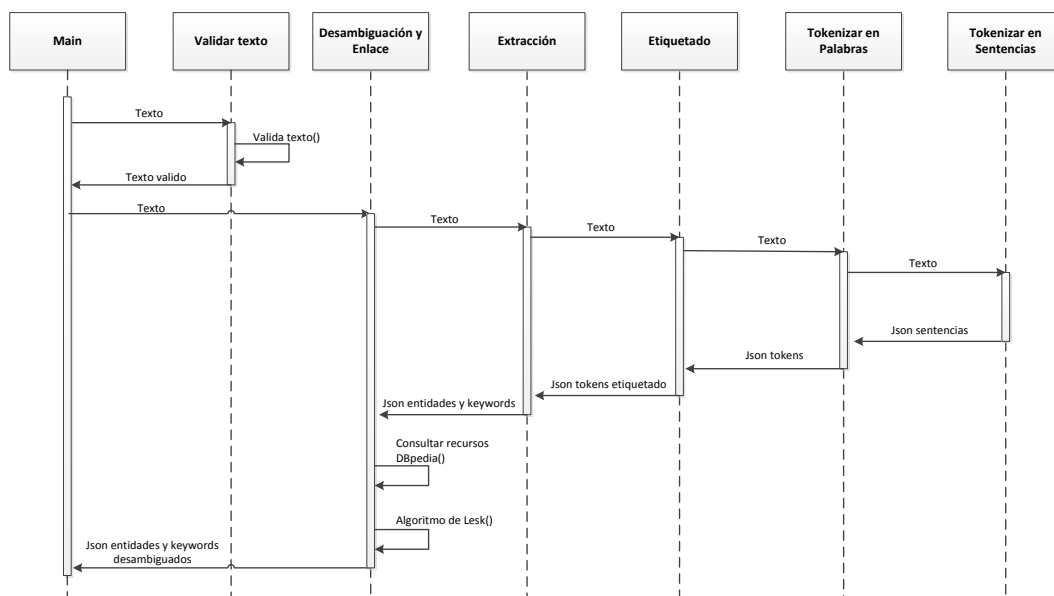


Figura 17. Diagrama de secuencia de desambiguación y enlace
Fuente: (propio)

3.4. Implementación

3.4.1. Servidor

Producto del análisis realizado en el diseño de la software, Python es el lenguaje de programación en el que se desarrollan los componentes de la lógica de la aplicación, gracias a las facilidades que proporciona la librería NLTK en el procesamiento del lenguaje natural, pero debido a que esta herramienta presenta limitaciones al trabajar con idiomas distinto al inglés, las tareas relacionadas con POS-Tagging se las realiza con TreeTagger permitiendo la posibilidad de trabajos futuros relacionados con otros idiomas distintos a inglés.

El código generado para el desarrollo del servidor del producto de software se encuentra en anexos.

Para extraer todos los recursos de DBpedia que a través de uso de sus atributos sea referenciado como el término que se ha encontrado del texto se realiza la consulta SPARQL capturada en la figura 18.

```
SELECT distinct ?x ?amb ?redir ?amb1
WHERE {
  ?x ?predicado '""'+term+""'@en.
  OPTIONAL { ?x dbpedia-owl:wikiPageDisambiguates ?amb }
  OPTIONAL { ?x dbpedia-owl:wikiPageRedirects ?redir }
  OPTIONAL { ?redir dbpedia-owl:wikiPageDisambiguates ?amb1 }
}
```

Figura 18. Captura de consulta de recursos de DBpedia
Fuente: (propio)

De los recursos de DBpedia resultantes, se requiere descripciones de estos para los procesos de desambiguación, para esto a todos los recursos se les realiza la consulta de la figura 19.

```
select ?abstract
where
{
  <""'+uri+""> <http://www.w3.org/2000/01/rdf-schema#comment> ?abstract.
  FILTER (lang (?abstract)="en")
}
```

Figura 19. Consulta de abstracts de recursos de DBpedia
Fuente: (propio)

Para determinar el tipo de recurso según la ontología de DBpedia que se ha encontrado se realiza el siguiente consulta:

```

select distinct *
where {
  <"""+term+"""> rdf:type ?typeself.
}

```

Figura 20. Consulta para extraer el tipo de recurso
Fuente: (propio)

3.4.1.1. Servicios

Los servicios levantados en base a los requerimientos funcionales del producto de software, corresponden a las funcionalidades del sistema disponibles por separado pero integrados entre sí para satisfacer las solicitudes por parte del usuario, quien va a poder elegir a que servicio quiere acceder.

Para la construcción de los servicios web se siguieron los principios de arquitectura REST, separando de esta forma las operaciones del servidor con las del cliente, los cuales se podrían construir utilizando tecnologías distintas.

La respuesta de los servicios se encuentra en formato JSON, la estructura del objeto cambia de acuerdo al servicio llamado, mediante un ejemplo en la figura 21 se muestra el resultado del servicio **Tokenización en Sentencias** mientras que le figura 22, el resultado de servicio **Desambiguación y Enlace**, en la figura 12 se da una idea clara por qué la estructura del objeto JSON aumenta notablemente conforme los servicios que intervienen para dar respuesta a la solicitud del usuario.

```
{
  "result":{
    "NumSentencias":3,
    "TokensSentencias":[3]
  }
}
```

Figura 21. Resultado de servicio de tokenizacion en sentencias
Fuente: (propio)

```
{
  "result":{
    "Entidades":[8],
    "EntidadesDesambiguadas":[10],
    "EtiquetadoPalabras":[3],
    "KeywordsCompuestas":[2],
    "KeywordsSimples":[9],
    "NumEntidades":8,
    "NumEntidadesDesambiguadas":10,
    "NumKeywordsCompuestas":2,
    "NumKeywordsSimples":9,
    "NumSentencias":3,
    "NumTokensPalabras":46,
    "TokensPalabras":[3],
    "TokensSentencias":[3]
  }
}
```

Figura 22. Resultado del servicio web de desambiguación y enlace
Fuente: (propio)

Cada servicio adiciona una nueva propiedad al resultado, en la tabla 17 se presenta una descripción de las propiedades que cada servicio suma al JSON resultante.

Tabla 17. Propiedades del JSON resultado de los servicios web

Servicios Web	Propiedad adicionada	Descripción propiedad
Tokenización en Sentencias	"TokensSentencias"	Contiene en una lista las sentencias en las que divide en texto.
	"NumSentencias"	En número de sentencias encontradas

Tokenización en Palabras	"TokensPalabras"	Las palabras y signos de puntuación que componen en texto, denominado tokens
	"NumTokensPalabras"	La cantidad de tokens encontradas
Etiquetado	"EtiquetadoPalabras"	Los token con las etiquetas de acuerdo a la función que realizan en la sentencias, del mismo número de "NumTokensPalabras"
Extracción	"NumKeywordsSimples"	La cantidad de palabras claves extraídas del texto
	"NumKeywordsCompuestas"	El número de palabras compuestas
	"KeywordsCompuestas"	Contiene las palabras claves que se descubrieron en el texto.
	"KeywordsSimples"	Las palabras claves simples del texto.
	"Entidades"	Las entidades extraídas del texto.
	"NumEntidades"	La cantidad de entidades extraídas.
Desambiguación y Enlace	"EntidadesDesambiguadas"	Los términos enlazados a DBpedia
	"NumEntidadesDesambiguadas"	Numero de términos Enlazados

Fuente: (Propio)

3.4.2. Servidor Dataset DBpedia Local

Archivos necesarios para la implantación de un repositorio local de recursos DBpedia para funcionalidad de desambiguación y enlace.

- Label de recursos:
 - labels_en.nt.bz2

- Datos personales de los recursos tipos Persona:
 - persondata_en.nt.bz2
- Resúmenes Corto de los recursos
 - short_abstracts_en.nt.bz2
- Links de Desambiguación de Wikipedia
 - disambiguations_en.nt.bz2
- Redirecciones entre Recursos
 - redirects_en.nt.bz2

El tamaño total de las importaciones es de 3.8 GB, de espacio en disco.

3.4.3. Cliente web

Presenta el resultado transparente para el usuario, cumple con el requerimiento de integrar los servicios levantados, además permite la selección de los servicios a los que se desea acceder de forma individual, respetando las dependencias que establecidas entre ellos para su funcionamiento.

Desarrollado en HTML, JavaScript y CSS interpreta el JSON recibido por parte de los servicios o el servicio invocado y procesa para que sea agradable al usuario además permite visualizar el JSON tal como se lo recibe desde el servicio, esto para usuarios interesados puedan analizar el resultado. Para introducir el texto a ser procesado se dispone a de un área donde colocarlo visible y amigable para el usuario.

Para seleccionar el servicio al que se requiere acceder existe un menú que los expone como funcionalidades del sistema, en este menú se encuentran:

- Segmentación en sentencias
- Tokenización
- Etiquetado
- Extracción de Etidades
- Desambiguación y Enlace

Las cuales puede ser seleccionados o deseleccionados por el usuario, una captura de la interfaz se puede visualizar en la figura 23.

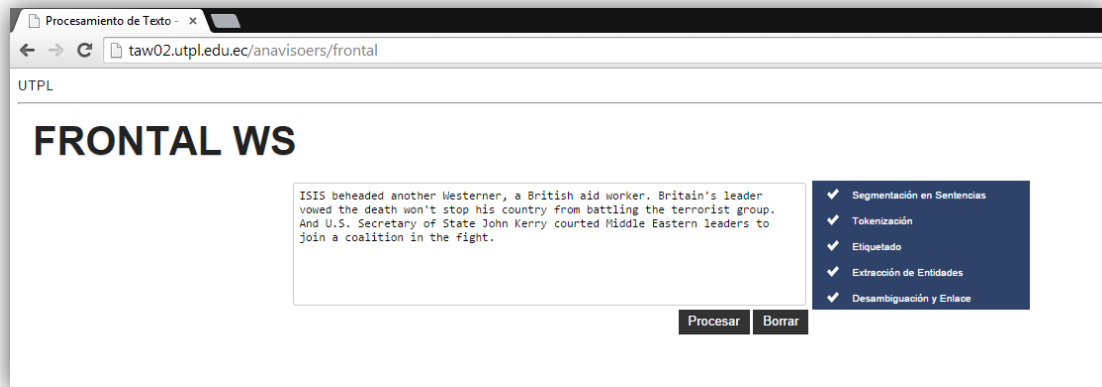


Figura 23. Captura de la interfaz es su estado inicial.

Fuente: (propio)

Cada uno de los ítems del menú de funcionalidad accede a un servicio diferente, pero en vista de que los servicios dependen entre ellos para su funcionamiento, cuando un usuario seleccionara un servicio se seleccionaran automáticamente las funcionalidad que acceden a los servicios de los cuales depende para realizar su funcionamiento, por ejemplo, en determinado momento todas las funcionalidades deseleccionadas y el usuario decide seleccionar la funcionalidad de **Etiquetado** (este caso esta capturado en la figura 24), en este caso de forma automática se seleccionará las funciones de **Segmentación en Sentencias** y **Tokenización** que acceden a los servicios necesario para poder realizar al etiquetado (el resultado de esta interacción se presenta en la figura 25).

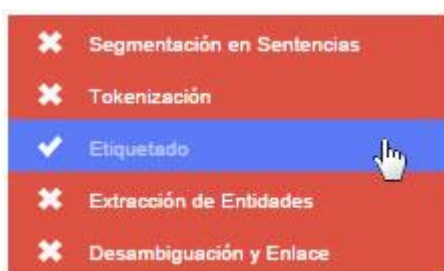


Figura 24. Momento previo a la selección de la funcionalidad de etiquetado

Fuente: (propio)

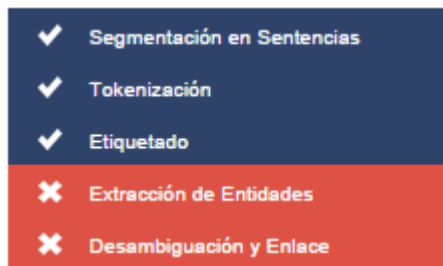


Figura 25. Función etiquetado seleccionado junto a las funcionalidades dependientes
Fuente: (propio)

De igual manera al deseleccionar un servicio al cual no se desea acceder se deseleccionaran las funcionalidades que no son necesarios para la resolución de esta petición, un ejemplo de esto lo se puede observar en la figura 26 donde se tiene todos los servicios seleccionados y se decide deseleccionar la funcionalidad de Tokenización y el resultado se observa en la figura 27.

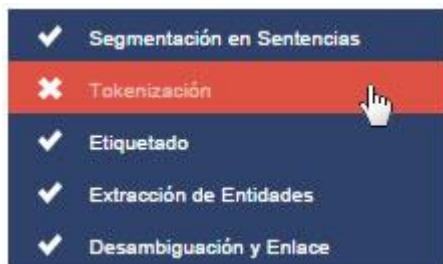


Figura 26. Momento previo a la deselección de la funcionalidad de tokenización
Fuente: (propio)

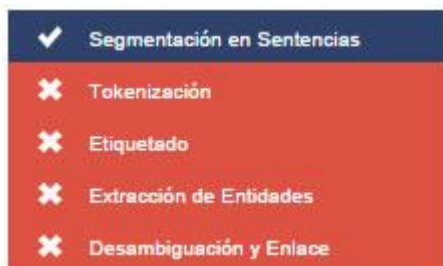


Figura 27. Función de tokenización deseleccionada junto las funciones que depende de esta.
Fuente: (propio)

Para iniciar el procesamiento se envía el texto base al presionar el botón **Procesar** disponible en la interfaz y se espera mientras se devuelve y procesa una respuesta. La respuesta es interpretado y se coloca en par te inferior de la interfaz, esto se puede visualizar en la figura 28.

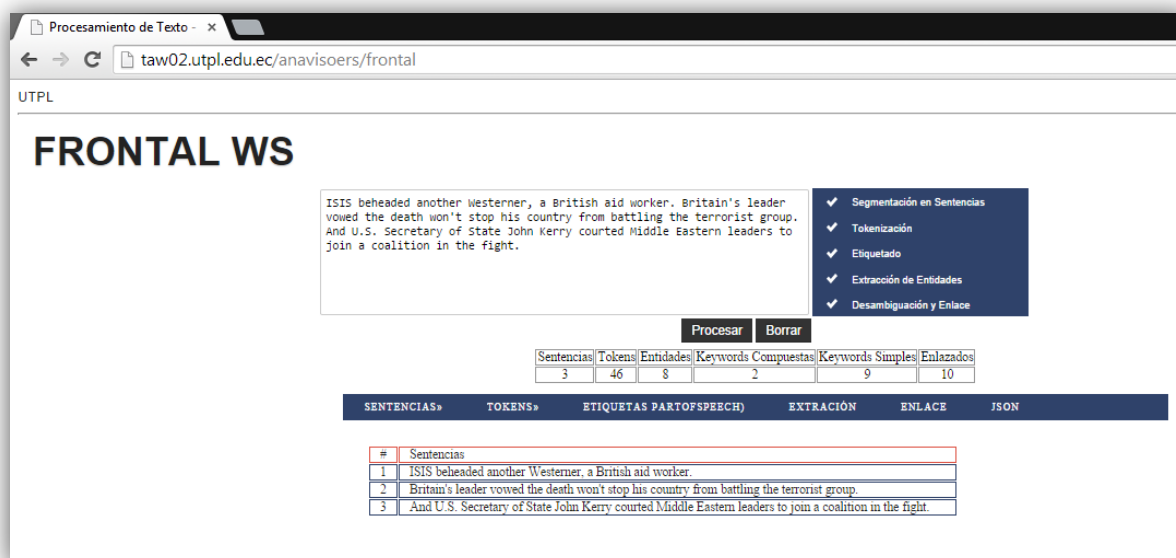


Figura 28. resultado del procesamiento del texto

Fuente: (propio)

En forma de resumen un conteo de lo encontrado en el texto mediante los procesos realizados sobre el texto y dependiendo de los servicios invocados se presenta en forma de tabla, la cual se encuentra capturada en la figura 29.

Sentencias	Tokens	Entidades	Keywords Compuestas	Keywords Simples	Enlazados
3	46	8	2	9	10

Figura 29. Captura de la tabla con datos cuantitativos de los servicios invocados.

Fuente: (propio)

Con la finalidad de una fácil interacción entre el usuario y la interfaz se presenta un menú correspondiente a cada funcionalidad del sistema invocado, esto se visualiza en la figura 30, este dependerá de los servicios accedidos para su construcción, es decir, cada ítem dentro de este menú corresponde a cada una de respuestas de los servicio, interpretados por el cliente y presentados por separado.

SENTENCIAS»	TOKENS»	ETIQUETAS PARTOFSPEECH)	EXTRACCIÓN	ENLACE	JSON
-------------	---------	-------------------------	------------	--------	------

Figura 30. Menú construido con todos los servicios

Fuente: (propio)

Al igual que el menú con las funcionalidades, este menú de resultado se encuentra en orden de forma que se puede evidenciar cómo evoluciona el tratamiento del texto a través de los servicios invocados. En la figura 31 se puede visualizar las sentencias separadas resultado del servicio de Tokenización en Sentencias, las cuales se encuentran numeradas.

[SENTENCIAS»]
TOKENS»
ETIQUETAS PARTOFSPEECH)
EXTRACCIÓN
ENLACE
JSON

#	Sentencias
1	ISIS beheaded another Westerner, a British aid worker.
2	Britain's leader vowed the death won't stop his country from battling the terrorist group.
3	And U.S. Secretary of State John Kerry courted Middle Eastern leaders to join a coalition in the fight.

Figura 31. Resultado de la función de tokenización
Fuente: (propio)

Separadas igualmente por sentencias y numeradas se encuentran los resultados de la Tokenización en la figura 32 y del Etiquetado, en las figura 33, donde se diferencias las funciones de las palabras dentro de las sentencias base para el servicio de Extracción.

SENTENCIAS »		[TOKENS »]	ETIQUETAS PARTOFSPEECH)	EXTRACCIÓN	ENLACE	JSON
SENTENCIA #1						
#	Token	#	Token	#	Token	
1	ISIS	2	beheaded	3	another	
4	Westerner	5	,	6	a	
7	British	8	aid	9	worker	
10	.					
SENTENCIA #2						
#	Token	#	Token	#	Token	
11	Britain	12	's	13	leader	
14	vowed	15	the	16	death	
17	wo	18	n't	19	stop	
20	his	21	country	22	from	
23	battling	24	the	25	terrorist	
26	group	27	.			
SENTENCIA #3						
#	Token	#	Token	#	Token	
28	And	29	U.S.	30	Secretary	
31	of	32	State	33	John	
34	Kerry	35	courted	36	Middle	
37	Eastern	38	leaders	39	to	
40	join	41	a	42	coalition	
43	in	44	the	45	fight	
46	.					

Figura 32. Resultado de la funcionalidad de tokenización.
Fuente: (propio)

SENTENCIAS»			TOKENS»			[ETIQUETAS PARTOFSPEECH]			EXTRACCIÓN			ENLACE			JSON		
-------------	--	--	---------	--	--	----------------------------	--	--	------------	--	--	--------	--	--	------	--	--

SENTENCIA #1

#	Token	Etiqueta	#	Token	Etiqueta	#	Token	Etiqueta
1	ISIS	NP	2	beheaded	VVD	3	another	DT
4	Westerner	NP	5	.	.	6	a	DT
7	British	JJ	8	aid	NN	9	worker	NN
10	.	SENT						

SENTENCIA #2

#	Token	Etiqueta	#	Token	Etiqueta	#	Token	Etiqueta
11	Britain	NP	12	's	POS	13	leader	NN
14	vowed	VVD	15	the	DT	16	death	NN
17	wo	MD	18	n't	RB	19	stop	VV
20	his	PP\$	21	country	NN	22	from	IN
23	battling	VVG	24	the	DT	25	terrorist	JJ
26	group	NN	27	.	SENT			

SENTENCIA #3

#	Token	Etiqueta	#	Token	Etiqueta	#	Token	Etiqueta
28	And	CC	29	U.S.	NP	30	Secretary	NP
31	of	IN	32	State	NP	33	John	NP
34	Kerry	NP	35	courted	VVD	36	Middle	NP
37	Eastern	NP	38	leaders	NNS	39	to	TO
40	join	VV	41	a	DT	42	coalition	NN
43	in	IN	44	the	DT	45	fight	NN
46	.	SENT						

Figura 33. Resultado de la funcionalidad de etiquetado.
Fuente: (propio)

Los resultados del servicio de extracción se presentan numeradas y no separados por sentencias, divididos en entidades keywords simple y keywords compuestas (keyword, en español palabra claves). Como lo visualización en la figura 34.

SENTENCIAS»			TOKENS»			ETIQUETAS PARTOFSPEECH)			[EXTRACCIÓN]			ENLACE			JSON		
-------------	--	--	---------	--	--	-------------------------	--	--	----------------	--	--	--------	--	--	------	--	--

ENTIDADES

#	Entidades	#	Entidades	#	Entidades
1	ISIS	2	Westerner	3	Britain
4	U.S. Secretary of State John Kerry	5	U.S. Secretary	6	State John Kerry
7	Middle Eastern leaders	8	Middle Eastern		

KEYWORDS COMPUESTAS

#	Keywords Compuestas	#	Keywords Compuestas	#	Keywords Compuestas
1	British aid worker	2	terrorist group		

KEYWORDS SIMPLES

#	Keywords Simples	#	Keywords Simples	#	Keywords Simples
1	aid	2	worker	3	leader
4	death	5	country	6	group
7	leaders	8	coalition	9	fight

Figura 34. Resultado del servicio de extracción.
Fuente: (propio)

Los enlaces a los recursos de DBpedia enlazados, después de ser desambiguados se muestran en una tabla como se observan en la figura 35.

SENTENCIAS» TOKENS» ETIQUETAS PARTOFSPEECH) EXTRACCIÓN [ENLACE] JSON			
ENTIDADES			
#	Entidad	Tipo	Enlace
1	ISIS		DBpedia
2	Westerner	http://www.w3.org/2002/07/owl#Thing http://dbpedia.org/ontology/TelevisionShow http://dbpedia.org/ontology/Work http://schema.org/CreativeWork	DBpedia
3	British aid worker	http://xmins.com/foaf0.1/Person	DBpedia
4	Britain		
5	terrorist group		
6	U.S. Secretary of State John Kerry		
7	U.S. Secretary		
8	State John Kerry		
9	Middle Eastern leaders		
10	Middle Eastern		DBpedia

Figura 35. Resultado del servicio de enlace.
Fuente: (propio)

Se puede visualizar el JSON resultante de los servicios invocados, antes de ser procesado por el cliente para esto se accede al último de los ítems del menú como lo se puede ver en la figura 36.

Procesamiento de Texto - x

← → ↻

taw02.utpl.edu.ec/anavisors/frontal#

UTPL

FRONTAL WS

ISIS beheaded another westerner, a British aid worker. Britain's leader vowed the death won't stop his country from battling the terrorist group. And U.S. Secretary of State John Kerry courted Middle Eastern leaders to join a coalition in the fight.

✓ Segmentación en Sentencias

✓ Tokenización

✓ Etiquetado

✓ Extracción de Entidades

✓ Desambiguación y Enlace

Procesar

Borrar

Sentencias	Tokens	Entidades	Keywords Compuestas	Keywords Simples	Enlazados
3	46	8	2	9	10

SENTENCIAS» TOKENS» ETIQUETAS PARTOFSPEECH) EXTRACCIÓN ENLACE JSON

{
 "result": {
 "Entidades": ["ISIS", "Westerner", "Britain", "U.S. Secretary of State John Kerry", "U.S. Secretary", "State John Kerry", "Middle Eastern leaders", "Middle Eastern"],
 "EntidadesDesambiguadas": [{"
 "dbpediaList": [{"
 "abstract": ["ISIS", "Goddess", "Ancient", "Egyptian", "Religious", "Beliefs", "Worship", "Spread", "Grec O-roman", "World"],
 "uri": "http://dbpedia.org/resource/ISIS",
 "abstract": ["Tamara", "Diane", "Wimer", "Born", "October", "Stage", "Name", "Isis", "Gee", "American", "Po p", "Singer-songwriter", "Arranger", "Programer", "Producer", "Own", "Music"],
 "uri": "http://dbpedia.org/resource/Isis_Gee",
 "abstract": ["Isis", "King", "Born", "October", "American", "Fashion", "Model", "Fashion", "Designer"],
 "uri": "http://dbpedia.org/resource/Isis_King"},
]
 }
 }
}

Figura 36. Captura de la visualización del JSON.
Fuente: (propio)

58

A continuación se presenta un ejemplo en el cual no se accede a todas las funcionalidades del sistema, sino solo a la función de etiquetado y funciones de las cuales depende, una captura del resultado se visualiza en la figura 37.

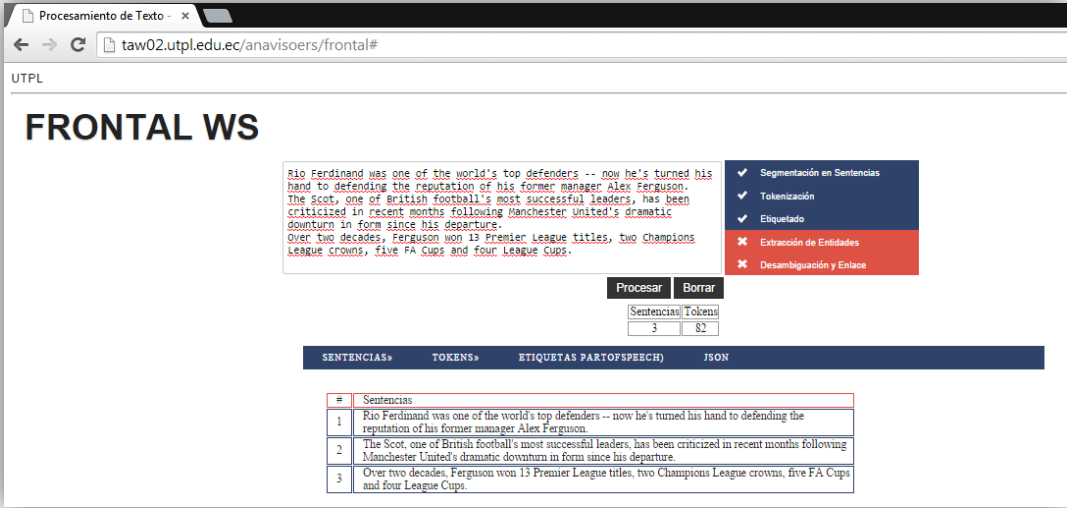


Figura 37. Captura del resultado del servicio de etiquetado de palabra
Fuente: (propio)

La presentación de la tabla resumen de los procesos así como el menú de que permite navegar por los resultados se ven afectados puesto que solo muestran los concernientes a los resultados invocados además del menú para la visualización de del JSON, una ampliación de esto se observa en la figura 38.



Figura 38. Tabla y menú generado de la llamada al servicio de etiquetado.
Fuente: (propio)

3.4.4. Resumen de prototipos

Prototipo 1

En inicio de se trató de utilizar herramientas abiertas disponibles en la web para el proceso de enlace y desambiguación de los elementos trascendentes en un texto, así se consumió el servicio ofrecido por DBpedia, spotlight el cual produce el resultado esperado

pero con problemas de conexión debido a que el servidor no se encontraba siempre disponible, para lo cual se construyó un cliente REST en Python. A partir de este intento y conforme se fue avanzando en la investigación de soluciones para los requerimientos presentados, se decidió desarrollar una propuesta de software propia.

Tabla 18. Tabla resumen del prototipo 1

Funcionalidad específica	Aprendizaje	Herramienta	Tiempo
Desambiguación y enlace	Consumir servicios libres existentes para procesos de Tiempo excesivos de estará para espera respuesta	DBpedia Spotlight	3 días
Construcción de cliente para servicios	Conexión con servidores REST	Python	2 días

Fuente: (propio)

Prototipo 2

A través de Procesamiento de Lenguaje Natural (PLN), se trabaja con el texto de entrada para obtener las entidades de quien se habla en las oraciones, así como las palabras claves que los acompañan. Este procesamiento se lo realiza especializado en el idioma inglés utilizando la librería *Natural Language Toolkit* (NLTK) de Python.

Tabla 19. Tabla resumen del prototipo 2

Funcionalidad específica	Aprendizaje	Herramienta	Tiempo
Tokenización en sentencias	▪ Procesamiento de lenguaje natural o PLN	▪ NLTK (Natural Language Toolkit) librería de Python.	2 semanas
Tokenización en palabras	▪ PLN	▪ NLTK	2 horas
Etiquetado (Part of speech) en idioma inglés	▪ PLN	▪ NLTK	2 horas
Extracción de entidades y palabras claves idioma inglés	▪ PLN	▪ NLTK	1 día

Fuente: (propio)

Prototipo 3

Todos los esfuerzos realizados se concentran en el idioma inglés, a través de NLTK para su procesamiento, para en otros idiomas como español por ejemplo no se puede realizar el mismo proceso que en el idioma inglés debido a lo diferente de su estructura. Para poder trabajar en otros idiomas se debe partir de un etiquetado (Part of Speech) propio del idioma, por lo cual se decidió trabajar con TreeTagger, que permite etiquetado en diferentes idiomas.

Tabla 20. Tabla resumen del prototipo 3

Funcionalidad específica	Aprendizaje	Herramienta	Tiempo
Etiquetado (Part of speech) en idioma inglés	▪ PLN	▪ TreeTagger	8 días
Extracción de entidades y palabras claves idioma inglés	▪ PLN	▪ NLTK	2 día

Fuente: (propio)

Prototipo 4

Una vez que se ha logrado extraer los datos importantes de un texto, se procede a enlazar esto con recursos disponibles en DBpedia.org, para esto se realizan consultas Sparql para obtener los recursos que coinciden con los datos extraídos. Una vez obtenidos los recursos de DBpedia, pueden existir más de un recurso que pueda coincidir para un dato del texto, este caso el que este término es ambiguo y es necesario romper esta ambigüedad para enlazarlo con un solo recurso de DBpedia, para lo cual se implementa el algoritmo de Lesk introducido por Michael E. Lesk en 1986.

Tabla 21. Tabla resumen del prototipo 4

Funcionalidad específica	Aprendizaje	Herramienta	Tiempo
Enlazar entidades y palabras claves encontrados con recursos similares en DBpedia	▪ Consultas en Sparql	▪ Python ▪ Sparql	4 días
Desambiguación entre los recursos encontrados	▪ Algoritmo de desambiguación de Lesk	▪ Python ▪ Sparql	5 días

Fuente: (propio)

Prototipo 5

Se levantan servicio **Rest** para ofrecer la propuesta y una interfaz para que pueda interactuar con usuarios.

Tabla 22. Tabla resumen del prototipo 5

Funcionalidad especifica	Aprendizaje	Herramienta	Tiempo
Servicio Web REST	<ul style="list-style-type: none">▪ Construcción Servicios Web▪ REST	<ul style="list-style-type: none">▪ REST▪ Python	2 semanas
Interfaz de usuario (Web) prototipo inicial	<ul style="list-style-type: none">▪ Construcción de cliente Web▪ JavaScript	<ul style="list-style-type: none">▪ JavaScript▪ HTML	3 días

Fuente: (propio)

Prototipo 6

Se levantan servicios diferenciados por los procesos realizado de esta forma el resultado de un servicio o es la entrada de toro servicio y así puedes ser consumido por otros servicios externos a esta propuesta.

Para poder visualizar el resultado de la interacción de los servicios Web de forma agradable para usuarios se mejora la interfaz gráfica en donde se diferencia el resultado de cada servicio.

Tabla 23. Tabla resumen del prototipo 6

Funcionalidad especifica	Aprendizaje	Herramienta	Tiempo
Servicios web diferenciados por Procesos	<ul style="list-style-type: none">▪ Levantar servicios Web Rest	<ul style="list-style-type: none">▪ REST▪ Python	3 semanas
Interfaz de usuario final	<ul style="list-style-type: none">▪ Mejoramiento de interfaz Web	<ul style="list-style-type: none">▪ Css▪ JavaScript	7 días

Fuente: (propio)

CAPITULO 4: VALIDACIÓN Y PRUEBAS

4. Validación de resultados

4.1. Objetivo.

Contabilizar las sentencias, tokens, entidades, keywords compuestas, keywords simples, resultantes de los servicios web de Segmentación, Tokenización, Etiquetado, Extracción y el servicio de Desambiguación y Enlace.

Medir los resultados del software desarrollado en base a los contenidos de publicaciones y proyectos reales a fin de verificar que lo obtenido se alinea y nutre los datos ya dados por los autores como keywords.

4.2. Contexto.

Las pruebas se realizarán en base a publicaciones desarrolladas en el departamento de Tecnologías Avanzadas de la Web y Sistemas Basados en Conocimiento (TAWSBC), de la Universidad Técnica Particular de Loja, así como el proyecto SMARTLAND (smartland.utpl.edu.ec) que actualmente se encuentra en desarrollo, a fin de extraer datos relevantes. Las pruebas se desarrollarán teniendo en cuenta los siguientes puntos:

- Las pruebas se desarrollarán en un servidor local
- No se modificarán los textos de las fuentes base, para la ejecución de las pruebas
- Los contenidos se encontrarán disponibles en la web

4.3. Pruebas sobre el abstract de publicación.

4.3.1. Descripción de publicaciones.

Consuming and producing linked open data: The case of OpenCourseWare

Autores: Piedra, Nelson; Tovar, Edmundo; Colomo-Palacios, Ricardo; Lopez-Vargas, Jorge; Chicaiza, Janneth Alexandra.

URL: <http://www.emeraldinsight.com/doi/full/10.1108/PROG-07-2012-0045>

Abstract: Purpose: The aim of this paper is to present an initiative to apply the principles of Linked Data to enhance the search and discovery of OpenCourseWare (OCW) contents created and shared by the universities. Design/methodology/approach: This paper is a case study of how linked data technologies can be applied for the enhancement

of open learning contents. Findings: Results presented under the umbrella of OCW-Universia consortium, as the integration and access to content from different repositories OCW and the development of a query method to access these data, reveal that linked data would offer a solution to filter and select semantically those open educational contents, and automatically are linked to the linked open data cloud. Originality/value: The new OCW-Universia integration with linked data adds new features to the initial framework including improved query mechanisms and interoperability.

4.3.2. Resultados de los servicios web

La tabla 24 resume los resultados de la ejecución prueba muestra los datos extraídos cuantificados, iniciando con las sentencias que componen el *abstract* analizado, los tokens y las etiquetas dentro de estas, los elementos extraídos distinguido entre entidades, keywords compuestas y simples, finalmente cuántos de estos elementos han sido enlazados con recursos de DBpedia. Para visualizar todos los datos ampliados dirigirse a Anexos 9.

Tabla 24: Resumen de resultados del procesamiento de la publicación

Sentencias	Tokens	Extracción			Desambiguación y Enlace		
		Entidades	Keywords Compuestas	Keywords Simples	Entidades	Keywords Compuestas	Keywords Simples
4	154	5	10	45	3	0	18

Fuente: (propio)

La figura 39 permite visualizar el contraste entre los resultados del servicio de Extracción y el servicio de Desambiguación y Enlace, en razón de elementos extraídos y de estos cuantos se han logrado enlazar, distinguiendo por tipos de elementos. En porcentajes en este caso en específico se puede decir que 60% de las entidades extraídas han sido enlazadas, el 0% de keywords compuestas y el 40 % de keywords simples.

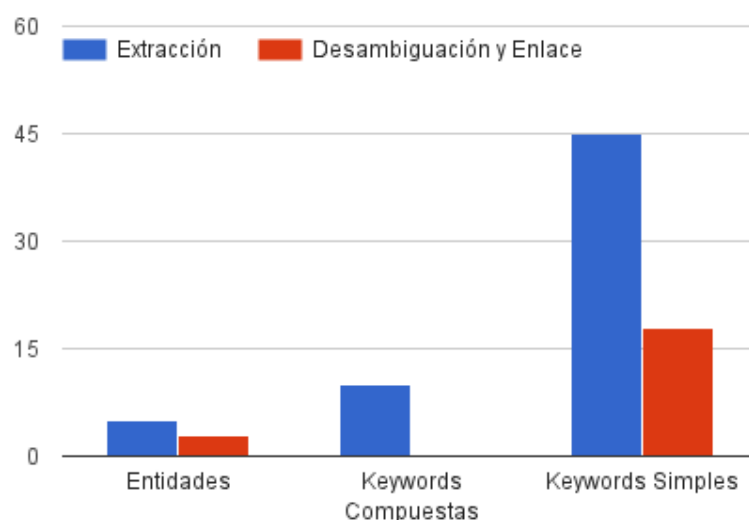


Figura 39: Comparación grafica entre elementos extraídos y enlazados
Fuente: (propio)

4.3.3. Comparación y enriquecimientos de datos

Las keywords dadas por el equipo de autores del artículo que has sido seleccionadas como las más representativas al momento comunicar a los lectores sobre el tema principal y los secundarios de la publicación, si a estos los comparamos con las entidades y keywords compuestas resultantes del proceso de servicios web tal como en la tabla XXX, se puede constatar la similitud que existen entre las palabras que las componen y las áreas relacionadas a las cuales se refieren. Y aún más se identifica una entidad como *OCW-Universa* y keywords compuestas extraídas tales como *open educational contents*, *open data*, *open learning contents*, *query mechanisms*, entro otros, que pueden dar nuevas ideas a los lectores del tema que aborda la publicación.

Tabla 25: comparación de resultados y keywords dados por los autores

Author Keywords	Entidades	Keywords Compuestas
Knowledge management	Linked Data	case study
Linked open data	OCW-Universia	data technologies
OpenCourseWare	OpenCourseWare	open learning contents
Retrieval and consumption of linked data	new OCW-Universia	different repositories OCW
	OCW	query method
		open educational contents
		open data
		new features
		initial framework
		query mechanisms

Fuente: (propio)

A través del servicio web de desambiguación y enlace, se vincula los datos extraídos con DBpedia el resultado se visualiza en el tabla XXX.

Tabla 26: Elementos enlazados con recursos de DBpedia

ENTIDADES			
#	Entidad	Tipo	Enlace
1	Linked Data		http://dbpedia.org/resource/Linked_data
2	OCW		http://dbpedia.org/resource/OpenCourseWare
3	OpenCourseWare		http://dbpedia.org/resource/OpenCourseWare
KEYWORD COMPUESTAS			
#	Keywords	Tipo	Enlace
KEYWORDS SIMPLES			
#	Keyword	Tipo	Enlace
1	universities		http://dbpedia.org/resource/Universities
2	principles		http://dbpedia.org/resource/Principles
3	paper		http://dbpedia.org/resource/Paper
4	initiative		http://dbpedia.org/resource/Initiative
5	Purpose		http://dbpedia.org/resource/Purpose
6	data		http://dbpedia.org/resource/Data
7	methodology		http://dbpedia.org/resource/Methodology
8	Design		http://dbpedia.org/resource/Design
9	learning		http://dbpedia.org/resource/Learning
10	technologies		http://dbpedia.org/resource/Technologies
11	umbrella		http://dbpedia.org/resource/Umbrella
12	Findings		http://dbpedia.org/resource/Findings
13	consortium		http://dbpedia.org/resource/Consortium
14	solution		http://dbpedia.org/resource/Solution
15	Results	http://www.w3.org/2002/07/owl#Thing http://dbpedia.org/ontology/Album http://dbpedia.org/ontology/MusicalWork http://dbpedia.org/ontology/Work http://schema.org/CreativeWork http://schema.org/MusicAlbum	http://dbpedia.org/resource/Results
16	method		http://dbpedia.org/resource/Methodology
17	Originality		http://dbpedia.org/resource/Originality
18	interoperability		http://dbpedia.org/resource/Interoperability

Fuente: (propio)

4.4. Pruebas en base a los contenido del proyecto SMARTLAND

4.4.1. Datos del proyecto

Proyecto: SMARTLAND - Intelligent Land Management

URL: <http://smartland.utpl.edu.ec/en/node/613>

Descripción: By “Smart land” we are referring to an area in which ICT is used intensively in order to improve both the quality of life of the citizens, as well as environmental management. Smart land broadens the concept of "Smart Cities" which focus on ensuring the intelligent development of cities within a certain territory. This initiative uses techniques including, advances in digital preservation, the representation and retrieval of information, the processing of large volumes of data, variable analysis techniques, sensor technologies, geographical information systems, information visualisation, and technologies emerging from the Semantic Web.

Smart Land is a UTPL initiative, which, through research projects (with the participation of more than 300 researchers), collects, manages and creates social, biological, environmental, cultural, and infrastructural data and indicator models, in order to stimulate innovative management systems for the area. Its purpose is to contribute to the optimal use of natural and cultural resources, reassessing their wealth so that in the medium term, they can be a support tool for decision-making, aimed at improving the quality of life of its inhabitants. This is achieved in collaboration with the regional, local and national governments, as well as public and private companies who wish to join the Smart Land project. This initiative's most significant Primary Research Area is the Province of Zamora-Chinchipe.

The immeasurable biodiversity, the ethnic groups, the rich archaeological heritage, the mineral wealth, and the tourist attractions, which include rivers, waterfalls, lagoons and orography, make it a quite unique area in the world due to the prevalence of the Amazon Rainforest alongside the foothills of the Andes. Several studies have shown that the biodiversity in southern Ecuador maintains a priceless natural heritage, and is considered as being a "hot spot" of biological diversity, resulting in UNESCO declaring it the “Podocarpus Biosphere Reserve - El Condor” (PBRC). This ecological and cultural wonder consists of an area of approximately 1.14 million acres.

4.4.2. Contabilización de Resultados.

Los contabilización de los resultados de la descripción del proyecto SMARTLAND se visualizan en la Tabla 27. Se trata de un texto más extenso que la primera prueba sobre el abstract de una publicación, queda evidente en los números procesados. La comparación entre las elementos extraídos y enlazados se presenta en la Figura 40.

Tabla 27: Resultados de servicios del procesamiento de la descripción del proyecto SMARTLAND

Sentencias	Tokens	Extracción			Desambiguación y Enlace		
		Entidades	Keywords Compuestas	Keywords Simples	Entidades	Keywords Compuestas	Keywords Simples
10	367	17	33	86	10	4	42

Fuente: (propio)

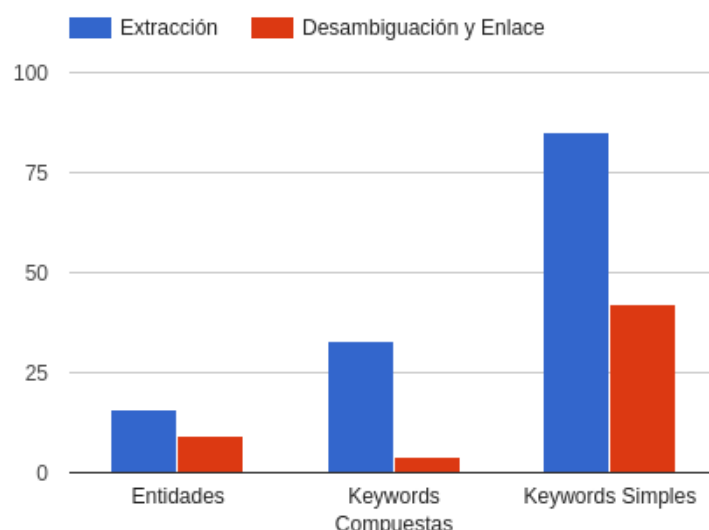


Figura 40: Comparación de datos extraídos y enlazados proyecto SMARTLAND
Fuente: (propio)

4.4.3. Comparación y enriquecimiento de datos.

Para la comparación en este caso se contrasta con principales paquetes de trabajo que se están desarrollando a fin con el proyecto, se evidencia que estas no están claramente definidos dentro de la descripción de proyecto, pero estos con los elementos extraídos guardan una clara relación, además de extraer elementos que no se evidencia en los paquetes de trabajo cumpliendo con el objetivos de descubrir nuevos datos.

Tabla 28: Comparación de elementos extraídos y los paquetes de trabajo del proyecto SMARTLAND

Work Packages	Entidades	Keywords Compuestas	
Patrimonial, cultural, touristic and recreational assets.	ICT		national governments
Biodiversity and ecosystem integrity.	Smart Cities	environmental management	immeasurable biodiversity
Cartography and geomorphology.	Semantic Web	Smart land	mineral wealth
Climate.	Smart Land	intelligent development	Several studies
Education: indicators of quality and coverage.	UTPL	large volumes	biological diversity
Energy and telecommunications.	significant Primary Research Area	geographical information systems	digital preservation

Infrastructure and transport.	Province of Zamora-Chinchipe	infrastructural data	sensor technologies
Water resources and water quality.	Province	optimal use	research projects
Public health.	Zamora-Chinchipe	support tool	innovative management systems
Production systems, entrepreneurship, innovation and economic indicators.	Amazon Rainforest	initiative ´ s	medium term
Society, human mobility and values.	Andes	rich archaeological heritage	private companies
Sustainable biodiversity use.	southern Ecuador	unique area	ethnic groups
	PBRC	hot spot	tourist attractions
	UNESCO	certain territory	priceless natural heritage
	Podocarpus Biosphere Reserve-EI Condor	variable analysis techniques	cultural wonder
		information visualisation	cultural resources
		indicator models	

Fuente: (propio)

El enriquecimiento de datos se hace evidente en la tabla 29 donde se muestran los recursos de DBpedia a los que se hace referencia en la descripción del proyecto.

Tabla 29: Elementos enlazado hacia recursos de DBpedia del análisis del proyecto SMARTLAND

ENTIDADES			
#	Entidad	Tipo	Enlace
1	ICT		http://dbpedia.org/resource/Information_and_communication_technologies_for_environmental_sustainability
2	Smart Cities		http://dbpedia.org/resource/Smart_city
3	Semantic Web		http://dbpedia.org/resource/Semantic_Web
4	Zamora-Chinchipe	http://www.w3.org/2002/07/owl#Thing http://dbpedia.org/ontology/Place http://dbpedia.org/ontology/PopulatedPlace http://dbpedia.org/ontology/Settlement http://schema.org/Place	http://dbpedia.org/resource/Zamora-Chinchipe_Province
5	Province		http://dbpedia.org/resource/Province
6	Amazon Rainforest	http://www.w3.org/2002/07/owl#Thing http://dbpedia.org/ontology/Place http://schema.org/Place	http://dbpedia.org/resource/Amazon_rainforest

7	Andes		http://dbpedia.org/resource/Andes
8	UNESCO		http://dbpedia.org/resource/UNESCO
9	PBRC		http://dbpedia.org/resource/Packed_red_blood_cells

KEYWORD COMPUESTAS

#	Keywords	Tipo	Enlace
1	environmental management		http://dbpedia.org/resource/Environmental_resources_management
2	digital preservation		http://dbpedia.org/resource/Digital_preservation
3	geographical information systems		http://dbpedia.org/resource/Geographic_information_system
4	ethnic groups		http://dbpedia.org/resource/Ethnic_group

KEYWORDS SIMPLES

#	Keyword	Tipo	Enlace
1	life		http://dbpedia.org/resource/Life
2	management		http://dbpedia.org/resource/Management
3	area		http://dbpedia.org/resource/Area
4	citizens		http://dbpedia.org/resource/Citizens
5	concept		http://dbpedia.org/resource/Concept
6	cities		http://dbpedia.org/resource/Cities
7	information		http://dbpedia.org/resource/Information
8	sensor		http://dbpedia.org/resource/Sensor
9	systems		http://dbpedia.org/resource/Systems
10	initiative		http://dbpedia.org/resource/Initiative
11	analysis		http://dbpedia.org/resource/Analysis
12	volumes		http://dbpedia.org/resource/Volumes
13	technologies		http://dbpedia.org/resource/Technologies
14	data		http://dbpedia.org/resource/Data
15	research		http://dbpedia.org/resource/Research
16	infrastructural		http://dbpedia.org/resource/Infrastructural
17	projects		http://dbpedia.org/resource/Projects

18	researchers		http://dbpedia.org/resource/Researchers
19	wealth		http://dbpedia.org/resource/Wealth
20	inhabitants		http://dbpedia.org/resource/Inhabitants
21	decision-making		http://dbpedia.org/resource/Decision-Making
22	purpose		http://dbpedia.org/resource/Purpose
23	tool		http://dbpedia.org/resource/Tool
24	resources		http://dbpedia.org/resource/Resources
25	project		http://dbpedia.org/resource/Projects
26	collaboration		http://dbpedia.org/resource/Collaboration
27	companies		http://dbpedia.org/resource/Companies
28	governments		http://dbpedia.org/resource/Governments
29	s		http://dbpedia.org/resource/Citizens
30	biodiversity		http://dbpedia.org/resource/Biodiversity
31	lagoons		http://dbpedia.org/resource/Lagoons
32	tourist		http://dbpedia.org/resource/Tourist
33	attractions		http://dbpedia.org/resource/Attractions
34	waterfalls		http://dbpedia.org/resource/Waterfalls
35	foothills		http://dbpedia.org/resource/Foothills
36	orography		http://dbpedia.org/resource/Orography
37	rivers		http://dbpedia.org/resource/Rivers
38	prevalence		http://dbpedia.org/resource/Prevalence
39	world		http://dbpedia.org/resource/World
40	mineral		http://dbpedia.org/resource/Mineral
41	diversity		http://dbpedia.org/resource/Biodiversity
42	acres		http://dbpedia.org/resource/Acres

Fuente: (propio)

5. Pruebas funcionales

5.1. Objetivo

Se pretende comprobar el buen funcionamiento de los servicios en las posibles interacciones con los usuarios, en específico se verificará:

- La aplicación responde de forma adecuada a las solicitudes.
- Se obtiene la respuesta específica solicitada a través de los servicios web.
- No presenta efectos secundarios sobre otras funcionalidades.
- La aplicación no se paraliza.
- No se presentan errores.

5.2. Escenario

Las pruebas se realizarán en un servidor local y se accederá a través de la interfaz web que permite la interacción entre los servicios construidos.

5.3. Pruebas sobre el servicio web de Segmentación en Sentencias

Tabla 30: Pruebas de funcionalidad sobre servicio web de segmentación en sentencias.

Entrada	WS - Segmentación en Sentencias			Observ.
	Resultado Esperado	Resultado Obtenido	Error (S/N)	
Texto vacío	json con mensaje de "error: error no text"	json con mensaje de "error: error no text"	N	
Caracteres especiales "#\$%&/()=?_!*- :; @½·~¬{[()]}	json con mensaje de "error: Text unsupported or language unsupportedd"	json con mensaje de "error: Text unsupported or language unsupportedd"	N	
Texto solo con caracteres en blanco	json con mensaje de "error: Text unsupported or language unsupportedd"	json con mensaje de "error: Text unsupported or language unsupportedd"	N	
Texto en español	json con mensaje de "error: Text unsupported or language unsupportedd"	json con mensaje de "error: Text unsupported or language unsupported"	N	
Texto en francés	json con mensaje de "error: Text unsupported or language unsupported"	json con mensaje de "error: Text unsupported or language unsupported"	N	
Texto en otro idioma (italiano)	json con mensaje de "error: Text unsupported or language unsupported"	json con mensaje de "error: Text unsupported or language unsupported"	N	
Texto con entidades en español	json son sentencias detectadas.	json son sentencias detectadas.	N	
Texto con entidades en francés	json son sentencias detectadas.	json son sentencias detectadas.	N	
Texto solo con números	json con mensaje de "error: Text unsupported or language unsupported"	json con mensaje de "error: Text unsupported or language unsupported"	N	

Fuente: (propio)

5.4. Prueba sobre el servicio web de Tokenización

Tabla 31: Pruebas del servicio web de Tokenización

Entrada	WS - Tokenización			Observ.
	Resultado Esperado	Resultado Obtenido	Error (S/N)	
Texto vacío	json con mensaje de "error: error no text"	json con mensaje de "error: error no text"	N	
Caracteres especiales "#\$%&/()=?_!*- :; @½·~¬{[]}"	json con tokens por cada caracter encontrado	json con tokens por cada caracter encontrado	N	
texto solo con caracteres en blanco	json vacío	json vacío	N	
Texto en español	json con mensaje de "error: Text unsupported or language unsupported"	json con mensaje de "error: Text is not English. spanish unsupported"	N	
texto en francés	json con mensaje de "error: Text unsupported or language unsupported"	json con mensaje de "error: Text is not English. french unsupported"	N	
Texto en otro idioma (italiano)	json con mensaje de "error: Text unsupported or language unsupported"	json con mensaje de "error: Text is not English. french unsupported"	N	
Texto con entidades en español	json con tokens por cada caracter encontrado	json con tokens por cada caracter encontrado	N	
texto con entidades en frances	json con tokens por cada caracter encontrado	json con tokens por cada caracter encontrado	N	
Texto solo con números	json con mensaje de "error: Text unsupported or language unsupported"	json con mensaje de "error: Text unsupported or language unsupported"	N	

Fuente: (propio)

5.5. Servicio web de Etiquetado

Tabla 32: Pruebas del servicio web de Etiquetado

Entrada	WS - Etiquetado			Observ.
	Resultado Esperado	Resultado Obtenido	Error (S/N)	
Texto vacío	json con mensaje de "error: error no text"	json con mensaje de "error: error no text"	N	
Caracteres especiales "#\$%&/()=?_!*- :; @½·~¬{[]}"	Json con etiquetas de los caracteres	Json con etiquetas de los caracteres	N	

:: @½·~~~{[()]}				
texto solo con caracteres en blanco	json vacío	json vacío	N	
Texto en español	json con mensaje de "error: Text unsupported or language unsupported"	json con mensaje de "error: Text unsupported or language unsupported"	N	
Texto en francés	json con mensaje de "error: Text unsupported or language unsupported"	json con mensaje de "error: Text unsupported or language unsupported"	N	
Texto en otro idioma (italiano)	Json con etiquetas de los tokens	Json con etiquetas de los tokens	N	
Texto con entidades en español	Json con etiquetas de los tokens	Json con etiquetas de los tokens	N	
texto con entidades en frances	Json con etiquetas de los tokens	Json con etiquetas de los tokens	N	
Texto solo con números	json con mensaje de "error: Text unsupported or language unsupported"	json con mensaje de "error: Text unsupported or language unsupported"	N	

Fuente: (propio)

5.6. Servicio web de Extracción

Tabla 33: Pruebas del servicio web de Extracción

Entrada	WS - Extracción			Observ.
	Resultado Esperado	Resultado Obtenido	Error (S/N)	
Texto vacío	json con mensaje de "error: error no text"	json con mensaje de "error: error no text"	N	
Caracteres especiales "#\$%&/()=? *~ :: @½·~~~{[()]}	json vacío	json vacío	N	
texto solo con caracteres en blanco	json vacío	json vacío	N	
Texto en español	json con mensaje de "error: Text is not English. spanish unsupported"	json con mensaje de "error: Text is not English. spanish unsupported"	N	
Texto en francés	json con mensaje de "error: Text is not English. french unsupported"	json con mensaje de "error: Text is not English. french unsupported"	N	
Texto en otro idioma (italiano)	json con mensaje de "error: Text is not English. french unsupported"	json con mensaje de "error: Text is not English. french unsupported"	N	
Texto con entidades en español	json con entidades y keyword extraídos	json con entidades y keyword extraídos	N	
Texto con entidades	json con entidades y keyword	json con entidades y keyword	N	

en frances	extraídos	extraídos		
Texto solo con números	json con mensaje de "error: Text unsupported or language unsupported"	json con mensaje de "error: Text unsupported or language unsupported"	N	

Fuente: (propio)

5.7. Servicio web de Desambiguación y Enlace

Tabla 34: Pruebas del servicio web de Desambiguación y Enlace

Entrada	WS - Desambiguación y Enlace			Observ.
	Resultado Esperado	Resultado Obtenido	Error (S/N)	
Texto vacío	json con mensaje de "error: error no text"	json con mensaje de "error: error no text"	N	
Caracteres especiales "#\$%&/()=? *~ :; @½~→{[()]}	json vacío	json vacío	N	
texto solo con caracteres en blanco	json vacío	json vacío	N	
Texto en español	json con mensaje de "error: Text is not English. spanish unsupported"	json con mensaje de "error: Text is not English. spanish unsupported"	N	
Texto en francés	json con mensaje de "error: Text is not English. french unsupported"	json con mensaje de "error: Text is not English. french unsupported"	N	
Texto en otro idioma (italiano)	json con mensaje de "error: Text is not English. french unsupported"	json con mensaje de "error: Text is not English. french unsupported"	N	
Texto con entidades en español	json con entidades y keywords enlazados	json con entidades y keywords enlazados	N	
texto con entidades en frances	json con entidades y keywords enlazados	json con entidades y keywords enlazados	N	
Texto solo con números	json con mensaje de "error: Text unsupported or language unsupported"	json con mensaje de "error: Text unsupported or language unsupported"	N	

Fuente: (propio)

DISCUSIÓN

Partiendo de la problemática planteada y valiéndonos de las tecnologías y técnicas existentes, se procede a la construcción de servicios web cumpliendo así también los objetivos propuestos. Estos servicios están estrechamente relacionados y son dependientes entre sí en forma secuencial para completar con sus funciones, a continuación se los menciona:

- Servicio web de Segmentación de Sentencias,
- Servicio web de Tokenización,
- Servicio web de Etiquetado,
- Servicio web de Extracción,
- Servicio web de Desambiguación y Enlace.

Uno de los puntos más fuertes y críticos para el desarrollo de la aplicación final es el Procesamiento de Lenguaje Natural (PLN), tanto es así que es la función principal de los cuatro primeros servicios, es el de dedicarse al tratamiento del texto de entrada que es base e inicio de los procesados. Si bien es cierto se pudo crear un solo servicio web que devuelva todo el resultado del PLN, pero su división en diferentes servicios ha permitido que estos puedan pasar a ser parte de los procesos de otro sistema, la utilización de diferentes herramientas de PNL, mejorar el diseño de la aplicación entre otros beneficios.

El quinto servicio web cumple con dos pases fundamentales como su nombre lo describe el de desambiguar que mantiene relación con PNL y enlazar las entidades y keywords que el cuarto servicio le provee, para lo cual intervienen componentes y tecnologías alineadas con la web semántica, como son RDF y SPARQL. A través de estas se puede obtener y manipular recursos de DBpedia (que es uno de los datasets más grandes de web con datos estructurados en tripletas y el nodo central de LOD cloud) que coincidan con las entidades y keyword extraídos, y si una de entidades o keywords llegase a tener más de una coincidencia para un recurso de DBpedia se ejecutarían procesos de desambiguación que tienen como eje principal el algoritmo de Lesk, a fin de descubrir el significado más acorde al contexto.

La implementación de un Dataset local con los recursos de DBpedia en el idioma inglés, ha permitido evitar posibles problemas ocasionados por conexión perdidas o lentas con el servidor de DBpedia y su implementación es posible gracias a que facilita sus recursos, y

las tecnologías libres para levantar un Triplestore con SPARQL endpoint como VIRTUOSO, estos ha utilizado en disco un espacio poco mayor en disco a 4GB que considerablemente bajo si tiene en cuenta los beneficios incluido una base propia para futuros proyecto.

Los resultados de los servicios web son visibles gracias a la construcción de una interfaz web que permite visualizar el resultado de cada uno y así la entrada del siguiente, además de gestionar el servicio que se desea ejecutar respetando la dependencia obligatoria entre ellos.

CONCLUSIONES

Una vez implementada la aplicación en cada una de sus funcionalidades, haberlas puesto a prueba y corregido errores, se puede concluir:

- En cuanto al servicio de extracción en las pruebas en promedio del total de elementos extraídos y de acuerdo con la clasificación establecida para los resultados de esta aplicación, el 8.14% son entidades, 24.14% keywords compuestos y 67.7% son keywords simples.
- Del total de elementos extraídos en las pruebas realizadas, en promedio solo 36.4% es enlazado hacia recursos de DBpedia, dejando un porcentaje 63.6% de elementos extraídos sin ser enlazados
- De todos los elementos enlazados 16.7% son entidades, 4.2% son keywords compuestas y el 79.1% son keywords simples.
- Del total de entidades extraídas el 59% se han enlazado con un recurso de DBpedia (es decir, que de 10 entidades extraídas 6 serían enlazadas) lo cual por tratarse de los elementos principales en los textos analizados, deja una buena impresión de los recursos disponibles en el dataset de DBpedia.
- Las keywords compuestas con 9.3% extracciones enlazadas son las que menor porcentaje de éxito al enlazar con recursos de DBpedia.
- Las keywords simples con 40.5% de enlaces realizados sobre las extracción logradas, y al ser más numerosas son las que mayormente se han enlazado con recursos de DBpedia.

RECOMENDACIONES

La creación de recursos propios dentro del Dataset local con los recursos de DBpedia, permitiría poder enlazar entidades y palabras claves que no tengan recursos en DBpedia, y que posiblemente pertenezcan a un entorno académico local, como un profesor por ejemplo.

Realizar un análisis previo de la estructura ontológica de los recursos de DBpedia antes de trabajar con estos, para facilitar el “moverse” mediante consultas SPARQL a través de sus propiedades y relaciones.

La utilización de lenguaje de programación de alto nivel, Python ya que está provisto de las librerías *Natural Language Toolkit (NLTK)* y *Treetagger* que permiten el procesamiento de lenguaje natural de forma fácil y potente, y que al igual que Python su curva de aprendizaje relativamente corta.

La utilización de la librería TreeTagger para el etiquetado (POS tagging) en el procesamiento de lenguaje natural permite tener la opción de procesar textos en diferentes al idioma ingles que una limitación que la librería NLTK no puede romper.

Utilizar las recomendaciones de la W3C para la publicación de datos enlazados, y enlazar los contenidos actuales publicados en la web a fuentes de datos estructurados con los principios de Datos Enlazados en pro de la difusión de la web semántica como estructura de la web.

Si bien es cierto se sentó la base para trabajar con otros idiomas distintos al inglés, esto no se encuentra desarrollado, así que se lo puede considerar como trabajo futuro.

Como trabajo futuro el desarrollo de un cliente que permita gestionar mejor los resultados y agregar funcionalidades como guardarlos con fines específicos como parte de una ontología.

BIBLIOGRAFÍA

- Albahari, J., & Albahari, B. (2012). C# 5.0 IN A NUTSHELL. En 656. O'Reilly.
- Beckett, D., Berners-Lee, T., Prud'hommeaux, E., Carothers, G., & Machina, L. (25 de 02 de 2014). *RDF 1.1 Turtle*. Obtenido de W3C Recommendation: <http://www.w3.org/TR/2014/REC-turtle-20140225/>
- Berners-Lee, T. (23 de Julio de 2006). Linked Data - Design Issues.
- Berners-Lee, T., Fielding, R., & Masinter, L. (01 de 2005). *Uniform Resource Identifier (URI): Generic Syntax*. Recuperado el 24 de 06 de 2014, de <http://tools.ietf.org/html/rfc3986>
- Bizer, C. (09 de 11 de 2009). *Dbpedia*. Recuperado el 10 de 06 de 2014, de The DBpedia Data Provision Architecture: <http://wiki.dbpedia.org/Architecture>
- Bizer, C., & Cyganiak, R. (25 de 02 de 2014). *RDF 1.1 TriG*. Obtenido de W3C Recommendation: <http://www.w3.org/TR/2014/REC-trig-20140225/>
- Carothers, G., & Seaborne, A. (25 de 02 de 2014). *RDF 1.1 N-Triples*. Recuperado el 25 de 06 de 2014, de W3C Recommendations: <http://www.w3.org/TR/2014/REC-n-triples-20140225/>
- Clark, K. G., Feigenbaum, L., & Torres, E. (01 de 15 de 2008). *SPARQL Protocol for RDF*. Recuperado el 24 de 06 de 2014, de W3C Recommendation 15 January 2008: <http://www.w3.org/TR/rdf-sparql-protocol/>
- Cutting, D., Kupiec, J., Pedersen, J., & Sibun, P. (Marzo de 1992). A Practical Part of Speech Tagger.
- Cyganiak, R., Wood, D., & Lanthaler, M. (25 de 02 de 2014). *RDF 1.1 Concepts and Abstract Syntax*. Recuperado el 25 de 06 de 2014, de W3C Recommendation: <http://www.w3.org/TR/2014/REC-rdf11-concepts-20140225/>
- Lapiente, M. J. (08 de 12 de 2013). *HIPERTEXTO: EL NUEVO CONCEPTO DE DOCUMENTO EN LA CULTURA DE LA* . Recuperado el 24 de 06 de 2014, de Tesis doctoral. Universidad Complutense de Madrid.: <http://www.hipertexto.info/>
- Lehmann, J., Isele, R., Jakob, M., Jentzsch, A., Kontokostas, D., Mendes, P. N., y otros. (2012). *DBpedia - A Large-scale, Multilingual Knowledge Base Extracted from Wikipedia*. Obtenido de <http://semantic-web-journal.net/system/files/swj499.pdf>
- McBride, B. (10 de 02 de 2004). *W3C Recommendation*. Recuperado el 26 de 06 de 2014, de RDF Primer: <http://www.w3.org/TR/2004/REC-rdf-primer-20040210/>
- Mihaela, I. N. (2004). *EL CONOCIMIENTO LINGÜÍSTICO EN LA DESAMBIGUACIÓN SEMÁNTICA AUTOMÁTICA*.
- Miller, E. (Mayo de 1998). *An Introduction to the Resource Description Framework*. Recuperado el 27 de 06 de 2014, de <http://www.dlib.org/dlib/may98/miller/05miller.html>

- Miller, E. (1998). Wiley Online Library. *Bulletin of the American Society for Information Science and Technology*, 15-19.
- Morocho, J. C., Piedra, N., & Valverde, M. F. (2012). *Estudio sobre la aplicación de Linked Data a la Legislación de Educación Superior en Latinoamérica*.
- Pautasso, C., Zimmermann, O., & Leymann, F. (Abril de 2008). RESTful Web Services vs. "Big" Web Services: Making the Right Architectural Decision.
- Peláez, A. R., Morocho, J. C., & Malla, P. (2012). *Desambiguación de URI's en el Contexto de Linked Open Data para Linked Universities Data*.
- Pérez, S. V. (2009). *Resolución de la ambigüedad semántica mediante métodos basados en conocimiento y su aportación a tareas de PLN*. Alicante.
- Piedra, N., Chicaiza, J. A., & López, J. (junio de 2014). An Architecture based on Linked Data technologies for the Integration and reuse of OER in MOOCs Context.
- Piedra, N., Tovar, E., Colomo-Palacios, R., Lopez-Varga, s. J., & Chicaiza, A. J. (2014). Consuming and producing linked open data: the case of OpenCourseWare.
- Prud'hommeaux, E., & Seaborne, A. (15 de 01 de 2008). *SPARQL Lenguaje de consulta para RDF*. Recuperado el 25 de 06 de 2014, de Recomendación del W3C de 15 de enero de 2008 : <http://skos.um.es/TR/rdf-sparql-query/>
- Richardson, L., & Amundsen, M. (2013). *RESTful Web APIs*. O'REILLY.
- Richardson, L., & Ruby, S. (2007). RESTful Web Services. En L. Richardson, & S. Ruby, *RESTful Web Services* (pág. 299). O'Reilly.
- Ruby, L. R. (2007). *RESTful Web Services*.
- San Martín Oliva, C. R. (s.f.). *COMPLEJO UNIVERSITARIO ISLAS MALVINAS*. Recuperado el 07 de 2014, de COMPLEJO UNIVERSITARIO ISLAS MALVINAS: <http://www.unsj-cuim.edu.ar/portalezonda/seminario08/archivos/MetodologiaICONIX.pdf>
- Sandeep Chatterjee, j. W. (2004). *Developing Enterprise Web Services: An Architect's Guide*. Person Education Inc.
- Satanjeev, B. (2002). *Adapting the Lesk Algorithm for Word Sense Disambiguation to WordNet*.
- Schreiber, G., & Raimond, Y. (24 de 06 de 2014). *RDF 1.1 Primer*. Obtenido de W3C Working Group: <http://www.w3.org/TR/2014/NOTE-rdf11-primer-20140624/>
- Sören Auer, Christian Bizer, Georgi Kobilarov, Jens Lehmann, Richard Cyganiak, and Zachary Ives. (2007). *In The semantic web*.

- Tello Leal, E. (17 de abril de 2009). *La Desambiguación del Sentido de las Palabras: revisión metodológica*. Obtenido de <http://www.nosolousabilidad.com/articulos/desambiguacion.htm>
- Tjong Kim Sang, E. F., & Buchholz, S. (Septiembre de 2000). Introduction to the CoNLL-2000 shared task: Chunking.
- W3C. (01 de 2005). *Uniform Resource Identifier (URI): Generic Syntax*. Recuperado el 22 de 02 de 2014, de January 2005
- W3C. (2013). *W3C*. Obtenido de <http://www.w3.org/standards/semanticweb/ontology>
- Wood, D. (25 de 02 de 2014). *What's New in RDF 1.1*. Obtenido de W3C: <http://www.w3.org/TR/2014/NOTE-rdf11-new-20140225/>

ANEXOS

1. Anexo 1: Especificación de Requerimientos de Software (ERS)

Especificación de Requerimientos de Software (ERS) *WS para descubrimiento, desambiguación y enlace en Datos Enlazados*

Versión [1.0]

Información de Documento

TÍTULO:	Especificación de casos de Uso
SUBTÍTULO:	Servicio Web de Extracción de Entidades
VERSIÓN:	[1.0]
ARCHIVO:	ECS_SW_ExtracciónEntidades
AUTOR:	Fabricio Montaña
ESTADO:	Borrador

Lista de Cambios

VERSIÓN	FECHA	AUTOR	DESCRIPCIÓN
1.0	17/07/2014	Fabricio Montaña	Emisión inicial

Firmas y Aprobaciones

ELABORADO POR:	Fabricio Montaña		
FECHA:	17/07/2014	FIRMA:	
REVISADO POR:	Ing. Nelson Piedra		
FECHA:		FIRMA:	

Especificación de Requerimientos de Software (ERS)

Introducción

Descripción

El presente documento tiene como finalidad redactar las funcionalidades con las que debe contar el sistema de descubrimiento, desambiguación y enlace en datos enlazados.

Problemas Conocidos

Después de un análisis inicial se detectan los siguientes problemas:

- Los *abstract* de las publicaciones universitarias que contienen datos relevantes a los que se desea acceder se encuentran en texto plano entendible solo para humanos.
- Los datos en texto no son referenciados hacia fuentes externas.
- Los datos en texto plano pueden ser ambiguos y tener más de un enlace posible en LOD Cloud específicamente DBpedia.
- No existe un proceso estándar para procesos de enlace y desambiguación lingüística.
- Las conexiones y/o consultas hacia el DataSet de DBpedia pueden demorar o fallar.

Referencias

ANSI/IEEE Std. 830-1984 Guía del IEEE para la Especificación de Requerimientos Software.³¹

Descripción General

³¹ <http://ieeexplore.ieee.org/xpl/mostRecentIssue.jsp?punumber=2228>

El fin del sistema es el descubrimiento de datos relevantes dentro del texto plano en los *abstracts* en las publicaciones universitarias, y si en caso un término extraído es ambiguo determinar el significado usado, para luego ser enlazado a DBpedia (LOD Cloud) si en caso existiera un recurso al cual referencie.

Para acceder al sistema se levantara servicios web diferentes para cada proceso relevante dentro del sistema que interactuaran entre sí. Se desarrollara un interfaz web para usuarios que permitirá visualizar los resultados individuales e integrales de los servicios.

Perspectiva del Proyecto

Esto se pretende desarrollar en base a la relevancia que toman los datos en la web semántica, buscando enlazar las publicaciones científicas a las fuentes de Datos Enlazados, donde se ubican los recursos a los cuales hacen referencia y permitiendo de esta forma ampliar la información y descubrir nuevos enlaces.

El sistema que permita extraer datos relevantes dentro del texto de las publicaciones científicas, desambiguar estos términos de ser necesario y enlazarlos a LOD Cloud. Se construirá separando e integrando los procesos relevantes mediante servicio web.

Características del Producto

El producto a desarrollar constara de las siguientes características:

- **Servidor**
 - Servicios web REST
 - Separados en los procesos importantes:
 - Procesamiento de lenguaje natural,
 - Desambiguación y
 - Enlace
 - Integrados entre si
 - Servidor local con DataSet de DBpedia (SPARQL EndPoint)
- **Cliente**
 - REST web
 - Visualizar resultados individuales de los servicios
 - Visualizar resultado integrado de los servicios
 - Permitir ver el JSON resultante del servicio consumido

Características del Usuario

Usuarios anónimos: a través de la construcción del cliente Rest web cualquier usuario podrá interactuar con el sistema.

Usuarios clientes Rest: los servicios web implementados podrán responder al cualquier servicio que se pueda construir a partir de esos.

Limitaciones Generales

A continuación se detallan limitaciones en cuanto al software:

- Todos los enlaces que se puedan realizar se los hará con recursos disponibles en el DataSet de DBpedia, esto significa que pueden existir recursos en otro u otros repositorios a los cuales no se los enlazara directamente.
- De no existir el término extraído en DBpedia, no podrá ser enlazado.
- La desambiguación de un “término” extraído de una publicación se realizara en base a los recursos disponibles en la DBpedia que son nombrado mediante este “termino”.

Asunciones y Dependencias

Asunciones

- Posibles errores en extracción de términos del texto de las publicaciones debido a faltar tipeado o error humano en la escritura del texto.
- Posibles errores el en enlace de a DBpedia producto de no existir recurso o error de desambiguación.

Dependencias

- Se desarrollara en lenguaje de programación de alto nivel Python 2.7 y algunas de sus librerías especializadas en procesamiento de lenguaje natural, levantamiento de servicios, consultas SPARQL, etc.
- Cliente desarrollara en base a HTML, CSS, JavaScript, etc.
- De navegadores web que soporte tecnologías en las que se construirá en cliente para poder acceder a este.

Requerimientos Funcionales.

REQ001 Extraer entidades y palabra relevantes

Descripción

Descubrir datos relevantes en el texto, a quien se describe y las palabra relevantes que lo acompañan

Entrada

- Texto

Proceso

1. Tokenización del texto en sentencias (oraciones), separa todas las sentencias.
2. Tokenización de las sentencias en palabras.
3. Etiquetar (Part of Speech).
4. Extracción en base a etiquetas.

Salida

- JSON con entidades y palabras claves extraídas

REQ002 Enlazar entidades y palabra relevantes con LOD Cloud

Descripción

Se enlazara los términos encontrados en caso de que sea posible con la LOD Cloud

Entrada

- JSON estructurado por procesos anteriores con entidades a enazar

Proceso

1. Consultar a DBpedia por recursos que sean nombrados con las entidades y palabras relevantes extraídas del texto de entrada.

Salida

- JSON estructurado con los enlaces de los recursos de DBpedia.

REQ002 Desambiguar entidades y palabra relevantes

Descripción

Se determinara el sentido con que las palabras estas siendo usadas en caso de que estas sean ambiguas

Entrada

- JSON procesos anteriores

Proceso

1. Aplicar algoritmo s de desambiguación en base a contexto.
2. Determinar sentido utilizado en base a mejor resultado de coincidencia.

Salida

- JSON con términos desambiguados.

REQ004 Levantar servicios REST separados para los procesos relevantes.

Descripción

Para que los procesos relevantes dentro del sistema puedan ser reutilizados se levantarán servicios individuales.

Entrada

- JSON.
- Texto.

Proceso

1. Se invoca las funciones necesarias del sistema para resolver la petición del servicio llamado.

2. Se estructurará la data en JSON.
3. Se devuelve el JSON ya sea al servicio que lo invocó o al cliente si este fuese el origen de la invocación del servicio.

Salida

- JSON estructurado de acuerdo al servicio invocado.

REQ005 Frontal UI Web

Descripción

Construir una interfaz web que permita visualizar el comportamiento del sistema, es decir, la integración de los servicios y su funcionamiento individual

Entrada

- Texto

Proceso

1. Introducir texto a ser procesado
2. Seleccionar servicios a ser invocados
3. Esperar resultado
4. Procesar resultado
5. Presentar resultado procesado

Salida

- Resultado gráfico de servicios invocados.

2. Anexo 2: Especificación de Caso de Uso (ECS) - Tokenización en Sentencias

Especificación de Caso de Uso (ECS) *Tokenización en Sentencias*

Versión [1.0]

Información de Documento

TÍTULO:	Especificación de casos de Uso
SUBTÍTULO:	Servicio Web - Tokenización en Sentencias
VERSIÓN:	[1.0]
ARCHIVO:	ECS_SW_TokenizaciónSentencias
AUTOR:	Fabricio Montaña
ESTADO:	Borrador

Lista de Cambios

VERSIÓN	FECHA	AUTOR	DESCRIPCIÓN
1.0	17/07/2014	Fabricio Montaña	Emisión inicial

Firmas y Aprobaciones

ELABORADO POR:	Fabricio Montaña		
FECHA:	17/07/2014	FIRMA:	
REVISADO POR:	Ing. Nelson Piedra		
FECHA:		FIRMA:	

Número	ECS-01	
Nombre	Tokenización en Sentencias	
Actores	Usuario, Cliente	
Descripción	Divide el texto de entrada en sentencias cortas separadas por un punto y parte, la salida es una lista de estas sentencias.	
Precondición	<ul style="list-style-type: none"> ▪ Ingreso texto como parámetro de la aplicación ▪ Texto ha sido validado y procesado ▪ Texto segmentado en sentencias 	
Secuencia Normal	Paso	Acción
	1	Entrada del texto validado, como parámetro para el servicio web. SA1
	2	Verifica el número de sentencias que comenten al texto, que estén separadas por un punto seguido (.)
	3	Divide cada una teniendo en cuenta la terminación con punto (.) estructura las sentencias dentro de una lista. SA1
	4	Estructura la lista de elementos formato JSON.
	5	Devuelve el JSON resultante.
Poscondición	<ul style="list-style-type: none"> ▪ El texto dividido en sentencias. 	
Secuencia alternativo	SA1 el número de sentencias es 1 Se estructura una lista de un solo elemento con la sentencia.	
Prioridad	Media	
Requerimientos Especiales		
Asunciones y Dependencias		
Notas adicionales		

3. Anexo 3: Especificación de Caso de Uso (ECS) - Tokenización en Palabras

Especificación de Caso de Uso (ECS) *Tokenización en Palabras*

Versión [1.0]

Información de Documento

TÍTULO:	Especificación de casos de Uso
SUBTÍTULO:	Servicio Web de Tokenización
VERSIÓN:	[1.0]
ARCHIVO:	ECS_SW_Tokenización
AUTOR:	Fabrizio Montaña
ESTADO:	Borrador

Lista de Cambios

VERSIÓN	FECHA	AUTOR	DESCRIPCIÓN
1.0	17/07/2014	Fabrizio Montaña	Emisión inicial

Firmas y Aprobaciones

ELABORADO POR:	Fabrizio Montaña		
FECHA:	17/07/2014	FIRMA:	
REVISADO POR:	Ing. Nelson Piedra		
FECHA:		FIRMA:	

Número	ECS-02	
Nombre	Tokenización en palabras	
Actores	Cliente, Servicio Web	
Descripción	Divide cada sentencia en palabras validas, tokens.	
Precondición	<ul style="list-style-type: none"> ▪ Ingreso texto como parámetro de la aplicación ▪ Texto ha sido validado y procesado 	
Secuencia Normal	Paso	Acción
	1	Entrada del texto validado, como parámetro para el servicio web.
	2	Segmentación del texto en sentencias. ECS-01
	3	Se recorre la lista de sentencias segmentadas.
	4	Se divide palabra por palabra de la sentencia en una lista, se obtiene una lista de listas.
	5	Se estructura en formato JSON
	6	Retorna el JSON con las sentencias divididas en "tokens"
Postcondición	<ul style="list-style-type: none"> ▪ Texto tokenizado por sentencias y estos a la vez tokenizados en palabras 	
Secuencia alternativo		
Prioridad	Baja	
Requerimientos Especiales	Del funcionamiento del Servicio web de Tokenización en Sentencias	
Asunciones y Dependencias		
Notas adicionales		

4. Anexo 4: Especificación de Caso de Uso (ECS) - Etiquetado

Fabrizio Montaña
Analista - Desarrollador

Ing. Nelson Piedra

Especificación de Caso de Uso (ECS)
Etiquetado

Versión [1.0]

Información de Documento

TÍTULO:	Especificación de casos de Uso
SUBTÍTULO:	Servicio Web de Etiquetado
VERSIÓN:	[1.0]
ARCHIVO:	ECS_SW_Etiquetado
AUTOR:	Fabricio Montaña
ESTADO:	Borrador

Lista de Cambios

VERSIÓN	FECHA	AUTOR	DESCRIPCIÓN
1.0	17/07/2014	Fabricio Montaña	Emisión inicial

Firmas y Aprobaciones

ELABORADO POR:	Fabricio Montaña		
FECHA:	17/07/2014	FIRMA:	
REVISADO POR:	Ing. Nelson Piedra		
FECHA:		FIRMA:	

Número	ECS-03	
Nombre	SW-Etiquetado	
Actores	Cliente, Servicio Web	
Descripción	Este servicio permite la tokenización de cada palabra y etiquetación de las mismas de acuerdo a la función que cumplen en el contexto que se encuentra, para hacerlo se apoya en el servicio web de tokenización en sentencias	
Precondición	<ul style="list-style-type: none"> ▪ Ingreso texto como parámetro de la aplicación ▪ Texto ha sido validado y procesado 	
Secuencia Normal	Paso	Acción
	1	Entrada del texto validado, como parámetro para el servicio web.
	2	Tokenización del texto en sentencias. ECS-01
	3	Recorrido de la lista de sentencias
	4	Etiquetado de las palabras que conforman cada sentencia
	5	Estructura y retorna data en JSON
Poscondición	<ul style="list-style-type: none"> ▪ Texto tokenizado a nivel de palabras y etiquetado. 	
Secuencia alternativo		
Prioridad	Alta	
Requerimientos Especiales		
Asunciones y Dependencias		
Notas adicionales	Depende del funcionamiento del servicio web de Etiquetado en Sentencias (ECS-01)	

5. Anexo 5: Especificación de Caso de Uso (ECS) - Extracción de Entidades

Especificación de Caso de Uso (ECS) *Extracción de Entidades*

Versión [1.0]

Información de Documento

TÍTULO:	Especificación de casos de Uso
SUBTÍTULO:	Servicio Web de Extracción de Entidades
VERSIÓN:	[1.0]
ARCHIVO:	ECS_SW_ExtracciónEntidades
AUTOR:	Fabricio Montaña
ESTADO:	Borrador

Lista de Cambios

VERSIÓN	FECHA	AUTOR	DESCRIPCIÓN
1.0	17/07/2014	Fabricio Montaña	Emisión inicial

Firmas y Aprobaciones

ELABORADO POR:	Fabricio Montaña		
FECHA:	17/07/2014	FIRMA:	
REVISADO POR:	Ing. Nelson Piedra		
FECHA:		FIRMA:	

Número	ECS-04	
Nombre	Extracción de Entidades	
Actores	Cliente, Servicio Web	
Descripción	Permite reconocer y extraer, las entidades y palabras relevantes o claves (keywords) que se encuentran dentro del texto, para lograr se apoya en el servicio web de Etiquetado (y en los que este a su vez , servicio web de tokenización de sentencias)	
Precondición	<ul style="list-style-type: none"> ▪ Ingreso texto como parámetro de la aplicación ▪ Texto ha sido validado y procesado 	
Secuencia Normal	Paso	Acción
	1	Entrada del texto validado, como parámetro para el servicio web.
	2	Tokenización del texto en sentencias. ECS-01
	3	Tokenización y Etiquetado de palabra ECS-03
	4	Reconocimiento de estructuras de Entidades y Keywords
	5	Extracción de Entidades y Keywords
	6	Estructuración de retorno de resultado en formato JSON
Poscondición	<ul style="list-style-type: none"> ▪ Entidades y palabra importantes que las acompañan reconocidos y extraídos 	
Secuencia alternativo		
Prioridad	Alta	
Requerimientos Especiales		
Asunciones y Dependencias		
Notas adicionales	Este servicio depende del funcionamiento del servicio web de Tokenización en Entidades (ECS-01) y Servicio web de Etiquetado (ECS-03)	

6. Anexo 6: Especificación de Caso de Uso (ECS) - Desambiguación y Enlace

Especificación de Caso de Uso (ECS) *Desambiguación y Enlace*

Versión [1.0]

Información de Documento

TÍTULO:	Especificación de casos de Uso
SUBTÍTULO:	Servicio Web – Desambiguación y Enlace
VERSIÓN:	[1.0]
ARCHIVO:	ECS_SW_DesambiguaciónEnlace
AUTOR:	Fabricio Montaña
ESTADO:	Borrador

Lista de Cambios

VERSIÓN	FECHA	AUTOR	DESCRIPCIÓN
1.0	17/07/2014	Fabricio Montaña	Emisión inicial

Firmas y Aprobaciones

ELABORADO POR:	Fabricio Montaña		
FECHA:	17/07/2014	FIRMA:	
REVISADO POR:	Ing. Nelson Piedra		
FECHA:		FIRMA:	

Número	ECS-05	
Nombre	Desambiguación y Enlace	
Actores	Cliente, Servicio Web	
Descripción	Enlaza las entidades y palabras relevantes (keywords) hacia LOD Cloud, más específicamente DBpedia, esto de existir un recurso al cual vincular, en caso de que una entidad o keyword tuviese más de uno posible recurso al cual enlazar, se realizara un proceso de desambiguación y luego de enlace.	
Precondición	<ul style="list-style-type: none"> ▪ Ingreso texto como parámetro de la aplicación ▪ Texto ha sido validado y procesado 	
Secuencia Normal	Paso	Acción
	1	Entrada del texto validado, como parámetro para el servicio web.
	2	Tokenización del texto en sentencias. ECS-01
	3	Tokenización y Etiquetado de palabra. ECS-03
	4	Extracción de Entidades y keywords. ECS-04
	5	Consulta de recursos a DBpedia.
	6	Consulta de “Abstract” de recurso a DBpedia
	7	Verificar si existen Entidades o keywords ambiguas
	8	Desambiguar Entidades y keywords ambiguos. SA1
	9	Estructurara resultado
	10	Retornar resultado
Poscondición	<ul style="list-style-type: none"> ▪ Entidades y palabra importantes que las acompañan reconocidos y extraídos 	
Secuencia alternativo	SA1 Entidades y keywords no ambiguos Se enlaza con los recursos únicos encontrados a las entidades y keywords del texto.	
Prioridad	Alta	
Requerimientos Especiales		

Notas adicionales	Este servicio depende de los servicios web de tokenización en sentencias (ECS-01), etiquetado (ECS-03), extracción de entidades (ECS-04).
--------------------------	---

7. Anexo 7: Pseudocódigo del algoritmo de Lesk encontrado en (Satanjeev, 2002)

```
for every word w[i] in the phrase
  let BEST_SCORE = 0
  let BEST_SENSE = null
  for every sense sense[j] of w[i]
    let SCORE = 0
    for every other word w[k] in the phrase, k != i
      for every sense sense[l] of w[k]
        SCORE = SCORE + number of words that occur in the gloss of
                           both sense[j] and sense[l]
      end for
    end for
    if SCORE > BEST_SCORE
      BEST_SCORE = SCORE
      BEST_SENSE = w[i]
    end if
  end for
if BEST_SCORE > 0
  output BEST_SENSE
else
```

8. Anexo 8: Pseudocódigo del algoritmo de Lesk encontrado en (Satanjeev, 2002)

1.	CC	Coordinating conjunction
2.	CD	Cardinal number
3.	DT	Determiner
4.	EX	Existential there
5.	FW	Foreign word
6.	IN	Preposition or subordinating conjunction
7.	JJ	Adjective
8.	JJR	Adjective, comparative
9.	JJS	Adjective, superlative
10.	LS	List item marker
11.	MD	Modal
12.	NN	Noun, singular or mass
13.	NNS	Noun, plural
14.	NP	Proper noun, singular
15.	NPS	Proper noun, plural
16.	PDT	Predeterminer
17.	POS	Possessive ending
18.	PP	Personal pronoun
19.	PP\$	Possessive pronoun
20.	RB	Adverb
21.	RBR	Adverb, comparative
22.	RBS	Adverb, superlative
23.	RP	Particle
24.	SYM	Symbol
25.	TO	to
26.	UH	Interjection
27.	VB	Verb, base form
28.	VBD	Verb, past tense
29.	VBG	Verb, gerund or present participle
30.	VCN	Verb, past participle
31.	VBP	Verb, non-3rd person singular present
32.	VBZ	Verb, 3rd person singular present
33.	WDT	Wh-determiner
34.	WP	Wh-pronoun
35.	WP\$	Possessive wh-pronoun
36.	WRB	Wh-adverb

9. Anexo 9: Resultado completo de la prueba sobre el artículo: Consuming and producing linked open data: The case of OpenCourseWare

9.1. Sentencias

#	Sentencias
1	Purpose: The aim of this paper is to present an initiative to apply the principles of Linked Data to enhance the search and discovery of OpenCourseWare (OCW) contents created and shared by the universities.
2	Design/methodology/approach: This paper is a case study of how linked data technologies can be applied for the enhancement of open learning contents.
3	Findings: Results presented under the umbrella of OCW-Universia consortium, as the integration and access to content from different repositories OCW and the development of a query method to access these data, reveal that linked data would offer a solution to filter and select semantically those open educational contents, and automatically are linked to the linked open data cloud.
4	Originality/value: The new OCW-Universia integration with linked data adds new features to the initial framework including improved query mechanisms and interoperability.

9.2. Tokens y Etiquetado

SENTENCIA #1

#	Token	Etiqueta
1	Purpose	NN
4	aim	NN
7	paper	NN
10	present	VV
13	to	TO
16	principles	NNS
19	Data	NP
22	the	DT
25	discovery	NN
28	((
31	contents	NNS
34	shared	VVN
37	universities	NNS

#	Token	Etiqueta
2	:	:
5	of	IN
8	is	VBZ
11	an	DT
14	apply	VV
17	of	IN
20	to	TO
23	search	NN
26	of	IN
29	OCW	NP
32	created	VVN
35	by	IN
38	.	SENT

#	Token	Etiqueta
3	The	DT
6	this	DT
9	to	TO
12	initiative	NN
15	the	DT
18	Linked	NP
21	enhance	VV
24	and	CC
27	OpenCourseWare	NP
30))
33	and	CC
36	the	DT

SENTENCIA #2

#	Token	Etiqueta
39	Design	NN
42	/	SYM
45	This	DT
48	a	DT
51	of	IN
54	data	NNS
57	be	VB
60	the	DT
63	open	JJ
66	.	SENT

#	Token	Etiqueta
40	/	SYM
43	approach	NN
46	paper	NN
49	case	NN
52	how	WRB
55	technologies	NNS
58	applied	VVN
61	enhancement	NN
64	learning	NN

#	Token	Etiqueta
41	methodology	NN
44	:	:
47	is	VBZ
50	study	NN
53	linked	VVN
56	can	MD
59	for	IN
62	of	IN
65	contents	NNS

SENTENCIA #3

#	Token	Etiqueta
67	Findings	NNS
70	presented	VVN
73	umbrella	NN
76	consortium	NN
79	the	DT
82	access	NN
85	from	IN
88	OCW	NN
91	development	NN
94	query	NN
97	access	VV
100	,	,
103	linked	VVN

#	Token	Etiqueta
68	:	:
71	under	IN
74	of	IN
77	,	,
80	integration	NN
83	to	TO
86	different	JJ
89	and	CC
92	of	IN
95	method	NN
98	these	DT
101	reveal	VVP
104	data	NNS

#	Token	Etiqueta
69	Results	NNS
72	the	DT
75	OCW-Universia	NP
78	as	IN
81	and	CC
84	content	VV
87	repositories	NNS
90	the	DT
93	a	DT
96	to	TO
99	data	NNS
102	that	IN/that
105	would	MD

106	offer	VV
109	to	TO
112	select	VV
115	open	JJ
118	,	,
121	are	VBP
124	the	DT
127	data	NNS

107	a	DT
110	filter	VV
113	semantically	RB
116	educational	JJ
119	and	CC
122	linked	VVN
125	linked	VVN
128	cloud	VVP

108	solution	NN
111	and	CC
114	those	DT
117	contents	NNS
120	automatically	RB
123	to	TO
126	open	JJ
129	.	SENT

SENTENCIA #4

#	Token	Etiqueta
130	Originality	NN
133	:	:
136	OCW-Universia	NP
139	linked	VVN
142	new	JJ
145	the	DT
148	including	VVG
151	mechanisms	NNS
154	.	SENT

#	Token	Etiqueta
131	/	SYM
134	The	DT
137	integration	NN
140	data	NN
143	features	NNS
146	initial	JJ
149	improved	VVN
152	and	CC

#	Token	Etiqueta
132	value	NN
135	new	JJ
138	with	IN
141	adds	VVZ
144	to	TO
147	framework	NN
150	query	NN
153	interoperability	NN

SENTENCIA #1

#	Token	Etiqueta
1	Purpose	NN
4	aim	NN
7	paper	NN
10	present	VV
13	to	TO
16	principles	NNS
19	Data	NP

#	Token	Etiqueta
2	:	:
5	of	IN
8	is	VBZ
11	an	DT
14	apply	VV
17	of	IN
20	to	TO

#	Token	Etiqueta
3	The	DT
6	this	DT
9	to	TO
12	initiative	NN
15	the	DT
18	Linked	NP
21	enhance	VV

22	the	DT
25	discovery	NN
28	((
31	contents	NNS
34	shared	VVN
37	universities	NNS

23	search	NN
26	of	IN
29	OCW	NP
32	created	VVN
35	by	IN
38	.	SENT

24	and	CC
27	OpenCourseWare	NP
30))
33	and	CC
36	the	DT

SENTENCIA #2

#	Token	Etiqueta
39	Design	NN
42	/	SYM
45	This	DT
48	a	DT
51	of	IN
54	data	NNS
57	be	VB
60	the	DT
63	open	JJ
66	.	SENT

#	Token	Etiqueta
40	/	SYM
43	approach	NN
46	paper	NN
49	case	NN
52	how	WRB
55	technologies	NNS
58	applied	VVN
61	enhancement	NN
64	learning	NN

#	Token	Etiqueta
41	methodology	NN
44	:	:
47	is	VBZ
50	study	NN
53	linked	VVN
56	can	MD
59	for	IN
62	of	IN
65	contents	NNS

SENTENCIA #3

#	Token	Etiqueta
67	Findings	NNS
70	presented	VVN
73	umbrella	NN
76	consortium	NN
79	the	DT
82	access	NN
85	from	IN
88	OCW	NN

#	Token	Etiqueta
68	:	:
71	under	IN
74	of	IN
77	,	,
80	integration	NN
83	to	TO
86	different	JJ
89	and	CC

#	Token	Etiqueta
69	Results	NNS
72	the	DT
75	OCW-Universia	NP
78	as	IN
81	and	CC
84	content	VV
87	repositories	NNS
90	the	DT

91	development	NN
94	query	NN
97	access	VV
100	,	,
103	linked	VVN
106	offer	VV
109	to	TO
112	select	VV
115	open	JJ
118	,	,
121	are	VBP
124	the	DT
127	data	NNS

92	of	IN
95	method	NN
98	these	DT
101	reveal	VVP
104	data	NNS
107	a	DT
110	filter	VV
113	semantically	RB
116	educational	JJ
119	and	CC
122	linked	VVN
125	linked	VVN
128	cloud	VVP

93	a	DT
96	to	TO
99	data	NNS
102	that	IN/that
105	would	MD
108	solution	NN
111	and	CC
114	those	DT
117	contents	NNS
120	automatically	RB
123	to	TO
126	open	JJ
129	.	SENT

SENTENCIA #4

#	Token	Etiqueta
130	Originality	NN
133	:	:
136	OCW-Universia	NP
139	linked	VVN
142	new	JJ
145	the	DT
148	including	VVG
151	mechanisms	NNS
154	.	SENT

#	Token	Etiqueta
131	/	SYM
134	The	DT
137	integration	NN
140	data	NN
143	features	NNS
146	initial	JJ
149	improved	VVN
152	and	CC

#	Token	Etiqueta
132	value	NN
135	new	JJ
138	with	IN
141	adds	VVZ
144	to	TO
147	framework	NN
150	query	NN
153	interoperability	NN

9.3. Extracción

ENTIDADES

#	Entidades
1	Linked Data
4	OCW-Universia

#	Entidades
2	OpenCourseWare
5	new OCW-Universia

#	Entidades
3	OCW

KEYWORDS COMPUESTAS

#	Keywords Compuestas
1	case study
4	different repositories OCW
7	open data
10	query mechanisms

#	Keywords Compuestas
2	data technologies
5	query method
8	new features

#	Keywords Compuestas
3	open learning contents
6	open educational contents
9	initial framework

KEYWORDS SIMPLES

#	Keywords Simples
1	Purpose
4	initiative
7	discovery
10	Design
13	paper
16	data
19	learning
22	Results
25	integration
28	OCW
31	method
34	solution
37	Originality
40	data
43	query

#	Keywords Simples
2	aim
5	principles
8	contents
11	methodology
14	case
17	technologies
20	contents
23	umbrella
26	access
29	development
32	data
35	contents
38	value
41	features
44	mechanisms

#	Keywords Simples
3	paper
6	search
9	universities
12	approach
15	study
18	enhancement
21	Findings
24	consortium
27	repositories
30	query
33	data
36	data
39	integration
42	framework
45	interoperability

9.4. Desambiguación y Enlace

ENTIDADES

#	Entidad	Tipo	Enlace
1	Linked Data		http://dbpedia.org/resource/Linked_data
2	OCW		http://dbpedia.org/resource/OpenCourseWare
3	OpenCourseWare		http://dbpedia.org/resource/OpenCourseWare

KEYWORD COMPUESTAS

KEYWORDS SIMPLES

#	Keyword	Tipo	Enlace
1	universities		http://dbpedia.org/resource/Universities
2	principles		http://dbpedia.org/resource/Principles
3	paper		http://dbpedia.org/resource/Paper
4	initiative		http://dbpedia.org/resource/Initiative
5	Purpose		http://dbpedia.org/resource/Purpose
6	data		http://dbpedia.org/resource/Data
7	methodology		http://dbpedia.org/resource/Methodology
8	Design		http://dbpedia.org/resource/Design
9	learning		http://dbpedia.org/resource/Learning
10	technologies		http://dbpedia.org/resource/Technologies
11	umbrella		http://dbpedia.org/resource/Umbrella
12	Findings		http://dbpedia.org/resource/Findings
13	consortium		http://dbpedia.org/resource/Consortium
14	solution		http://dbpedia.org/resource/Solution
15	Results	http://www.w3.org/2002/07/owl#Thing http://dbpedia.org/ontology/Album http://dbpedia.org/ontology/MusicalWork http://dbpedia.org/ontology/Work http://schema.org/CreativeWork http://schema.org/MusicAlbum	http://dbpedia.org/resource/Results
16	method		http://dbpedia.org/resource/Methodology
17	Originality		http://dbpedia.org/resource/Originality
18	interoperability		http://dbpedia.org/resource/Interoperability

10. Anexo 10: Descripción del Proyecto SMARTLAND

10.1. Sentencias

#	Sentencias
1	By "Smart land" we are referring to an area in which ICT is used intensively in order to improve both the quality of life of the citizens, as well as environmental management.
2	Smart land broadens the concept of "Smart Cities" which focus on ensuring the intelligent development of cities within a certain territory.
3	This initiative uses techniques including, advances in digital preservation, the representation and retrieval of information, the processing of large volumes of data, variable analysis techniques, sensor technologies, geographical information systems, information visualisation, and technologies emerging from the Semantic Web.
4	Smart Land is a UTPL initiative, which, through research projects (with the participation of more than 300 researchers), collects, manages and creates social, biological, environmental, cultural, and infrastructural data and indicator models, in order to stimulate innovative management systems for the area.
5	Its purpose is to contribute to the optimal use of natural and cultural resources, reassessing their wealth so that in the medium term, they can be a support tool for decision-making, aimed at improving the quality of life of its inhabitants.
6	This is achieved in collaboration with the regional, local and national governments, as well as public and private companies who wish to join the Smart Land project.
7	This initiative's most significant Primary Research Area is the Province of Zamora-Chinchipe.
8	The immeasurable biodiversity, the ethnic groups, the rich archaeological heritage, the mineral wealth, and the tourist attractions, which include rivers, waterfalls, lagoons and orography, make it a quite unique area in the world due to the prevalence of the Amazon Rainforest alongside the foothills of the Andes.
9	Several studies have shown that the biodiversity in southern Ecuador maintains a priceless natural heritage, and is considered as being a "hot spot" of biological diversity, resulting in UNESCO declaring it the "Podocarpus Biosphere Reserve-EI Condor" (PBRC).
10	This ecological and cultural wonder consists of an area of approximately 1.14 million acres.

10.2. Tokens y etiquetas

SENTENCIA #1

#	Token	Etiqueta	#	Token	Etiqueta	#	Token	Etiqueta
1	By	IN	2	?Smart	NN	3	land?	NN
4	we	PP	5	are	VBP	6	referring	VVG
7	to	TO	8	an	DT	9	area	NN
10	in	IN	11	which	WDT	12	ICT	NP

13	is	VBZ
16	in	IN
19	improve	VV
22	quality	NN
25	of	IN
28	,	,
31	as	IN
34	.	SENT

14	used	VVN
17	order	NN
20	both	CC
23	of	IN
26	the	DT
29	as	RB
32	environmental	JJ

15	intensively	RB
18	to	TO
21	the	DT
24	life	NN
27	citizens	NNS
30	well	RB
33	management	NN

SENTENCIA #2

#	Token	Etiqueta
35	Smart	JJ
38	the	DT
41	"	``
44	"	"
47	on	IN
50	intelligent	JJ
53	cities	NNS
56	certain	JJ

#	Token	Etiqueta
36	land	NN
39	concept	NN
42	Smart	NP
45	which	WDT
48	ensuring	VVG
51	development	NN
54	within	IN
57	territory	NN

#	Token	Etiqueta
37	broadens	VVZ
40	of	IN
43	Cities	NPS
46	focus	VVP
49	the	DT
52	of	IN
55	a	DT
58	.	SENT

SENTENCIA #3

#	Token	Etiqueta
59	This	DT
62	techniques	NNS
65	advances	NNS
68	preservation	NN
71	representation	NN
74	of	IN
77	the	DT
80	large	JJ
83	data	NNS

#	Token	Etiqueta
60	initiative	NN
63	including	VVG
66	in	IN
69	,	,
72	and	CC
75	information	NN
78	processing	NN
81	volumes	NNS
84	,	,

#	Token	Etiqueta
61	uses	VVZ
64	,	,
67	digital	JJ
70	the	DT
73	retrieval	NN
76	,	,
79	of	IN
82	of	IN
85	variable	JJ

86	analysis	NN
89	sensor	NN
92	geographical	JJ
95	,	,
98	,	,
101	emerging	VVG
104	Semantic	NP

87	techniques	NNS
90	technologies	NNS
93	information	NN
96	information	NN
99	and	CC
102	from	IN
105	Web	NP

88	,	,
91	,	,
94	systems	NNS
97	visualisation	NN
100	technologies	NNS
103	the	DT
106	.	SENT

SENTENCIA #4

#	Token	Etiqueta
107	Smart	NP
110	a	DT
113	,	,
116	through	IN
119	((
122	participation	NN
125	than	IN
128))
131	,	,
134	creates	VVZ
137	biological	JJ
140	,	,
143	and	CC
146	and	CC
149	,	,
152	to	TO
155	management	NN
158	the	DT

#	Token	Etiqueta
108	Land	NP
111	UTPL	NP
114	which	WDT
117	research	NN
120	with	IN
123	of	IN
126	300	CD
129	,	,
132	manages	VVZ
135	social	JJ
138	,	,
141	cultural	JJ
144	infrastructural	NN
147	indicator	NN
150	in	IN
153	stimulate	VV
156	systems	NNS
159	area	NN

#	Token	Etiqueta
109	is	VBZ
112	initiative	NN
115	,	,
118	projects	NNS
121	the	DT
124	more	JJR
127	researchers	NNS
130	collects	VVZ
133	and	CC
136	,	,
139	environmental	JJ
142	,	,
145	data	NNS
148	models	NNS
151	order	NN
154	innovative	JJ
157	for	IN
160	.	SENT

SENTENCIA #5

#	Token	Etiqueta
161	Its	PP\$
164	to	TO
167	the	DT
170	of	IN
173	cultural	JJ
176	reassessing	VVG
179	so	RB
182	the	DT
185	,	,
188	be	VB
191	tool	NN
194	,	,
197	improving	VVG
200	of	IN
203	its	PP\$

#	Token	Etiqueta
162	purpose	NN
165	contribute	VV
168	optimal	JJ
171	natural	JJ
174	resources	NNS
177	their	PP\$
180	that	IN/that
183	medium	JJ
186	they	PP
189	a	DT
192	for	IN
195	aimed	VVN
198	the	DT
201	life	NN
204	inhabitants	NNS

#	Token	Etiqueta
163	is	VBZ
166	to	TO
169	use	NN
172	and	CC
175	,	,
178	wealth	NN
181	in	IN
184	term	NN
187	can	MD
190	support	NN
193	decision-making	NN
196	at	IN
199	quality	NN
202	of	IN
205	.	SENT

SENTENCIA #6

#	Token	Etiqueta
206	This	DT
209	in	IN
212	the	DT
215	local	JJ
218	governments	NNS
221	well	RB
224	and	CC
227	who	WP
230	join	VV
233	Land	NP

#	Token	Etiqueta
207	is	VBZ
210	collaboration	NN
213	regional	JJ
216	and	CC
219	,	,
222	as	IN
225	private	JJ
228	wish	VVP
231	the	DT
234	project	NN

#	Token	Etiqueta
208	achieved	VVN
211	with	IN
214	,	,
217	national	JJ
220	as	RB
223	public	JJ
226	companies	NNS
229	to	TO
232	Smart	NP
235	.	SENT

SENTENCIA #7

#	Token	Etiqueta
236	This	DT
239	s	NNS
242	Primary	NP
245	is	VBZ
248	of	IN

#	Token	Etiqueta
237	initiative	NN
240	most	RBS
243	Research	NP
246	the	DT
249	Zamora-Chinchi	NP

#	Token	Etiqueta
238	'	NN
241	significant	JJ
244	Area	NP
247	Province	NP
250	.	SENT

SENTENCIA #8

#	Token	Etiqueta
251	The	DT
254	,	,
257	groups	NNS
260	rich	JJ
263	,	,
266	wealth	NN
269	the	DT
272	,	,
275	rivers	NNS
278	,	,
281	orography	NN
284	it	PP
287	unique	JJ
290	the	DT
293	to	TO
296	of	IN
299	Rainforest	NP
302	foothills	NNS
305	Andes	NP

#	Token	Etiqueta
252	immeasurable	JJ
255	the	DT
258	,	,
261	archaeological	JJ
264	the	DT
267	,	,
270	tourist	NN
273	which	WDT
276	,	,
279	lagoons	NNS
282	,	,
285	a	DT
288	area	NN
291	world	NN
294	the	DT
297	the	DT
300	alongside	IN
303	of	IN
306	.	SENT

#	Token	Etiqueta
253	biodiversity	NN
256	ethnic	JJ
259	the	DT
262	heritage	NN
265	mineral	NN
268	and	CC
271	attractions	NNS
274	include	VVP
277	waterfalls	NNS
280	and	CC
283	make	VV
286	quite	RB
289	in	IN
292	due	JJ
295	prevalence	NN
298	Amazon	NP
301	the	DT
304	the	DT

SENTENCIA #9

#	Token	Etiqueta
307	Several	JJ
310	shown	VVN
313	biodiversity	NN
316	Ecuador	NP
319	priceless	JJ
322	,	,
325	considered	VVN
328	a	DT
331	spot	NN
334	biological	JJ
337	resulting	VVG
340	declaring	VVG
343	"	``
346	Reserve-EI	NP
349	((
352	.	SENT

#	Token	Etiqueta
308	studies	NNS
311	that	IN/that
314	in	IN
317	maintains	VVZ
320	natural	JJ
323	and	CC
326	as	IN
329	"	``
332	"	"
335	diversity	NN
338	in	IN
341	it	PP
344	Podocarpus	NP
347	Condor	NP
350	PBRC	NP

#	Token	Etiqueta
309	have	VHP
312	the	DT
315	southern	JJ
318	a	DT
321	heritage	NN
324	is	VBZ
327	being	VBG
330	hot	JJ
333	of	IN
336	,	,
339	UNESCO	NP
342	the	DT
345	Biosphere	NP
348	"	"
351))

SENTENCIA #10

#	Token	Etiqueta
353	This	DT
356	cultural	JJ
359	of	IN
362	of	IN
365	million	CD

#	Token	Etiqueta
354	ecological	JJ
357	wonder	NN
360	an	DT
363	??approximately	RB
366	acres	NNS

#	Token	Etiqueta
355	and	CC
358	consists	VVZ
361	area	NN
364	1.14	CD
367	.	SENT

10.3. Extracción

ENTIDADES

#	Entidades
1	ICT
4	Smart Land
7	significant Primary Research Area
10	Zamora-Chinchi
13	southern Ecuador
16	PBRC

#	Entidades
2	Smart Cities
5	UTPL
8	Province of Zamora-Chinchi
11	Amazon Rainforest
14	UNESCO

#	Entidades
3	Semantic Web
6	Smart Land
9	Province
1 2	Andes
1 5	Podocarpus Biosphere Reserve-El Condor

KEYWORDS COMPUESTAS

#	Keywords Compuestas
1	?Smart land?
4	intelligent development
7	large volumes
10	geographical information systems
13	infrastructural data
16	optimal use
19	support tool
22	initiative ´ s
25	rich archaeological heritage
28	unique area
31	hot spot

#	Keywords Compuestas
2	environmental management
5	certain territory
8	variable analysis techniques
11	information visualisation
14	indicator models
17	cultural resources
20	national governments
23	immeasurable biodiversity
26	mineral wealth
29	Several studies
32	biological diversity

#	Keywords Compuestas
3	Smart land
6	digital preservation
9	sensor technologies
12	research projects
15	innovative management systems
18	medium term
21	private companies
24	ethnic groups
27	tourist attractions
30	priceless natural heritage
33	cultural wonder

KEYWORDS SIMPLES

#	Keywords Simples
1	?Smart

#	Keywords Simples
2	land?

#	Keywords Simples
3	area

4	order
7	citizens
10	concept
13	territory
16	advances
19	retrieval
22	volumes
25	techniques
28	information
31	visualisation
34	research
37	researchers
40	indicator
43	management
46	purpose
49	wealth
52	tool
55	life
58	governments
61	initiative
64	groups
67	wealth
70	rivers
73	orography
76	prevalence
79	biodiversity
82	diversity
85	acres

5	quality
8	management
11	development
14	initiative
17	preservation
20	information
23	data
26	sensor
29	systems
32	technologies
35	projects
38	infrastructural
41	models
44	systems
47	use
50	term
53	decision-making
56	inhabitants
59	companies
62	s
65	heritage
68	tourist
71	waterfalls
74	area
77	foothills
80	heritage
83	wonder

6	life
9	land
12	cities
15	techniques
18	representation
21	processing
24	analysis
27	technologies
30	information
33	initiative
36	participation
39	data
42	order
45	area
48	resources
51	support
54	quality
57	collaboration
60	project
63	biodiversity
66	mineral
69	attractions
72	lagoons
75	world
78	studies
81	spot
84	area

10.4. Desambiguación y Enlace

ENTIDADES

#	Entidad	Tipo	Enlace
1	ICT		http://dbpedia.org/resource/Information_and_communication_technologies_for_environmental_sustainability
2	Smart Cities		http://dbpedia.org/resource/Smart_city
3	Semantic Web		http://dbpedia.org/resource/Semantic_Web
4	Zamora-Chinchipe	http://www.w3.org/2002/07/owl#Thing http://dbpedia.org/ontology/Place http://dbpedia.org/ontology/PopulatedPlace http://dbpedia.org/ontology/Settlement http://schema.org/Place	http://dbpedia.org/resource/Zamora-Chinchipe_Province
5	Province		http://dbpedia.org/resource/Province
6	Amazon Rainforest	http://www.w3.org/2002/07/owl#Thing http://dbpedia.org/ontology/Place http://schema.org/Place	http://dbpedia.org/resource/Amazon_rainforest
7	Andes		http://dbpedia.org/resource/Andes
8	UNESCO		http://dbpedia.org/resource/UNESCO
9	PBRC		http://dbpedia.org/resource/Packed_red_blood_cells

KEYWORD COMPUESTAS

#	Keywords	Tipo	Enlace
1	environmental management		http://dbpedia.org/resource/Environmental_resources_management
2	digital preservation		http://dbpedia.org/resource/Digital_preservation
3	geographical information systems		http://dbpedia.org/resource/Geographic_information_system
4	ethnic groups		http://dbpedia.org/resource/Ethnic_group

KEYWORDS SIMPLES

#	Keyword	Tipo	Enlace
1	life		http://dbpedia.org/resource/Life
2	management		http://dbpedia.org/resource/Management
3	area		http://dbpedia.org/resource/Area

4	citizens		http://dbpedia.org/resource/Citizens
5	concept		http://dbpedia.org/resource/Concept
6	cities		http://dbpedia.org/resource/Cities
7	information		http://dbpedia.org/resource/Information
8	sensor		http://dbpedia.org/resource/Sensor
9	systems		http://dbpedia.org/resource/Systems
10	initiative		http://dbpedia.org/resource/Initiative
11	analysis		http://dbpedia.org/resource/Analysis
12	volumes		http://dbpedia.org/resource/Volumes
13	technologies		http://dbpedia.org/resource/Technologies
14	data		http://dbpedia.org/resource/Data
15	research		http://dbpedia.org/resource/Research
16	infrastructural		http://dbpedia.org/resource/Infrastructural
17	projects		http://dbpedia.org/resource/Projects
18	researchers		http://dbpedia.org/resource/Researchers
19	wealth		http://dbpedia.org/resource/Wealth
20	inhabitants		http://dbpedia.org/resource/Inhabitants
21	decision-making		http://dbpedia.org/resource/Decision-Making
22	purpose		http://dbpedia.org/resource/Purpose
23	tool		http://dbpedia.org/resource/Tool
24	resources		http://dbpedia.org/resource/Resources
25	project		http://dbpedia.org/resource/Projects
26	collaboration		http://dbpedia.org/resource/Collaboration
27	companies		http://dbpedia.org/resource/Companies
28	governments		http://dbpedia.org/resource/Governments
29	s		http://dbpedia.org/resource/Citizens
30	biodiversity		http://dbpedia.org/resource/Biodiversity
31	lagoons		http://dbpedia.org/resource/Lagoons
32	tourist		http://dbpedia.org/resource/Tourist

33	attractions		http://dbpedia.org/resource/Attractions
34	waterfalls		http://dbpedia.org/resource/Waterfalls
35	foothills		http://dbpedia.org/resource/Foothills
36	orography		http://dbpedia.org/resource/Orography
37	rivers		http://dbpedia.org/resource/Rivers
38	prevalence		http://dbpedia.org/resource/Prevalence
39	world		http://dbpedia.org/resource/World
40	mineral		http://dbpedia.org/resource/Mineral
41	diversity		http://dbpedia.org/resource/Biodiversity
42	acres		http://dbpedia.org/resource/Acres