

Presentación de Trabajo de Fin de Máster
**PROPUESTA DE BÚSQUEDA SEMÁNTICA:
APLICACIÓN AL CATÁLOGO DE MAPAS,
PLANOS Y DIBUJOS DEL ARCHIVO GENERAL
DE SIMANCAS**

Máster en Lenguajes y Sistemas Informáticos: Tecnologías del Lenguaje en la Web

Universidad de Educación a Distancia

Marzo 2013

Autor: José Alberto Benítez Andrades

Directora: Ana M^a García Serrano



Introducción

► Objetivos

- Plantear un trabajo de investigación aplicada y una metodología de tratamiento de información estructurada
- Comparar tres aproximaciones en un dominio concreto

► Estado del arte, 2 perspectivas

► Tecnológica

► Teórico-práctica

- Uso de recursos de la ingeniería lingüística y RI.
- Uso de ontologías en el área de la Ingeniería del Conocimiento.



Introducción – Trabajo teórico práctico

- La metodología abordada del trabajo realizado es
 - Análisis principal del problema a resolver
 - Estudio de las posibles soluciones existentes para resolver dicho problema
 - Selección de la solución o soluciones
 - Elección de las herramientas necesarias para aplicar dicha solución
 - Desarrollo de herramientas necesarias para abarcar el problema
 - Fase de aplicación de la solución: Experimentación y análisis de resultados.



Parte 1: Panorama Tecnológico

- Web Semántica
 - Concepto, herramientas, estándares
 - Definición de ontologías: Protégé
 - Estándares para descripción de contenidos: RDF DC y OWL
 - Método para extraer relaciones semánticas desde la Wikipedia
- Propuesta de un buscador ontológico en Protégé y traductor de RDF DC a OWL necesario para el conjunto de documentos (fichas) estructurados.



Parte 1: Panorama Tecnológico

- Recuperación de información
 - Modelos clásicos de RI
 - Estándares y herramientas relacionadas con Lucene, Solr, Apache y Sparql.
 - Estudio de trabajos relacionados con
 - Los dirigidos por P.Castells sobre búsqueda semántica
 - Sistemas pregunta-respuesta (Question Answering o QA) dirigidos por P. Martínez Barco.
- **Propuesta: Metodología para construir una versión de un conjunto de documentos (fichas) estructurados y comparación de un modelo RI facetado y otro textual**



Parte 1: Análisis del dominio

- ▶ Estructura de las fichas del catálogo **mapas, planos y dibujos del Archivo General de Simancas (AGS)**
 - ▶ Poseen 8 campos, facetas:
 - ▶ Fecha, Referencias, Creador, Tipo, Idiomas, Temática, Técnica Utilizada
 - ▶ El primer paso: descargar las fichas en dos formatos
 - ▶ Texto Plano
 - ▶ RDF DC

Parte 1: Análisis del dominio



The screenshot shows the website interface for the 'Catálogo Colectivo de la Red de Bibliotecas de los Archivos Estatales'. The header includes the Spanish government logo and navigation links in various languages. The main content area displays search results for 'Mapas, Planos y Dibujos'. A sidebar on the left shows a 'Nube de materias' (Subject Cloud) with categories like 'Lengua' (Language) and 'Pertenece a' (Belongs to). The main results list includes items such as 'Planos y Perfiles de la obra que se construye antes de empezarse a fundir para recipientes del Mineral y que...' (1762) and 'Pamplona. Hospitales. Planos. 1721. (1ª 2ª y 3ª Planos según el proyecto del Ingeniero)'.

➤ http://www.mcu.es/ccbae/es/consulta/resultados_busqueda.cmd?tipo_busqueda=mapas_planos_dibujos&posicion=1&id=30485


Parte 1: Análisis del dominio

Formato: **Ficha** [MARC XML](#) [MODS](#) [BartIX](#) [ISAC](#) [Volver a resultados](#)
[Enlace persistente](#)

1 de 7792


 [...] Planos y Perfiles de la obra que se construye antes de empezarse a fundir para recipientes del Mineral y que... (1762) 

Disponibles

Sección:	Material gráfico AGS
Número de control:	BA020100036921
Título:	[...] Planos y Perfiles de la obra que se construye antes de empezarse a fundir para recipientes del Mineral y que se deshace despues de evacuada la fundición que es la que se supone llamarse Crisol [Material gráfico no proyectable]
Publicación:	[1762]
Descripción física:	1 Dibujo: ms., col.; 33 x 39 cm
Notas:	Referencias: Mapas, planos y dibujos (Años 1503-1805). Volumen I : p. 405 Tinta y colores. Con explicación Manuscrito sobre papel AGS. Secretaría de Marina, 00679. Acompaña a carta de don Maximino de la Croix a don Ricardo Wall, Chaves, 16 de julio de 1762.
Materia / geográfico:	Hornos metalúrgicos-S.XVIII-Dibujos 
Tipo de publicación:	 Ilustraciones y Fotos
Ejemplares:	Archivo General de Simancas. Signatura: MPD, 95, 212. Ubicación Anterior: SMA, 00679. Préstamo: <input checked="" type="checkbox"/> Disponible <input type="checkbox"/> Copia Digital

1 de 7792

© Ministerio de Educación, Cultura y Deporte
Nota legal [Accesibilidad](#)



Parte 2: Propuesta

- ▶ Necesidad de realizar un buscador semántico, para mejorar la búsqueda en las fichas.
 - ▶ Usuarios no informáticos.
 - ▶ Desconocedores de la estructura.
- ▶ Planteamiento del trabajo de investigación para conseguirlo:
 1. Desarrollo del entorno para la experimentación
 2. Extracción de información desde el AGS en el formato deseado
 3. Realización de tres aproximaciones de almacenamiento del catálogo
 4. Generación de conjunto de preguntas y comparación entre buscador textual (facetado y no) y ontológico.
 5. Comparación de aproximaciones con TRECEval. No se puede realizar de forma completa debido a la falta de disponibilidad de juicios de relevancia. Sin embargo se ha realizado una evaluación manual.

Parte 2: RDF DC a OWL

► Ficha en formato RDF DC

```
<?xml version="1.0" encoding="UTF-8" standalone="no"?>
<rdf:RDF xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#" xmlns:dc=
"http://purl.org/dc/elements/1.1/" xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance">
  <rdf:Description>
    <dc:relation>Referencias: Mapas, planos y dibujos (Años 1503-1805). Volumen I : p. 405
    </dc:relation>
    <dc:coverage>-S.XVIII</dc:coverage>
    <dc:title>[ ...] Planos y Perfiles de la obra que se construye antes de empezarse a
fundir para recipientes del Mineral y que se deshace despues de evacuada la fundición
que es la que se supone llamarse Crisol [Material gráfico no proyectable]</dc:title>
    <dc:description>AGS. Secretaria de Marina, 00679. Acompaña a carta de don Maximino de
la Croix a don Ricardo Wall, Chaves, 16 de julio de 1762</dc:description>
    <dc:description>Tinta y colores. Con explicación</dc:description>
    <dc:description>Manuscrito sobre papel.</dc:description>
    <dc:type>Ilustraciones y Fotos</dc:type>
    <dc:language>spa</dc:language>
    <dc:date>1762</dc:date>
    <dc:identifier>http://www.mcu.es/ccbae/es/consulta/registro.cmd?id=178725</dc:identifier>
    <dc:format>image/jpeg</dc:format>
    <dc:subject>Hornos metalúrgicos-Dibujos</dc:subject>
  </rdf:Description>
```

Parte 2: RDF DC a OWL

- Se selecciona el fichero origen
- Se selecciona el fichero de destino
- Se define el nombre del objeto
- Se definen los campos identificadores
- Se agregan las propiedades
- Pulsamos convertir

The screenshot shows the 'ParserRDFtoOWL' application window. The title bar reads 'ParserRDFtoOWL por @jabenitez88'. The interface includes a menu bar with 'Archivo' and 'Ayuda'. Below the menu bar, there are two text input fields: 'Fichero Origen' and 'Fichero Destino', each with a 'Seleccionar fichero' button. A third text input field is labeled 'Objeto en Fichero Origen (por ejemplo rdf.Description en RDF Dublin Core)'. Below this, there are two more text input fields: 'Campo de ID en Fichero Origen' (containing 'dc:identifier') and 'Nombre de Objeto en OWL' (containing 'Ficha'). At the bottom, there is a table with two columns: 'Propiedad Anterior' and 'Propiedad Nueva'. The table lists several mappings: 'dc:type' to 'Tipo', 'dc:format' to 'Formato', 'dc:coverage' to 'Materia', 'dc:relation' to 'Referencias', 'dc:titulo' to 'Titulo', and 'dc:description' to 'Notas'. Below the table are two buttons: 'Añadir Propiedades' and 'Borrar Props Seleccionadas'. At the very bottom, there is a 'Convertir' button and a checkbox labeled 'Debugger activado (texto en consola)'.

Propiedad Anterior	Propiedad Nueva
dc:type	Tipo
dc:format	Formato
dc:coverage	Materia
dc:relation	Referencias
dc:titulo	Titulo
dc:description	Notas

Parte 2: RDF DC a OWL

- Ficha en formato final OWL

```
<Ficha rdf:about="http://www.mcu.es/ccbae/es/consulta/registro.cmd?id=178725">
  <Titulo> Planos y Perfiles de la obra que se construye antes de empezarse a fundir para
  recipientes del Mineral y que se deshace despues de evacuada la fundición que es la que se supone
  llamarse Crisol [Material gráfico no proyectable]</Titulo>
  <Referencias>Referencias: Mapas, planos y dibujos (Años 1503-1805). Volumen I : p.
  405</Referencias>
  <Materia>-S.XVIII</Materia>
  <Notas>AGS. Secretaría de Marina, 00679. Acompaña a carta de don Maximino de la Croix a
  don Ricardo Wall, Chaves, 16 de julio de 1762</Notas>
  <Notas>Tinta y colores. Con explicación</Notas>
  <Notas>Manuscrito sobre papel.</Notas>
  <Tipo>Ilustraciones y Fotos</Tipo>
  <Idioma>spa</Idioma>
  <Publicacion>1762</Publicacion>
  <Formato>image/jpeg</Formato>
  <Tematica>Hornos metalúrgicos-Dibujos</Tematica>
</Ficha>
```

Parte 2: Modelo Ontológico con Protégé

- Información de las entidades

The screenshot displays the Protégé ontology editor interface. The main window shows the ontology 'datos_total' with a search bar and various toolbars. The left pane displays the 'Class hierarchy' for 'Ficha', showing a tree structure with 'Thing' as the root, followed by 'DIRECTED-BINARY-RELATION', 'Ficha', and 'PAL-CONSTRAINT'. The bottom-left pane shows 'Individuals by type: registro.cmd?id=181571', listing 15 instances of the 'Ficha' class, each with a unique ID. The right pane shows 'Individual Annotations' for 'registro.cmd?id=181571', listing properties such as 'Formato', 'image/jpeg', 'Idioma', 'spa', 'Materia', 'España-Pais Vasco-Guipúzcoa-Hondarribia', 'Notas', 'Manuscrito sobre papel', 'Tinta y colores. Con rotulación', and 'Publicacion'. The bottom-right pane shows 'Property assertions' for the same individual, including 'Object property assertions', 'Data property assertions', 'Negative object property assertions', and 'Negative data property assertions'.



Parte 2: Segundo modelo de búsqueda semántica

- ▶ ¿Ventajas de búsqueda ontológica con Protégè vs otro tipo de buscador?
 - ▶ Realización importación de fichas en Apache SOLR
- ▶ La fase de aplicación a la solución se subdividió en:
 - ▶ Indexación
 - ▶ Generación de preguntas de tipo QA
 - ▶ Realización de búsquedas de elementos de forma facetada y de forma textual
 - ▶ Comparación de los resultados obtenidos

Parte 2: Clasificación de consultas (ej)


- Clases: Qué, quién, dónde, cuándo, cómo

DÓNDE	se realizó / se utilizó / se realizó	El/la/los/las Algún/alguno/ algunos/ algunas un/una/ unos/unas <i>planos/ dibujos mapas</i>	por el autor \$AUTOR con el título \$TÍTULO con los colores \$COLOR entre las fechas \$FECHA y \$FECHA en el idioma \$IDIOMA en los idiomas \$IDIOMA y \$IDIOMA en la época \$ÉPOCA con la temática \$TEMÁTICA sobre el soporte \$SOPORTE	[Ciudad] [País] [Continente]
CUÁNDO	se realizó / se utilizó / se realizó	El/la/los/las Algún/alguno/ algunos/ algunas un/una/ unos/unas <i>planos/ dibujos mapas</i>	por el autor \$AUTOR con el título \$TÍTULO con los colores \$COLOR entre las fechas \$FECHA y \$FECHA en el idioma \$IDIOMA en los idiomas \$IDIOMA y \$IDIOMA en la época \$ÉPOCA con la temática \$TEMÁTICA sobre el soporte \$SOPORTE sobre la ciudad \$CIUDAD sobre el país de \$PAÍS	[Época] [Año]

Parte 2: Consultas en Solr

Solr Admin (example core one)

licantropo.Benired10B:80
cwd=/var/lib/tomcat6 SolrHome=/var/solr/fichasFacetadas/
HTTP caching is ON



Request Handler	<input type="text" value="/select"/>
Query String	<input type="text" value="titulo:*1950*"/>
Filter Query	<input type="text"/>
Start Row	<input type="text" value="0"/>
Maximum Rows Returned	<input type="text" value="10"/>
Fields to Return	<input type="text" value="titulo,score,id"/>
Output Type	<input type="text" value="xml"/>
Debug: enable	<input type="checkbox"/> <small>Note: you may need to "view source" in your browser to see explain() correctly indented.</small>
Debug: explain others	<input type="text"/> <small>Apply original query scoring to matches of this query to see how they compare.</small>
Enable Highlighting	<input checked="" type="checkbox"/>
Fields to Highlight	<input type="text" value="titulo"/>
<input type="button" value="Search"/>	

This form demonstrates the most common query options available for the built in Query Types. Please consult the Solr Wiki for additional Query Parameters.

http://casa.jabenitez.com/solr/fichasFacetadas/select?indent=on&version=2.2&q=titulo%3A*1952*&fq=&start=0&rows=10&fl=titulo%2Cscore&wt=&explainOther=&hl=on&hl.fl=titulo

Parte 2: Consultas SPARQL en Protégé

- ▶ Ejemplo de una de las consultas
- ▶ Resto, ver PDF páginas 111 a 115

Consulta en Protégé para **Qué fichas datan sobre una obra realizada en 1950:**

```
SELECT distinct ?s ?p ?o
WHERE {
    ?s ?p ?o .
    ?s a :Ficha .
    FILTER (regex(?o, "^1950"))}
```

Resultados obtenidos 4:

Results		
s	p	o
◆ http://www.mcu.es/ccbae/es/consulta/registro.cmd?id=182452	Publicacion	1950
◆ http://www.mcu.es/ccbae/es/consulta/registro.cmd?id=182451	Publicacion	1950
◆ http://www.mcu.es/ccbae/es/consulta/registro.cmd?id=182449	Publicacion	1950
◆ http://www.mcu.es/ccbae/es/consulta/registro.cmd?id=182450	Publicacion	1950

Parte 2: Buscador textual

► **buscadorTextual.js**

- Se crea un objeto AJAX: realizará la consulta a SOLR y recibirá la respuesta en JSON
- Almacenamos los datos en una variable llamada **datos_formulario**
- Obtenida la cadena, asignamos la propiedad URL a nuestro objeto AJAX con el siguiente valor:
 - http://casa.jabenitez.com/solr/fichas/select?indent=on&version=2.2&q='+datos_formulario+'&fq=&start=0&rows=7500&fl=id%2Cscore&wt=json&explainOther=&hl=on&hl.fl=titulo
- Un ejemplo práctico: Supongamos que el usuario desea tener la siguiente información: ¿Qué mapas fueron realizados en el año 1950 por el autor Juan Baptista?
- En el caso de nuestro buscador semántico, recogerá la cadena de texto introducida por el usuario y nos devolverá los resultados de buscar esa cadena completa, no sabiendo interpretar qué tipo de datos necesitaría ni de dónde debe obtenerlos (y tampoco se lo indica el usuario, que no lo conoce).

Parte 2: Buscador facetado

► **buscadorSemantico.js**

1. Teniendo en cuenta el sistema QA creado, sabemos que existen 5 tipos de pregunta (Qué, Quién, Cómo, Dónde y Cuándo).
El parser detecta en un primer análisis qué pregunta se encuentra en la cadena insertada por el usuario.
2. En un segundo procesado de la cadena, se analiza el objeto al que se hace referencia: **autor, mapa, obra, dibujo u otro campo**.
3. En una tercera fase, se detecta el verbo de la frase, siendo capaz de detectar **a qué campo** de las fichas facetadas se debe consultar.
4. Finalmente, el objeto AJAX que habíamos creado, recibe la información en formato JSON y la muestra por pantalla.

Parte 2: Parser semántico que traduce consultas a Solr

➤ Interfaz gráfica del buscador:

➤ Buscador semántico

➤ Buscador textual

The screenshot shows a web browser window with the URL `casa.jabenitez.com/index.php?modo=0`. The page has a dark blue header with two tabs: "BUSCADOR SEMÁNTICO" (selected) and "BUSCADOR TEXTUAL". Below the header, the text "TIPO DE BUSCADOR ACTIVADO: BUSCADOR SEMÁNTICO" is displayed. A search input field contains the text "Que mapas se hicieron en el año 1639 por Galvarreta" and a "Buscar" button. Below the search bar, the URL executed is shown: `http://casa.jabenitez.com/solr/fichasFacetadas/select?indent=on&version=2.2&q=tipo%3aMapas+AND+idioma%3a*+AND+creador%3aGalvarreta+OR+publicacion%3a*16*+OR+publicacion%3a*1639*+OR+titulo%3a*16*+OR+titfq=&start=0&rows=10&fl=titulo%2Ccreador%2Ctematica+AND+%2Cscore&wt=json&explainOther=&hl=on&hl.fl=titulo`. The results section states "Se encontraron 1 resultados" and shows a single result in a dashed box: "Obra con id: <http://www.mcu.es/ccbae/es/consulta/registro.cmd?id=184605>
Titulo: Bayona (Francia). Fortificaciones. Planos. 1639,[Planta de la fortificación de Bayona] [Material cartográfico]
Autor: Galvarreta



Parte 2: Experimentación y comparación

- ▶ Teniendo en cuenta los diferentes buscadores que hemos obtenido tras realizar un análisis específico de nuestro caso particular, a continuación detallaré en una tabla los distintos resultados obtenidos de las consultas realizadas en la base de datos de las fichas sin facetar y facetadas y en el ontológico.
- ▶ Por falta de juicios de relevancia, solo se pudo realizar una comparación de forma manual.
- ▶ Los análisis realizados, aunque de forma manual, tienen buenas expectativas, tanto en la precisión como en la cobertura.

Parte 2: Experimentación (ejm)

- **Consulta 1** en Solr para ¿Qué fichas datan sobre una obra realizada en 1950?:

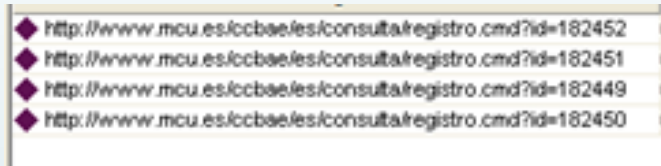
SOLR FACETADA – 4 RESULTADOS

```
<response>
  <lst name="responseHeader">
    <int name="status">0</int>
    <int name="QTime">238</int>
  </lst>
  <result name="response" numFound="4" start="0" maxScore="0.033296965">...</result>
  <lst name="highlighting">
    <lst name="http://www.mcu.es/ccbae/es/consulta/registro.cmd?id=182452"/>
    <lst name="http://www.mcu.es/ccbae/es/consulta/registro.cmd?id=182451"/>
    <lst name="http://www.mcu.es/ccbae/es/consulta/registro.cmd?id=182450"/>
    <lst name="http://www.mcu.es/ccbae/es/consulta/registro.cmd?id=182449"/>
  </lst>
</response>
```

SOLR SIN FACETAR – 9 RESULTADOS

```
<response>
  <lst name="responseHeader">
    <int name="status">0</int>
    <int name="QTime">932</int>
  </lst>
  <result name="response" numFound="9" start="0" maxScore="1.0">...</result>
  <lst name="highlighting">...</lst>
</response>
```

PROTÈGÈ – 4 RESULTADOS



◆ <http://www.mcu.es/ccbae/es/consulta/registro.cmd?id=182452>

◆ <http://www.mcu.es/ccbae/es/consulta/registro.cmd?id=182451>

◆ <http://www.mcu.es/ccbae/es/consulta/registro.cmd?id=182449>

◆ <http://www.mcu.es/ccbae/es/consulta/registro.cmd?id=182450>

Conclusiones:

- De los 9 resultados en SOLRsin facetar, 4 comunes.
- 5 Que poseen otro campo distinto a la fecha de creación que coincide con la búsqueda 1950

Parte 2: Experimentación (ej)

► **Consulta 2** en Solr para ¿Qué fichas están fabricadas en pergamino?:

SOLR FACETADA – 21 RESULTADOS

```
<response>
  <lst name="responseHeader">
    <int name="status">0</int>
    <int name="QTime">235</int>
  </lst>
  <result name="response" numFound="21" start="0" maxScore="1.0">...</result>
  <lst name="highlighting">
```

SOLR SIN FACETAR – 21 RESULTADOS

```
<response>
  <lst name="responseHeader">
    <int name="status">0</int>
    <int name="QTime">843</int>
  </lst>
  <result name="response" numFound="21" start="0" maxScore="1.0">...</result>
  <lst name="highlighting">...</lst>
</response>
```

PROTÈGÈ – 21 RESULTADOS

id
http://www.mcu.es/fcb/aeles/consulta/registro.cmd?id=177769
http://www.mcu.es/fcb/aeles/consulta/registro.cmd?id=176258
http://www.mcu.es/fcb/aeles/consulta/registro.cmd?id=181544
http://www.mcu.es/fcb/aeles/consulta/registro.cmd?id=176653
http://www.mcu.es/fcb/aeles/consulta/registro.cmd?id=180724
http://www.mcu.es/fcb/aeles/consulta/registro.cmd?id=177489
http://www.mcu.es/fcb/aeles/consulta/registro.cmd?id=176266
http://www.mcu.es/fcb/aeles/consulta/registro.cmd?id=177490
http://www.mcu.es/fcb/aeles/consulta/registro.cmd?id=176475
http://www.mcu.es/fcb/aeles/consulta/registro.cmd?id=177761

Conclusiones:

- En este caso coinciden los resultados porque la información se obtiene del campo descripción, es una búsqueda muy similar en ambas bases de datos.

Parte 2: Experimentación (ej)

- **Consulta 4** en Solr para ¿Qué fichas fueron creadas por Juan Baptista?:

SOLR FACETADA – 1 RESULTADO

```
▼<response>
  ▼<lst name="responseHeader">
    <int name="status">0</int>
    <int name="QTime">13</int>
  </lst>
  ▶<result name="response" numFound="1" start="0" maxScore="1.4142135">...</result>
  ▼<lst name="highlighting">
    <lst name="http://www.mcu.es/ccbae/es/consulta/registro.cmd?id=178297"/>
  </lst>
</response>
```

SOLR SIN FACETAR – 3 RESULTADOS

```
▼<response>
  ▼<lst name="responseHeader">
    <int name="status">0</int>
    <int name="QTime">1828</int>
  </lst>
  ▶<result name="response" numFound="3" start="0" maxScore="1.4142135">...</result>
</response>
```

PROTÈGÈ – 3 RESULTADO

Results	
	id
◆	http://www.mcu.es/ccbae/es/consulta/registro.cmd?id=180501
◆	http://www.mcu.es/ccbae/es/consulta/registro.cmd?id=178297
◆	http://www.mcu.es/ccbae/es/consulta/registro.cmd?id=176716

Conclusiones:

- De los 3 resultados, 1 de ellos común.
- Los otros dos no son obras de Juan Baptista, pero aparece en las “Notas” de la obra.

Parte 2: Experimentación (ej)

- **Consulta 6** en Solr para ¿Qué fichas poseen como temática ballenas?:

SOLR FACETADA – 1 RESULTADO

```
▼<response>
  ▼<lst name="responseHeader">
    <int name="status">0</int>
    <int name="QTime">13</int>
  </lst>
  ▶<result name="response" numFound="1" start="0">...</result>
  ▼<lst name="highlighting">
    <lst name="http://www.mcu.es/ccbae/es/consulta/registro.cmd?id=183373"/>
  </lst>
</response>
```

PROTÈGÈ – 1 RESULTADO

Results	
	id
◆	http://www.mcu.es/ccbae/es/consulta/registro.cmd?id=183373

SOLR SIN FACETAR – 1 RESULTADO

```
▼<response>
  ▼<lst name="responseHeader">
    <int name="status">0</int>
    <int name="QTime">1117</int>
  </lst>
  ▶<result name="response" numFound="1" start="0" maxScore="0.5">...</result>
</response>
```

Conclusiones:

- El resultado es el mismo que en los otros casos.

Parte 2: Experimentación (ej)

- **Consulta 7** en Solr para ¿Qué obras tienen las medidas en varas?:

SOLR FACETADA – 467 RESULTADOS

```
▼<response>
  ▼<lst name="responseHeader">
    <int name="status">0</int>
    <int name="QTime">544</int>
  </lst>
  ▶<result name="response" numFound="467" start="0" maxScore="1.0">...</result>
  ▶<lst name="highlighting">...</lst>
</response>
```

PROTÈGÈ – 467 RESULTADOS

Results	
	id
◆	http://www.mcu.es/ccbae/es/consulta/registro.cmd?id=180666
◆	http://www.mcu.es/ccbae/es/consulta/registro.cmd?id=177247
◆	http://www.mcu.es/ccbae/es/consulta/registro.cmd?id=176667
◆	http://www.mcu.es/ccbae/es/consulta/registro.cmd?id=178574
◆	http://www.mcu.es/ccbae/es/consulta/registro.cmd?id=184043
◆	http://www.mcu.es/ccbae/es/consulta/registro.cmd?id=176947
◆	http://www.mcu.es/ccbae/es/consulta/registro.cmd?id=176922
◆	http://www.mcu.es/ccbae/es/consulta/registro.cmd?id=179169
◆	http://www.mcu.es/ccbae/es/consulta/registro.cmd?id=178424
◆	http://www.mcu.es/ccbae/es/consulta/registro.cmd?id=180908

SOLR NO FACETADA – 1683 RESULTADOS

```
▼<response>
  ▼<lst name="responseHeader">
    <int name="status">0</int>
    <int name="QTime">980</int>
  </lst>
  ▶<result name="response" numFound="1683" start="0" maxScore="1.0">...</result>
</response>
```

Conclusiones:

- En SOLR no facetada detecta la palabra varas en el campo, pero no están las medidas realmente realizadas en esa unidad.

Parte 2: Experimentación (ej)

- **Consulta 8** en Solr para ¿Qué obras representan un monumento de ceuta?:

SOLR FACETADA – 257 RESULTADOS

```
▼<response>
  ▼<lst name="responseHeader">
    <int name="status">0</int>
    <int name="QTime">18</int>
  </lst>
  ▶<result name="response" numFound="257" start="0" maxScore="1.0">...</result>
```

SOLR NO FACETADA – 3 RESULTADOS

```
▼<response>
  ▼<lst name="responseHeader">
    <int name="status">0</int>
    <int name="QTime">446</int>
  </lst>
  ▶<result name="response" numFound="3" start="0" maxScore="1.0">...</result>
</response>
```

Conclusiones

- Los resultados en este caso son menos debido a que el campo "materia" en la exportación de fichas en texto plano, no se incluía en dicha exportación.



6. Conclusiones

- ▶ El problema de la búsqueda semántica a partir de texto en lenguaje natural.
- ▶ La comparación entre el modelo de búsqueda textual no facetado, facetado y ontológico (o clasificación con la ontología) nos permite observar que para alcanzar conclusiones relevantes es necesario seguir investigando en modelos y técnicas semánticas.




6. Futuras líneas de trabajo

- Se proponen dos metodologías para el enriquecimiento semi-automático:
 - Metodología basada en la distancia contextual y ganancia de conocimiento
 - Metodología basada en roles semánticos.
- Ambas abordarían la instanciación automática desde la combinación del análisis lingüístico tradicional y las tecnologías de extracción de conocimiento textual



6. Otras futuras líneas de trabajo

- ▶ Otras líneas de trabajo futuro
 - ▶ Importancia en el tiempo de respuesta y los resultados obtenidos en las consultas realizadas a las fichas facetadas.
 - ▶ Buscador semántico vs buscador textual
 - ▶ Análisis de utilización de recursos lingüísticos:
 - ▶ Automatización de eliminación de stopwords y otros.
 - ▶ Mejora en el algoritmo de funcionamiento del parser semántico.
 - ▶ Enriquecedor semántico automatizado



Presentación de Trabajo de Fin de Máster
**PROPUESTA DE BÚSQUEDA SEMÁNTICA:
APLICACIÓN AL CATÁLOGO DE MAPAS,
PLANOS Y DIBUJOS DEL ARCHIVO GENERAL
DE SIMANCAS**

Máster en Lenguajes y Sistemas Informáticos: Tecnologías del Lenguaje en la Web

Universidad de Educación a Distancia

Marzo 2013

Autor: José Alberto Benítez Andrades

Directora: Ana M^a García Serrano