

Recuperación de Información en Web Semántica

Crawler: ¿Qué hay de comer?

Adrián Blanco Blanco
adrian.blanco@udc.es

Brais Freire
brais.freire@udc.es

1	Introducción	1
2	Herramientas	2
3	Configuración	2
4	Implementación de Plugins	5
4.1	Configuración del plugin	5
4.2	Implementación	6
4.2.1	parse-justeat	6
4.2.2	index-justeat	6
5	Interfaz Web	6
5.1	¿Cómo está hecho?	7
6	Conclusiones y Trabajo Futuro	11
7	Anexo I: Manual de usuario	11

1 Introducción

En esta memoria relataremos cómo hemos obtenido la solución para el problema que planteamos. Este se basa en situaciones típicas donde un grupo de personas no saben donde ir a comer, pero tienen ganas de distintos tipos de platos. Se introducen estos en nuestra aplicación y nos devolverá qué restaurantes cumplen con estos requisitos. Para ello, nos basaremos en la web de reparto de comida *JustEat*, donde podemos hacer búsquedas por zonas, pero no por restaurante o por plato concreto. Elegimos esta página por ser la que nos ofrece una información más clara y homogénea de la carta de los lugares de comida de la ciudad. Para conseguir la información implementaremos un plugin de *Apache Nutch* y se guardará y organizará la información con *hbase* y *Solr*. Después esto se mostrará en una web para hacer más cómoda

su utilización.

Nos centraremos en la zona de A Coruña, y la información a guardar para cada restaurante será el *nombre, tipo, dirección, logo, lista de platos con sus precios y posiblemente una descripción*

2 Herramientas

Vamos a ver un poco más en profundidad las distintas herramientas que utilizamos para este proyecto.

- **Apache Nutch 2.3.1:** Web Crawler Open Source escalable donde el almacenamiento es independiente del mismo, por lo que ofrece gran flexibilidad según las necesidades de cada uno. Parte de su fuerza viene en la utilización de plugins configurables y desarrollados desde cero mediante la implementación de librerías específicas. Puede utilizarse en una máquina solamente, pero podría utilizarse en varias con la ayuda de *Hadoop*. Lo utilizaremos para obtener direcciones de restaurantes y para crear dos plugins para parsear y obtener información de las páginas.
- **Solr 4.8.1:** Herramienta para el indexado, replicado y balanceado de búsquedas sobre la información, mostrándola en un formato más comprensible y fácilmente manejable para enseñar al usuario según los requisitos pedidos.
- **HBase 0.98.19-hadoop2:** Base de datos open source, no relacional y distribuida que sirve como soporte para guardar los datos obtenidos con las otras herramientas. Ofrece un servicio de clustering que hace posible la lectura y escritura de Big Data.

3 Configuración

Vamos a ver, paso a paso, cómo configurar un entorno desde cero para poder ejecutar nuestro crawler. Empezaremos configurando el archivo **apache-nutch-2.3.1/conf/nutch-site.xml**, donde pondremos el siguiente texto:

```
1 <configuration>
2   <property>
3     <name>storage.data.store.class</name>
4     <value>org.apache.gora.hbase.store.HBaseStore</value>
5     <description>Default class for storing data</description>
6   </property>
7
8   <property>
9     <name>http.agent.name</name>
10    <value>My Nutch Spider</value>
11  </property>
12 </configuration>
```

Con esto especificaremos el tipo de sistema para la base de datos y le daremos un nombre a la *araña* que crawleará las webs. Lo siguiente es descomentar las dos líneas que nos indican para hbase¹. Una vez hecho esto, podremos compilar nutch ejecutando **ant runtime** desde la carpeta raíz del mismo.

Para configurar hbase, lo primero será descargarlo y descomprimirlo:

```
1 wget http://archive.apache.org/dist/hbase/hbase-0.94.14/hbase-0.94.14.tar.gz
2 tar xzvf hbase-0.94.14.tar.gz
```

Ahora solo tendremos que abrir el fichero **conf/hbase-site.xml** e introducir el siguiente texto:

¹En nuestro caso, 116 y 117

```

1 <configuration>
2   <property>
3     <name>hbase.rootdir</name>
4     <value>file:///home/brais/crawler/hbasedata</value>
5   </property>
6   <property>
7     <name>hbase.zookeeper.property.dataDir</name>
8     <value>/home/brais/crawler/zookeeper</value>
9   </property>
10 </configuration>

```

Ahora vamos a configurar los campos que obtendremos e indexaremos de la web, para ello añadimos el siguiente código entre las etiquetas de `<fields>` del archivo **apache-nutch-2.3.1/runtime/local/conf/schema.xml**:

```

1 <fields>
2   [...]
3   <field name="name" type="string" stored="true" indexed="true"/>
4   <field name="address" type="string" stored="true" indexed="true"/>
5   <field name="city" type="string" stored="true" indexed="true"/>
6   <field name="typeRest" type="string" stored="true" indexed="true"/>
7   <field name="imgUrl" type="string" stored="true" indexed="false"/>
8   <field name="dishes" type="string" stored="true" indexed="true"/>
9   <field name="prizes" type="string" stored="true" indexed="true"/>
10 </fields>

```

Una vez copiados, sobrescribimos el archivo **solr-4.8.1/example/solr/collection1/conf/schema.xml**. A continuación, añadimos el siguiente código a **apache-nutch-2.3.1/runtime/local/conf/solrindex-mapping.xml**:

```

1 <fields>
2   [...]
3   <!-- Our fields -->
4   <field dest="name" source="name" />
5   <field dest="address" source="address"/>
6   <field dest="typeRest" source="typeRest"/>
7   <field dest="imgUrl" source="imgUrl"/>
8   <field dest="dishes" source="dishes"/>
9   <field dest="city" source="city"/>
10  <field dest="prizes" source="prizes"/>
11 </fields>

```

Ahora ya tendremos solr y hbase listos, por lo que vamos a terminar con Nutch. Primero, en la ruta **apache-nutch-2.3.1/runtime/local** creamos la carpeta *urls*, y dentro de ella, un archivo llamado *seed.txt*, donde pondremos las urls iniciales, que son estas:

```

1 https://www.just-eat.es/area/15001-a-coru%c3%b1a
2 https://www.just-eat.es/area/15002-a-coru%c3%b1a
3 https://www.just-eat.es/area/15003-a-coru%c3%b1a
4 https://www.just-eat.es/area/15004-a-coru%c3%b1a
5 https://www.just-eat.es/area/15005-a-coru%c3%b1a
6 https://www.just-eat.es/area/15006-a-coru%c3%b1a
7 https://www.just-eat.es/area/15007-a-coru%c3%b1a
8 https://www.just-eat.es/area/15008-a-coru%c3%b1a
9 https://www.just-eat.es/area/15009-a-coru%c3%b1a
10 https://www.just-eat.es/area/15010-a-coru%c3%b1a
11 https://www.just-eat.es/area/15011-a-coru%c3%b1a
12
13 https://www.just-eat.es/area/27001-lugo
14 https://www.just-eat.es/area/27002-lugo
15 https://www.just-eat.es/area/27003-lugo
16 https://www.just-eat.es/area/27004-lugo
17
18 https://www.just-eat.es/area/32001-ourense
19 https://www.just-eat.es/area/32002-ourense
20 https://www.just-eat.es/area/32003-ourense
21 https://www.just-eat.es/area/32004-ourense
22 https://www.just-eat.es/area/32005-ourense
23
24 https://www.just-eat.es/area/36001-pontevedra
25 https://www.just-eat.es/area/36002-pontevedra
26 https://www.just-eat.es/area/36003-pontevedra

```

```
27 https://www.just-eat.es/area/36004-pontevedra
28 https://www.just-eat.es/area/36005-pontevedra
```

Estas son las principales áreas de cada una de las capitales de provincia gallegas. A partir de aquí, se sacarán las distintas url's que contenga cada una para obtener la de los distintos restaurantes de cada zona. Además, tendremos que aplicar un filtro a las urls para que estas no se alejen de la página y no indexe las que sabemos que no nos interesan, como los menús, assets, pie de página, y los menús de los laterales donde se clasifican los restaurantes por tipos. Al eliminarlas, conseguiremos reducir de manera importante el tiempo de ejecución, al tener que parsear muchas menos cosas. Añadiremos esto al final del archivo **apache-nutch-2.3.1/runtime/local/conf/regex-urlfilter.txt**

```
1  ~^https://www.just-eat.es/assets/. *
2  ~^https://www.just-eat.es/cookiespolicy
3  ~^https://www.just-eat.es/faq
4  ~^https://www.just-eat.es/blog/. *
5  ~^https://www.just-eat.es/blog
6  ~^https://www.just-eat.es/contact
7  ~^https://www.just-eat.es/adomicilio/. *
8  ~^https://www.just-eat.es/privacy-policy
9
10 ~^https://www.just-eat.es/area/15001-a-coru%c3%b1a/. *
11 ~^https://www.just-eat.es/area/15002-a-coru%c3%b1a/. *
12 ~^https://www.just-eat.es/area/15003-a-coru%c3%b1a/. *
13 ~^https://www.just-eat.es/area/15004-a-coru%c3%b1a/. *
14 ~^https://www.just-eat.es/area/15005-a-coru%c3%b1a/. *
15 ~^https://www.just-eat.es/area/15006-a-coru%c3%b1a/. *
16 ~^https://www.just-eat.es/area/15007-a-coru%c3%b1a/. *
17 ~^https://www.just-eat.es/area/15008-a-coru%c3%b1a/. *
18 ~^https://www.just-eat.es/area/15009-a-coru%c3%b1a/. *
19 ~^https://www.just-eat.es/area/15010-a-coru%c3%b1a/. *
20 ~^https://www.just-eat.es/area/15011-a-coru%c3%b1a/. *
21
22 ~^https://www.just-eat.es/area/27001-lugo/. *
23 ~^https://www.just-eat.es/area/27002-lugo/. *
24 ~^https://www.just-eat.es/area/27003-lugo/. *
25 ~^https://www.just-eat.es/area/27004-lugo/. *
26
27 ~^https://www.just-eat.es/area/32001-ourense/. *
28 ~^https://www.just-eat.es/area/32002-ourense/. *
29 ~^https://www.just-eat.es/area/32003-ourense/. *
30 ~^https://www.just-eat.es/area/32004-ourense/. *
31 ~^https://www.just-eat.es/area/32005-ourense/. *
32
33 ~^https://www.just-eat.es/area/36001-pontevedra/. *
34 ~^https://www.just-eat.es/area/36002-pontevedra/. *
35 ~^https://www.just-eat.es/area/36003-pontevedra/. *
36 ~^https://www.just-eat.es/area/36004-pontevedra/. *
37 ~^https://www.just-eat.es/area/36005-pontevedra/. *
38
39 # accept anything else
40 +^https://www.just-eat.es/. *
```

Para poder utilizar la aplicación web con jquery, fue necesario modificar Solr para evitar el error "Access-Control-Allow-Origin". Para ello, deberemos añadir en **solr-4.8.1/example/solr-webapp/webapp/WEB-INF/web.xml** lo siguiente:

```
1 <filter>
2   <filter-name>cross-origin</filter-name>
3   <filter-class>org.eclipse.jetty.servlets.CrossOriginFilter</filter-class>
4   <init-param>
5     <param-name>allowedOrigins</param-name>
6     <param-value>*</param-value>
7   </init-param>
8   <init-param>
9     <param-name>allowedMethods</param-name>
10    <param-value>*</param-value>
11  </init-param>
12  <init-param>
13    <param-name>allowedHeaders</param-name>
14    <param-value>*</param-value>
15  </init-param>
```

```

16 </filter>
17 <filter-mapping>
18   <filter-name>cross-origin</filter-name>
19   <url-pattern>/*</url-pattern>
20 </filter-mapping>

```

Esto actuará de filtro en las conexiones, permitiendo conexiones del tipo GET para nuestra página sin error.

4 Implementación de Plugins

Lo primero que haremos, será decirle a Nutch que queremos que trabaje con los dos plugins que vamos a hacer. Para ello, modificaremos el campo *plugin.includes* del archivo `runtime/local/conf/nutch-site.xml`, poniendo el siguiente:

```

1 <property>
2   <name> plugin.includes </name>
3   <value> protocol-http|urlfilter-regex|parse-(html|tika)|index-(basic|anchor)|indexer-
4     solr|urlnormalizer-(pass|regex|basic)|scoring-opic|parse-justeat|index-justeat|
      nutch-extensionpoints|protocol-httpclient</value>
5 </property>

```

En nuestro caso hemos creado dos, *parse-justeat*, el cual recibe la información de un restaurante y saca la información relevante del mismo, *index-justeat* para indexar la información obtenida por el parser y poder mandarla a solr. Primero explicaremos como crear la estructura que rodea al plugin con el ejemplo del indexer, ya que el del parser se haría de manera análoga, para después explicar ambos plugins.

4.1 Configuración del plugin

Comenzaremos creando una carpeta con el nombre del plugin en `apache-nutch-2.3.1/src/plugin`, y dentro de ella los siguientes tres archivos:

build.xml

```

1 <project name="index-justeat" default="jar-core">
2   <import file="../../build-plugin.xml"/>
3 </project>

```

ivy.xml

```

1 <ivy-module version="1.0">
2   <info organisation="org.apache.nutch" module="${ant.project.name}">
3     <license name="Apache 2.0"/>
4     <ivyauthor name="Apache Nutch Team" url="http://nutch.apache.org"/>
5     <description>
6       Apache Nutch
7     </description>
8   </info>
9
10  <configurations>
11    <include file="../../ivy/ivy-configurations.xml"/>
12  </configurations>
13
14  <publications>
15    <!--get the artifact from our module name-->
16    <artifact conf="master"/>
17  </publications>
18
19  <dependencies>
20  </dependencies>
21
22 </ivy-module>

```

plugin.xml

```
1 <?xml version="1.0" encoding="UTF-8" ?>
2 <plugin
3   id="index-justeat"
4   name="JustEat Indexing Filter"
5   version="1.0.0"
6   provider-name="nutch.org">
7   <runtime>
8     <library name="index-justeat.jar">
9       <export name="*" />
10    </library>
11  </runtime>
12  <requires>
13    <import plugin="nutch-extensionpoints" />
14  </requires>
15  <extension id="org.apache.nutch.indexer.justeat"
16    name="Nutch JustEat Indexing Filter"
17    point="org.apache.nutch.indexer.IndexingFilter">
18    <implementation id="JustEatIndexer"
19      class="org.apache.nutch.indexer.justeat.JustEatIndexer" />
20  </extension>
21 </plugin>
```

Para terminar la configuración, añadiremos las siguientes líneas en **apache-nutch-2.3.1/src/plugin**.

```
1 <target name="deploy">
2   [...]
3   <ant dir="parse-justeat" target="deploy" />
4   <ant dir="index-justeat" target="deploy" />
5 </target>
```

4.2 Implementación

4.2.1 parse-justeat

Dentro de su carpeta, creamos la siguiente ristra, **src/java/org/apache/nutch/parse/justeat/JustEatParser.java**. Dentro de la misma tendremos tres funciones importantes, que se llaman en este orden:

1. **filter:** Implementación de *ParseFilter*. Recibe el contenido de las distintas páginas para pasarlas a la siguiente función. Si resulta tener el contenido esperado, osea, es un restaurante, se mandan al *pushIntoMetadata* para que llegue al index-justeat.
2. **getRestaurant:** Recibe el contenido de una página y recorre sus etiquetas para ver si alguna contiene un nombre igual al que buscamos. Por ejemplo, para obtener el nombre del restaurante buscamos la etiqueta *h1* con un nombre de clase *infoTextBlock-item-title*.
3. **pushIntoMetadata:** Para cerrar, recibimos los atributos del restaurante, y los ponemos en la *metadata* de la página, para que pueda ser indexada posteriormente por el otro plugin.

4.2.2 index-justeat

En esta ocasión, crearemos **src/java/org/apache/nutch/indexer/justeat/JustEatIndexer.java**, donde implementaremos *IndexingFilter*. Para ello solo tendremos una función *filter*, la cual extrae el metadata de la página para prepararlo en un *NutchDocument* para que solr sea capaz de interpretarlo.

5 Interfaz Web

Para hacer esto accesible a un usuario medio, vamos a preparar una interfaz web que lea los datos de Solr. Tendremos una caja de texto donde el usuario pondrá los distintos platos que quiere que contenga

el restaurante al que comer. Además, tiene un desplegable donde puede elegir si quiere buscar en todos los restaurantes o filtrar por una ciudad específica.

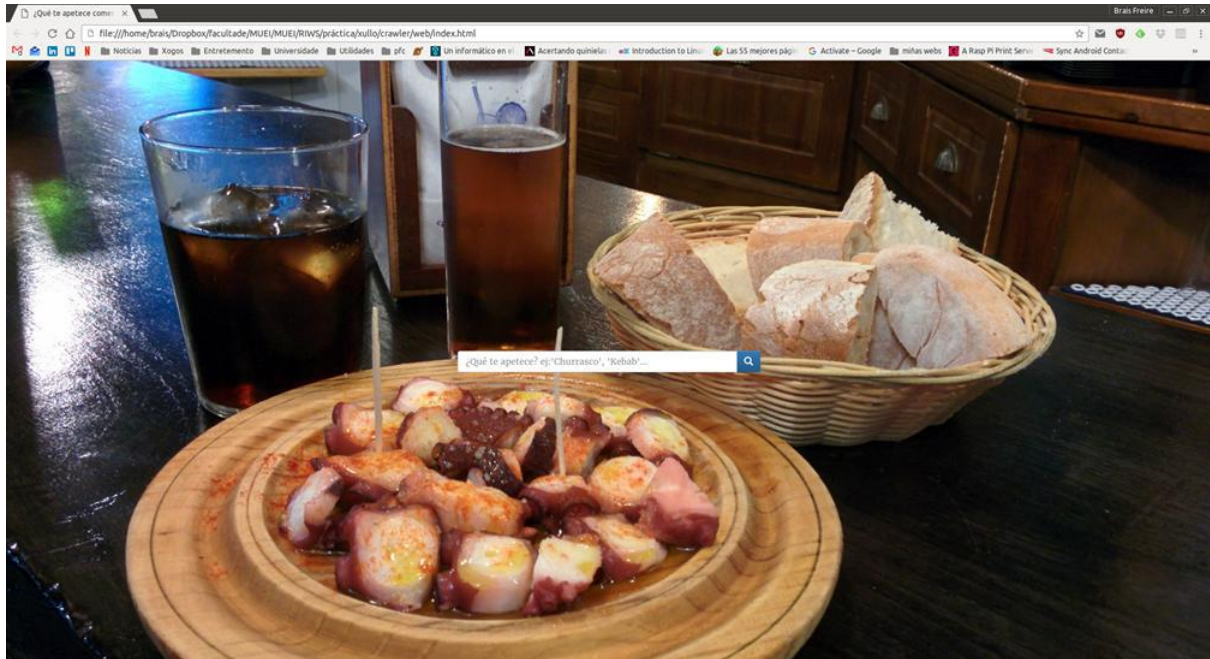


Figure 1: Pantalla inicial de la web

La consulta se realiza poniendo, en la caja de texto, los nombres de los platos separados por AND. Por ejemplo:

```
1 Pulpo AND Churrasco
```

O, simplemente, un único plato:

```
1 Ensalada
```

Una vez dentro, tendremos los distintos restaurantes que cumplen los requisitos, donde podremos ver un botón que posibilita ver el menú que ofrecen a través de un modal junto a sus precios. Además, podemos acceder directamente a la página de JustEat por si prefiere hacer el pedido a domicilio y ahorrar en los gastos de envío evitando tener que pedir a más de un sitio.

También podemos filtrar la búsqueda por tipo de restaurante, pudiendo así afinar todavía más la misma y adaptar las preferencias del usuario al local que más le apetece.

5.1 ¿Cómo está hecho?

La aplicación web está hecha sobre html+css+jquery. Consta de 2 páginas unicamente: index.html y results.html.

Index.html

Desde index.html se introduce el plato a buscar, que se envía a results.html a través de GET. También se valida el formulario para evitar enviar datos vacíos.

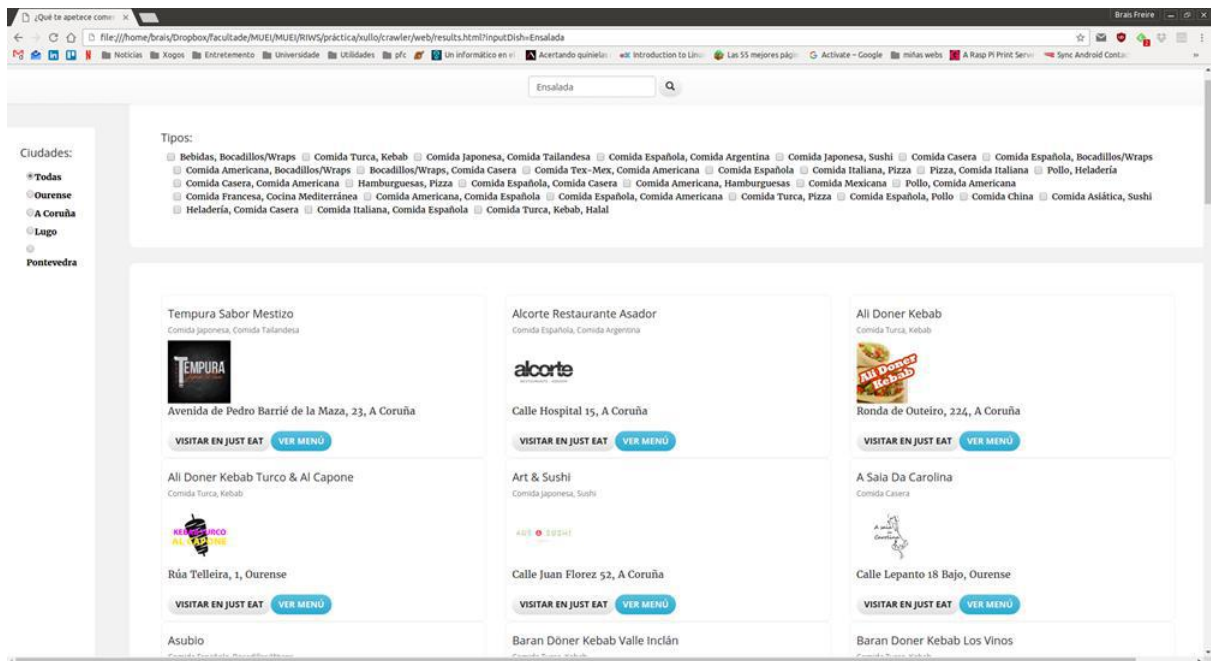


Figure 2: Búsqueda por platos

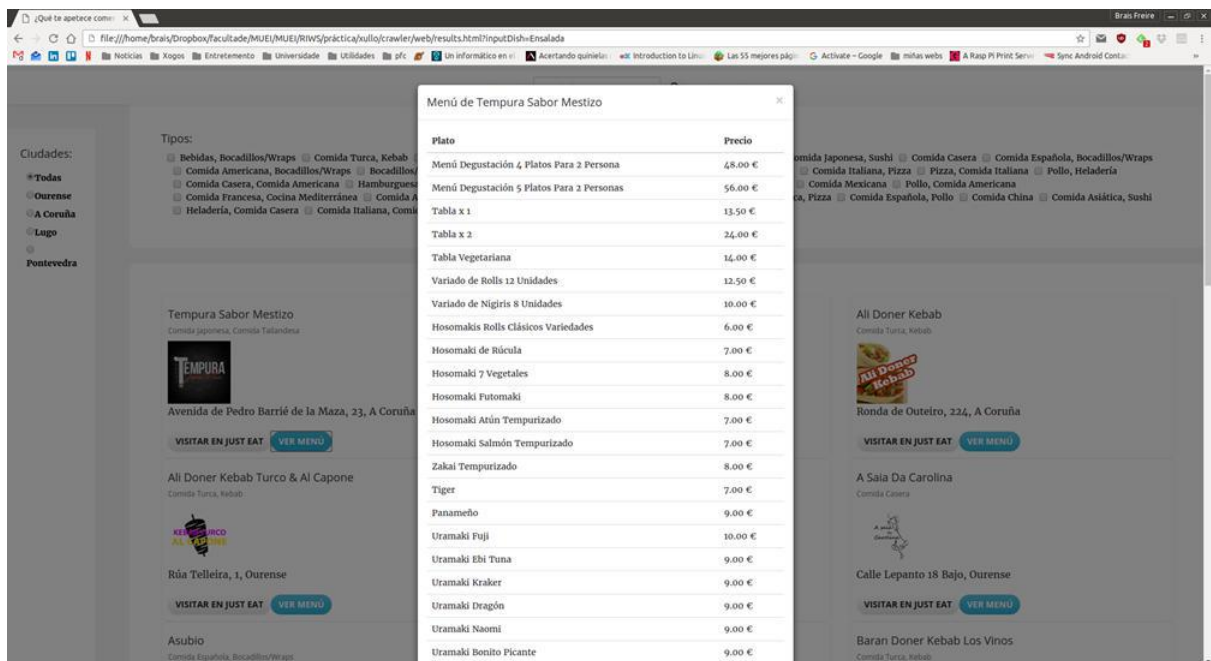


Figure 3: Ejemplo de menú de un restaurante

Results.html

En esta página es donde se realiza practicamente todo. Consta de 4 partes: una zona superior con otro componente de texto para buscar, inmediatamente debajo tenemos la zona donde se muestran los distintos tipos de restaurantes, a la izquierda el filtro por ciudades y, por último, el área donde se muestran los resultados.

El core de esta página es un archivo llamado search.js. Aquí se construyen las distintas url al endpoint de solr. Al cargar la página por primera vez, se generan los filtros (tanto de restaurantes como de ciudades) y se imprimen a la página. Esto se hace en la función llamada printMenus().

```
1 function printMenus(){
2   //Getting the diferent cities of the restaurants in Solr
3
4   $.getJSON(cities, function(json){ //Connecting to Solr and getting data
5
6   }).done(function(json){
7     $('#loadingDivC').hide();
8     var cities = parseElements(json.response.docs); //We transform items for usability
9     $('<label class="checkbox"><input class="box" type="radio" name="city" value="All"
10      onclick="filterCity(this);" checked="checked">Todas</label>').appendTo('#
11      cityCheckboxes');
12     for (var i=0; i<cities.length; i++){
13       $('<label class="checkbox"><input class="box" type="radio" name="city" value="'+
14        cities[i]+'>'+cities[i]+'</label>').appendTo('#
15        cityCheckboxes');
16     }
17   });
18
19   //Getting the diferent types of the restaurants in Solr
20   $.getJSON(typeRestaurants, function(json){
21
22   }).done(function(json){
23     $('#loadingDivT').hide();
24     var types = parseTypeElements(json.response.docs);
25
26     for (var i=0; i<types.length; i++){
27       $('<label class="checkbox-inline"><input type="checkbox" onclick="filterType(this)
28        ;" value="'+types[i]+'>'+types[i]+'</label>').appendTo('#typeCheckboxes');
29     }
30   });
31 }
```

Cabe destacar que hemos escogido radiobuttons para el filtro de las ciudades porque entendemos que interesa buscar en un lugar concreto. No se trata de una elección del tipo “comer en Coruña o en Lugo”.

También se produce una primera conexión para obtener los restaurantes del plato buscado utilizando la función connectSolr(). Cada vez que se selecciona un filtro, se produce una actualización del área de resultados, llamando de nuevo a la función connectSolr(), que recibe un plato como parámetro.

```
1 function connectSolr(dish){ //Connecting the Solr endpoints
2   var connect = solrAPI+dish+"&wt=json&indent=true&rows=50"
3
4   if (city != null){ //City Filter
5     if (city.length > 0){
6       if (city != "All"){
7         connect = connect+"&fq=city%3A*"+city;
8       }
9     }
10  } else {
11  }
12
13  if (filterTypeRest.length>0){ //Type Restaurant Filter
14    for (var i=0; i<filterTypeRest.length; i++){
15      var filterRest;
16      var res = filterTypeRest[i].split(" ");
17      for (var j=0; j<res.length; j++){
```

```

18         if (j==0){
19             filterRest=res[j];
20         } else {
21             filterRest=filterRest+" "+res[j];
22         }
23     }
24
25     if (i==0){
26         connect = connect+"&fq=typeRest%3A%22"+filterRest+"%22";
27     } else {
28         connect = connect+"OR+typeRest%3A%22"+filterRest+"%22";
29     }
30 }
31 }
32
33 $('#inner').empty(); //Clear previous results
34
35 $.getJSON( connect, function(json){ //Getting new data
36     //console.log(json);
37 })
38 .done(function(json){
39     $('#loadingDiv').hide();
40
41     //Parsing restaurants to show the result in the page
42     console.log(json.response.docs)
43     var restaurants = parseRestaurants(json.response.docs);
44     if (restaurants.length > 0){
45         printRestaurants(restaurants);
46     }else{
47         $('#<div class="alert alert-warning" role="alert"><center><strong>Lo sentimos!</strong> No disponemos de restaurantes que cumplan el criterio de búsqueda :(</center></div>').appendTo('#inner');
48     }
49 })
50 .fail(function (){
51     $('#loadingDiv').hide();
52     $('#<div class="alert alert-warning" role="alert"><center><strong>Lo sentimos!</strong> Se ha perdido la conexion con Solr </center></div>').appendTo('#inner');
53 });
54 }

```

De esta función, cabe destacar varias cosas. La generación de las URL de endpoints se crean de la siguiente manera. Tenemos generadas, en variables globales, 3 urls que son:

```

1 var cities = 'http://localhost:8983/solr/collection1/select?q=%3A*&wt=json&group=true&group.field=city&group.main=true&rows=500' //distintas ciudades
2 var typeRestaurants = 'http://localhost:8983/solr/collection1/select?q=%3A*&wt=json&group=true&group.field=typeRest&group.main=true&rows=500' //distintos tipos de restaurante
3 var solrAPI = "http://localhost:8983/solr/collection1/select?q=";

```

Las dos primeras se utilizan para recoger las distintas ciudades y tipos de restaurante que están almacenadas, haciendo uso del group.main y group.field. Se utilizan para generar los filtros que va a tener la página. La última de ellas es la base para generar consultas. A partir de la misma se generan las distintas consultas del usuario. Esto se hace de la siguiente manera:

```

1 var connect =solrAPI+dish+"&wt=json&indent=true&rows=50"
2 if (city != null){ //City Filter
3     if (city.length > 0){
4         if (city != "All"){
5             connect = connect+"&fq=city%3A*"+city;
6         }
7     }
8
9     if (filterTypeRest.length>0){ //Type Restaurant Filter
10         for (var i=0; i<filterTypeRest.length; i++){
11             var filterRest;
12             var res = filterTypeRest[i].split(" ");
13             for (var j=0; j<res.length; j++){
14                 if (j==0){
15                     filterRest=res[j];

```

```

16     } else {
17         filterRest=filterRest+" "+res[j];
18     }
19 }
20
21 if (i==0){
22     connect = connect+"&fq=typeRest%3A%22"+filterRest+"%22";
23 } else {
24     connect = connect+"OR+typeRest%3A%22"+filterRest+"%22";
25 }
26 }
27 }

```

6 Conclusiones y Trabajo Futuro

Como hemos visto, hemos conseguido *darle la vuelta* a JustEat, ofreciendo una nueva idea de negocio. Ofrecemos más formas de búsqueda, como por platos o restaurantes, donde en JustEat no se permite, dejando solo buscar por código postal. Además, y donde creemos que puede estar parte de la ventaja competitiva, ofrecemos un buscador para que un grupo de personas sea capaz de encontrar un restaurante a su medida, e incluso permitimos poder hacer el pedido en la propia página de JustEat, consiguiendo así un ahorro en posibles gastos de envío por tener que llamar a dos sitios distintos. También nos puede servir como una ayuda a encontrar nuevos sitios para visitar.

Aún así, este trabajo es solo un punto de partida. Muchos de los restaurantes de la ciudad no se ofertan en esta página, por lo que tendríamos que buscar más páginas que nos ofrezcan este tipo de servicio. La elección sigue siendo buena, no encontramos otra página que nos dé tanta información como JustEat, pero aún así, depender de un único proveedor no es bueno. Una forma *sencilla* de atajar esto puede ser abrir una forma de que los dueños de un restaurante fueran capaces de mandar sus datos a la página.

En el ámbito económico, podríamos ofertar *restaurantes patrocinados*, donde con búsquedas específicas, aparezcan siempre arriba. Por ejemplo, un restaurante italiano puede querer patrocinar su local y cree que su punto fuerte son los raviolis, por lo que puede pagar una suscripción para que, en caso de que el usuario busque este plato, aparezca siempre en los primeros puestos.

7 Anexo I: Manual de usuario

Para hacer más cómoda la ejecución de la aplicación, hemos creado una serie de scripts que automatizarán gran parte del trabajo. Comenzaremos iniciando solr y hbase con el siguiente script, que utilizaremos cada vez que reiniciemos el ordenador:

```

1 #!/bin/bash
2
3 cd solr-4.8.1/example/
4 java -Dsolr.solr.home=solr -jar start.jar &
5 cd ../../
6 sleep 5
7
8 printf "Loading HBASE...\n"
9 hbase-0.98.19-hadoop2/bin/start-hbase.sh &

```

Para crawlear, utilizaremos el siguiente, donde decidiremos si inyectar nuevas urls si pulsamos *s*. Para obtener los restaurantes, ejecutaremos dos veces este script, inyectando solo la primera vez.

```

1 #!/bin/bash
2
3 cd apache-nutch-2.3.1/runtime/local/
4

```

```

5 read -n1 -r -p "Press s to inject..." key
6 if [ "$key" = 's' ]; then
7     bin/nutch inject urls
8 fi
9 bin/nutch generate -topN 3
10 bin/nutch fetch -all
11 bin/nutch parse -all
12 bin/nutch updatedb -all

```

Una vez obtenido el contenido buscado, lo subiremos a solr con el siguiente comando:

```

1 #!/bin/bash
2
3 cd apache-nutch-2.3.1/runtime/local/
4 bin/nutch solrindex http://localhost:8983/solr/ -all

```

Como un script auxiliar, también hemos utilizado este para eliminar el contenido guardado en solr y en hbase.

```

1 #!/bin/bash
2
3 cd "solr-4.8.1/example/exampledocs"
4 java -Ddata=args -jar post.jar "<delete><query>*:*</query></delete>"
5
6 cd ../../..
7 rm -r zookeeper
8 mkdir zookeeper
9 rm -r hbasedata
10 mkdir hbasedata

```

Para visualizar los datos, abriremos un navegador e iremos a <http://localhost:8983/solr/>, donde entraremos en el panel de control. Si después elegimos *collection1* en *core selection*, podemos elegir la opción *Query* para hacer consultas sobre los datos.