

***Recuperación de Información
y Recomendaciones en la
Web 2016***

Grupo 11

Docente : Libertad Tansini

Integrantes:

Nicolas Fabre

Valentina Franchi

Índice

1.Introducción.....	2
2.Enfoque de la solución.....	2
3.Diseño e implementación.....	3
3.1.Arquitectura del sistema.....	4
3.2 Tecnologías utilizadas.....	4
4.Pruebas.....	5
4.1 Prueba 1.....	5
4.2 Prueba 2.....	6
5.Dificultades.....	8
6.Trabajo a futuro.....	9
7. Conclusiones.....	10
8. Manual de usuario.....	11
9.Referencias.....	13

1.Introducción

Esta aplicación planea brindar una solución a la búsqueda y localización de supermercados en el Uruguay que cumplan con ciertos requisitos que desee el cliente. Estos requisitos pueden ir desde características comunes como la ubicación, cadena de supermercado y el horario hasta servicios más específicos como Banred, Red Brou, estacionamiento, entre otros.

Se eligieron las cadenas de supermercados Devoto, Disco, Kinko y Multiahorro por ser de las cadenas más grandes en el Uruguay y con mayor y más clara información en su página web. A pesar del hecho que dichas páginas web proveen información de cada supermercado y su localización, la idea es proveer toda la información de los mismos junta, donde también se consideró importante realizar un buscador brindando filtros avanzados sobre sus características.

2.Enfoque de la solución

El sistema de obtención de las distintas noticias fue desarrollado utilizando el lenguaje Python y el framework Scrapy. También se utilizó una Base de Datos no relacional (MongoDB), y para el frontend se utilizaron AngularJS , HTML5 y CSS.

La solución que se planteó fue obtener la información de los supermercados mediante web scraping, utilizando la herramienta scrapy, sobre las páginas de Devoto, Disco, Kinko y Multiahorro, para luego procesarla y poder mostrarla en la aplicación, así como también localizarla en un mapa. Los datos obtenidos para cada supermercado fueron:

- dirección,
- barrio,
- teléfono,
- horario,
- coordenadas,
- servicios brindados.

Al obtener esta información se decidió persistir la misma en una base de datos, debido a que los datos de los mismos son rígidos y no cambian

constantemente, reduciendo así el costo en tiempo que implica conectarse con otro servicio externo, así como en ancho de banda.

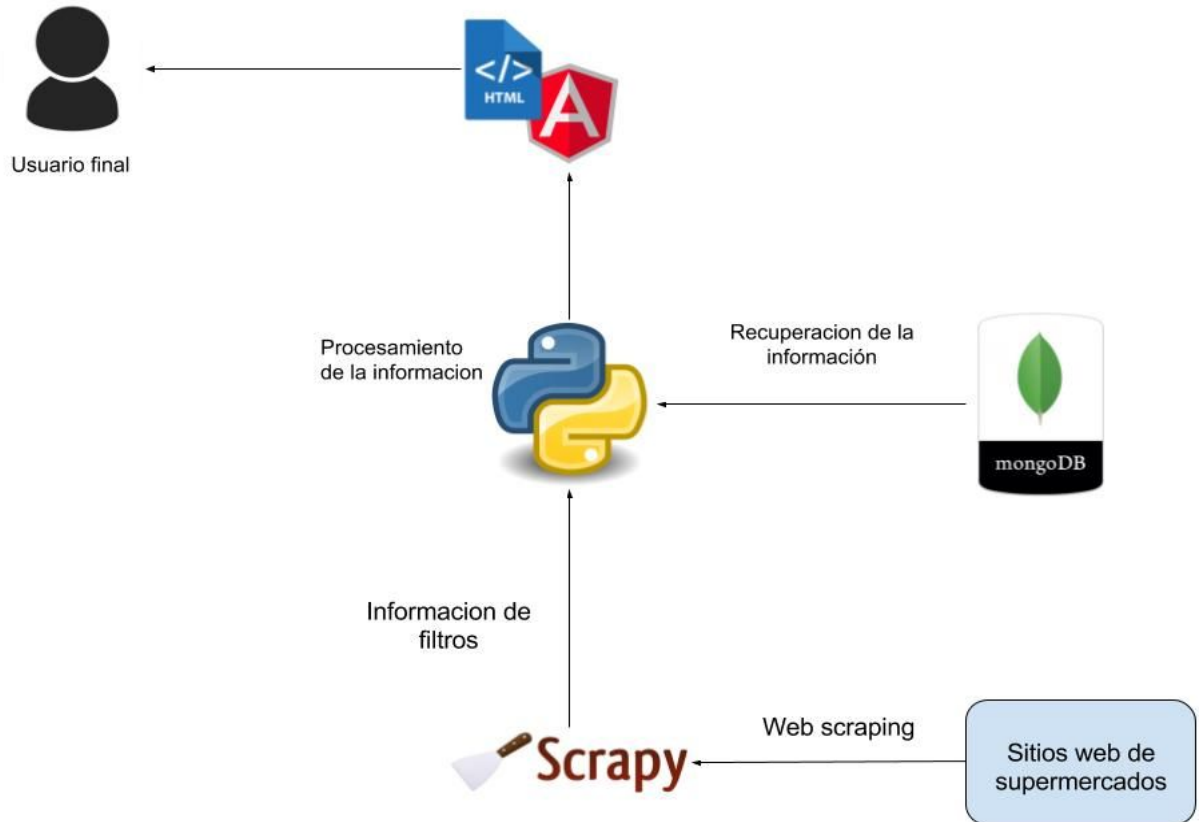
Una vez obtenida toda la información se procedió a procesarla para luego mostrarla en la aplicación, esto implicó entre otras cosas, obtener las coordenadas de algunos supermercados ya que en otros casos fueron obtenidas de los datos de la misma página.

Los filtros implementados fueron la búsqueda por barrio, cadena de supermercado, horario, si desea que tenga estacionamiento y todo tipo de servicios, permitiendo cualquier combinación entre ellos.

3.Diseño e implementación

El diseño del sistema corresponde a un proyecto Python, implementado con el framework Scrapy , el cual tiene como principales componentes a los Ítems y Spiders. Los Ítems son los contenedores de la información resultante de la extracción. Es un DSL mediante el cual se define la estructura de los datos a extraer. Los Spiders son clases, piezas de código y módulos escritos con el fin de descargar el código fuente de la URL objetivo, y recorrer la información del sitio web descargado, para luego depositarlo en los Ítems. Decidimos persistir los datos mencionados en una base de datos NoSQL, ya que persistimos en formato JSON, y las operaciones involucradas son get de un super (se requiere de la información completa del objeto) y get de todos los objetos super. Luego, la web accede a los JSON obtenidos como resultado del procesamiento de los datos obtenidos, para luego presentar la información en el mapa y como filtros, para hacerla visible al usuario.

3.1. Arquitectura del sistema



3.2 Tecnologías utilizadas

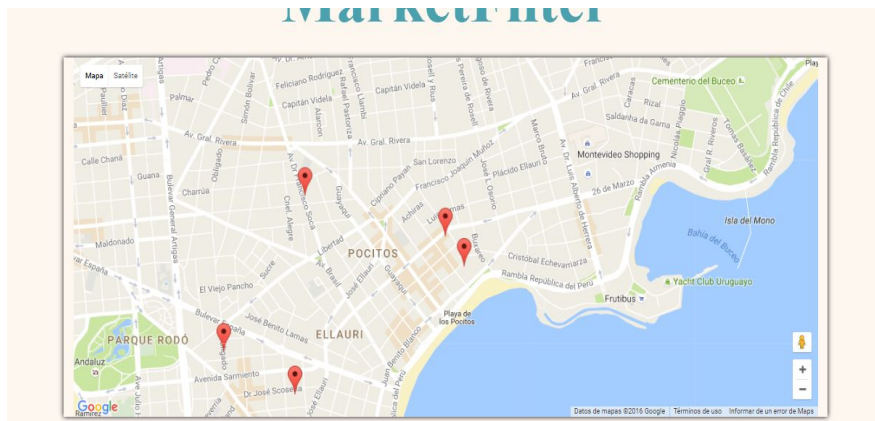
- Se utilizaron para el desarrollo de la lógica del programa el lenguaje python y la herramienta de Scrapy.
- Para la persistencia se utilizó mongoDB.
- En cuanto a Frontend se utilizaron las tecnologías HTML 5, CSS y AngularJS.

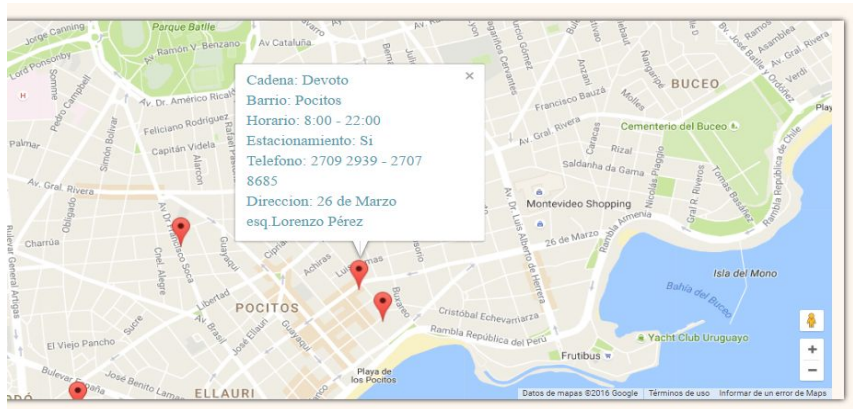


4.Pruebas

4.1 Prueba 1

Se toma como precondition que exista la información de algún supermercado en la base de datos del sistema. Al seleccionar cualquier filtro ya sea de cadena de supermercado, horario, estacionamiento o servicio y seleccionando un solo barrio, el mapa se centra sobre el barrio seleccionado mostrando aquellos supermercados que cumplan además del barrio elegido, con el resto de los filtros.



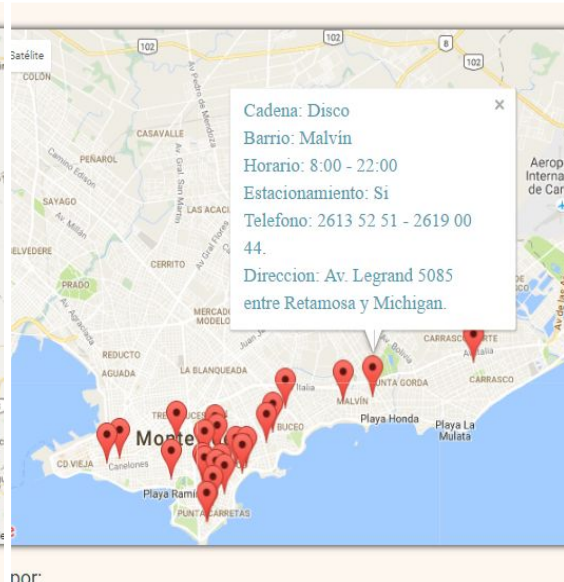
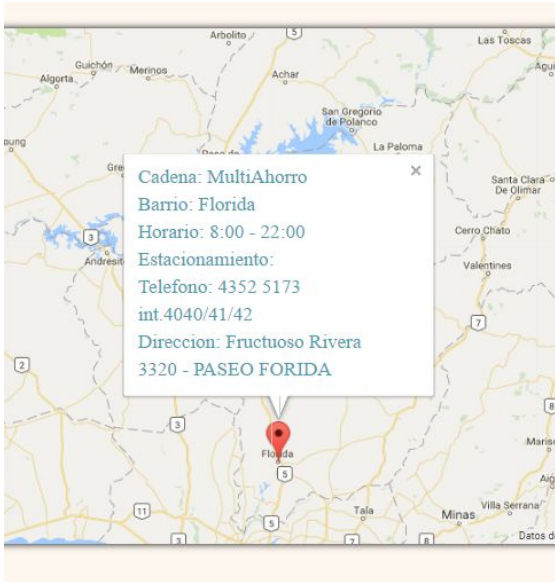
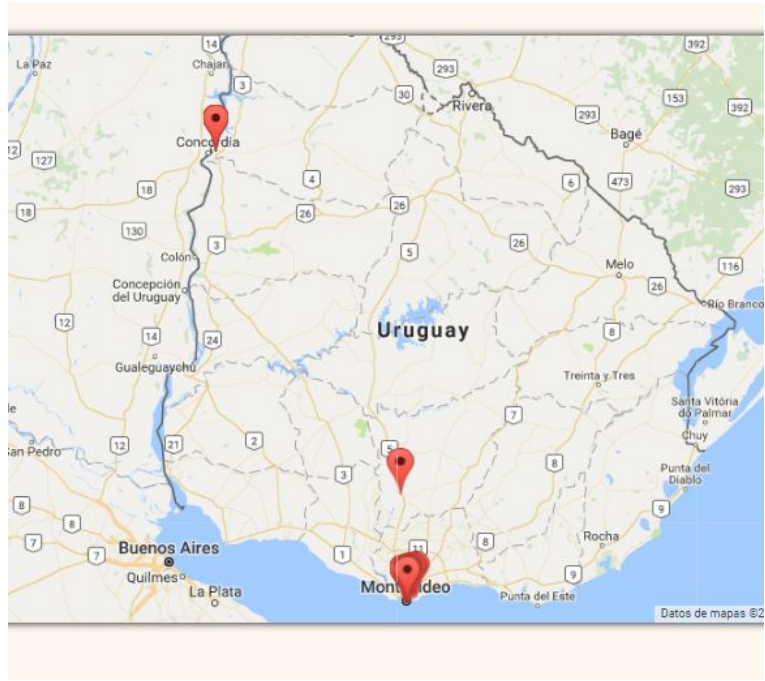


4.2 Prueba 2

Se toma como precondition que exista la información de algún supermercado en la base de datos del sistema. Al seleccionar varios filtros y hacer click en el botón "Consultar", se muestran en el mapa solo los supermercados que cumplan con los mismos, siendo cualquiera esa combinación de filtros. A su vez clickeando sobre los markers en el mapa se podrá ver toda la información de ese supermercado, y de esta forma también se podrá comprobar que los mismos cumplan con los filtros seleccionados.

Filtrar por:

<p>Cadena</p> <p><input type="checkbox"/> Devoto</p> <p><input checked="" type="checkbox"/> Kinko</p> <p><input checked="" type="checkbox"/> MultiAhorro</p> <p><input checked="" type="checkbox"/> Disco</p> <p>Horario</p> <p>21:48</p>	<p><input type="checkbox"/> Estacionamiento</p> <p>Servicio</p> <p><input type="checkbox"/> BanRed</p> <p><input type="checkbox"/> Peluquería</p> <p><input type="checkbox"/> Santander</p> <p><input type="checkbox"/> Western Union</p> <p><input type="checkbox"/> Optica Sfera</p> <p><input type="checkbox"/> Cerrajería</p> <p><input type="checkbox"/> Librería</p> <p><input type="checkbox"/> Tintorería</p> <p><input type="checkbox"/> Redpagos</p> <p><input type="checkbox"/> Marcel calzados</p> <p>Zapatería</p> <p><input type="checkbox"/> Visión Plus Óptica</p> <p><input type="checkbox"/> Game Stop</p> <p><input type="checkbox"/> Video Club</p> <p><input type="checkbox"/> Play Time Juegos</p> <p><input type="checkbox"/> Farmacia</p> <p><input type="checkbox"/> Abitab</p> <p><input type="checkbox"/> Red Brou</p> <p><input type="checkbox"/> Cambio Val</p> <p><input type="checkbox"/> CUTCSA</p> <p><input type="checkbox"/> Florería</p> <p><input type="checkbox"/> Bar</p> <p><input type="checkbox"/> Geant Travel</p> <p><input type="checkbox"/> Baúl Gitano Regalos</p> <p><input type="checkbox"/> Sitio del Lector librería</p> <p><input type="checkbox"/> Heladería Cherry 's</p> <p><input type="checkbox"/> Cellular Center</p> <p><input type="checkbox"/> Movistar</p>	<p>Barrio</p> <p><input checked="" type="checkbox"/> Pocitos</p> <p><input checked="" type="checkbox"/> Cordón</p> <p><input checked="" type="checkbox"/> Centro</p> <p><input type="checkbox"/> Piedras Blancas</p> <p><input type="checkbox"/> Tres Cruces</p> <p><input checked="" type="checkbox"/> Salto</p> <p><input checked="" type="checkbox"/> Florida</p> <p><input type="checkbox"/> Maroñas</p> <p><input type="checkbox"/> Bella Vista</p> <p><input type="checkbox"/> Parque Batlle</p> <p><input type="checkbox"/> Atlántida</p> <p><input type="checkbox"/> Sayago</p> <p><input type="checkbox"/> Carrasco Norte</p> <p><input type="checkbox"/> Piriapolis</p> <p><input type="checkbox"/> San Martín</p> <p><input checked="" type="checkbox"/> Buceo</p> <p><input checked="" type="checkbox"/> Carrasco</p> <p><input type="checkbox"/> Las Acacias</p> <p><input checked="" type="checkbox"/> Punta Carretas</p> <p><input type="checkbox"/> El Tanque</p> <p><input type="checkbox"/> Perez Castellanos</p> <p><input type="checkbox"/> Las Piedras</p> <p><input type="checkbox"/> Ciudad Vieja</p> <p><input type="checkbox"/> Barrio Sur</p> <p><input type="checkbox"/> Punta del Este</p> <p><input type="checkbox"/> El Pinar</p> <p><input type="checkbox"/> Shangrilá</p> <p><input type="checkbox"/> Balvedere</p> <p><input type="checkbox"/> Prado</p> <p><input checked="" type="checkbox"/> Parque Rodó</p> <p><input checked="" type="checkbox"/> Malvin</p> <p><input type="checkbox"/> Vista Linda</p> <p><input type="checkbox"/> La Blanqueada</p> <p><input type="checkbox"/> Goes</p> <p><input type="checkbox"/> Paysandu</p> <p><input type="checkbox"/> Solymar</p> <p><input checked="" type="checkbox"/> Melo</p> <p><input type="checkbox"/> Palermo</p> <p><input type="checkbox"/> Canelones</p> <p><input type="checkbox"/> Colón</p> <p><input type="checkbox"/> Punta Gorda</p> <p><input type="checkbox"/> Union</p> <p><input type="checkbox"/> Pando</p>	<p>Consultar</p>
---	--	---	-------------------------



por:

5.Dificultades

- Surgieron varios problemas al unificar datos de distintas fuentes debido a que la calidad de los datos contenidos no es la mejor y el formato no se mantiene en todos los casos. Se encontraron diferencias semánticas en las fuentes elegidas, por ejemplo en el caso de Devoto se maneja el concepto de estacionamiento como un servicio más, y por otro lado el Disco provee información adicional referenciada como Parking.
- También se encontraron muchos datos inconsistentes en la misma fuente así como direcciones mal ingresadas o abreviaciones de calles incorrectas, lo cual dificultó para la obtención de los datos para su Geolocalización.
- Se requiere de un gran trabajo de procesamiento de lenguaje natural para el desarrollo de la aplicación.
- El hecho de que no existan APIs predefinidas por parte de las cadenas importantes de supermercados, dificulta el objetivo de centralizar y normalizar la información. Dado que el resultado del web scraping depende del código HTML de la página web y este puede variar.

6.Trabajo a futuro

- Una de las características que puede mejorar la calidad del producto podría ser agregar más cadenas de supermercados, logrando centralizar la información sobre las mismas, para poder brindarle al usuario más opciones para poder elegir mejor.
- Por otro lado, a su vez como en el punto anterior se podría expandir la aplicación no solo a supermercados sino también a estaciones de servicio, farmacias, etc. La desventaja que podría tener esta sería que no todos tienen página web por lo que sería difícil la obtención de los datos.
- Otro punto que consideramos importante es crear una aplicación móvil, ya que permite el acceso inmediato a la información de manera sencilla. A su vez, otra característica de Geolocalización y filtrado que podría ser de interés para el usuario sería la sugerencia de supermercados que se encuentren cerca de él mediante la obtención de su ubicación.
- También la posibilidad de que la web se actualice cada cierto periodo de tiempo, realizando nuevamente el scraping de la web y actualizando los datos de los supermercados, y no que se realice en parte de forma manual como se encuentra actualmente.

7. Conclusiones

Se desarrolló una aplicación web con el fin de brindar a los usuarios información personalizada de supermercados de interés. Ya que se ofrece la posibilidad de consultar por los locales que cumplan con las preferencias y necesidades deseadas.

El producto final es un servicio simple e intuitivo que cumple con los requerimientos iniciales propuestos por el equipo, combinando diferentes tecnologías y técnicas de recuperación, procesamiento y presentación de información de la web.

8. Manual de usuario

MarketFilter



Filtrar por:

Filtrar por:

Cadena <ul style="list-style-type: none"><input type="checkbox"/> Devoto<input type="checkbox"/> Kinko<input type="checkbox"/> MultiAhorro<input type="checkbox"/> Disco	Estacionamiento	Barrio	Consultar
Horario <input type="text"/>	Servicio <ul style="list-style-type: none"><input type="checkbox"/> BanRed<input type="checkbox"/> Peluquería<input type="checkbox"/> Santander<input type="checkbox"/> Western Union<input type="checkbox"/> Optica Sfera<input type="checkbox"/> Cerrajería<input type="checkbox"/> Librería<input type="checkbox"/> Tintorería<input type="checkbox"/> Redpagos<input type="checkbox"/> Marcel calzados<input type="checkbox"/> Zapatería<input type="checkbox"/> Visión Plus Óptica<input type="checkbox"/> Game Stop<input type="checkbox"/> Video Club<input type="checkbox"/> Play Time Juegos<input type="checkbox"/> Farmacia<input type="checkbox"/> Abitab<input type="checkbox"/> Red Brou<input type="checkbox"/> Cambio Val<input type="checkbox"/> CUTCSA<input type="checkbox"/> Florería<input type="checkbox"/> Bar<input type="checkbox"/> Geant Travel<input type="checkbox"/> Baúl Gitano Regalos<input type="checkbox"/> Sitio del Lector librería<input type="checkbox"/> Heladería Cherry 's<input type="checkbox"/> Cellular Center<input type="checkbox"/> Movistar	<ul style="list-style-type: none"><input type="checkbox"/> Pocitos<input type="checkbox"/> Cordón<input type="checkbox"/> Centro<input type="checkbox"/> Piedras Blancas<input type="checkbox"/> Tres Cruces<input type="checkbox"/> Salto<input type="checkbox"/> Florida<input type="checkbox"/> Maroñas<input type="checkbox"/> Bella Vista<input type="checkbox"/> Parque Battle<input type="checkbox"/> Atlántida<input type="checkbox"/> Sayago<input type="checkbox"/> Carrasco Norte<input type="checkbox"/> Piriapolis<input type="checkbox"/> San Martín<input type="checkbox"/> Buceo<input type="checkbox"/> Carrasco<input type="checkbox"/> Las Acacias<input type="checkbox"/> Punta Carretas<input type="checkbox"/> El Tanque<input type="checkbox"/> Perez Castellanos<input type="checkbox"/> Las Piedras<input type="checkbox"/> Ciudad Vieja<input type="checkbox"/> Barrio Sur<input type="checkbox"/> Punta del Este<input type="checkbox"/> El Pinar<input type="checkbox"/> Shangrilá<input type="checkbox"/> Balvedere<input type="checkbox"/> Prado<input type="checkbox"/> Parque Rodó<input type="checkbox"/> Malvin<input type="checkbox"/> Vista Linda<input type="checkbox"/> La Blanqueada<input type="checkbox"/> Goes<input type="checkbox"/> Paysandu<input type="checkbox"/> Solymar<input type="checkbox"/> Melo<input type="checkbox"/> Palermo<input type="checkbox"/> Canelones<input type="checkbox"/> Colón<input type="checkbox"/> Punta Gorda<input type="checkbox"/> Union<input type="checkbox"/> Pando	Consultar

Cuando se inicia la aplicación, los datos ya están cargados en la base de datos y no es necesario realizar ningún procedimiento para cargar los mismos. Se mostrará en la pantalla el mapa sin cargar las ubicaciones aun, ya que las mismas son muchas y demoran en cargar. Una vez iniciado ya se podrán utilizar los filtros.

Los filtros posibles son los siguientes:

1. Por barrio
2. Por horario
3. Por cadena de supermercado
4. Si desea que tenga estacionamiento
5. Por servicios

- (1) Se utiliza para filtrar los supermercados por el barrio en que se encuentran.
- (2) Se filtra por horarios, el usuario ingresa una hora en que desea encontrar supermercados abiertos y se muestran los mismos en el mapa.
- (3) Se utiliza para filtrar por la cadena de supermercado, como lo son Devoto, Disco, Multiahorro y kinko.
- (4) Se utiliza para obtener aquellos supermercados que tienen estacionamiento.
- (5) Se filtra por todo tipo de servicios, por ejemplo si desea que el mismo tenga farmacia, peluquería, abitab, BanRed, Red Brou y otros.

Los filtros no son excluyentes entre sí, es decir, se pueden combinar los mismos para generar una búsqueda más específica.

Una vez que se marquen los filtros deseados por el usuario, seleccionando el botón "Consultar" se mostrarán en el mapa aquellos supermercados que cumplan con aquellos que fueron seleccionados.

9.Referencias

- [1] Scrapy 1.0 framework, scrapy.org.
- [2] Base de datos No SQL, Mongo DB, www.mongodb.org.
- [3] Angular JS, angularjs.org.
- [4] Tutorial proyecto en Scrapy,
<http://doc.scrapy.org/en/latest/intro/tutorial.html>.
- [5] Items, <http://doc.scrapy.org/en/latest/topics/items.html>.
- [6] Spiders, <http://doc.scrapy.org/en/latest/topics/spiders.html>.
- [7] Pagina web de Devoto, <http://www.devoto.com.uy/aindex.aspx>
- [8] Pagina web del Kinko, <http://kinko.com.uy/>
- [9] Pagina web del Disco, <http://www.disco.com.uy/>
- [10] Pagina web de Multiahorro, <http://www.multiahorro.com.uy/>