

REDES DE COMPUTADORES

Proyecto: ¿Cómo Funciona Google?

Felipe Fernández Pino
201121011-5

René Fredes Verdugo
201121031-K

Felipe Gil Silva
201121019-0

28 de septiembre de 2015

1. Resumen

El proyecto abarca el funcionamiento general del buscador web Google. Primero se estudia el proceso de rastreo e indexación permanente que realiza Google sobre todos los sitios web encontrados a partir de los ya identificados. Luego se mencionan los principales algoritmos de búsqueda que Google utiliza para mejorar los resultados y personalizarlos de acuerdo a los intereses de cada usuario. Finalmente se realizan unas demostraciones sobre las características de personalización de Google y las restricciones de búsqueda que imponen algunos sitios de importancia.

2. Introducción

En el presente informe se detalla sobre la exposición realizada, llamada: ¿Cómo funciona el buscador Google?. Para cualquier navegante de la red resulta casi imperativo encontrarse en algún punto con este buscador (o algún otro, pero en notable menor medida), dado que suele ser la primera opción para cualquier *plug-in* de búsqueda y además su gran cantidad de aplicaciones ya implementadas en los diferentes dispositivos computacionales existentes en la actualidad.

Este proyecto busca responder a varias preguntas que cualquier usuario puede hacerse muy rápidamente luego de realizar una búsqueda cualquiera, como por ejemplo: ¿De dónde Google obtiene la información buscada?, ¿cómo decide qué resultados entregar?, ¿qué tipo de resultados son discriminados?, etc.

El lector observará que para responder a las preguntas propuestas se llegará a los algoritmos de rastreo e indexación y algoritmos de búsqueda, destacándose las palabras claves: **Googlebot** y **PageRank** respectivamente.

3. Rastreo e Indexación

Seguramente, lo que cada usuario de Google piensa cuando realiza una pregunta es que el motor de búsqueda buscará, valga la redundancia, por toda web el contenido especificado, pero lo que en realidad busca en una especie de índice que se ha creado conforme los años y que ahora mismo sigue en crecimiento.

Este índice es rastreado mediante softwares llamados “Rastreadores Web”, también conocidos como *Spiders*, del inglés Arañas. Esta última acepción contribuye con una buena analogía del funcionamiento del programa, que consiste en rastrear a alguna página, para luego rastrear a todos los enlaces a los que ésta lleve, creando así toda una “telaraña” de información sobre las páginas existentes en la web.

Uno de los *spiders* más conocidos es “*Googlebot*”. El proceso de rastreo que éste realiza comienza con un listado de sitios web de rastreos anteriores y *sitemaps*. *Googlebot* recopila tanta información como puede de cada página a la que ingresa, y luego la almacena en un índice dependiendo del tipo de información: el título y los enlaces se almacenan en un índice especializado en búsquedas amplias; mientras que el resto del contenido de la página se almacena en otro índice dedicado a las búsquedas más complejas y específicas. Este comportamiento se ilustra en la Figura 1.

Los *webmasters* (o dueños de sitios web) pueden restringir el acceso a los rastreadores para evitar que algunas páginas de su sitio aparezcan en los resultados de las búsquedas. Esto se logra agregando un archivo de texto `robots.txt` en la raíz del sitio web y especificar en este archivo cuales páginas se deben indexar y cuales no. En la sección de demostración se observan ejemplos de su configuración en algunos

sitios web más reconocidos.

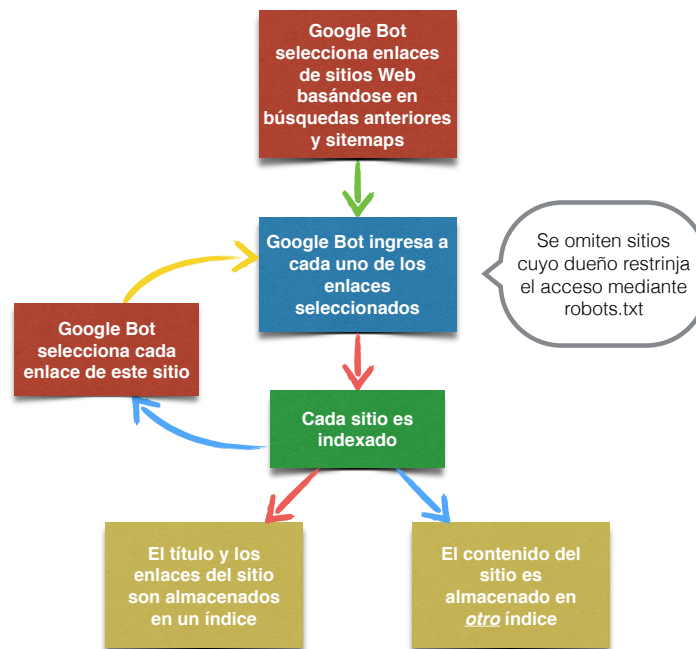


Figura 1: Rastreo e indexación de Google Bot.

Si se piensa con detenimiento, la idea de Google es bastante parecida a la conducta de algún investigador que recurre a un libro, esta persona sabe a qué libro consultar, pero no a qué capítulo ir, por lo tanto su primera idea será buscar en el índice, buscará el capítulo e inciso en cuestión para luego saber a qué página prestar atención.

Los índices son almacenados en los centros de datos (ó *Data Centers*), estos consisten en servidores de gran capacidad de almacenamiento. Google tiene 36 centros de datos activos a lo largo del mundo, 19 en Estados Unidos y 17 más en el resto del mundo, 1 de ellos en América Latina, específicamente en Quilicura, Chile. Entre ellos se tienen aproximadamente 2.000.000 de servidores y más de 10 Exabytes (10^{10} GB) de información. Cada centro de datos puede tener un tamaño de hasta 500.000 metros cuadrados y llegar a costar 600 millones de dolares (El ubicado en Chile costo 150 millones de dolares).

Estos tienen la siguiente arquitectura: existen racks de entre 40 y 80 equipos, cada rack tiene un switch de Ethernet, los racks están organizados en clusters conformados por 30 o más racks. El cluster se encarga del balance de carga, es decir, cuando alguien trata de conectarse a Google. Los servidores DNS traducen la dirección `www.google.com` a varias direcciones IP distintas permitiendo que se distribuya la carga a varios clusters, de forma que el hardware envía la consulta al servidor que se encuentre menos ocupado.

Los servidores dentro del centro de datos se dividen en distintas tareas, entre ellos se encuentran los servidores *Proxy Squid*, que aceptan la petición y devuelven el resultado desde el cache. Cuando la información no se encuentra en el cache esta es enviada al servidor web, este tiene la función de coordinar los envíos a los servidores índices, los cuales almacenan un conjunto de trozos del índice. Los servidores

índices son llenados con la información recopilada por los servidores de recolección. Por último se tienen los servidores de documentos que como su nombre lo dice almacena documentos y los servidores de anuncios que gestionan la publicidad.

4. Algoritmos de Búsqueda

El objetivo del buscador de Google es conseguir que el usuario obtenga la respuesta que esta buscando lo más rápido posible, creando una conexión casi perfecta entre el usuario y el conocimiento que este busca. Para ello, Google no entrega las millones de paginas que mencionen lo que se esta buscando, si no que filtra principalmente según la actualidad del contenido buscado, la región y el PageRank.

Para realizar estas tarea las base de datos debe actualizarse constantemente, con la finalidad que se entreguen datos de contingencia a la búsqueda. Además se rastrea la IP del usuario para entregar resultados de importancia regional. Luego de filtrar los resultados se les da un valor numérico a las paginas web, esto lo hace mediante el PageRank, con la finalidad de mostrar primero las paginas con el PageRank más alto.

4.1. PageRank

Se le llama PageRank a la gran conglomeración de diferentes protocolos y algoritmos que le dan una connotación numérica a cada sitio que sea indexado. El algoritmo original es una función matemática definida de la siguiente manera:

$$PR(A) = (1 - d) + d \sum_{i=1}^n \frac{PR(i)}{C(i)} \quad (1)$$

Con:

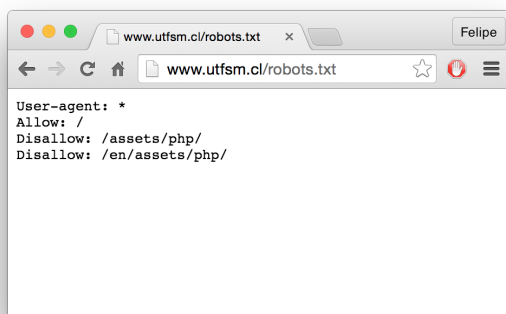
- $PR(A)$: PageRank asociado a un sitio A.
- d : Factor de Amortiguamiento.
- $PR(i)$: PageRank asociado a un sitio i -ésima que enlaza a la página A.
- $C(i)$: Número total de enlaces salientes de la página i -ésima.

Se puede concluir rápidamente a partir de la Ecuación (1) que se trata de un algoritmo recursivo, y relativamente simple. Se observa que el PageRank de un sitio A depende del de los sitios que enlazan a A, ponderado por el número de sitios enlazados de cada una de aquellas páginas y que es un valor numérico acotado entre 1 y 10. Además existe el factor de amortiguamiento “ d ” que representa la probabilidad de que se llegue al sitio A a partir de escribir la `url` o acceder a un hiperenlace que guíe a la página. Por lo tanto, este protocolo se enfoca en la cantidad de enlaces que citan a una página, y qué página las cite, puesto que una página con PageRank alto le ayudará a una página a que su PageRank sea más alto. Se podría decir que PageRank es una especie de concurso de popularidad donde cada voto es un hiperenlace y, bueno, la diferencia clara es que cada voto tiene un valor numérico diferenciado.

5. Demostraciones

5.1. Restricciones de indexación

A modo de ejemplo se ilustran en las Figuras 2 y 3 los archivos `robots.txt` de los sitios web `www.utfsm.cl` y `www.google.com`.

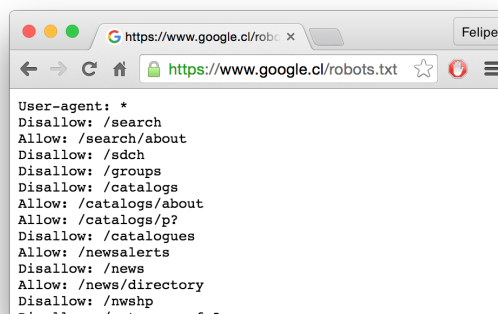


```

User-agent: *
Allow: /
Disallow: /assets/php/
Disallow: /en/assets/php/

```

Figura 2: Archivo `robots.txt` del sitio web `www.utfsm.cl`.



```

User-agent: *
Disallow: /search
Allow: /search/about
Disallow: /sdch
Disallow: /groups
Disallow: /catalogs
Allow: /catalogs/about
Allow: /catalogs/p?
Disallow: /catalogues
Allow: /newsalerts
Disallow: /news
Allow: /news/directory
Disallow: /nwshp
Disallow: /getnewsprofs?

```

Figura 3: Archivo `robots.txt` del sitio web `www.google.cl`.

Se puede observar a partir de la Figura 2 que se configuran restricciones para todos los rastreadores al utilizar la línea `User-agent: *`. En este caso se permite la indexación de todo el sitio (`Allow: /`) pero se agregan dos subdirectorios como excepción: `/assets/php/` y `/en/assets/php/`.

Por otro lado se observa en la Figura 3 que el mismo sitio de Google tiene configurado un archivo `robots.txt` para restringir su propio rastreador. Esto tiene lógica al darse cuenta que no es conveniente entregar como resultado de una búsqueda un enlace a otra búsqueda; esto se evita al agregar la línea `Disallow: /search`, restringiendo la indexación del subdirectorio que almacena las búsquedas.

5.2. Personalización de los resultados de búsqueda

Para demostrar los resultados realizó la búsqueda de una estación de gas y de un hotel utilizando una IP de Viña del Mar, Chile (Figuras 7 y 9) y otra IP de Ámsterdam, Holanda (Figuras 8 y 10) (Anexo 1).

En las figuras se aprecia que al realizar la búsqueda del mismo texto, se obtienen distintos resultados según la ubicación del usuario que realiza la búsqueda (particularmente es la localización de la IP que realiza la búsqueda). En las figuras 7 y 9 se obtienen resultados de estaciones de bencina y hoteles en las proximidades del plan de Viña del Mar, en cambio en las figuras 8 y 10 se tienen resultados de las proximidades de Ámsterdam, Holanda (Anexo 1).

En caso de que el usuario posea una cuenta Google, el algoritmo utilizará información recopilada creando un historial de búsqueda para que en futuras búsquedas arrojar el resultado óptimo según cada usuario.

En las siguientes figuras 4 5 se demuestra la capacidad de almacenamiento de Google.



Figura 4: Historial de búsqueda de Google toda la historia.



Figura 5: Historial de búsqueda de Google en pasado mes.

5.3. Comprobación del rastreador

Como medio de seguridad para los dueños de sitios web se puede verificar si alguno de los accesos al sitio corresponde o no a un rastreador. Aquí se analiza como comprobar que el acceso lo realizó Googlebot.

En un sistema con Windows se puede utilizar el comando `nslookup x.x.x.x` y en sistemas Unix se puede utilizar el comando `host x.x.x.x`; donde `x.x.x.x` corresponde a la dirección IP con la que se desea realizar el requerimiento. En la Figura 6 se observa que el resultado del requerimiento utilizando la IP `66.249.66.1` es el dominio `googlebot.com`, verificándose así que el acceso lo realizó el rastreador Googlebot.

```

FelipeFernandez — bash — 65x10
Last login: Sun Sep 27 23:07:04 on ttys000
MacBook-Pro-de-Pineit0r:~ FelipeFernandez$ host 66.249.66.1
1.66.249.66.in-addr.arpa domain name pointer crawl-66-249-66-1.googlebot.com.
MacBook-Pro-de-Pineit0r:~ FelipeFernandez$ host crawl-66-249-66-1.googlebot.com.
crawl-66-249-66-1.googlebot.com has address 66.249.66.1
MacBook-Pro-de-Pineit0r:~ FelipeFernandez$

```

Figura 6: Requerimiento a DNS inverso.

6. Conclusiones

A lo largo de este informe se respondieron cada una de las preguntas propuestas en la etapa introductoria. Se observa que al realizar una búsqueda Google no irá revisando página por página y enviará los resultados de búsqueda, sino que buscará en su propio índice de sitios web, que están jerarquizados mediante el protocolo PageRank. A los datos que entregue les dará particular importancia a la región en donde se realiza la búsqueda, además de la información reciente que se encuentre sobre ella.

El rastreo de páginas web continua actualizándose día a día dando especial cuidado a sitios de contingencia, noticias y avisos publicitarios relevantes.

En un comienzo, el conjunto de algoritmos llamado PageRank contemplaba solamente la ecuación (1), donde se evalúa la importancia de cada sitio indexado, destacando numéricamente el número de enlaces que dirigen al sitio en cuestión, y quienes lo citan.

Mediante la creación de un archivo `robots.txt` en la raíz de un sitio web se puede especificar las páginas de este sitio que se desean indexar y las que se quiere evitar su indexación.

Se puede corroborar la autenticidad de un rastreador mediante un requerimiento a DNS inverso con la IP del rastreador.


Toda la información indexada es guardada en Data Centers, que tienen gran capacidad de almacenamiento, estos distribuyen las tareas a los diversos servidores que lo conforman.






7. Referencias

1. <http://www.google.com/patents/US6285999>
2. <http://infolab.stanford.edu/~backrub/google.html>
3. <http://www.google.com/intl/es/insidesearch/howsearchworks/index.html>

8. Anexo 1

Filling station - Wikipedia, the free encyclopedia
https://en.wikipedia.org/wiki/Filling_station Traducir esta página
 Pasar a **Service stations** - [edit]. A service station or "servo" is the terminology predominantly used in Australia and New Zealand. In Australia, a servo is ...



A Shell Libertad Viña del Mar	 Sitio web	 Ruta
B Copec Quilota Viña del Mar		 Ruta
C Shell Camino Real Viña del Mar	 Sitio web	 Ruta

☰ Más gas station

Figura 7: Búsqueda utilizando IP perteneciente a Viña del Mar.

Mapa de gas station



Find BP Gas Stations Located Near Me | My BP Station

<https://mybpstation.com/station-finder> - En caché

Find BP gas stations near you quickly and easily with My BP Station. All stations locations across the U.S. are listed, so BP's got you covered.

 3150 W Chicago Ave, Chicago, IL 60608, Estados Unidos
 +1 773-722-9604

2 reseñas

Gas Station Locator | Find a Gas Station or APlus Convenience ...

<https://www.sunoco.com/gas-station-locator/> - En caché

Finding the nearest gas station is easy - enter a city and state or zip code to find a Sunoco gas station or APlus convenience store near you. You can narrow your ...

Filling station - Wikipedia, the free encyclopedia

https://en.wikipedia.org/wiki/Filling_station - En caché

[edit]. A service station or "servo" is the terminology predominantly used in Australia and New Zealand. In Australia, a servo is ...

Driver tries to scare off spider using lighter, sets gas station on fire ...

www.theguardian.com/.../driver-spider-petrol-gas-station-fire-lighter - En caché

2 hours ago ... A motorist scared of spiders gave himself a double shock when he accidentally set a gas station on fire trying to get rid of one using a cigarette ...

Shell Station Locator and Route Planner | Shell

www.shell.us/motorist/shell-station-locator-and-route-planner.html - En caché

Quickly find details of your nearest station or route. ... most advanced road fuel, actively cleans for better performance, creating a more exciting drive—staring ...

Shell Station Locator - Shell Global

www.shell.com/global/products-services/.../shell-station-locator.html - En caché -

Figura 8: Búsqueda utilizando un proxy con IP de Ámsterdam.

Hoteles, alojamientos baratos | Despegar.com
www.despegar.cl/hoteles/
 Encuentre ofertas en hoteles y alojamientos. Reserve online y descubra hoteles baratos, céntricos, económicos o lujosos en un solo lugar.

hotel

dom., 11 oct. lun., 12 oct.

Queen Royal Hotel
 3,7 ★★★★★ 10 opiniones
 Calle Cinco Norte

Hotel O'Higgins
 3,8 ★★★★★ 25 opiniones · Hotel de 3 estrellas
 Plaza Vergara

Mar Poniente
 4 opiniones
 Calle Uno Poniente

Más hotel

Hotel Plaza San Francisco - Santiago centro - Chile
www.plazasanfrancisco.cl/
 Hotel Plaza San Francisco ubicado en el centro de Santiago te invita a disfrutar una estadia de lujo a pasos de los atractivos turísticos de Santiago.
 Ubicación - Restaurant Bristol - Habitaciones - Capacitación-Seminarios

Hotel Torremayor Lyon - Providencia - Chile - Santiago
www.hotelorremayor.cl/
 Considerado uno de los Hoteles preferidos en la zona de Providencia, Hotel Torremayor destaca por su ubicación, agradable arquitectura, cálida decoración y ...

Imágenes de hotel Notificar imágenes





Figura 9: Búsqueda utilizando IP perteneciente a Viña del Mar.

Resultados de negocios locales que coinciden con hotel



A Hotel Okura Amsterdam
 020 678 7111

B Intel Hotels Rotterdam Centre
 010 413 4139

C Grand Hotel Karel V
 030 233 7555

D Intel Hotels Amsterdam Zaandam
 075 631 1711

E Hotel Arena
 020 850 2400

F Lloyd Hotel & Cultural Embassy
 020 581 3836

G Hampshire Hotel - Theatre District Amsterdam
 020 607 7900

Hoteles.com - Encuentra y reserva hotel entre más de 257.000 ...
es.hotels.com/ - En caché - Similares
 En Hoteles.com encuentra miles de ofertas en mas de 240 000 hoteles, desde hoteles de lujo hasta hoteles mas económicos.

Hotels.com - Encuentra y reserva hotel entre más de 257.000 ...
es.hotels.com/ - En caché - Similares
 En Hoteles.com encuentra millares de ofertas de hoteles para reserva, desde hotel con todo incluido hasta los más baratos.
 Hotels.com® Rewards - Ofertas de hotel - Hoteles en Las Vegas - Hoteles en Miami

Hotels.com | Cheap Hotels, Discounts, Hotel Deals and Offers
www.hotels.com/ - En caché - Similares

Figura 10: Búsqueda utilizando un proxy con IP de Ámsterdam.