

Capítulo 6

Retorno de la información en la web

Retorno de información en la web

- Las máquinas de búsqueda en la web procesan los sitios web y colecciones de documentos.
 - Los documentos son páginas web.
 - Cubren una parte de la web.
 - Se mantiene un repositorio indexado de páginas web, usualmente usando índices invertidos.
 - Más aun, se deben actualizar los índices regularmente.
 - Los resultados a consultas son las páginas web más relevantes para el usuario ordenadas en orden descendente de relevancia.
- Las máquinas de búsqueda verticales son personalizadas para tópicos específicos; cubren una colección específica de documentos en la web.

Rastreadores web

- **Problema:** En la web la colección de documentos donde hacer las búsquedas no viene dada de antemano.
 - Hay que encontrar y construir la colección
 - **Idea:** Para ello se pueden aprovechar los **hiperenlaces**.
- **Solución:** Los **rastreadores web** son **programas que localizan** y **recolectan** información en la web.
 - Se **siguen los hiperenlaces** presentes en **documentos conocidos** para **encontrar otros documentos**.
 - Se puede comenzar por un **conjunto semilla** de documentos.
 - **Toma meses** realizar un rastreo.

Rastreadores web

- **Problema:** Como la web es inmensa, la etapa de rastreo puede tomar demasiado (años o meses).
- **Solución:** El rastreo **se hace por varios procesos en varias máquinas ejecutando en paralelo.**
 - El conjunto de enlaces a ser rastreados se puede almacenar en una base de datos.
 - Los nuevos enlaces encontrados en las páginas rastreadas se pueden añadir a este conjunto para ser rastreados a continuación.

Indexado

- Los documentos recolectados por los rastreadores son procesados por un **sistema de indexado**.
- Los documentos recolectados pueden ser descartados luego del indexado o almacenados como una **copia en caché**.
- **Problema:** cuando se está haciendo el indexado se debería poder contestar consultas al mismo tiempo.
 - El indexado toma tiempo.
 - Las máquinas de búsqueda en la web no pueden parar.
- **Solución:**
 - El proceso de indexado **se ejecuta** en varias máquinas.
 - **Se crea** una **nueva copia** del **índice** en lugar de modificar el índice viejo.
 - El **índice viejo** **se usa** para **contestar consultas**.
 - **Luego de completar la fase de rastreo** y justa antes de iniciar un nuevo indexado, el **índice nuevo** **se convierte** en el **índice viejo**.

Búsquedas

- **Problema:** puede haber demasiadas consultas simultáneas para una máquina de búsqueda.
- **Solución:** Usar múltiples máquinas para contestar consultas.
 - Se pueden mantener los índices en memoria
 - Las consultas pueden ser enrutadas a diferentes máquinas para balanceo de carga.

Relevancia de documentos en la web

- Las **máquinas de búsqueda en la web tempranas** ordenaban las respuestas basándose solo en relevancia TF-IDF.
- **Pero usar este enfoque para la web tiene sus problemas:**
 - El número de documentos relevantes a una consulta puede ser enorme si solo frecuencia de términos es usada.
 - Muchas páginas que los usuarios quieren ver pueden tener frecuencia de términos baja y no van a obtener valor TF-IDF alto.
 - Puede haber **páginas spam**: Es solo agregar ciertas palabras en una página sin valor para que aparezca en las búsquedas.
 - La mayoría de los usuarios se interesa por las páginas de los **sitios populares**.

Relevancia de documentos en la web

- **Solución:** usar la **popularidad de un sitio web** (por ejemplo, cuánta gente lo visita) para dar rango a sus páginas web en el resultado de las consultas.
 - Pero es difícil encontrar la popularidad real de un sitio.
 - Porque para obtener esta información es necesaria la cooperación del sitio.
 - Algunos sitios podrían mentir sobre esto para obtener un rango mayor.
- **Solución más refinada:** medidas tradicionales de relevancia de una página como **TF-IDF** pueden ser combinadas con la **popularidad** del sitio de la página
 - para obtener una **medición global de la relevancia** de la página para una consulta.
 - Las páginas con mayor valor de relevancia pueden retornarse como las respuestas top de la consulta.

Popularidad de sitios web

- Ahora vemos como medir la popularidad de un sitio web.
- **Idea 1:** usar los **archivos de marcadores** de cada usuario (para sus navegadores) para saber qué páginas son populares:
 - Los sitios que aparecen en una gran cantidad de archivos de marcadores se los puede considerar como muy populares.
 - Sin embargo, los archivos de marcadores usualmente son almacenados privadamente y no son accesibles en la web.

Popularidad de sitios web

- Las páginas web tienen información muy importante de la cual carece el texto plano: **hiperenlaces**.
- **Idea 2:** usar el **número de enlaces a un sitio S** como una medida de la **popularidad o prestigio del sitio**.
 - Contar un enlace a S desde cada sitio que lo enlaza a S .
 - La popularidad es para sitios, no para páginas web (la mayoría de los enlaces son a la página principal del sitio).
- **Refinamiento de idea 2:** Cuando se computa el prestigio basado en enlaces a un sitio,
 - se le puede dar más peso a los **enlaces de sitios** que tienen **mayor prestigio**.

Popularidad de sitios web

- **Idea 3:** Las máquinas de búsqueda llevan la pista de qué fracción de veces los usuarios seleccionan una página retornada como resultado de una búsqueda.
 - Esta medida puede ser usada como una medición de popularidad del sitio.

PageRank

- **Algoritmo de ordenamiento PageRank**

- Es un algoritmo definido por Google como una medida de popularidad de una página basada en la popularidad de las páginas que la enlazan.
- Analiza los enlaces hacia fuera y los enlaces hacia dentro.
- Se considera a las páginas altamente enlazadas por otras páginas como más importantes (con mayor autoridad) que las páginas con menos enlaces hacia ellas.
- Se dice que una página P tiene un rango alto si la suma de los rangos de las páginas que apuntan a P tiene un valor alto.

PageRank

- **Algoritmo de ordenamiento PageRank (continuación)**

- Supongamos que una persona navegando en la web realiza una caminata aleatoria de páginas web de la siguiente manera:

- El primer paso comienza en una **página web aleatoria**.
 - En los pasos siguientes la persona hace una de las siguientes:
 - Con una probabilidad δ la persona salta a una página web elegida aleatoriamente.,
 - Con una probabilidad $1 - \delta$ la persona elige aleatoriamente uno de los enlaces hacia afuera de la página actual y sigue ese enlace.

PageRank

- **Algoritmo de ordenamiento PageRank (continuación)**

- Asumiendo la caminata aleatoria anterior, el **PageRank de una página** web es la **probabilidad** de que esa **persona** **visite la página** en un **determinado punto** del tiempo.

- Las páginas apuntadas por muchas páginas web es más probable que sean visitadas y van a tener un PageRank más alto.
 - Las páginas apuntadas por páginas web con alto PageRank van a tener una mayor probabilidad de ser visitadas y entonces van a tener un mayor PageRank.

PageRank

- **Algoritmo de ordenamiento PageRank (continuación)**

- PageRank puede definirse por un **conjunto de ecuaciones lineales**.
- Primero se les da **identificadores enteros** a las páginas web.
- La **matriz de probabilidad de salto** T : $T[i, j]$ es la probabilidad que un caminante aleatorio que sigue un enlace afuera de la página i siga el enlace hacia la página j .
- Suponiendo que cada enlace de i tiene la misma probabilidad de ser seguido, $T[i, j] = 1/N_i$, donde N_i es el número de enlaces afuera de la página i .

PageRank

- **Algoritmo de ordenamiento PageRank (continuación)**

- El PageRank de la página j puede definirse como:

$$P[j] = \delta/N + (1 - \delta) * \sum_{i=1}^N (T[i, j] * P[i])$$

- Donde δ es una constante entre 0 y 1 y N es el número de páginas.
 - δ representa la probabilidad de un paso en la caminata aleatoria sea un salto aleatorio.

PageRank

- **Algoritmo de ordenamiento PageRank (continuación)**

- El **conjunto de ecuaciones** se resuelve usando una técnica iterativa.
- Se comienza inicializando cada $P[i]$ a $1/N$.
- Cada paso de la iteración computa nuevos valores para $P[i]$ usando los valores de P de la iteración previa.
- La iteración acaba cuando el valor máximo de cambio en un valor de $P[i]$ en la iteración va por debajo de un cierto valor de corte.

Popularidad de sitios web

- **Problema:** PageRank asigna una medida de popularidad que no considera los términos de la consulta.
- **Solución:** usar las palabras clave en el texto del áncora de los enlaces a una página para juzgar para qué tópicos la página es altamente relevante.
 - La *popularidad basada en texto de áncoras* puede ser usada en combinación con otras medidas de popularidad y con TF-IDF para obtener un ranking de los resultados de una consulta.
- **Implementación 1:** Si se considera el texto de esas áncoras como parte de la página apuntada, entonces TF-IDF toma texto de áncoras en cuenta.

Popularidad de sitios web

- **Implementación 2:** se computa una medida de popularidad usando solo páginas que contienen los términos de la consulta en lugar de computar popularidad usando todas las páginas web disponibles.
 - Este enfoque es más costoso porque el computo del ranking de popularidad debe ser hecho dinámicamente cuando se recibe una consulta,
 - mientras que PageRank se computa estáticamente una vez y se reutiliza para todas las consultas.
 - El algoritmo Hits se basa en esta idea de implementación.

Relevancia de documentos en la web

- Los **valores TF-IDF** de una página **deben ser combinados** con el **ranking** de **popularidad** (p.ej. PageRank).
- **Problema:** ¿cómo combinar TF-IDF con otras medidas de popularidad?
 - Este es un secreto mantenido por muchas empresas.
 - Esto ayuda a protegerse de la competencia y de los que quieren producir páginas spam.
- **Idea de solución:** una **fórmula** que combina puntajes es fija y **toma como parámetros** pesos para cada factor considerado.
 - Se usa un **conjunto de entrenamiento** de resultados de consultas cuyos rangos son fijados por humanos.
 - Usando el mismo, un **algoritmo de entrenamiento automático** puede calcular valores para esos parámetros.

Buscador de Google

- **Características principales del buscador de Google para la web:**

- **Rastreo:** Google utiliza **arañas web** (Googlebot) que **recorren la web** para encontrar nuevas y actualizadas páginas.
- **Indexación:** Las **páginas encontradas** son almacenadas en el **índice** de Google, una base de datos masiva que contiene la copia de todas las páginas web visitadas por los rastreadores.
- **Algoritmos de búsqueda:** Google emplea **complejos algoritmos** para **interpretar las consultas** de los usuarios y encontrar los resultados más relevantes. Estos algoritmos consideran cientos de factores, incluyendo palabras clave, sinónimos y la ubicación del usuario.
- **Relevancia de respuestas:** Los algoritmos de Google no sólo **encuentran coincidencias** de **palabras clave**, sino que también **intentan entender el contexto** y la **intención** de la **búsqueda** para ofrecer los resultados más útiles y precisos.

Buscador de Google

- **Características del buscador de Google (Cont):**

- **Rankings y PageRank:** Google clasifica los resultados de búsqueda basándose en la autoridad y relevancia de las páginas, utilizando medidas como el PageRank, que evalúa la cantidad y calidad de los enlaces hacia una página.
- **Personalización:** Las búsquedas se personalizan en función de la historia del usuario, su ubicación y otros datos disponibles.
- **Calidad del contenido:** Google valora contenido original, de alta calidad y relevante, penalizando a sitios con prácticas de spam o contenido duplicado.
- **Lenguaje de consulta:** ya hablamos de ello en el archivo de filminas anterior.

Buscador de Google

- **Datos que contiene un enlace devuelto por el buscador de Google:**
 1. **Título de la página:** El encabezado principal, generalmente en azul, que resume de qué trata la página.
 2. **URL:** La dirección web del sitio, a menudo en verde. Esto te da una idea de la fuente del contenido.
 3. **Extracto o snippet:** Un breve fragmento de texto que Google extrae de la página para mostrar una vista previa del contenido. Esto suele estar en gris y puede incluir las palabras clave de tu búsqueda resaltadas en negrita.
 4. **Breadcrumb:** Una ruta de navegación que muestra dónde se encuentra la página dentro de la jerarquía del sitio web.
 5. **Fecha:** En algunos casos, especialmente en artículos y noticias, Google incluye la fecha de publicación o de la última actualización.

Buscador de Google

- El buscador de Google ofrece mucho más que solo enlaces como resultado de las búsquedas:
 - 1. Fragmentos destacados:** Resúmenes directos de respuestas que se muestran en la parte superior de la página de resultados, a menudo en formato de párrafo (p.ej. definiciones, descripciones), lista (p.ej. pasos de instrucciones o listas de elementos) o tabla (p.ej. datos comparativos en forma tabular).
 - 2. Panel de conocimiento:** Información detallada sobre un tema, persona, lugar o cosa, normalmente a la derecha de los resultados.
 - 3. Videos:** Recomendaciones de videos de plataformas como YouTube que son relevantes para tu consulta.
 - 4. Imágenes:** Una selección de imágenes relacionadas con tu búsqueda.
 - 5. Preguntas relacionadas:** Una lista de preguntas adicionales que otros usuarios han hecho sobre el mismo tema o uno similar, junto con respuestas breves que se pueden expandir para obtener más información.
 - 6. Gráficos y estadísticas:** En algunos casos, como para datos financieros o demográficos, Google puede mostrar gráficos interactivos.