

## Resumen 9

---

Edgar André Araya Vargas 2020142856

En este documento se nos introduce nuevamente pero desde una vista diferente el concepto de índices invertidos utilizados en motores de búsqueda web e intranet para ejecutar consultas de texto. Estos índices invertidos almacenan, para cada palabra, la lista de documentos en los que aparece, en contraposición a la forma tradicional de almacenar documentos como listas de palabras. La actualización de los índices invertidos es costosa, por lo que los motores de búsqueda suelen reconstruirlos desde cero periódicamente. Cuanto más frecuentemente se reconstruya el índice, más rápido se reflejarán las actualizaciones en los resultados de búsqueda, lo que mejora la calidad de la búsqueda. Los índices invertidos de alto rendimiento son ampliamente utilizados en motores de búsqueda web y de búsqueda.

Entonces, en pocas palabras, se puede decir que los documentos se almacenan como listas de palabras, pero los índices invertidos invierten esta estructura, almacenando para cada palabra la lista de documentos en los que aparece, de ahí su nombre. Existen diferentes variaciones en los índices invertidos pero se comparte que lo que requiere como mínimo es almacenar para cada palabra la lista de documentos en los que aparece. Si se desean admitir consultas de frases y proximidad, es necesario almacenar las posiciones de las palabras en cada documento. La granularidad de las posiciones puede variar desde desplazamientos de bytes hasta palabras, párrafos o secciones, pero generalmente se almacena a nivel de posición de palabra. De alguna manera también se puede apoyar con el posible almacenamiento de la frecuencia de las palabras en cada documento en lugar de las posiciones de las palabras.

Almacenar la frecuencia total de cada palabra puede ser útil para optimizar los planes de ejecución de las consultas. Algunas implementaciones almacenan dos listas invertidas: una que contiene solo las listas de documentos (y frecuencias de palabras) y otra que almacena las listas completas de posiciones de palabras. De esta manera, las consultas simples se pueden responder consultando solo las listas de documentos, que son mucho más cortas. Algunas implementaciones incluso van más allá y almacenan metainformación sobre cada "hit" o posición de palabra, utilizando uno o dos bytes que contienen información como el tamaño de fuente y el tipo de texto. Esta información se puede utilizar para mejorar la clasificación de los resultados de búsqueda, ya que las palabras con formato especial suelen tener mayor importancia.

Otra variación posible es si se almacena el léxico por separado o no. El léxico contiene todos los tokens indexados en la colección y, generalmente, también almacena información estadística, como el número de documentos en los que aparece cada token. El léxico puede ser útil en varias formas, a las que se hará referencia más adelante.

El tamaño del índice invertido varía entre el 5% y el 100% del tamaño total de los documentos indexados. Esta amplia variación se debe a las diferentes implementaciones de los índices invertidos. Algunas almacenan posiciones de palabras, mientras que otras no lo hacen. Algunas realizan un procesamiento agresivo de los documentos para reducir el tamaño del índice, otras no lo hacen. Algunas admiten actualizaciones dinámicas, lo que puede causar fragmentación y requiere espacio adicional para futuras actualizaciones. Además, algunas implementaciones utilizan métodos de compresión más potentes pero más lentos que otras.

En resumen, los índices invertidos son una herramienta fundamental en los motores de búsqueda para realizar consultas de texto. Almacenar la lista de documentos en los que aparece cada palabra permite una búsqueda eficiente y la posibilidad de realizar consultas más complejas, como las de frases y proximidad. Sin

embargo, la implementación y configuración de los índices invertidos puede variar ampliamente, lo que afecta tanto al rendimiento como al tamaño del índice resultante.

Es importante en el documento como se describen las técnicas necesarias para implementar un motor de búsqueda utilizando un índice invertido. Se mencionan tres etapas principales en el procesamiento de documentos: la tokenización, el stemming y la eliminación de palabras irrelevantes. En la tokenización, los documentos se convierten en una lista de tokens, donde cada token es una palabra alfanumérica. El stemming implica transformar cada palabra en su raíz morfológica para reducir el tamaño del índice y aumentar los resultados relevantes. Por último, se eliminan las palabras irrelevantes, como artículos y preposiciones comunes.

Algo que también podemos encontrar expresado en el documento son los diferentes tipos de consultas que se pueden realizar en un motor de búsqueda. Las consultas normales buscan documentos que contengan un solo término, mientras que las consultas booleanas permiten combinar términos utilizando operadores lógicos como AND, OR y NOT. Las consultas de frase se utilizan para encontrar documentos que contengan palabras específicas en un orden determinado. Las consultas de proximidad buscan términos que estén cerca uno del otro en un número específico de palabras. Por último, las consultas con comodines permiten realizar búsquedas aproximadas mediante el uso de caracteres comodín.

El ranking de resultados es crucial para las aplicaciones de búsqueda, ya que permite mostrar los documentos más relevantes primero. Tradicionalmente, se ha utilizado una medida de similitud entre la consulta y el documento para clasificar los resultados. Esta medida se basa en factores como la frecuencia de términos en el documento y en la colección, la longitud del documento y de la consulta, y el factor de inversa de la frecuencia de documentos (IDF). Sin embargo, las medidas de similitud tradicionales no funcionan bien para consultas cortas, que son comunes en los motores de búsqueda web. En estas consultas, los documentos con varias instancias del término más raro tienden a clasificarse primero, incluso si no contienen los otros términos de la consulta... Esto no cumple con las expectativas de los usuarios, que esperan que los documentos que contengan todos los términos de la consulta se clasifiquen en primer lugar.

Para optimizar la evaluación de consultas, se han propuesto estrategias de poda dinámica que permiten omitir o procesar parcialmente las listas invertidas de los términos de la consulta. Estas estrategias se dividen en grupos seguros y no seguros, dependiendo de si producen los mismos resultados que las consultas no optimizadas. También se han investigado diferentes métodos de evaluación de consultas, como la evaluación término por término o documento por documento.

Para concluir este resumen del documento de índices invertidos me parece importante recalcar nueva y ultimamente como la construcción de un índice invertido es un aspecto fundamental en los motores de búsqueda web. A través de este artículo, se ha explorado cómo construir eficientemente dicho índice, considerando diversas estrategias y técnicas. Se ha destacado la importancia del ranking de resultados y se ha señalado la necesidad de optimizar la evaluación de consultas para mejorar la calidad de los resultados y reducir el tiempo de procesamiento. En definitiva, se han presentado avances significativos en la investigación sobre este tema, aunque existen desafíos continuos debido a la constante evolución de los motores de búsqueda y a la protección de los resultados por parte de los actores del mercado.