

Hashing

Look-up tables

Una tabla de look-up es un diccionario con velocidad de acceso muy eficiente. Suele ser $O(1)$. Se implementan con un arreglo en el que la clave del diccionario se usa como índice al arreglo. Esta eficiencia la tienen tanto para lectura como para escritura.

Esta eficiencia se logra calculando el índice en el cual colocar un dato a partir de su clave. Entonces tengo una forma de corresponder una clave con un índice. El rendimiento en memoria de una look-up table, entonces, será mejor para tablas más densas.

Hash Tables

Para mejorar este problema de la eficiencia de memoria de las tablas look-up se hicieron las tablas hash. En una tabla hash, la clave se convierte mediante una **función hash**, que devuelve un **valor hash**, que es el índice en el arreglo de la tabla.

En una tabla hash, las entradas no están ordenadas. Esto quiere decir que no se puede hacer una búsqueda dentro de un rango de la manera que se podría en las tablas look-up comunes.

La ventaja de las tablas hash es que permiten mapear entradas más dispersas a un espacio denso, para mejorar la eficiencia en memoria.

Cosas que las tablas hash deberían hacer:

- producir, para cada clave, un valor determinístico
 - si no sucede así entonces una clave puede corresponder a más de una entrada. Esto es desastroso
- aprovechar todos los bits
- los valores resultantes deberían estar distribuidos uniformemente en todas las posibles entradas de la tabla hash
 - no queremos que valores se concentren en una región de la tabla hash
- produzcan valores hash muy diferentes para entradas muy similares
- ser eficientes en cálculo

Funciones Hash

Estas hay que tomarlas como si fueran un camino solamente de ida. Este es el motivo por el cual funcionan. Una función hash tiene como objetivo tomar una clave y devolver un valor hash. Sirven para seguridad y autenticación porque es prácticamente imposible partir del valor hash y obtener la clave. Según el número de bits que tenga la función hash, se puede tardar más o menos en hallar la clave. Una función hash de 256 bits devuelve hasta 2^{256} valores hash.

Ejemplos de funciones hash:

- **Folding:** esta separa la clave en secciones de N bits y luego las suma o combina con XOR
- **Mid-squares:** esta eleva la clave al cuadrado y devuelve los N bits centrales de ese resultado (esta no da los mejores resultados)
- **Multiplicativa:** $h(n) = w * n$, donde w es un coprimo del tamaño de la tabla hash.
- **Criptográficas:** MD5, SHA-1, SHA-2, SHA-3

Puede suceder que dos claves distintas devuelvan hashes iguales. Esto se conoce como una colisión. Para esto se escriben **políticas de resolución de colisiones**.

Una política **abierta** inserta una segunda estructura dentro de la tabla. Cada entrada sería una lista o un AVL. Entonces varias entradas pueden vivir con distintas claves y valores en el mismo índice en la tabla. Esto permite crear una nueva función hash a partir de la anterior. $h_2(n) = h(n+K)$.

Una política **cerrada** guarda el valor en la siguiente posición libre en la memoria. Es menos usado.

Factor de Carga

Es la relación entre la cantidad de entradas ocupadas en la tabla hash y el tamaño de la tabla hash. Las tablas hash suelen ser óptimas cuando el factor de carga está entre 0,6 y 0,75. Si es superior, tendremos muchas colisiones y si es menor se está desperdiciando mucha memoria. Cuando nos estamos yendo del rango de valores eficientes, es conveniente redimensionar dinámicamente la tabla hash.

El problema con esto es que van a cambiar los índices de la tabla. Para subsanar este efecto es que se lleva a cabo el proceso de **rehashing**. Este consiste en volver a calcular los valores hash e índices para todas las claves previamente existentes. Puede ser un proceso largo ya que no hay más remedio que iterar sobre todos los elementos existentes en la tabla hash original.

En C++

En STL, las tablas hash estan implementadas en `unordered_set`, `unordered_map`, `unordered_multiset` y `unordered_multimap`.