

PRÁCTICA HASH: ANÁLISIS DE DISTINTAS IMPLEMENTACIONES

Existe en la aplicación del vacunódromo una base de datos de 10.000 pacientes, almacenada en un archivo de texto (`pacientes_vacunodromo.txt`). Dicho fichero contiene en cada línea los datos de contacto de un/a paciente: su nombre y apellidos, su nickname o nombre de usuario para la plataforma y su correo electrónico asociado. Se deben cargar estos/as pacientes en memoria utilizando una tabla hash de forma que la inserción, búsqueda y eliminación de pacientes sea lo más eficiente posible. Como **CLAVE** para almacenar los/as pacientes en la tabla se utilizará en primera instancia el campo del nickname de usuario/a, que consiste en el carácter @ seguido del nombre y las iniciales de los apellidos de la persona, con un dígito a continuación en caso de que haya usuarios/as que repitan el mismo nombre e iniciales que alguno/a ya existente. Las cuentas de correo son de la forma `nombre.apellido1.apellido2@vacunas.gal`, con un dígito al final del segundo apellido en caso de que haya usuarios/as coincidentes en nombre y apellidos.

El objetivo de este ejercicio es analizar el comportamiento de las distintas implementaciones de tabla hash estudiadas en cuanto a eficiencia y recursos consumidos (número de operaciones necesarias, memoria ocupada). Se proporcionan dos implementaciones diferentes del tipo de dato `TablaHash`, una con resolución de colisiones por encadenamiento, y la otra con resolución de colisiones por recolocación:

- La implementación con encadenamiento (`TablaHashEncadenamiento.zip`) incluye el tipo lista de pacientes para las listas enlazadas en cada posición de la tabla.
- La implementación con recolocación (`TablaHashRecolocacion.zip`) incluye en el código dos estrategias diferentes para la recolocación (lineal y cuadrática).

Ambas implementaciones incluyen 3 tipos de funciones hash, con objeto de poder estudiar qué función distribuye mejor las claves en la tabla.

OBJETIVOS ESPECÍFICOS

Analizar el comportamiento de distintas implementaciones prestando atención al número de operaciones requeridas a la hora de insertar o buscar datos en la tabla. Para ello debéis implementar un mecanismo que cuente estas operaciones, modificando para ello el TAD `TablaHash` de manera que se muestren las “colisiones” producidas durante una operación sobre los datos (inserción o búsqueda). Esto no es una operación que esté contemplada en el TAD, sino que la añadimos para poder realizar la experimentación.

(0) Definición de colisión en inserción y búsqueda. Implementación de contadores de operaciones necesarias para insertar o buscar un elemento.

En una tabla hash se define una colisión como la situación en la que una función hash devuelve el mismo valor para dos claves diferentes. Esto implica que haya que realizar más pasos de los esperados teóricamente para insertar o buscar un elemento en la tabla. Debéis explicar cómo modificáis el código proporcionado para poder calcular los datos siguientes para inserción y búsqueda:

- **Inserción:**
 - **nColisiones:** colisiones producidas: en el momento en que se obtiene el hash de un elemento. Debe comprobarse si ese hash ya fue obtenido anteriormente (en esa posición ya hay otro dato). Este dato debe calcularse tanto en encadenamiento como en recolocación.
 - **nPasosExtra:** Los pasos adicionales que requiere la inserción del dato: en la estrategia de recolocación, hay que buscar una posición libre, y puede no encontrarse a la primera. Hemos de contar cuántos intentos son necesarios para ubicar el dato a insertar en caso de colisión hash. Esto sólo sucede en recolocación, ya que en

encadenamiento el dato simplemente se inserta al principio de la lista correspondiente a su código hash.

- **Búsqueda:**

- **nPasosExtraB:** Hay que calcular los **pasos adicionales que requiere la búsqueda de un dato** cuando se produce una colisión hash: porque el dato no está de primero en la lista de esa posición (encadenamiento), o porque no está en la posición indicada (fue recolocado).

Una vez realizadas las modificaciones, y disponibles estos contadores de colisiones, y de operaciones “extra”, se procederá con los siguientes experimentos.

(1) Influencia del tamaño (N) elegido para la tabla hash en el proceso de inserción

El primer experimento consiste en observar el comportamiento de las implementaciones con encadenamiento y con recolocación simple (lineal con constante $a=1$), en cuanto a **número de colisiones y operaciones extra requeridas en la inserción** para diferentes tamaños de la tabla hash con la función **hash2**. Para ello leeremos todos los datos del fichero y los insertaremos en la tabla correspondiente.

Pasos a seguir:

1. **Encadenamiento:** modificar el programa base para añadir un contador de colisiones (**nColisionesI**) en el proceso de inserción.
2. **Recolocación:** modificar el programa base para añadir un contador de colisiones (**nColisionesI**) y de operaciones extra (**nPasosExtraI**) en el proceso de inserción utilizando **recolocación simple con $a=1$** .
3. Para ambos casos, debéis probar con la función **hash2** y con **distintos tamaños N de tabla**, considerando lo visto en clase respecto al factor de carga y a los números primos. Seleccionad los tamaños a probar en función de los valores recomendados para cada tipo de tabla para poder comprobar las propiedades teóricas. Se trata de elegir tamaños que se espera sean los adecuados, y alguno más que se espera que no lo sea, para observar las diferencias. **Debéis justificar la elección de los valores con los que probáis recordando que para encadenamiento se recomienda un factor de carga $L \leq 0,75$ y que para recolocación se recomienda un factor de carga $L \leq 0,5$.**

Tabla 1. Influencia de N en el proceso de inserción usando la función hash2

Encadenamiento	N=n1	N=n2	N=n3	N=n4	N=n5	N=n6
nColisionesI						
Recolocación simple ($a=1$)	N=n1	N=n2	N=n3	N=n4	N=n5	N=n6
nColisionesI						
nOperacionesExtraI						

CONCLUSIÓN: ¿Cuál es el tamaño N que tiene mejor comportamiento para Encadenamiento y Recolocación?

(2) Influencia de la función hash (1-2-3) y de la clave en el proceso de inserción (nColisionesI)

Para los tamaños N seleccionados en el apartado anterior, se debe comprobar y explicar cómo afecta la selección de la función hash en cuanto al número de colisiones producidas en el proceso de inserción. Probad con las 3 funciones hash (1-2-3) y, para la última, analizad al menos dos valores distintos de la constante K. Debéis explicar además el motivo por el que la función hash tipo 1 es especialmente mala para este problema concreto.

Tabla 2. Influencia de la función hash en el proceso de inserción en Encadenamiento

Clave=nickname

	N=n1	N=n2	N=n3	N=n4	N=n5	N=n6
nColisionesI Hash1						
nColisionesI Hash2						
nColisionesI Hash3 K=500						
nColisionesI Hash 3 K=otro valor						

Tabla 3. Influencia de la función hash en el proceso de inserción en Recolocación Lineal (a=1)

Clave=nickname

	N=n1	N=n2	N=n3	N=n4	N=n5	N=n6
nColisionesI Hash1						
nPasosExtral Hash1						
nColisionesI Hash2						
nPasosExtral Hash2						
nColisionesI Hash3 K=500						
nPasosExtral Hash3, K=500						
nColisionesI Hash 3 K=otro valor						
nPasosExtral Hash3, K=otro valor						

Además, para ver cómo influye la selección de la clave utilizada para indexar los datos de entrada, haremos los mismos experimentos pero utilizando como campo clave el correo electrónico del usuario, en lugar del nickname, por lo que debéis modificar el tipo de elemento en la tabla adecuadamente.

Tabla 4. Influencia de la función hash en el proceso de inserción en Encadenamiento
clave=correo

	N=n1	N=n2	N=n3	N=n4	N=n5	N=n6
nColisionesl Hash1						
nColisionesl Hash2						
nColisionesl Hash3 K=500						
nColisionesl Hash 3 K=otro valor						

Tabla 5. Influencia de la función hash en el proceso de inserción en Recolocación Lineal (a=1)
clave=correo

	N=n1	N=n2	N=n3	N=n4	N=n5	N=n6
nColisionesl Hash1						
nPasosExtral Hash1						
nColisionesl Hash2						
nPasosExtral Hash2						
nColisionesl Hash3 K=500						
nPasosExtral Hash3, K=500						
nColisionesl Hash 3 K=otro valor						
nPasosExtral Hash3, K=otro valor						

(3) Influencia de estrategia de recolocación en el proceso de inserción

Después de la experimentación anterior, debéis fijar el tamaño y la función hash que combinados funcionen mejor (produzcan menos colisiones). Elegid un par de combinaciones que hayan tenido buen comportamiento.

- Caso 1: Tamaño/funciónHash?
- Caso 2: Tamaño/funciónHash?

En esta experimentación, para los casos 1 y 2, estudiaremos el comportamiento del número de colisiones en la inserción ($n_{\text{ColisionesI}}$) y el número de pasos extra ($n_{\text{PasosExtral}}$) al variar la estrategia de recolocación: simple (lineal con $a=1$), lineal (con al menos dos valores distintos de la constante a) y cuadrática. **Intenta justificar los resultados obtenidos.**

Tabla 6. Influencia de la estrategia de recolocación en el proceso de inserción

Casos	Variables	Simple (lineal con $a=1$)	Lineal ($a=\text{valor1}$)	Lineal ($a=\text{valor2}$)	Cuadrática
Caso 1	$n_{\text{ColisionesI}}$				
	$n_{\text{PasosExtral}}$				
Caso 2	$n_{\text{ColisionesI}}$				
	$n_{\text{PasosExtral}}$				

(4) Estimación de la eficiencia en el acceso a los datos (búsqueda)

Compararemos en este apartado el número de pasos requeridos en la búsqueda de los elementos de la tabla. Vamos a buscar todos los elementos una vez insertados, y comprobar el número total de operaciones extra necesarias para encontrarlos (**nOperacionesExtraB**).

Comprobad los resultados para las implementaciones con encadenamiento, recolocación simple, y recolocación cuadrática, con distintos valores de tamaño de tabla y función hash.

Debéis por tanto, repetir la experimentación de las tablas 1, 2, 3 y 6 para analizar el **nOperacionesExtraB**, lo que se puede resumir en rellenar la tabla siguiente:

Tabla 7. Influencia de N en el proceso de búsqueda (nOperacionesExtraB)

	N=n1	N=n2	N=n3	N=n4	N=n5	N=n6
Encadenamiento Hash 1						
Encadenamiento Hash 2						
Encadenamiento Hash 3, k=500						
Encadenamiento Hash 3, k=otro valor						
Recolocación simple (a=1) Hash 1						
Recolocación simple (a=1) Hash 2						
Recolocación simple (a=1) Hash 3, k=500						
Recolocación simple (a=1) Hash 3, k=otro valor						
Recolocación cuadrática Hash 1						
Recolocación cuadrática Hash 2						
Recolocación cuadrática Hash 3, k=500						
Recolocación cuadrática Hash 3, k=otro valor						

Conclusión: ¿Cuál es la mejor estrategia de implementación para los datos proporcionados?

FORMATO DE ENTREGA:

Para este trabajo debéis entregar a través del aula virtual un fichero comprimido, con nombre **ApellidosNombre_5.zip**, incluyendo el proyecto o proyectos completos que hayáis implementado, así como un documento pdf analizando y razonando los resultados. **NO SE REVISARÁ NINGUNA PRÁCTICA EN LA QUE NO SE ADJUNTE EL INFORME EN PDF DEL ANÁLISIS DE LA EXPERIMENTACIÓN.**