

HASH TABLES

1. INTRODUCTION

2. DICTIONARY DATA TYPE

Definition 2.1. Let \mathcal{X} be a set of *items* and let \mathcal{U} be a set of *keys*. Consider an abstract data type D which dynamically stores a collection of pairs (k, x) where $k \in \mathcal{U}$ and $x \in \mathcal{X}$ in such a way that D does not store two pairs having the same key at the same time. Moreover, we assume that D supports the following operations.

INSERT($D, (k, x)$)

Adds pair (k, x) into D if there is no other pair stored in D with k as a first entry.

DELETE(D, k)

Removes a pair with k as a first entry from D if such pair is stored in D .

SEARCH(D, k)

Returns x if a pair (k, x) is stored in D . Otherwise returns *nil*.

An abstract data type with these properties and interface is called *an associative array* or *a dictionary*.

Definition 2.2. Let \mathcal{X} and \mathcal{U} be sets. *Dictionary problem* for \mathcal{X} and \mathcal{U} is the task of designing a dictionary with \mathcal{X} as the set of items and \mathcal{U} as the set of keys.

3. HASH FUNCTIONS

In this section we introduce the important notion of a hash function and we discuss some probabilistic properties of it.

Definition 3.1. Let \mathcal{U} be a set. A *hash function* is a mapping $h : \mathcal{U} \rightarrow \{0, 1, \dots, m-1\}$ where $m \in \mathbb{N}_+$. Given a hash function h a *collision* is a pair of keys $k_1, k_2 \in \mathcal{U}$ such that $h(k_1) = h(k_2)$.

Definition 3.2. Let X be a set and let $n \in \mathbb{N}_+$. Then a set

$$X^{\wedge n} = \{(x_1, \dots, x_n) \in X^n \mid \forall 1 \leq i < j \leq n, x_i \neq x_j\}$$

is called *antisymmetric cartesian power* of X .

Definition 3.3. Let \mathcal{U} be a measurable space. We consider $\mathcal{U}^{\wedge n}$ as the measurable subspace of a product space \mathcal{U}^n . Suppose that P is a probability distribution on $\mathcal{U}^{\wedge n}$. Let $h : \mathcal{U} \rightarrow \{0, 1, \dots, m-1\}$ be a hash function for some $m \in \mathbb{N}_+$. Suppose that

$$P(h(k_i) = h(k_j) \mid (k_1, \dots, k_n) \in \mathcal{U}^{\wedge n}) = \frac{1}{m}$$

for every pair of distinct elements $i, j \in \{1, \dots, n\}$. Then h is a *simple uniform hashing with respect to P* .

Example 3.4. Let $\mathcal{U} = [0, m]$ for some $m \in \mathbb{N}_+$. Then \mathcal{U} is a measurable space with respect to Borel algebra $\mathcal{B}([0, m])$. We define a hash function $h : \mathcal{U} \rightarrow \{0, 1, \dots, m-1\}$ by formula

$$h(x) = \lfloor x \rfloor$$

Then h is a simple uniform hashing with respect to the normalization of n -dimensional Lebesgue measure on $[0, m]^{\wedge n}$.

Definition 3.5. Let \mathcal{U} be a measurable space. We consider $\mathcal{U}^{\wedge n}$ as the measurable subspace of a product space \mathcal{U}^n . Suppose that P is a probability distribution on $\mathcal{U}^{\wedge n}$. Let $h : \mathcal{U} \rightarrow \{0, 1, \dots, m-1\}$ be a hash function for some $m \in \mathbb{N}_+$. Fix a real number $\epsilon > 0$ and suppose that

$$P(h(k_i) = h(k_j) \mid (k_1, \dots, k_n) \in \mathcal{U}^{\wedge n}) = \frac{1}{m} + \epsilon$$

for every pair of distinct elements $i, j \in \{1, \dots, n\}$. Then h is a simple ϵ -uniform hashing with respect to P .

Example 3.6. Let $\mathcal{U} = \{0, 1, \dots, m^2 - 1\}$ for some $m \in \mathbb{N}_+$. Then \mathcal{U} is a measurable space with respect to the power algebra $\mathcal{P}(\{0, 1, \dots, m^2 - 1\})$. Consider $\mathcal{U}^{\wedge n}$ as a probability space with respect to distribution describing random sampling of n -elements without replacement from \mathcal{U} . We define a hash function $h : \mathcal{U} \rightarrow \{0, 1, \dots, m-1\}$ by formula

$$h(x) = x \bmod m$$

where m is a divisor of N . Fix distinct $i, j \in \{1, \dots, n\}$ and note that

$$P(h(k_i) = h(k_j) \mid (k_1, \dots, k_n) \in \mathcal{U}^{\wedge n}) = \frac{m}{m^2} = \frac{1}{m}$$

4. HASH TABLES WITH CHAINING AS A SOLUTION TO DICTIONARY PROBLEM

In this section we present the solution to the dictionary problem and discuss its efficiency.

Definition 4.1. Let \mathcal{U} and \mathcal{X} be sets. Let $h : \mathcal{U} \rightarrow \{0, 1, \dots, m-1\}$ be a hash function for some $m \in \mathbb{N}_+$. We consider an m -element array D_h such that $D_h[i]$ is a linked list storing values from $\mathcal{U} \times \mathcal{X}$ for every $i \in \{0, 1, \dots, m-1\}$. We describe dictionary operations.

INSERT($D_h, (k, x)$)

Inserts pair (k, x) to the linked list $D_h[h(k)]$ as its new head.

DELETE(D_h, k)

Deletes a pair with first entry k from the linked list $D_h[h(k)]$.

SEARCH(D_h, k)

Searches for the pair with the first entry k in the list $D_h[h(k)]$. If such pair is found, then returns its second entry. Otherwise returns *nil*.

Then D_h together with these operations is a solution of dictionary problem for \mathcal{U} and \mathcal{X} . We call it the hash table with collisions resolved by chaining.

For the sequel we need the notion of sampling without replacement.

Definition 4.2. Let (Ω, \mathcal{F}, P) be a probability space and let \mathcal{U} be a measurable space. We fix a random variable $\mathcal{K} : \Omega \rightarrow \mathcal{U}$ and consider independent and identically distributed random variables $\mathcal{K}_1, \dots, \mathcal{K}_n : \Omega \rightarrow \mathcal{U}$ for $n \in \mathbb{N}_+$ with distributions identical to \mathcal{K} . Define

$$\Delta = \{(k_1, \dots, k_n) \in \mathcal{U}^n \mid \forall_{i \neq j} k_i \neq k_j\}$$

Then $\langle \mathcal{K}_1, \dots, \mathcal{K}_n \rangle : \Omega \rightarrow \mathcal{U}^n$ is a random variable which gives rise to a probability distribution on a measurable space \mathcal{U}^n . Hence it also gives rise to a probability distribution $\mu_{\mathcal{K}, n}$ on its subspace $\mathcal{U}^n \setminus \Delta$. Then $\mu_{\mathcal{K}, n}$ is the distribution of sampling without replacement with respect to \mathcal{K} .

Definition 4.3. Let $n \in \mathbb{N}_+$ be the number of elements stored in D . Then $\alpha = \frac{n}{m}$ is called load factor.

Theorem 4.4. Let \mathcal{U} be a measurable space and let \mathcal{X} be a set. Let $h : \mathcal{U} \rightarrow \{0, 1, \dots, m-1\}$ be a hash function for some $m \in \mathbb{N}_+$. Suppose that \mathcal{K} is a random variable with target \mathcal{U} such that h is a simple uniform hashing with respect to \mathcal{K} . For a sequence of keys chosen randomly according to distribution $\mu_{\mathcal{K}, n}$ and stored in the hash table D_h the expected time of an unsuccessful search is $\Theta(\frac{n}{m})$.

Proof.

□