

## 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 such functions.

**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  $m \in \mathbb{N}_+$ .

**Definition 3.2.** Let  $h : \mathcal{U} \rightarrow \{0, 1, \dots, m-1\}$  be a hash function. A *collision* is a pair of keys  $k_1, k_2 \in \mathcal{U}$  such that  $h(k_1) = h(k_2)$ .

**Definition 3.3.** 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 *the antisymmetric cartesian power* of  $X$ .

**Definition 3.4.** Let  $\mathcal{U}$  be a measurable space. We consider  $\mathcal{U}^{\wedge n}$  as the measurable subspace of the 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 measurable hash function for some  $m \in \mathbb{N}_+$ . Assume that the following assertions hold.

(1)

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

for every element  $i \in \{1, \dots, n\}$  and every  $l \in \{0, 1, \dots, m-1\}$ .

(2)

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

for every pair of distinct elements  $i, j \in \{1, \dots, n\}$ .

Then  $h$  is a *simple uniform hashing with respect to  $P$* .

**Example 3.5.** 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  $[0, m]^n$ .

**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  $\mathcal{U}^n$  as a probability space with respect to the uniform distribution  $P$ . We define a hash function  $h : \mathcal{U} \rightarrow \{0, 1, \dots, m-1\}$  by formula

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

For  $i \in \{1, \dots, n\}$  and  $l \in \{0, 1, \dots, m-1\}$  we have

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

Fix distinct  $i, j \in \{1, \dots, n\}$  and  $l \in \{0, 1, \dots, m-1\}$ . Note that

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

Hence

$$\begin{aligned} P((k_1, \dots, k_n) \in \mathcal{U}^n \mid h(k_i) = h(k_j)) &= \sum_{l=0}^{m-1} P(h(k_i) = l, h(k_j) = l \mid (k_1, \dots, k_n) \in \mathcal{U}^n) = \\ &= \frac{m}{m \cdot (m+1)} = \frac{1}{m+1} \leq \frac{1}{m} \end{aligned}$$

Thus  $h$  is a simple uniform hashing with respect to  $P$ .

#### 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  $D_h[l]$  is a linked list storing values from  $\mathcal{U} \times \mathcal{X}$  for every  $l \in \{0, 1, \dots, m-1\}$ . 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  $h$* .

Suppose that  $\mathcal{U}$  and  $\mathcal{X}$  are sets. Let  $h : \mathcal{U} \rightarrow \{0, 1, \dots, m-1\}$  be a hash function. Consider the hash table  $D_h$ . Fix  $l \in \{0, 1, \dots, m-1\}$  and  $n \in \mathbb{N}_+$ . Suppose that pairs  $(k_1, x_1), \dots, (k_n, x_n)$  for  $(k_1, \dots, k_n) \in \mathcal{U}^n$  and  $x_1, \dots, x_n \in \mathcal{X}$  are consecutively inserted to initially empty  $D_h$ . After these sequence of insertions is performed the length of the linked list stored in  $D_h[l]$  is equal to the cardinality of the set  $\{i \mid h(k_i) = l\}$ . We denote the function

$$\mathcal{U}^n \ni (k_1, \dots, k_n) \rightarrow |\{i \mid h(k_i) = l\}| \in \mathbb{N}$$

by  $coll_l$ .

**Theorem 4.2.** 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 measurable hash function and fix  $n \in \mathbb{N}_+$ . Then the following assertions hold.

(1) The function  $\text{coll}_l : \mathcal{U}^n \rightarrow \mathbb{N}$  is measurable for every  $l \in \{0, 1, \dots, m-1\}$ .

(2) If  $h$  is a simple uniform hashing with respect to some probability distribution  $P$  on  $\mathcal{U}^n$ , then

$$\mathbb{E} \text{coll}_l = \int_{\mathcal{U}^n} \text{coll}_l dP = \frac{n}{m}$$

for every  $l \in \{0, 1, \dots, m-1\}$ .

*Proof.* Suppose that  $X_i$  is the indicator function of the measurable set  $\{(k_1, \dots, k_n) \in \mathcal{U}^n \mid h(k_i) = l\}$ . Then

$$\text{coll}_l = \sum_{i=1}^n X_i$$

This proves that  $\text{slot}_l$  is measurable. If in addition  $h$  is a simple uniform hashing with respect to some probability distribution  $P$  on  $\mathcal{U}^n$ , then

$$\mathbb{E} \text{coll}_l = \mathbb{E} \left( \sum_{i=1}^n X_i \right) = \sum_{i=1}^n \mathbb{E} X_i = \sum_{i=1}^n P((k_1, \dots, k_n) \in \mathcal{U}^n \mid h(k_i) = l) = \frac{n}{m}$$

□

Suppose that  $\mathcal{U}$  and  $\mathcal{X}$  are sets. Let  $h : \mathcal{U} \rightarrow \{0, 1, \dots, m-1\}$  be a hash function. Consider the hash table  $D_h$ . Fix  $n \in \mathbb{N}_+$  and  $i \in \{1, \dots, n\}$ . Suppose that pairs  $(k_1, x_1), \dots, (k_n, x_n)$  for  $(k_1, \dots, k_n) \in \mathcal{U}^n$  and  $x_1, \dots, x_n \in \mathcal{X}$  are consecutively inserted to initially empty  $D_h$ . After these sequence of insertions is performed the number of elements in the linked list stored in  $D_h[k(k_i)]$  which precede  $(k_i, x_i)$  is equal to the cardinality of the set  $\{j \in \{i+1, \dots, n\} \mid h(k_j) = h(k_i)\}$ . We denote the function

$$\mathcal{U}^n \ni (k_1, \dots, k_n) \rightarrow |\{i \mid h(k_i) = l\}| \in \mathbb{N}$$

by  $\text{coll}_{<i}$ .

**Theorem 4.3.** 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 measurable hash function and fix  $n \in \mathbb{N}_+$ . Fix  $i \in \{1, \dots, n\}$ . Then the following assertions hold.

(1) The function

$$\mathcal{U}^n \ni (k_1, \dots, k_n) \mapsto \#\{j \mid i \leq j \leq n \text{ and } h(k_i) = h(k_j)\} \in \mathbb{N}$$

is measurable.

(2) If  $h$  is a simple uniform hashing with respect to some probability distribution  $P$  on  $\mathcal{U}^n$ , then

$$\mathbb{E} \#\{j \mid i \leq j \leq n \text{ and } h(k_i) = h(k_j)\} = \int_{\mathcal{U}^n} \#D_h[l] dP = \frac{n}{m}$$

for every  $l \in \{0, 1, \dots, m-1\}$ .

*Proof.* Suppose that  $X_i$  is the indicator function of the measurable set  $\{(k_1, \dots, k_n) \in \mathcal{U}^n \mid h(k_i) = l\}$ . Then

$$\#D_h[l] = \sum_{i=1}^n X_i$$

This proves that  $\#D_h[l]$  is measurable. If in addition  $h$  is a simple uniform hashing with respect to some probability distribution  $P$  on  $\mathcal{U}^n$ , then

$$\mathbb{E} \#D_h[l] = \mathbb{E} \left( \sum_{i=1}^n X_i \right) = \sum_{i=1}^n \mathbb{E} X_i = \sum_{i=1}^n P((k_1, \dots, k_n) \in \mathcal{U}^n \mid h(k_i) = l) = \frac{n}{m}$$

□