

Weekplan: Hashing

Philip Bille

References and Reading

[1] Notes from Aarhus, Peter Bro Miltersen.

[2] Scribe notes from MIT.

[3] Universal Classes of Hash Functions, J. Carter and M. Wegman, J. Comp. Sys. Sci., 1977.

[4] Storing a Sparse Table with $O(1)$ Worst Case Access Time, M. Fredman, J. Komlos and E. Szemerédi, J. ACM., 1984.

[5] Notes on Discrete Probability, Jeff Erickson.

We recommend reading [1] and [2] in detail. [3] and [4] provide background on universal and perfect hashing. [5] provides a concise refresh of basic discrete probability.

Exercises

1 [w] Streaming Statistics An IT-security friend of yours wants a high-speed algorithm to count the number of *distinct* incoming IP-addresses in his router to help detect denial of service attacks. Can you help him?

2 [w] Dense Set Hashing A set $S \subseteq U = \{0, \dots, u-1\}$ is called *dense* if $|S| = \Theta(u)$. Suggest a simple and efficient dictionary data structure for dense sets.

3 [w] Multi-Set Hashing A multi-set is a set M , where each element may occur multiple times. Design an efficient data structure supporting the following operations:

- $\text{add}(x)$: Add an(other) occurrence of x to M .
- $\text{remove}(x)$: Remove an occurrence of x from M . If x does not occur in M do nothing.
- $\text{report}(x)$: Return the number of occurrences of x .

4 Properties of Universal Hashing Let $h \in H$ be a hash function from a universal family mapping $U = \{0, \dots, u-1\}$ to $M = \{0, \dots, m-1\}$. Solve the following exercises.

4.1 If h has no collision on U , how large must m be?

4.2 Suppose $m \geq u$. Is the identity function $f(x) = x$ a universal hash function?

4.3 A family G of hash functions mapping U to M is *family of pair-wise independent hash function* if for any $g \in G$,

$$\Pr(g(x) = \alpha \wedge g(y) = \beta) = 1/m^2 \quad \forall x \neq y \in U, \quad \forall \alpha, \beta \in M.$$

Show that any family of pairwise independent hash functions is a family of universal hash functions.

5 Linear Space Hashing The chained hashing solution for the dynamic dictionary problem presented assume that $m = \Theta(n)$. Solve the following exercises.

5.1 What is the space complexity of chained hashing without this assumption?

5.2 Give a solution that achieves $O(n)$ space and the same time complexities without assuming $m = \Theta(n)$. *Hint:* Think dynamic arrays.

6 Graph Adjacency Let G be a graph with n vertices and m edges. We want to represent G efficiently and support the following operation.

- $\text{adjacent}(v, w)$: Return true if nodes v and w are adjacent and false otherwise.

Solve the following exercises:

6.1 Analyse the space and query time in terms of n and m for the classic adjacency matrix and adjacency list representation.

6.2 Design a data structure that improves both the adjacency matrix and adjacency list.

7 Perfect Hashing Analysis Consider the 2-level FKS perfect hashing scheme. A friend suggest the following two "optimizations" to the data structure. What happens to the performance of the data structure for each of these?

7.1 Modify level 1 of the data structure to map U to an array of size $n\sqrt{n}$ instead of n to further decrease the probability of collisions.

7.2 Replace the universal hash function with a faster *near-universal hash function* on both levels. Near-universal hashing is the same as universal hashing except that $\leq 1/m$ guarantee on the probability is changed to $\leq 2/m$.

8 Lost Integer Puzzles Suppose that you receive a stream of $n - 1$ distinct integers from the set $\{1, \dots, n\}$, i.e., the stream consists of all of $\{1, \dots, n\}$ except a single missing integer. We want a space-efficient algorithm that efficiently computes this integer during a single pass over the input stream. Solve the following exercises:

8.1 Show how to find the lost integer using $O(n)$ space.

8.2 [*] Show how to find the lost integer using $O(1)$ space.

8.3 [**] Suppose there are now two lost integers. Show how to find them using $O(1)$ space.

9 Basic Probability Theory Refresh Bonus In case your knowledge of probability theory is rusty. Solve the following self-help exercises.

9.1 Prove linearity of expectation.

9.2 Prove that the expectation of the *indicator function* for $h(x) = h(y)$ (1 if $h(x) = h(y)$ and 0 otherwise) is equal to the probability that $h(x) = h(y)$.

9.3 Show that the expected number of trials to get a perfect hashing function using an array of size n^2 is ≤ 2 .