



アルゴリズムとデータ 構造

第7回探索のためのデータ構造(3) その1: ハッシュ表

ハッシュ表

■ 今日の内容:

- 抽象データ型の「辞書」
- 探索 (search) と, 挿入 (Insert), (delete) を $O(1)$ 平均時間で実現するデータ構造.
- 実装が簡単で, 実際の性能も良いため, 広く使われる

■ ポイント

- ハッシュ関数のアイデア
- 平均計算時間量が $O(1)$ のデータ構造

Q. 検索・挿入・削除の時間計算量は、 $O(\log n)$ の壁を破れないか？

平均時間計算量なら可能！

- **ハッシュ表**を用いれば、平均時間計算量を $O(1)$ （定数時間）で行える。
- ただし、ハッシュ表のサイズ m を格納要素数 n のオーダー $O(n)$ にとった場合に、 $O(1)$ を実現可能

ハッシュ表は最も実用的な辞書に適した構造

ハッシング

ハッシュ表(hash table)とは

ハッシュ関数 h を使って、要素 x があらかじめ準備した m 個の場所のうち、位置 $h(x)$ に格納される表

hash は英語で「ごた混ぜにする」という意味。



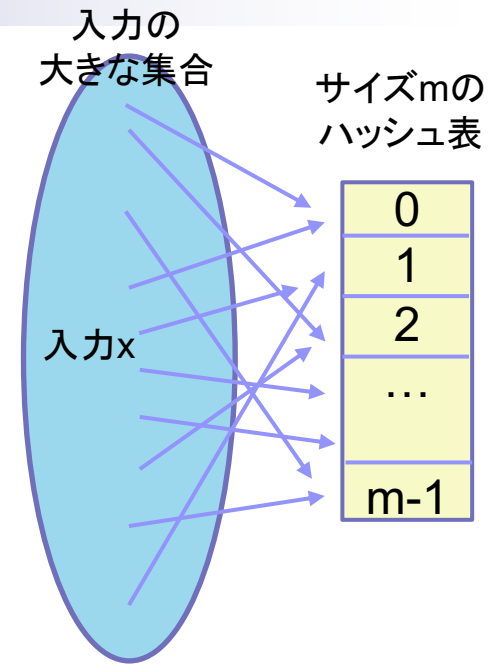
Texas hash with cornbread and green beans
(Wikipedia, Hash (food), United States)

ハッシング

ハッシュ表(hash table)とは

ハッシュ関数 h を使って、要素 x があらかじめ準備した
 m 個の場所のうち、位置 $h(x)$ に格納される表

hash は英語で「ごた混ぜにする」という意味。



[ハッシュ関数 h のもつべき性質]

1. 関数値 $h(x)$ の高速な算出が可能である
2. 要素となりうる x に対して、ハッシュ値 $h(x)$ が m 個の位置にできるだけ偏りなく分布すること

[よく用いられるハッシュ関数]

・ x が整数の場合

$$h(x) = x \% m \quad (x \text{ を } m \text{ で割った余り})$$

・ x が文字列の場合

$$h(x) = \left(\sum_{i=0}^{k-1} \text{ord}(x[i]) \right) \% m$$

ただし、 $\text{ord}(a)$ は文字 a の整数コード(ASCII, JIS, EUC等)であり、 k は x の文字列長

一般にデータはいろいろな周期のものを含んでいるため、 m がデータの周期を約数にもつと関数値の衝突が起こりやすくなる。そのため m としては素数を選ぶことが多い。

しかし、どのようなハッシュ関数を選んでも衝突は避けられない！

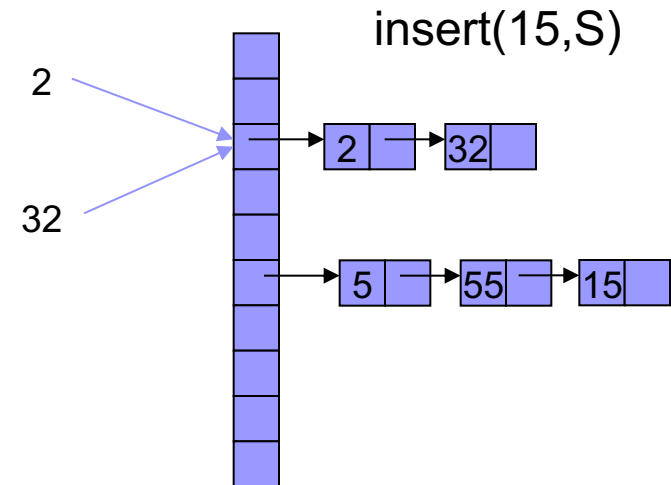
異なる要素 x, y に対し、ハッシュ値が等しくなる($h(x)=h(y)$)ことがある。
 そのような場合の対処法として次の2つがある。

1. 外部ハッシュ法、チェイニング (open hashing, chaining)

同じハッシュ値をもつ要素を、
 その値に対応するバケットに格納する。
 バケットは連結リストで実現できる。

32

$h(x)=x\%m$ の場合の

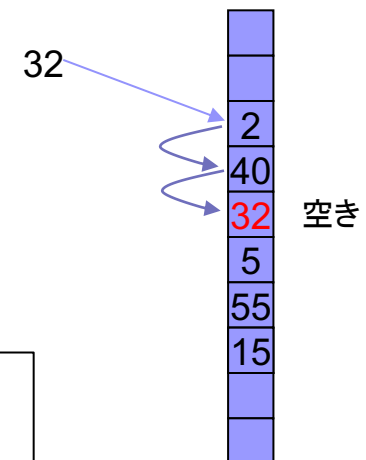


2. 内部ハッシュ法、オープンアドレッシング (closed hashing, open addressing)

x を格納しようとしたときに、位置 $h(x)$ がすでに
 使われている場合、新しいハッシュ値 $h_i(x)$ ($i=1, 2, \dots$)を
 次々と求め、最初に見つかった空いている位置
 $h_i(x)$ に格納する。 h_i としては、以下の関数などが使われる。

$$h_i(x) = (h(x) + i) \% m$$

コメント: この関数は、位置 $h(x)$ が使われ
 ていたら、添字が一つ大きな隣を順に見て
 いき、最初の空き位置に入れることと同じ。



外部ハッシュ法の基本操作

member(x,S) $h(x)$ に対応するバケット内で値がxのものを探し、見つければyes、見つからなければnoを返す。

Insert(x,S) $h(x)$ に対応するバケット内で値がxのものを探し、見つければ何もしない、見つからなければバケットに追加。

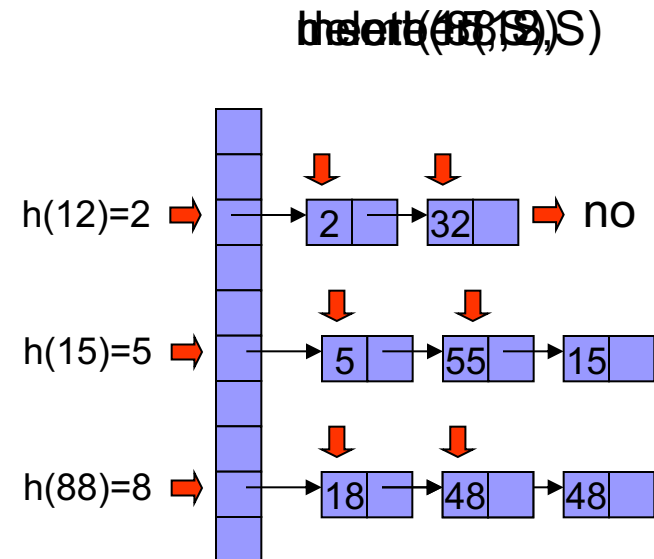
delete(x,S) $h(x)$ に対応するバケット内で値がxのものを探し、見つければ削除、見つからなければ何もしない。

基本操作の時間計算量

$\alpha = n/m$: 占有率 (n: 格納要素数, m: バケット数)

最悪時間計算量 $O(n)$
(全ての要素が同じハッシュ値をもつ場合)

平均時間計算量 $O(\alpha)$
($n = O(m)$ であれば $O(1)$)



member(x,S) 位置 $h(x)$ から探し始め、 $h(x), h_1(x), h_2(x), \dots$ の順でその位置の要素が x と等しいかをチェックする。等しいものがみつかったらyesを返す。位置 $h_i(x)$ が空("empty")となるまでチェックし、見つからなければnoを返す。

Insert(x,S) 位置 $h(x)$ から探し始め、 $h(x), h_1(x), h_2(x), \dots$ の順でその位置の要素が x と等しいかをチェックする。位置 $h_i(x)$ が空("empty")となるまでチェックし、 x と等しい要素が見つければ何もしない。みつからなければ、検索途中で見つけた空("empty" or "deleted")に x を格納する。

delete(x,S) 位置 $h(x)$ から探し始め、 $h(x), h_1(x), h_2(x), \dots$ の順でその位置の要素が x と等しいかをチェックする。位置 $h_i(x)$ が空("empty")となるまでチェックし、 x と等しい要素が見つければ削除し"deleted"のフラグを立てる。みつからなければ何もしない。

基本操作の時間計算量

$\alpha = n/m$: 占有率 (n : 格納要素数, m : バケット数)

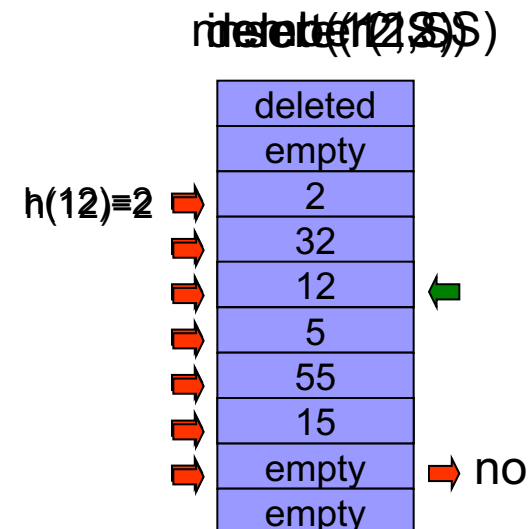
最悪時間計算量 $O(n)$

(全ての要素が同じハッシュ値をもつ場合)

平均時間計算量 $O(1/(1-\alpha))$ (x が表にない場合)

$O(-(1/\alpha)\log(1-\alpha))$ (x が表にある場合)

どちらの場合も $n=O(m)$ であれば、 $O(1)$



内部ハッシュ法の平均時間計算量の証明

・xが表にない場合

$h_0(x)(=h(x)), h_1(x), \dots$ の順にハッシュ値の計算を行うものとする。

位置 $h_i(x)$ で初めて空を見つける確率は

$$\frac{n(n-1)\cdots(n-i+1)(m-n)}{m(m-1)\cdots(m-i+1)(m-i)}$$

である。(ただし、 $i=0$ のとき、上式は $(m-n)/m$ をあらわすものとする。)

したがって、要素の比較回数の期待値は、

$$\begin{aligned} \sum_{i=0}^n (i+1) \frac{n(n-1)\cdots(n-i+1)(m-n)}{m(m-1)\cdots(m-i+1)(m-i)} &= \sum_{i=0}^n (i+1) \frac{n(n-1)\cdots(n-i+1)}{m(m-1)\cdots(m-i+1)} \left(1 - \frac{n-i}{m-i}\right) \\ &= \sum_{i=0}^n (i+1) \frac{n(n-1)\cdots(n-i+1)}{m(m-1)\cdots(m-i+1)} - \sum_{i=1}^{n+1} i \frac{n(n-1)\cdots(n-i+1)}{m(m-1)\cdots(m-i+1)} \\ &\leq \sum_{i=0}^n \frac{n(n-1)\cdots(n-i+1)}{m(m-1)\cdots(m-i+1)} \\ &\leq \sum_{i=0}^{\infty} \left(\frac{n}{m}\right)^i = \frac{1}{1 - \frac{n}{m}} = \frac{1}{1 - \alpha} \end{aligned}$$

比較回数の定数倍の時間で計算できるので平均時間計算量は $O(1/(1-\alpha))$

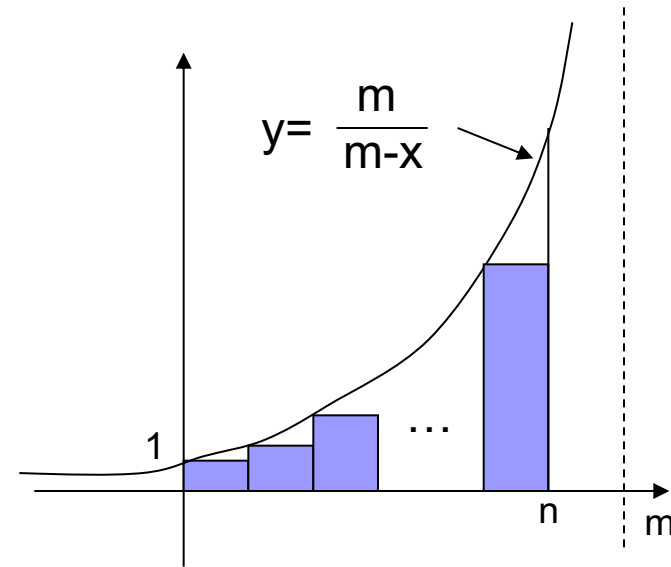
内部ハッシュ法の平均時間計算量の証明(続き)

・xが表にある場合

比較回数の期待値は、n個の異なる要素を空の表に格納するのに必要な比較回数の平均に等しい。したがって

$$\begin{aligned} \frac{1}{n} \sum_{i=0}^{n-1} \frac{m}{m-i} &\leq \frac{1}{n} \int_0^n \frac{m}{m-x} dx \\ &= \frac{m}{n} \ln \frac{m}{m-n} = -\frac{1}{\alpha} \ln(1-\alpha) \end{aligned}$$

比較回数の定数倍で計算できるので、平均時間計算量は $O\left(\frac{1}{\alpha} \log(1-\alpha)\right)$ (証明終わり)



確率pで生起する事象がk回目で初めて生起するとした場合、kの期待値は1/pであるという事実を使っている。つまり、確率(m-i)/mで生起する事象では、kの期待値はm/(m-i)である。

ハッシュ表

■ 今日の内容:

- 抽象データ型の「辞書」
- ハッシュ関数を用い, n 個の入力要素を長さ m の配列に格納するデータ構造
- ($m=O(n)$ のとき), 探索(search)と, 挿入(Insert), (delete)を $O(1)$ 平均時間で実現するデータ構造.
- 最悪時間は $O(n)$.
- 実装が簡単で, 実際の性能も良いため, 広く使われる

■ ポイント

- ハッシュ関数のもつべき性質と作り方
- 平均計算時間量が $O(1)$ のデータ構造