

روش درهم‌سازی

از فصل ۱۱ کتاب CLRS

مسئله‌ی جدول نمادها (Symbol Table Problem)

- جدول نمادهای T که n رکورد را ذخیره می‌کند.
- هر رکورد علاوه بر مولفه‌ی کلید (key) مولفه‌های دیگر هم دارد.
- اعمال: درج، حذف و جست‌وجو (اعمال فرهنگ داده‌ای)
- داده چه‌گونه باید ذخیره شود؟

جدول آدرس دهی مستقیم (Direct Access Table)

- اگر کلیدها متمایز و $K \subseteq \{0, 1, \dots, m-1\}$

$$\Leftarrow \text{آرایه‌ی } T[0..m-1]$$

$$T[k] = \begin{cases} x & \text{if } k \in K \text{ and } \text{key}[x] = k \\ \text{null} & \text{otherwise} \end{cases}$$

- اعمال در $O(1)$

جدول آدرس دهی مستقیم (اعمال)

$$\text{DIRECT-ADDRESS-SEARCH}(T, k)$$

$$1 \text{ return } T[k]$$

$$\text{DIRECT-ADDRESS-INSERT}(T, x)$$

$$1 \text{ } T[\text{key}(x)] \leftarrow x$$

$$\text{DIRECT-ADDRESS-DELETE}(T, x)$$

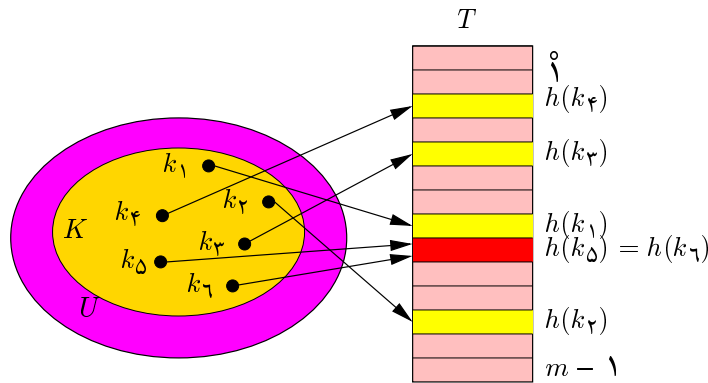
$$1 \text{ } T[\text{key}(x)] \leftarrow \text{null}$$

جدول آدرس دهی مستقیم (مشکل)

• m می‌تواند بسیار بزرگ باشد

- کلید ۶۴ بیتی می‌تواند $۱۸,۴۴۶,۷۴۴,۰۷۳,۷۰۹,۵۵۱,۶۱۶$ مقدار متفاوت داشته باشد!
- کلید اگر رشته باشد، تعداد حالت‌ها حتی بیش‌تر از این می‌تواند باشد.

جدول‌های درهم‌سازی (hashing tables)



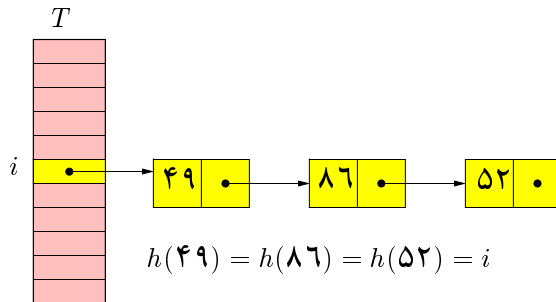
عنصر x با کلید $k = key(x)$ در درایه‌ی $h(k)$ ذخیره می‌شود. \Leftarrow تصادم (collision)

تابع درهم‌سازی (hashing function)

$$h : U \rightarrow \{0, 1, \dots, m - 1\}$$

تعریف: اگر دو عنصر توسط تابع درهم‌سازی به یک درایه از جدول متناظر شوند، می‌گوییم یک تصادم (collision) روی داده است.

روش زنجیره‌ای (chaining) برای حل تصادم



عناصری که به یک درایه نگاشت می‌شوند ($h(k)$ برابر) در لیست پیوندی $T[h(k)]$ لیست تهی (null)

روش زنجیره‌ای (اعمال)

CHAINED-HASH-SEARCH(T, k)1 search an element with key k in list $T[h(k)]$ CHAINED-HASH-INSERT(T, x)1 insert x at the head of list $T[h(key[x])]$ CHAINED-HASH-DELETE(T, x)1 delete x from the list $T[h(key[x])]$

روش زنجیره‌ای (تحلیل)

• فرض: جدول درهم‌سازی (T) به اندازه‌ی m عنصر را ذخیره می‌کند.• تعریف: α ضریب بارگذاری (load factor) T ، یعنی $\alpha = \frac{n}{m}$ • α متوسط تعداد عناصری است که در یک زنجیره قرار می‌گیرند.

تحلیل روش زنجیره‌ای (بدترین حالت)

• بدترین حالت: n عنصر در یک درایه قرار گیرند.

• بدترین زمان برای جست‌وجو:

 $O(n)$ به علاوه زمان لازم برای محاسبه تابع درهم‌سازی.

تحلیل روش زنجیره‌ای (در حالت متوسط)

• فرض: تابع درهم‌سازی ساده و یکنواخت است (simple uniform hashing).

• یعنی هر عنصر با احتمال مساوی در هر یک از درایه‌های جدول قرار می‌گیرد.

• n_j : اندازه‌ی لیست T_j ($j = 0, \dots, m-1$)• داریم: $n = n_0 + n_1 + \dots + n_{m-1}$ • مقدار متوسط n_j برابر است با $\alpha = \frac{n}{m}$ • فرض: مقدار $h(k)$ را در $O(1)$ حساب می‌کنیم.• زمان مورد نیاز برای جست‌وجوی یک عنصر با کلید k به صورت خطی به $h(k)$

تحلیل متوسط روش زنجیره‌ای (جست و جوی)

دو حالت برای یک عمل جست و جوی: جست و جوی موفق یا ناموفق ندارد.

قضیه:

در یک جدول درهم‌سازی که تصادم‌ها به روش زنجیره‌ای حل شده‌اند، با فرض استفاده از تابع درهم‌سازی ساده‌ی یک‌نواخت، یک جست و جوی ناموفق به زمان متوسط $O(1 + \alpha)$ نیاز دارد.

تحلیل متوسط روش زنجیره‌ای (جست و جوی ناموفق)

اثبات:

- طبق فرض هر کلید k که هنوز در جدول ذخیره نشده باشد، با احتمال یکسان می‌تواند در هر یک از درایه‌های جدول قرار گیرد.
- پس زمان لازم برای انجام یک جست و جوی ناموفق برابر زمان لازم برای جست و جوی تا انتهای لیست $T[h(k)]$ می‌باشد، که دارای طول متوسط α است.
- بنابراین زمان جست و جوی ناموفق برابر $O(1 + \alpha)$ است. (زمان محاسبه $h(k)$ ، $O(1)$ فرض شده است).

تحلیل متوسط روش زنجیره‌ای (جست و جوی موفق)

قضیه:

در یک جدول درهم‌سازی که تصادم‌ها به روش زنجیره‌ای حل شده‌اند، با فرض استفاده از تابع درهم‌سازی ساده‌ی یک‌نواخت، یک جست و جوی موفق زمان متوسط $O(1 + \alpha)$ نیاز دارد.

ادامه‌ی تحلیل متوسط روش زنجیره‌ای (جست و جوی موفق)

اثبات:

- فرض: x_i امین عنصری باشد که در جدول درج می‌شود
 $k_i = \text{key}[x_i]$ ، و $(i = 1, 2, \dots, n)$
- متغیر تصادفی X_{ij} : $X_{ij} = I(h(k_i) = h(k_j))$
- داریم: $pr(h(k_i) = h(k_j)) = \frac{1}{m}$
- بنابراین $E[X_{ij}] = \frac{1}{m}$

تعداد متوسط عناصر مورد بررسی برای جست‌وجوی موفق

$$\begin{aligned}
 E\left[\frac{1}{n} \sum_{i=1}^n \left(1 + \sum_{j=i+1}^n X_{ij}\right)\right] &= \\
 &= \frac{1}{n} \sum_{i=1}^n \left(1 + \sum_{j=i+1}^n E[X_{ij}]\right) \\
 &= \frac{1}{n} \sum_{i=1}^n \left(1 + \sum_{j=i+1}^n \frac{1}{m}\right) \\
 &= 1 + \frac{1}{nm} \sum_{i=1}^n (n-i) \\
 &= 1 + \frac{1}{nm} \left(n^2 - \frac{n(n+1)}{2}\right) \\
 &= 1 + \frac{n-1}{2m} = 1 + \frac{\alpha}{2} - \frac{\alpha}{2n}
 \end{aligned}$$

بنابراین

زمان جست‌وجوی موفق (با احتساب زمان لازم برای محاسبه تابع درهم‌سازی) برابر است با

$$\Theta(2 + \alpha/2 - \alpha/2n) = \Theta(1 + \alpha)$$

توابع درهم‌سازی

- یک تابع درهم‌سازی خوب باید یکنواخت و ساده simple uniform hashing
- بسیاری از توابع درهم‌سازی فرض می‌کنند کلید عناصر اعداد طبیعی هستند،
- بنابراین اگر کلید عناصری که می‌خواهیم ذخیره کنیم طبیعی نباشد باید روشی برای متناظر کردن آن‌ها با اعداد طبیعی بیابیم.
- برای مثال، اگر کلید عناصر رشته‌هایی از نویسه‌ها باشد، می‌توانیم کلید را مجموع کد اسکی نویسه‌های آن رشته در نظر بگیریم.

توابع درهم‌سازی (روش تقسیم)

- $h(k) = k \bmod m$.
- باید از انتخاب مقادیر خاصی برای m اجتناب کنیم.
- مثلاً m نباید توانی از ۲ باشد، زیرا اگر $m = 2^p$ معادل p بیت کم‌ارزش k است
- به‌تر است از مقادیر همگی بیت‌ها برای محاسبه‌ی $h(k)$ استفاده شود.
- به‌طور کلی یک عدد اول که نزدیک به توانی از ۲ نباشد انتخاب مناسبی برای m است.

توابع درهم‌سازی (روش تقسیم)

توزیع و مقدار کلیدها در خوبی یا بدی تابع درهم‌سازی نقش دارند:

مثال

• عناصر مجموعه، اعداد مجذور کامل باشند یعنی $X_i = i^2$ و $(i = 1..100)$

• $m = 7$

• می‌دانیم $xy \bmod m = (x \bmod m)(y \bmod m) \bmod m$

• $h(x) = i^2 \bmod 7 = (i \bmod 7)^2 \bmod 7$

i	$i \bmod 7$	$(i \bmod 7)^2 \bmod 7$
1	1	1
2	2	4
3	3	2
4	4	2
5	5	4
6	6	1
7	0	0

$14 * 7 + 2 = 100$

• تصادم در درایه 0 = 14

$14 * 2 + 1 = 29$ تصادم در درایه 1

$14 * 2 = 28$ تصادم در درایه 2

$14 * 2 + 1 = 29$ تصادم در درایه 4

توابع درهم‌سازی (روش ضرب)

• فرض: کلیدها اعداد صحیح، $m = 2^r$ و کلمه‌ی کامپیوتر w بیتی است.

$h(k) = (A.k \bmod 2^w) \text{ rsh } (w - r)$

• rsh یعنی شیفت راست بیتی

• A یک عدد صحیح فرد است به طوری که $2^{w-1} < A < 2^w$.

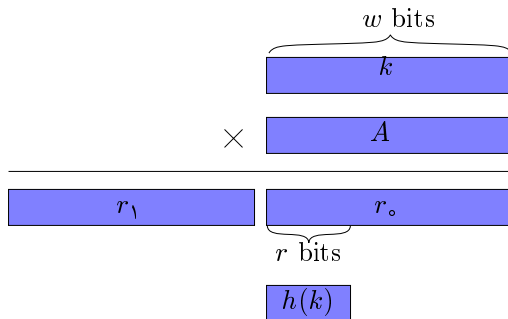
• توصیه‌ها و نکات:

-- A نزدیک به 2^w نباشد.

-- عمل ضرب در $2^w \bmod$ سریع است.

-- عمل rsh هم سریع است.

توابع درهم‌سازی به روش ضرب (مثال)



توابع درهم‌سازی به روش ضرب (مثال)

$$\begin{array}{r}
 1011001 \quad A \\
 1101011 \quad k \\
 \hline
 10010100110011
 \end{array}$$

مشکل توابع درهم‌سازی

- برای هر تابع درهم‌سازی h که دیده‌ایم، می‌توان مجموعه‌ای از کلیدها پیدا کرد که متوسط زمان جست‌وجو بسیار زیاد شود.
- مثل یک بازی: شما تابع و من کلید
- ایده‌ی راه‌حل: درهم‌سازی سراسری (universal hashing)
- از بین یک مجموعه از توابع درهم‌سازی یکی را به صورت تصادفی و مستقل انتخاب کن
- حتی اگر حریق همه‌ی توابع را از قبل بدانند، نمی‌توانند کلیدهای بد را انتخاب کنند.
- چون نمی‌دانند کدام تابع درهم‌سازی انتخاب خواهد شد.

آدرس‌دهی باز (open hashing)

در آدرس‌دهی باز هر عنصر در یک درایه جدول ذخیره می‌شود

درج

برای مشخص کردن این‌که کدام درایه‌ها باید بررسی شوند تابع درهم‌سازی را به این صورت تعریف می‌کنیم:

$$h : U * \{0, 1, \dots, m - 1\} \rightarrow \{0, 1, \dots, m - 1\}$$

در آدرس‌دهی باز برای هر کلید k نیاز به چک کردن متوالی درایه‌های زیر داریم:

$$\langle h(k, 0), h(k, 1), \dots, h(k, m - 1) \rangle$$

HASH-INSERT(T, k)

```

1  $i \leftarrow 0$ 
2 repeat  $j \leftarrow h(k, i)$ 
3     if  $T[j] = \text{null}$ 
4         then  $T[j] \leftarrow k$ 
5         return  $j$ 
6     else  $i \leftarrow i + 1$ 
7 until  $i = m$ 
8 error "hash table overflow"

```

HASH-SEARCH(T, k)

```

1  $i \leftarrow 0$ 
2 repeat  $j \leftarrow h(k, i)$ 
3     if  $T[j] = k$ 
4         then return  $j$ 
5      $i \leftarrow i + 1$ 
6 until  $T[j] = \text{null}$  or  $i = m$ 
7 return null

```

برای انجام عمل حذف، به جای این که در درایه عنصر حذف شده null قرار دهیم مقدار خاص "deleted" را قرار می‌دهیم (بنابراین الگوریتم درج هم نیاز به اندکی اصلاح دارد به این صورت که آن درایه می‌توان عنصر جدیدی درج کرد).

HASH-DELETE(T, k)

```

1  $i \leftarrow \text{HASH-SEARCH}(T, k)$ 
2 if  $i \neq \text{null}$ 
3     then  $T[i] \leftarrow \text{"deleted"}$ 

```

HASH-INSERT(T, k)

```

1  $i \leftarrow 0$ 
2 repeat  $j \leftarrow h(k, i)$ 
3     if  $T[j] = \text{null}$  or  $T[j] = \text{"deleted"}$ 
4         then  $T[j] \leftarrow k$ 
5         return  $j$ 
6     else  $i \leftarrow i + 1$ 
7 until  $i = m$ 
8 error "hash table overflow"

```


روش‌های معمول برای آدرس‌دهی باز

- واری خطی (linear probing)
- واری درجه ۲ و (quadratic probing)
- واری دوگانه (double hashing)

واری خطی

$$h(k, i) = (h'(k) + i) \bmod m$$

برای $i = 0, 1, \dots, m - 1$

واری درجه دو

$$h(k, i) = (h'(k) + c_1 i + c_2 i^2) \bmod m$$

h' یک تابع درهم‌سازی کمکی است و

c_1, c_2 ثابت‌های کمکی هستند ($c_2 \neq 0$)

($i = 0, 1, \dots, m - 1$)

اگر دو کلید واری اولیه مشابه داشته باشند رشته واری‌های آن‌ها مثل هم است:

$$h(k_1, 0) = h(k_2, 0) \implies h(k_1, i) = h(k_2, i)$$

درهم‌سازی دوگانه

درهم‌سازی دوگانه تابع درهم‌سازی به این شکل دارد:

$$h(k, i) = (h_1(k) + i h_2(k)) \bmod m$$

که h_1, h_2 توابع درهم‌سازی کمکی هستند.

تحلیل آدرس‌دهی باز

• در آدرس‌دهی باز در هر درایه جدول حداکثر یک عنصر قرار می‌گیرد،

$$\alpha = \frac{n}{m} \leq 1 \iff n \leq m$$

• فرض: درهم‌سازی یکنواخت

• در این صورت رشته واری: $\langle h(k, 0), h(k, 1), \dots, h(k, m-1) \rangle$

• درج یا جست‌وجوی کلید k : معادل یکی از جایگشت‌های $\langle 0, 1, \dots, m-1 \rangle$

تحلیل آدرس‌دهی باز (جست‌وجوی ناموفق)

قضیه:

حداکثر تعداد واری‌ها در یک جدول درهم‌سازی باز با $\alpha < 1$ برای یک جست‌وجوی ناموفق برابر است با $\frac{1}{1-\alpha}$

اثبات

X : یک متغیر تصادفی برابر با تعداد واری‌ها در یک جست‌وجوی ناموفق

A_i : رخداد این‌که i امین واری به درایه‌ی پر برود.

رخداد $\{X \geq i\}$ یعنی: یک جست‌وجوی ناموفق که حتماً $i-1$ واری اولش پر بوده است و پس از آن یکی از واری‌هایش (i ام یا بیش‌تر) تهی بوده است. پس

$$\{X \geq i\} = A_1 \cap A_2 \cap \dots \cap A_{i-1}$$

$$\begin{aligned} \Pr\{X \geq i\} &= \Pr\{A_1 \cap A_2 \cap \dots \cap A_{i-1}\} \\ &= \Pr\{A_1\} \times \Pr\{A_2|A_1\} \times \Pr\{A_3|A_1 \cap A_2\} \dots \\ &= \frac{n}{m} \times \frac{n-1}{m-1} \times \frac{n-2}{m-2} \times \dots \times \frac{n-i+2}{m-i+2} \\ &\leq \left(\frac{n}{m}\right)^{i-1} \\ &= \alpha^{i-1} \end{aligned}$$

$$\begin{aligned}
 E[X] &= \sum_{i=1}^8 i \Pr\{X = i\} \\
 &= \sum_{i=1}^8 i (\Pr\{X \geq i\} - \Pr\{X \geq i + 1\}) \\
 &= \sum_{i=1}^8 \Pr\{X \geq i\} \\
 &\leq \sum_{i=1}^8 \alpha^{i-1} \\
 &= \sum_{i=0}^8 \alpha^i \\
 &= \frac{1}{1 - \alpha}
 \end{aligned}$$

تحلیل آدرس دهی باز (درج)

نتیجه:

درج یک عنصر در یک جدول درهم‌سازی باز در حالت متوسط حداکثر $\frac{1}{1-\alpha}$ واریسی نیا دارد.

اثبات:

یک عنصر درج می‌شود اگر حداقل یک فضای خالی وجود داشته باشد. بنابراین $\alpha < 1$ در ضمن درج یک عنصر مستلزم انجام یک جست‌وجو و قرار دادن آن در اولین درایه‌ی خالی است. پس تعداد متوسط واریسی‌ها با توجه به قضیه‌ی قبل برابر $\frac{1}{1-\alpha}$ خواهد بود.

تحلیل آدرس دهی باز (جست‌وجوی موفق)

قضیه:

مقدار متوسط تعداد واریسی‌ها در یک جدول درهم‌سازی باز با $\alpha < 1$ برای یک جست‌وجوی موفق حداکثر برابر است با

$$\frac{1}{\alpha} \ln \frac{1}{1 - \alpha}$$

اثبات:

جست‌وجو برای k همان دنباله‌ی واریسی‌ها را دارد که اگر k را بخواهیم درج کنیم. اگر $k (i + 1)$ امین کلیدی باشد که درج می‌شود، متوسط واریسی‌ها برابر است با:

$$\frac{1}{1 - i/m} = m/(m - i)$$

پس

$$\begin{aligned}
 \frac{1}{n} \sum_{i=0}^{n-1} \frac{m}{m - i} &= \frac{m}{n} \sum_{i=0}^{n-1} \frac{1}{m - i} \\
 &= \frac{1}{\alpha} (H_m - H_{m-n})
 \end{aligned}$$

که $H_i = \sum_{j=1}^i 1/j$ تابع هارمونی i ام است.

$$\begin{aligned} \frac{\lambda^{n-1}}{n} \sum_{i=0}^{m-1} \frac{m}{m-i} &= \frac{m^{n-1}}{n} \sum_{i=0}^{m-1} \frac{\lambda}{m-i} \\ \frac{\lambda}{\alpha} (H_m - H_{m-n}) &= \frac{\lambda}{\alpha} \sum_{k=m-n+1}^m \lambda/k \\ &\leq \frac{\lambda}{\alpha} \int_{m-n}^m (\lambda/x) dx \\ &= \frac{\lambda}{\alpha} \ln \frac{m}{m-n} \\ &= \frac{\lambda}{\alpha} \ln \frac{\lambda}{\lambda - \alpha} \end{aligned}$$