



جلسه‌ی ۱۷: درهم‌سازی

نگارنده: مجتبی تفاق

مدرّس: دکتر شهرام خزائی

۱ مقدمه

در بسیاری از کاربردها، برای داده‌ساختار مورد نیاز کافی است تا کارکردهای یک فرهنگ لغت عادی را انجام دهد:

- اضافه کردن یک داده‌ی جدید،
- حذف یک داده‌ی موجود،
- جستجوی یک داده.

تا کنون داده‌ساختارهای زیر را جهت ذخیره اطلاعات دیده‌ایم:

حذف	جست‌وجو	درج	داده‌ساختار
n	n	۱	آرایه‌ی نامرتب
n	$\log(n)$	n	آرایه‌ی مرتب
n	n	۱	لیست پیوندی یک‌سویه
۱	n	۱	لیست پیوندی دوسویه
$\log(n)$	$\log(n)$	$\log(n)$	هرم بیشینه (کمینه)

در این جلسه می‌خواهیم داده‌ساختاری را معرفی کنیم، که این اعمال را در زمان $\Theta(1)$ انجام دهد.

۲ آدرس‌دهی مستقیم

فرض کنید می‌خواهیم اطلاعات ۱۰۰۰۰۰ دانشجوی را ذخیره کنیم. همچنین فرض کنید شماره‌های دانشجویی، اعدادی بین ۱ تا یک میلیون هستند. برای این کار می‌توانیم از یک "جدول آدرس‌دهی مستقیم"^۱ استفاده کنیم، به این ترتیب که آرایه‌ای مانند T با یک میلیون خانه و برای هر شماره‌ی دانشجویی یک خانه از T را در نظر می‌گیریم. بنابراین یک آرایه‌ی خیلی بزرگ خواهیم داشت، که در آن از شماره‌ی دانشجویی به عنوان اندیس آرایه استفاده کرده‌ایم و در هر

^۱direct addressing table

خانه‌ی متناظر با شماره دانشجویی یک دانشجو، اطلاعات و یا اشارگری به اطلاعات آن دانشجو وجود دارد. این روش وقتی که کلیدها (در این جا شماره‌ی دانشجویی) از یک مجموعه‌ی به طور نسبی کوچک آماده باشند، روشی کارآمد است و هدف ما را برآورده می‌کند زیرا که سه عمل بالا را با الگوریتم‌هایی که در پایین می‌آیند می‌توان در زمان ثابت انجام داد:

DIRECT-ADDRESS-SEARCH (T, k)

Return $T[k]$

DIRECT-ADDRESS-INSERT (T, x)

$T[x.key] \leftarrow x$

DIRECT-ADDRESS-DELETE (T, x)

$T[x.key] \leftarrow null$

۳ چند کاربرد

حال فرض کنید که به جای اطلاعات ۱۰۰۰۰۰ دانشجو، می‌خواستیم که اطلاعات یک کلاس ۱۰۰ نفره را ذخیره کنیم. دوباره مجبور بودیم آرایه‌ای با همان اندازه‌ی قبلی در نظر بگیریم. همان‌گونه که مشاهده می‌شود، در مواقعی که تعداد داده‌های مورد نظر نسبت به اندازه‌ی مجموعه‌ای که کلیدها از آن می‌آیند کوچک باشد، روش آدرس‌دهی مستقیم کارایی خود را به علت اتلاف حافظه از دست خواهد داد. دقت کنید که کاربردهای این چنینی بسیار زیاد رخ می‌دهند. به عنوان یک مثال دیگر فرض کنید می‌خواهیم خود یک فرهنگ لغت را پیاده‌سازی کنیم. با احتساب فرم‌های مختلف کلی حدود ۶۰۰۰۰۰ لغت در زبان انگلیسی وجود دارد، اما اگر بخواهیم رشته‌های حروف به طول حداکثر سی را به عنوان کلید برای آدرس‌دهی مستقیم در نظر بگیریم، آن‌گاه به آرایه‌ای به اندازه‌ی 26^{30} نیاز داریم! برای مثال آخر، رهیابی^۲ را در یک شبکه در نظر بگیرید که به طور متوسط در هر لحظه باید ۵۰۰ آدرس را ذخیره کند. تعداد IP‌های ممکن برابر ۲^{۳۲} است که همان اندازه‌ی مجموعه‌ی کلیدهای ماست. آیا راه‌حلی برای استفاده‌ی بهتر از حافظه در این موارد وجود دارد؟

^۲router

۴ درهم‌سازی

در این بخش، با داده‌ساختاری به نام جدول درهم‌سازی^۲ و روش درهم‌سازی آشنا می‌شویم. در جدول‌های درهم‌سازی سه عمل مورد نظر در بالا بسته به شرایط مختلف، زمان‌های متفاوتی می‌برند ولی با انتخاب مناسب روش، می‌توان میانگین این زمان را در حد $\Theta(1)$ پایین آورد.

ایده‌ی کلی بر این اساس است که به جای در نظر گرفتن خود کلیدها، برای آدرس‌دهی از تابعی از کلیدها استفاده و هر داده را در اندیس متناظر با مقدار آن تابع ذخیره کنیم. حال اگر برد این تابع از مرتبه‌ی تعداد داده‌ها باشد و این تابع روی کلیدهای مورد نظر به صورت "تقریباً" یک‌به‌یک عمل کند، تمامی اطلاعات را می‌توان در آرایه‌ای به اندازه‌ی داده‌ها ذخیره کرد که مطلوب‌ترین حالت ممکن است.

فرض می‌کنیم مجموعه‌ی داده‌ها حداکثر شامل m عنصر است. کلید عنصر x را $x.k$ می‌نامیم که کلیدها متمایز و عناصری از مجموعه‌ای به نام K هستند. همچنین می‌دانیم که $K \subseteq U$ در آن U مجموعه‌ی جهانی مقادیر ممکن کلیدهاست. همان‌طور که اشاره شد، تابع درهم‌ساز از فضای کلیدهای ممکن به فضای اندیس‌های موجود است. یعنی:

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

اگر تابع h روی مجموعه‌ی K یک‌به‌یک باشد آن‌گاه داده‌ی x را در $h(x.k)$ ذخیره می‌کنیم. اما اگر دو کلید متفاوت به یک اندیس مشابه نگاشته شوند، چه باید کرد؟ به این مشکل برخورد^۴ می‌گوییم که باید آن‌را به روشی برطرف کنیم.

۵ برخورد

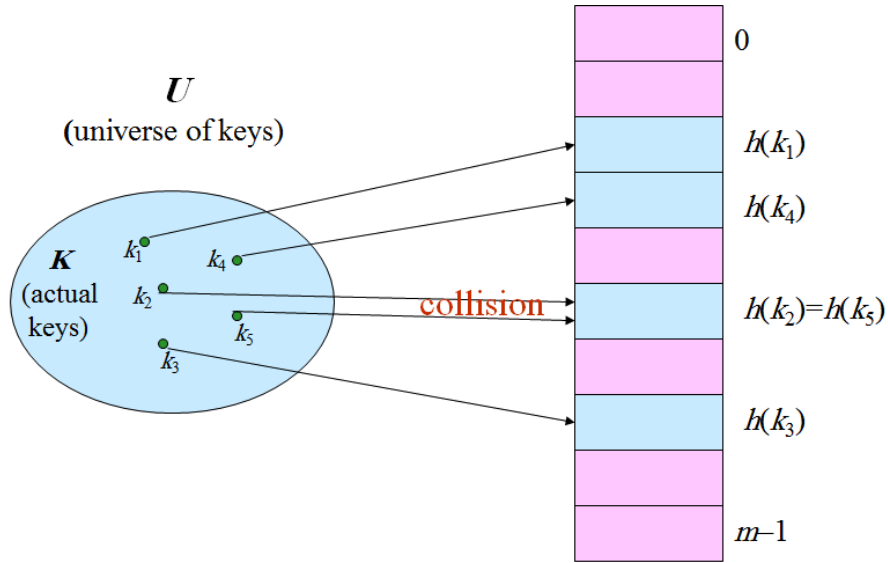
شکل ۱ یک جدول درهم‌سازی با مجموعه‌ی جهانی U ، مجموعه‌ی کلیدهای K و نیز تابع h را برای پنج کلید نشان می‌دهد. در این جا کلیدهای k_5 و k_2 برخورد دارند چون تحت h به یک درایه نگاشته می‌شوند. برای رفع مشکل برخورد به دو راه حل اشاره می‌کنیم. روش اول "روش زنجیره‌ای"^۵ و روش دوم "آدرس دهی باز"^۶ نامیده می‌شود. در ادامه به توصیف این دو روش خواهیم پرداخت.

^۲ hashing table

^۴ collosion

^۵ chaind hashing

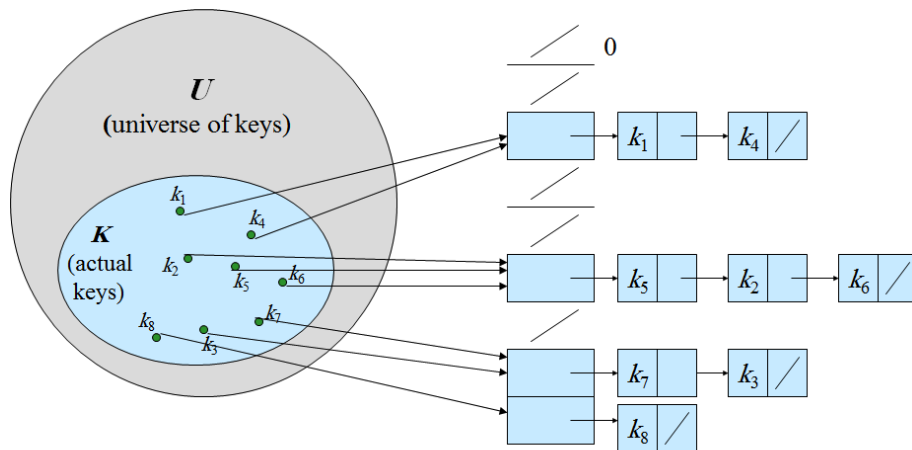
^۶ open addressing



شکل ۱: یک جدول درهم‌سازی با تابع درهم‌سازی h .

۶ روش زنجیره‌ای

یکی از راه‌حل‌های رفع مشکل برخورد، ذخیره‌ی اندیس‌هایی که برخورد کرده‌اند در یک لیست پیوندی است. درایه‌های خالی جدول هم لیست‌های تهی ($null$) هستند. به چنین جدولی "جدول درهم‌سازی با روش زنجیره‌ای" می‌گوییم. شکل ۲ مثالی از این روش است.



شکل ۲: روش زنجیره‌ای برای حل برخورد.

در زیر الگوریتم‌های متناظر با سه عمل مورد نظر در این روش آمده است:

CHAINED-HASH-INSERT (T, x)

insert x at the head of the list $T[h(x.k)]$

CHAINED-HASH-SEARCH (T, k)

search for an element with key k in list $T[h(k)]$

CHAINED-HASH-DELETE (T, x)

delete x from list $T[h(x.k)]$

اگر تابع درهم‌ساز تقریباً یک‌به‌یک باشد، تعداد برخوردها کم خواهد بود و هر سه عمل درج، حذف، و جست‌وجو در $\Theta(1)$ انجام خواهند شد. در غیر این صورت، بعضی از کلیدها یک لیست طولانی خواهند داشت. با این‌که زمان درج $\Theta(1)$ است (چون می‌توانیم به ابتدای لیست اضافه کنیم)، ولی حذف و جست‌وجو کند خواهند بود و در بدترین حالت ممکن است در $\Theta(n)$ انجام شوند که در این‌جا n طول لیست پیوندی مورد نظر است.

۷ عامل بار

میزان برخورد مستقیماً به تعداد اعداد ورودی بستگی ندارد و به “عامل بار”^۷ جدول مربوط است. عامل بار، که با نماد α نشان می‌دهیم، برابر است با نسبت تعداد داده‌های ورودی (تعداد کلیدها) به اندازه‌ی آرایه:

$$\alpha = \frac{n}{m}$$

عامل بار را می‌توان به این صورت تفسیر کرد که در m خانه‌ی حافظه‌ی با n داده‌ی ورودی به طور متوسط α داده در هر خانه قرار می‌گیرد.

برای به دست آوردن زمان متوسط سه عمل بالا فرض می‌کنیم تابع h برای قرار دادن هر کلید در حافظه، یکی از m خانه‌ی موجود را به صورت تصادفی انتخاب می‌کند، یعنی احتمال انتخاب شدن هر خانه برابر با $\frac{1}{m}$ است. در ضمن فرض می‌کنیم کلیدها مستقل از هم در خانه‌های جدول ذخیره می‌شوند. به این فرض “درهم‌سازی یک‌نواخت ساده”^۸ می‌گوییم، زیرا هر کلید شانس یکسانی برای ذخیره شدن در هر خانه‌ی جدول دارد. توجه کنید که این فرض صحیحی نیست، چون h یک تابع غیرتصادفی است و این فرض صرفاً جهت ساده‌سازی انجام می‌شود.

حال برای مثال میانگین زمان جست‌وجو را محاسبه می‌کنیم. بدترین حالت زمانی اتفاق می‌افتد که کلید k که به دنبال آن هستیم ذخیره نشده باشد. در این صورت باید کل لیست $h(k)$ را پیمایش کنیم. بنابر تعریف، طول این لیست به طور متوسط برابر است با عامل بار:

$$E[n_{h(k)}] = \alpha$$

پس زمان جست‌وجو در حالت ناموفق از مرتبه‌ی $\Theta(1 + \alpha)$ است. زمان ۱ برای این است که ببینیم کلید در کدام خانه‌ی

^۷load factor

^۸simple uniform hashing

حافظه است و در زمان α باید آرایه‌ی مربوط به آن خانه را جست‌وجو کنیم. در حالت جست‌وجوی موفق نیز با همین استدلال و با توجه به این که کلید k به طور متوسط در فاصله‌ی $\frac{\alpha}{2}$ از ابتدای لیست قرار دارد، زمان جست‌وجوی از مرتبه‌ی $\Theta(1 + \frac{\alpha}{2})$ است.

در مورد زمان درج و حذف اگر از لیست پیوندی دوسویه استفاده کنیم، هر دو عمل در $\Theta(1)$ انجام می‌شوند. بنابراین در مجموع اگر $\alpha = O(1)$ آن‌گاه همان‌طور که انتظار داشتیم هر سه عمل مورد نظر در $\Theta(1)$ انجام خواهند شد.

۸ تابع درهم‌ساز

در این قسمت می‌خواهیم چند نوع تابع درهم‌ساز مقدماتی را معرفی کنیم. فرض کنید می‌خواهیم برای مثال رهیاب در بالا، یک تابع درهم‌ساز بسازیم. یک پیشنهاد می‌تواند به صورت زیر باشد:

$$h(k) = k \pmod{256}$$

دقت کنید که عمل این تابع روی IP های مختلف را می‌توان به این صورت نمایش داد:

$$(x_1, x_2, x_3, x_4) \rightarrow x_4$$

آیا این تابع مناسب است؟ برای پاسخ دادن به این سوال باید فرض درهم‌سازی یک‌نواخت ساده را بررسی کنیم. در این حالت اگر بدانیم که هشت رقم آخر یک آدرس در مبنای دو توزیع یکنواخت دارند، این تابع در این فرض صدق می‌کند. البته در عمل موضوعات دیگری را نیز باید در نظر گرفت، مانند این که طراحی یک حمله برای چنین تابع ساده‌ای بسیار راحت است که در این جا به این مباحث نمی‌پردازیم. برای کاربردهای ساده، یک تابع مناسب در بسیاری از این موارد تابع زیر است:

$$h(k) = [m \{A \cdot k\}]$$

در این عبارت m مانند گذشته اندازه‌ی آرایه است. تابع $\{ \cdot \}$ همان تابع جز اعشار است، یعنی $\{x\} = x - [x]$. برای مقدار ثابت A ، انتخاب‌های زیادی می‌توان داشت ولی آن‌چه در عمل بسیار استفاده می‌شود،

$$A = \frac{\sqrt{5} - 1}{2} = 0,6180339887\dots$$

می‌باشد.

۹ آدرس‌دهی باز

همان‌گونه که گفته شد برای رفع مشکل برخورد، روش دیگری نیز وجود دارد که همان آدرس‌دهی باز است. در این روش برعکس روش قبل تمامی داده‌ها داخل خود آرایه ذخیره می‌شوند. شاید به نظر آید که می‌توان همان لیست‌های پیوندی را درون خود آرایه ذخیره کرد ولی نکته‌ی بعدی آن است که در این روش به اشاره‌گر نیز نیازی نیست، که این خود باعث صرفه‌جویی در حافظه است.

ایده‌ی کلی آن است که اگر برخورد رخ داد، به خانه‌ی بعدی برویم و همین روند را ادامه دهیم تا در نهایت به یک خانه‌ی خالی برسیم و یا این که کل آرایه پر شده باشد. ایده‌ی بعدی آن است که این ”وارسی“^۹ برای رسیدن به خانه‌ی خالی، به این صورت نیست که در هر مرحله به خانه‌ی بعدی نگاه کنیم و در عوض ”دنباله‌ی وارسی“^{۱۰} تابعی از کلید مورد نظر است.

برای این منظور باید به ازای هر کلید یک جایگشت از $\{0, 1, 2, \dots, m-1\}$ را به عنوان دنباله‌ی وارسی متناظر داشته باشیم، پس تابع درهم‌ساز را به این طریق گسترش می‌دهیم:

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

با در دست داشتن این تابع، دنباله‌ی وارسی متناظر با هر کلید به این صورت به دست می‌آید که:

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

حال برای درج داده با کلید k کافی است با همین ترتیب جدول را وارسی و داده را در اولین خانه‌ی خالی ذخیره کنیم. به طریق مشابه، برای جست‌وجوی داده با کلید k دوباره با همین ترتیب خانه‌ها را وارسی می‌کنیم. اگر کلید k را یافتیم که جست‌وجو موفق بوده است و اگر به خانه‌ای خالی رسیدیم، دیگر نیازی به ادامه‌ی وارسی نیست و جست‌وجو ناموفق بوده است. در زیر شبه‌کدهای متناظر با آنچه گفته شد، آورده شده است.

HASH-INSERT (T, k)

```

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

۱۰ الگوریتم حذف

در روش آدرس‌دهی باز الگوریتم حذف نیاز به توجه ویژه‌ای دارد. آنچه در ابتدا به ذهن می‌رسد این است که همان الگوریتم جست‌وجو را اجرا و پس از یافتن کلید مورد نظر، آن را حذف کنیم. اما توجه کنید اگر خانه‌ی پاک شده را با $null$ علامت‌گذاری کنید، در الگوریتم جست‌وجو و در حالت جست‌وجوی ناموفق دچار مشکل خواهید شد. بنابراین باید از نماد دیگری مانند $deleted$ استفاده کرد که در الگوریتم درج دقیقاً مانند $null$ عمل می‌کند، ولی در الگوریتم

^۹probe

^{۱۰}probe sequence

HASH-SEARCH (T, k)

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

جست‌وجو با رسیدن به آن جست‌وجو متوقف نمی‌شود. با اندکی دقت، باز مشکل دیگری وجود دارد و آن این که در حالتی که تعداد زیادی داده اضافه و پاک شوند، آن بخش از دنباله‌های واریسی که در الگوریتم جست‌وجو دنبال می‌شود به علت وجود تعداد زیادی *deleted* مستقل از عامل بار بلند خواهد شد و بنابراین اجرای الگوریتم جست‌وجو به صورت قابل ملاحظه‌ای کند می‌شود. پس با این که دیدیم روش آدرس‌دهی باز نسبت به استفاده از لیست‌های پیوندی مدیریت حافظه‌ی کارآمدتری دارد، در مواردی که داده‌ساختار به صورت مداوم تغییر می‌کند مناسب نیست.

۱۱ درهم‌سازی یک‌نواخت

برای تحلیل کارایی روش آدرس‌دهی باز به تعمیمی از فرض درهم‌سازی یک‌نواخت ساده نیاز داریم، که آن را فرض ”درهم‌سازی یک‌نواخت“^{۱۱} می‌نامیم. این فرض دقیقاً مانند فرض قبلی است، با این تفاوت که در این حالت تابع درهم‌ساز به ازای هر کلید یک جایگشت تصادفی تولید می‌کند، پس فرض درهم‌سازی یک‌نواخت به این صورت بیان می‌شود که دنباله‌ی واریسی متناظر با یک کلید مشخص، می‌تواند با احتمال برابری هر یک از $m!$ جایگشت ممکن باشد. با در نظر گرفتن این فرض، قضیه زیر در مورد زمان اجرای الگوریتم جست‌وجو برقرار است که اثبات آن در کتاب منبع درس آمده است و در این جا به آن اشاره نمی‌شود.

قضیه ۱ برای یک جدول درهم‌سازی با آدرس‌دهی باز که در آن ضریب بار $\alpha = \frac{n}{m} < 1$ ، امید ریاضی تعداد پیمایش‌ها در یک جستجوی ناموفق حداکثر برابر است با:

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

همچنین با فرض این که هر کلید در جدول با احتمال یکسانی جستجو می‌شود، امید ریاضی تعداد پیمایش‌ها در یک جستجوی موفق حداکثر برابر است با:

$$\frac{1}{\alpha} \ln\left(\frac{1}{1 - \alpha}\right)$$

^{۱۱}uniform hashing

۱۲ ساخت تابع درهم‌ساز از روی تابع درهم‌ساز ساده

در این قسمت می‌خواهیم به طور اجمالی به روش ساخت یک تابع درهم‌ساز برای آدرس‌دهی باز مانند

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

از روی یک تابع درهم‌ساز کمکی مانند

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

بپردازیم. به عنوان تلاش اول، تابع زیر را در نظر بگیرید:

$$h(k, i) = h'(k) + i \pmod{m}$$

این تابع در حقیقت همان ساده‌ترین روش موجود است که هر موقع خانه‌ای پر بود، به خانه‌ی بعدی آن برویم. مشکلی که برای این روش وجود دارد، زمان بالای الگوریتم جست‌وجو است. فرض کنید که i خانه‌ی متوالی پر باشند و خانه‌ی بعدی آن‌ها خالی باشد. آن‌گاه این خانه‌ی خالی در درج بعدی با احتمال $(i+1)/m$ پر می‌شود. به این ترتیب دنباله‌های طولانی از خانه‌های اشغال شده بلندتر و بلندتر می‌شوند و زمان جست‌وجو را کندتر و کندتر می‌کنند. به جای آن روشی را در نظر بگیرید که در هر مرحله به جای خانه‌ی بعد، به چند خانه بعد می‌پریم و طول این پرش برای هر کلید متفاوت و عددی تصادفی است. به این منظور می‌توان از یک تابع درهم‌ساز کمکی دیگری مانند h'' ، به صورت زیر استفاده نمود:

$$h(k, i) = h'(k) + i \times h''(k) \pmod{m}$$

به سادگی دیده می‌شود که مشکل تجمع خانه‌های اشغال شده که در بالا به آن اشاره شد، در این صورت اصلاح شده به وجود نمی‌آید.