

به نام خدا

ساختمان داده

آرش شفيعی



درهم سازی

- الگوریتم‌ها ممکن است نیاز داشته باشند عملیات متعددی بر روی مجموعه‌های پویا انجام دهند. بسیاری از الگوریتم‌ها نیاز دارند تنها عملیات درج، حذف و جستجو انجام دهند. مجموعه‌ای پویا که این عملیات را پشتیبانی کند دیکشنری¹ نامیده می‌شود.
- برای مثال در یک کامپایلر جدولی به نام جدول علائم وجود دارد که متغیرها و ویژگی‌های آنها را ذخیره می‌کند. برای ذخیره‌سازی متغیرهای یک برنامه و ویژگی‌های آن متغیرها از دیکشنری استفاده می‌کنیم.

¹ dictionary

- جدول درهم سازی¹ یک ساختار داده کارا برای پیاده سازی دیکشنری است.
- اگرچه جستجوی یک عنصر در جدول درهم سازی می تواند در بدترین حالت مانند جستجو در یک لیست پیوندی یعنی $\Theta(n)$ باشد، اما با تعدادی پیش فرض می توان جستجوی یک عنصر در جدول درهم سازی را در زمان $O(1)$ انجام داد.
- نوع داده ای دیکشنری در زبان پایتون توسط جداول درهم سازی پیاده سازی شده است.
- درواقع جدول درهم سازی مفهوم آرایه ها را تعمیم می دهد تا جستجو در مجموعه پویا همانند دسترسی به یک عنصر در آرایه در زمان $O(1)$ انجام شود.
- ایده اصلی جدول درهم سازی این است که اندیس یک آرایه حاوی یک مقدار کلید از طریق مقدار کلید آن با استفاده از روشی که توضیح خواهیم داد، محاسبه می شود.

¹ hash table

– آدرس دهی مستقیم¹ یک روش ساده است که وقتی مجموعه مرجع U کوچک باشد می تواند مورد استفاده قرار بگیرد.

¹ direct addressing

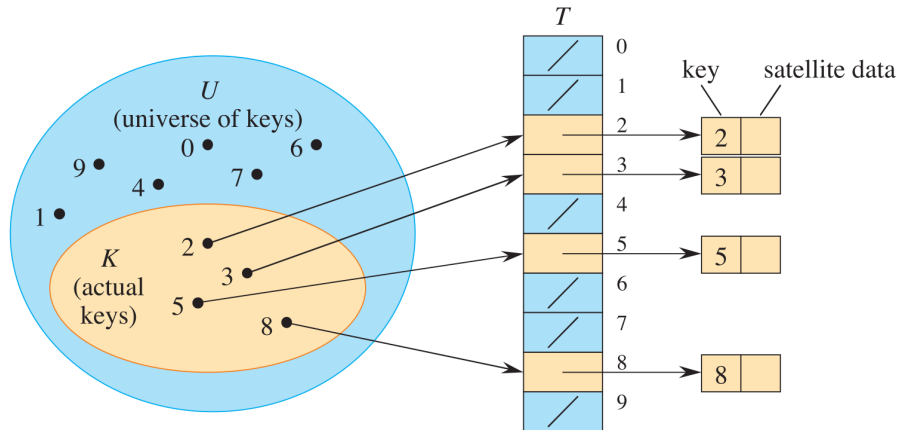
- فرض کنید برنامه‌ای می‌خواهد از یک مجموعهٔ پویا استفاده کند که کلیدهای عناصر آن یکتا هستند و این کلیدها زیر مجموعه‌ای از مجموعهٔ مرجع $U = \{0, 1, \dots, m-1\}$ ¹ هستند به طوری که m عدد بسیار بزرگی نیست.
- برای نمایش یک مجموعه پویا، می‌توانید از یک آرایه استفاده کنید که به آن جدول آدرس دهی مستقیم گفته می‌شود. این آرایه را با $T[0 : m-1]$ نمایش می‌دهیم هریک از مکان‌های آن (که به آن اسلات² نیز گفته گفته می‌شود) به یکی از کلیدها در مجموعهٔ مرجع U اختصاص دارد.

¹ universal set (universe)

² slot

جداول آدرس دهی مستقیم

- شکل زیر این روش را نمایش می دهد. مکان k به عنصری اشاره می کند که مقدار کلید آن k است. اگر مجموعه هیچ عنصری با کلید k نداشته باشد، آنگاه خواهیم داشت $T[k]=NIL$.



- همه توابع جستجو Direct-Address-Search و درج Direct-Address-Insert و حذف Direct-Address-Delete در زیر نشان داده شده اند که هرکدام در زمان $O(1)$ انجام می شوند.

Algorithm Direct Address Search

```
function DIRECT-ADDRESS-SEARCH(T,k)
1: return T[k]
```

Algorithm Direct Address Insert

```
function DIRECT-ADDRESS-INSERT(T,x)
1: T[x.key] = x
```

Algorithm Direct Address Delete

– **function** DIRECT-ADDRESS-DELETE(T, x)

1: $T[x.key] = \text{NIL}$

جداول درهم‌سازی

- اگر مجموعه مرجع U بسیار بزرگ یا نامحدود باشد، ذخیره‌سازی جدول T به اندازه $|U|$ غیر عملی یا حتی غیرممکن است.
- علاوه بر این، مجموعه k از کلیدهایی که واقعا ذخیره شده‌اند ممکن است نسبت به مجموعه مرجع U بسیار کوچک باشد و در نتیجه بیشتر فضایی که به T تخصیص داده شده است ممکن است هدر رود.
- وقتی مجموعه k یعنی کلیدهای ذخیره شده در دیکشنری بسیار کوچک‌تر از مجموعه مرجع U یعنی تمام کلیدهای ممکن باشد، جداول درهم‌سازی کمک می‌کنند که فضای بسیار کوچک‌تری نسبت به جدول آدرس‌دهی مستقیم اشغال شود.
- پیچیدگی حافظه با استفاده از جداول درهم‌سازی به $\Theta(|k|)$ کاهش پیدا می‌کند و پیچیدگی زمانی جستجو همچنان $O(1)$ باقی می‌ماند.

جداول درهم‌سازی

- در آدرس‌دهی مستقیم یک عنصر با کلید k ذخیره می‌شود، اما در جداول درهم‌سازی از یک تابع درهم‌سازی¹ h استفاده می‌کنیم تا مکان کلید k را محاسبه کنیم و بنابراین کلید k در مکان $h(k)$ قرار می‌گیرد. تابع درهم‌سازی h مجموعه مرجع U را به مکان‌های جدول درهم‌سازی $T[0 : m - 1]$ نگاشت می‌کند.

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

- در اینجا m اندازه جدول درهم‌سازی است که معمولاً بسیار کوچک‌تر از U است.

- می‌گوییم عنصر با کلید k به مکان $h(k)$ نگاشت می‌شود².

- همچنین می‌گوییم $h(k)$ مقدار نگاشت شده یا مقدار هش³ برای کلید k است.

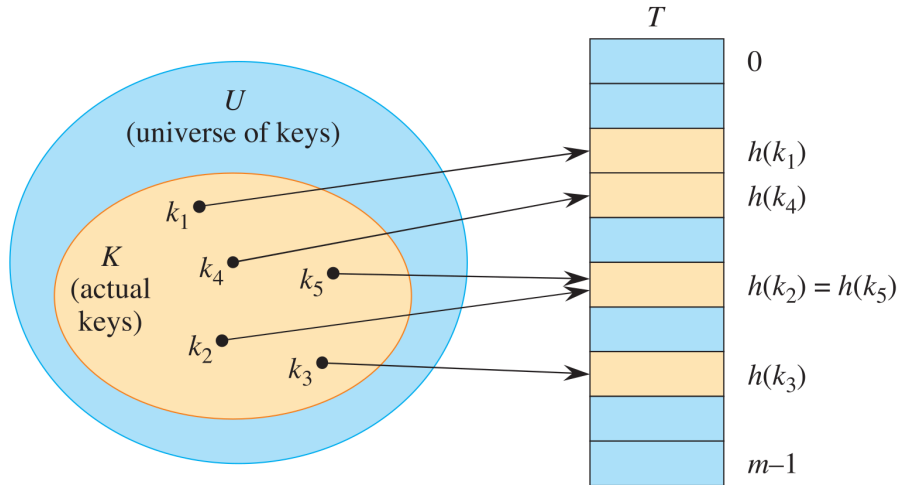
¹ hash function

² k hashes to slot $h(k)$

³ hash value

جداول درهم سازی

- ایده اصلی جدول درهم سازی در زیر نشان داده شده است.



- تابع درهم‌سازی محدوده اندیس‌های آرایه را کاهش می‌دهد و بنابراین اندازه آرایه مورد نیاز برای ذخیره را کاهش می‌دهد، در نتیجه به جای استفاده از جدولی به اندازه $|L|$ از جدولی به اندازه m استفاده می‌کنیم.
- یک مثال ساده از تابع درهم‌سازی که البته تابع کاربردی بسیار مناسبی نیست می‌تواند $h(k) = k \bmod m$ باشد.
- مشکلی که ممکن است به وجود آید این است که دو کلید به یک مکان حافظه نگاشت شوند. به چنین وضعیتی برخورد یا تصادم¹ گفته می‌شود.

¹ collision

- روش های کارایی برای حل مشکل تصادم وجود دارد.
- در شرایط ایده آل هیچ تصادمی نباید رخ دهد و تلاش می کنیم تابع درهم سازی را به نحوی انتخاب کنیم که تصادم حداقل شود.
- یک ایده برای درهم سازی این است که از یک تابع تولید کننده عدد تصادفی در فرایند نگاشت استفاده کنیم تا تعداد تصادم ها را کاهش دهیم.

- البته تابع درهم‌سازی باید قطعی¹ باشد به طوری که به ازای ورودی k همیشه خروجی $h(k)$ را تولید کند.
- از آنجایی که $|U| > m$ بنابراین امکان تصادم وجود دارد. بنابراین اگرچه یک تابع درهم‌سازی خوب تعداد تصادم‌ها را کاهش می‌دهد، اما همچنان امکان تصادم وجود دارد و باید روشی برای برطرف کردن تصادم به کار ببریم.
- در ادامه در مورد انتخاب تابع درهم‌سازی خوب و همچنین روش‌های برخورد با تصادم صحبت خواهیم کرد.

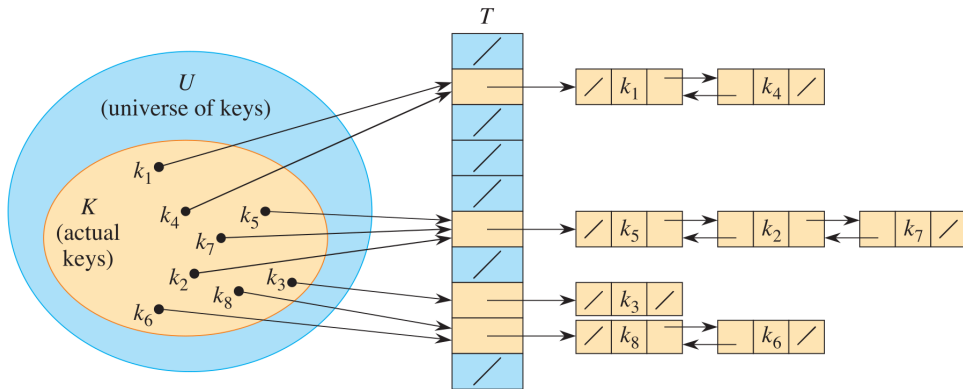
¹ deterministic

- درهم‌سازی مستقل یکنواخت¹ : یک تابع درهم‌سازی ایده‌آل h به ازای هر ورودی k در دامنه U خروجی $h(k)$ را مستقلاً به طور یکنواخت از محدود $\{0, 1, \dots, m-1\}$ انتخاب می‌کند.
- وقتی مقدار $h(k)$ انتخاب شد، هر فراخوانی بعدی h با ورودی k همان مقدار $h(k)$ را تولید می‌کند.
- چنین تابعی را یک تابع درهم‌سازی مستقل یکنواخت² می‌نامیم.

¹ independent uniform hashing

² independent uniform hash function

- رفع تصادم با روش زنجیره‌ای : شکل زیر ایده اصلی روش زنجیره‌ای¹ برای رفع تصادم را نشان می‌دهد.



¹ chaining

- هر مکان غیرتهی در جدول به یک لیست پیوندی اشاره می‌کند و همه عناصری که به یک مکان نگاشت می‌شوند وارد لیست پیوندی آن مکان می‌شوند.
- مکان z یک اشاره‌گر به ابتدای یک لیست پیوندی دارد که همه عناصری که به مکان z نگاشت می‌شوند را نگهداری می‌کند. اگر هیچ عنصری به مکان z نگاشت نشده باشد، مکان z تهی است و مقدار NIL را نگهداری می‌کند.
- اگر از چنین روشی برای رفع تصادم استفاده کنیم عملیات دیکشنری به صورت زیر پیاده‌سازی خواهند شد.

Algorithm Chained Hash Insert

```
function CHAINED-HASH-INSERT(T,x)
1: List-Prepend (T[h(x.key)],x)
```

Algorithm Chained Hash Search

— **function** CHAINED-HASH-SEARCH(T, k)
1: **return** List-Search ($T[h(k), k]$)

Algorithm Chained Hash Delete

function CHAINED-HASH-DELETE(T, x)
1: List-Delete ($T[h(x.key)], x$)

- زمان اجرای عملیات درج در بدترین حالت $O(1)$ است. عملیات درج سریع است زیرا فرض می کنیم عنصر x که می خواهیم درج کنیم در جدول وجود ندارد. بدون وجود این فرض، قبل از درج باید جستجو کنیم که آیا x در لیست پیوندی وجود دارد یا خیر.
- زمان اجرای عملیات جستجو در بدترین حالت متناسب با طول لیست است که این زمان اجرا را دقیق تر مورد بررسی قرار خواهیم داد.
- زمان اجرای عملیات حذف $O(1)$ است اگر لیست های پیوندی دو طرفه باشند. از آنجایی که تابع حذف یک اشاره گر x به عنصر مورد نظر دریافت می کند و نه مقدار کلید آن را، بنابراین نیازی به جستجو نیست.

- تحلیل درهم سازی با روش زنجیره ای : به ازای جدول درهم سازی T شامل m مکان که n عنصر را ذخیره می کند، ضریب بار α^1 برای T را به صورت n/m تعریف می کنیم که تعداد میانگین عناصر ذخیره شده در یک زنجیر است. تحلیلی که ارائه می کنیم براساس α است که می تواند کمتر، مساوی یا بیشتر از یک باشد.
- در بدترین حالت همه n کلید در یک مکان جدول ذخیره می شوند که درواقع لیستی به طول n می سازند. زمان اجرا در بدترین حالت برای جستجو $\Theta(n)$ است به علاوه زمان مورد نیاز برای محاسبه تابع درهم سازی.
- کارایی درهم سازی در حالت میانگین بستگی به این دارد که تابع درهم سازی h چقد خوب مجموعه n کلید را در m مکان توزیع می کند.

¹ load factor

- فرض کنید به ازای $j = 0, 1, \dots, m-1$ طول لیست $T[j]$ برابر با n_j باشد. بنابراین خواهیم داشت $n = n_0 + n_1 + \dots + n_{m-1}$ و مقدار مورد انتظار¹ برای n_j برابر با $E[n_j] = \alpha = n/m$ است.
- فرض می‌کنیم زمان $O(1)$ برای محاسبه مقدار هش $h(k)$ کافی است و بنابراین زمان لازم برای جستجوی یک مقدار با کلید k به طول $n_{h(k)}$ از لیست $T[h(k)]$ بستگی دارد.
- دو حالت را بررسی خواهیم کرد: حالتی که جستجوی یک کلید ناموفق باشد و هیچ عنصری در جدولی با کلید k وجود نداشته باشد و حالتی که جستجو موفق باشد و عنصری با کلید k یافت شود.

¹ expected value

- قضیه : در یک جدول درهم سازی که تصادم با زنجیر سازی رفع می شود، یک جستجوی ناموفق به طور متوسط در زمان $\Theta(1 + \alpha)$ انجام می شود.

- اثبات : با فرض بر این که درهم‌سازی مستقل و یکنواخت است، هر کلید k که هنوز در جدول ذخیره نشده است، می‌تواند با احتمال یکنواخت در هر یک از m مکان جدول ذخیره شود. زمان مورد نیاز برای یافتن کلیدی که در جدول وجود ندارد یعنی برای جستجوی ناموفق کلید k برابر با زمان مورد انتظار آن $E[n_{h(k)}] = \alpha$ است.
- بنابراین تعداد عناصر بررسی شده در یک جستجوی ناموفق برابر با α است و کل زمان مورد نیاز (شامل زمان محاسبه $h(k)$) برابر با $\Theta(1 + \alpha)$ است.

- وضعیت برای یک جستجوی موفق متفاوت است. یک جستجوی ناموفق با احتمال یکسان همه مکان های جدول درهم سازی را بررسی می کند. یک جستجوی موفق اما همه مکان ها را بررسی نمی کند، زیرا کلید مورد نظر در یکی از لیست های پیوندی موجود است. فرض می کنیم عنصر مورد بررسی با احتمال یکنواخت یکی از عناصر جدول باشد، بنابراین هرچه یک لیست بزرگ تر باشد احتمال این که کلید در آن لیست باشد بیشتر است. زمان مورد انتظار همچنان $\Theta(1 + \alpha)$ است.

- قضیه : در یک جدول درهم سازی که تصادم ها توسط زنجیر سازی رفع شده اند، یک جستجوی موفق به طور میانگین در زمان $\Theta(1 + \alpha)$ انجام می شود.

- اثبات : فرض می کنیم عنصر مورد جستجو با احتمال یکسان هریک از n عنصر ذخیره شده در جدول باشد. تعداد عناصر مورد بررسی در یک جستجوی موفق برای عنصر x برابر است با تعداد عناصر در لیستی که x متعلق به آن است به علاوه یک. از آنجایی که عناصر جدید در ابتدای لیست قرار می گیرند، عناصر قبل از x همه بعد از این که x وارد لیست شده درج شده اند.
- فرض کنید x_i عنصر i ام درج شده در جدول باشد به ازای $i = 1, 2, \dots, n$ و $k_i = x_i.key$ باشد.

- به ازای هر مکان q در جدول و برای هر زوج از کلیدهای متمایز k_i و k_j ، متغیر تصادفی زیر را تعریف می کنیم.

$$X_{ijq} = I \{h(k_i) = q, h(k_j) = q, x_i \text{ جستجو برای است}\}$$

- بنابراین وقتی k_i و k_j در مکان q برخورد کنند $X_{ijq} = 1$ است وقتی جستجو برای x_i باشد.

- از آنجایی که $\Pr\{x_i \text{ جستجو برای}\} = 1/n$ و $\Pr\{h(k_i) = q\} = 1/m$ و $\Pr\{h(k_j) = q\} = 1/m$ و این احتمال ها مستقل هستند، داریم $\Pr\{X_{ijq} = 1\} = 1/nm^2$

- حال برای هر عنصر x_i متغیر تصادفی زیر را تعریف می کنیم.

$$Y_j = I \{x_j \text{ در یک لیست قبل از عنصر مورد جستجو باشد}\}$$

$$= \sum_{q=0}^{m-1} \sum_{i=1}^{j-1} X_{ijq}$$

- زیرا حداکثر یکی از X_{ijq} ها برابر با یک است یعنی وقتی عنصر x_i که مورد جستجو است متعلق به همان لیست x_j باشد (که در مکان q قرار دارد) و $i < j$ است (یعنی x_i بعد از x_j در لیست ظاهر می شود).

- در نهایت متغیر تصادفی Z را تعریف می کنیم که تعداد عناصری که در لیست قبل از عنصر مورد جستجو قرار دارند را می شمارد.

$$Z = \sum_{j=1}^n Y_j$$

- از آنجایی که باید عنصر مورد جستجو و همه عناصر قبل از آن در لیست بشماریم، می‌خواهیم $E[Z + 1]$ را محاسبه کنیم.
- با محاسبه این مقدار به دست می‌آوریم $E[Z + 1] = 1 + \frac{\alpha}{2} + \frac{\alpha}{2n}$.
- بنابراین کل زمان مورد نیاز برای یک جستجوی موفق (شامل زمان مورد نیاز برای محاسبه تابع درهم‌سازی) برابر با $\Theta(1 + \alpha) = \Theta(2 + \alpha/2 - \alpha/2n)$ است.

- این تحلیل بدین معناست که اگر تعداد عناصر در جدول حداکثر متناسب با تعداد مکان‌های جدول درهم‌سازی باشد، داریم $n = O(m)$ و $\alpha = n/m = O(m)/m = O(1)$.
- بنابراین جستجو به طور میانگین در زمان ثابت انجام می‌شود. از آنجایی که درج بدترین حالت در زمان $O(1)$ و حذف در بدترین حالت در زمان $O(1)$ انجام می‌شود، وقتی لیست دو طرفه باشد، همه عملیات دیکشنری به طور متوسط در زمان $O(1)$ انجام می‌شوند.
- تحلیل ارائه شده تنها وقتی معتبر است که درهم‌سازی یکنواخت و مستقل باشد. یکنواخت بودن بدین معناست که کلید با احتمال برابر در هریک از m مکان نگاشت می‌شود و مستقل بودن بدین معناست که هر دو کلید متمایز با احتمال $1/m$ برخورد می‌کنند.

- برای اینکه درهم‌سازی خوبی داشته باشیم نیاز به یک تابع درهم‌سازی خوب داریم.
- علاوه بر اینکه تابع باید به طور کارا قابل محاسبه باشد، بررسی می‌کنیم یک تابع خوب چه ویژگی‌هایی باید داشته باشد و چگونه طراحی می‌شود.
- در این قسمت ابتدا دو درهم‌سازی با روش تقسیم و روش ضرب را معرفی می‌کنیم.
- یک روش کارا این است که خانواده‌ای از توابع درهم‌سازی داشته باشیم و در زمان اجرا به طور تصادفی از یکی از توابع خانواده استفاده کنیم. به این روش درهم‌سازی تصادفی گفته می‌شود. یکی از انواع درهم‌سازی تصادفی، درهم‌سازی سراسری که به معرفی آن خواهیم پرداخت.

توابع درهم‌سازی

- تابع درهم‌سازی خوب : یک تابع درهم‌سازی خوب بر اساس فرض درهم‌سازی یکنواخت مستقل عمل می‌کند بدین معنا که در آن هرکلید با احتمال یکسان به هرکدام از m مکان جدول نگاشت می‌شود مستقل از این که بقیه کلیدها به چه مکانی نگاشت نشده‌اند.
- احتمال یکسان در اینجا به چه معناست؟ اگر تابع درهم‌سازی ثابت باشد، هر احتمالی براساس توزیع احتمالی کلیدهای ورودی خواهد بود. اما معمولاً دسترسی به توزیع احتمالی کلیدهای ورودی نداریم.
- در برخی موارد ممکن است توزیع احتمالی را بدانیم. برای مثال، اگر بدانید کلیدها اعداد تصادفی حقیقی k هستند که به طور یکنواخت و مستقل در محدوده $0 \leq k < 1$ توزیع شده‌اند، آنگاه تابع درهم‌سازی $h(k) = [km]$ یک تابع یکنواخت و مستقل خواهد بود.
- یک تابع درهم‌سازی خوب مقدار را به گونه‌ای تولید می‌کند که مستقل از هرگونه الگویی در داده‌های ورودی است.

- در عمل، یک تابع درهم‌سازی برای نوع داده‌های زیر طراحی می‌شود.

- (۱) کلیدهای که اعداد صحیح غیر منفی کوچک هستند و در W بیت جای می‌گیرند. معمولاً مقدار W برابر با ۳۲ یا ۶۴ است.
- (۲) کلیدهایی که یک وکتور کوچک از اعداد صحیح غیر منفی هستند که هرکدام اندازه محدودی دارند. برای مثال، هر عنصر وکتور می‌تواند ۸ بیت باشد. یک وکتور را در این حالت یک رشته نیز می‌نامیم.

- درهم سازی ایستا¹ از یک تابع درهم سازی ثابت استفاده می کند.
- در اینجا در مورد دو روش استاندارد برای درهم سازی ایستا صحبت می کنیم که روش ضرب و تقسیم نام دارند.
- اگرچه درهم سازی ایستا در عمل استفاده نمی شود و از درهم سازی تصادفی استفاده می شود، اما درهم سازی ایستا پایه ای برای یادگیری مفاهیم درهم سازی است.

¹ static hashing

- روش تقسیم: در روش تقسیم¹ برای ساختن یک تابع درهم‌سازی کلید k را به یکی از m مکان با محاسبه باقیمانده k بر m نگاشت می‌کنیم. بنابراین تابع درهم‌سازی برابر است با $h(k) = k \bmod m$
- برای مثال، اگر اندازه جدول درهم‌سازی $m = 12$ باشد و کلید مورد نظر $k = 100$ باشد، آنگاه $h(k) = 4$. از آنجایی که تنها عملیات مورد نیاز تقسیم است، روش تقسیم روش نسبتاً سریعی است.
- روش تقسیم وقتی m یک عدد اول باشد که به اعداد توان ۲ نزدیک نباشد خوب عمل می‌کند. اما هیچ تضمینی وجود ندارد که این روش در حالت میانگین کارایی خوبی داشته باشد.

¹ division method

- روش ضرب : روش ضرب

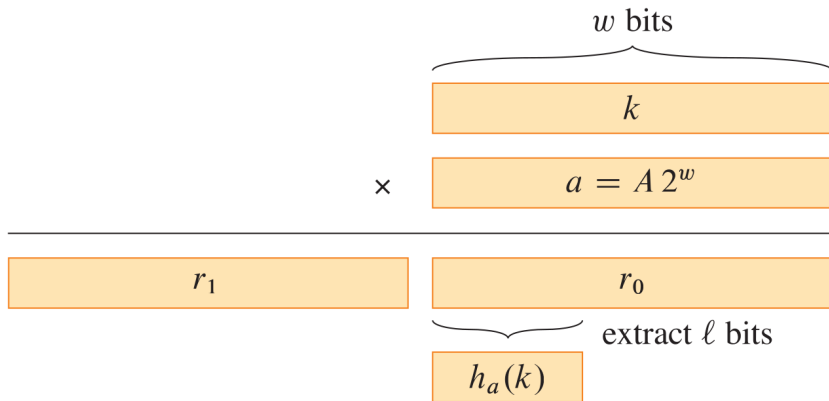
برای ساخت تابع درهم‌سازی در دو گام عمل می‌کند. ابتدا، مقدار کلید k را در یک ثابت A در محدوده $0 < A < 1$ ضرب می‌کند و قسمت اعشاری kA را جدا می‌کند. سپس مقدار به دست آمده را m ضرب می‌کند و کف آن را به عنوان نتیجه محاسبه می‌کند. بنابراین تابع درهم‌سازی برابر است با

$$h(k) = \lfloor m(kA \bmod 1) \rfloor$$

- در اینجا $kA \bmod 1$ به معنای قسمت اعشاری kA است یعنی $kA - \lfloor kA \rfloor$.

- روش ضرب-انتقال : در عمل، روش ضرب وقتی m توانی از عدد ۲ است یعنی $m = 2^l$ (به ازای عدد صحیح l) به طوری که $l \leq w$ و w تعداد بیت های یک کلمه در ماشین است، خوب عمل می کند.
- اگر یک عدد صحیح مثبت w بیتی $a = A2^w$ را انتخاب کنید به طوری که $0 < A < 1$ و $0 < a < 2^w$ ، می توانیم تابع درهم سازی را به صورت زیر پیاده سازی کنیم. فرض می کنیم کلید k در یک کلمه w بیتی جای می گیرد.

- این روش در شکل زیر نشان داده شده است.



درهم سازی ایستا

- ابتدا k را در عدد صحیح w بیتی a ضرب می کنیم. نتیجه مقدار $2w$ بیتی $r_0 + r_1 2^w$ است جایی که r_1 مقدار w بیت مرتبه بالا و r_0 مقدار w بیت مرتبه پایین نتیجه ضرب است.
- مقدار l بیتی هش تشکیل شده است از l بیت پر ارزش r_0 .
- از آنجایی که از r_1 چشم پوشی می شود، تابع درهم سازی می تواند بر روی یک کامپیوتر به گونه ای پیاده سازی شود که یک حاصل ضرب بر 2^w محاسبه کند.
- تابع درهم سازی در اینجا به صورت $h_a(k) = (ka \bmod 2^w) \gg (w - l)$ محاسبه می شود.
- از آنجایی که حاصل ضرب ka از دو عدد $2w$ فضا اشغال می کند، باقیمانده ضرب بر 2^w تنها w بیت پر ارزش را به دست می دهد (r_1) و w بیت کم ارزش را حذف می کند (r_0).

- عملگر \gg یک انتقال به چپ به اندازه $w-1$ بیت انجام می دهد، که معادل است با تقسیم عدد بر 2^{w-1} و محاسبه کف نتیجه.
- تابع درهم سازی h_a می تواند با سه دستور پیاده سازی شود، ضرب، تفریق و انتقال.
- به عنوان مثال، فرض کنید $k = 123456$ و $l = 14$ و $m = 2^{14} = 16384$ و $w = 32$.
- همچنین فرض کنید $a = 2654435769$ (که توسط دونالد کنوت پیشنهاد شده است) در این صورت $r_0 = 17612864$ و $r_1 = 76300$ و $ka = 327706022297664 = (76300 \cdot 2^{32}) + 17612864$.
- بنابراین خواهیم داشت $h_a(k) = 67$.

- اگرچه روش ضرب-انتقال روش سریعی است اما تضمین نمی کند که کارایی خوبی در حالت میانگین داشته باشد. برای درهم سازی سراسری که درمورد آن صحبت خواهیم کرد این تضمین وجود دارد.
- اگر a به صورت تصادفی از بین اعداد فرد انتخاب شود، روش ضرب-انتقال در حالت میانگین خوب عمل می کند.

درهم سازی تصادفی

- فرض کنید یک نفر کلیدها به گونه ای انتخاب کند که همه آنها به یک مکان جدول نگاشت شوند. در این حالت پیچیدگی زمانی جستجو $\Theta(n)$ خواهد بود. هرگونه درهم سازی ایستا دارای این نقطه ضعف است.
- یک روش برای مقابله با این مشکل استفاده از درهم سازی تصادفی¹ است.
- یکی از حالات خاص این روش درهم سازی جامع² است که کارایی خوبی دارد وقتی تصادم ها با زنجیر سازی مدیریت می شوند.
- برای استفاده از درهم سازی تصادفی، در ابتدای اجرای برنامه تابع درهم سازی را به صورت تصادفی از میان خانواده ای از توابع درهم سازی انتخاب می کنیم. از آنجایی که تابع درهم سازی را به صورت تصادفی انتخاب می کنید، الگوریتم در هر اجرا رفتار متفاوتی دارد.

¹ random hashing

² universal hashing

- فرض کنید H یک خانواده محدود از توابع درهم‌سازی باشد که یک مجموعه مرجع U از کلیدها را به بازه $\{0, 1, \dots, m-1\}$ نگاشت می‌کند. به این خانواده یک خانواده جامع گفته می‌شود اگر به ازای هرجفت از کلیدهای متمایز $k_1, k_2 \in U$ ، تعداد توابع $h \in H$ که در رابطه $h(k_1) = h(k_2)$ صدق می‌کنند، حداکثر $|H|/m$ باشد.
- به عبارت دیگر با انتخاب یک تابع درهم‌سازی به صورت تصادفی که از مجموعه H انتخاب شده باشد، احتمال تصادم بین دو کلید متمایز k_1 و k_2 بیشتر از احتمال تصادم $1/m$ وقتی که $h(k_1)$ و $h(k_2)$ به صورت تصادفی و مستقل از مجموعه $\{0, 1, \dots, m-1\}$ انتخاب شده باشند نیست.
- درهم‌سازی یکنواخت مستقل مانند انتخاب تابع درهم‌سازی به صورت یکنواخت و تصادفی از مجموعه m^n تابع درهم‌سازی است، به طوری که هر عضو از آن خانواده n کلید را به m مقدار نگاشت می‌کنند.

- هر خانواده تصادفی یکنواخت مستقل از توابع درهم سازی جامع است، اما برعکس آن صحیح نیست.
- مجموعه $U = \{0, 1, \dots, m-1\}$ را در نظر بگیرید و فرض کنید تنها یک تابع همانی در خانواده به عنوان تابع درهم سازی در نظر گرفته شده است. احتمال این که دو کلید متمایز برخورد کنند صفر است اگرچه هر کلید به یک مقدار ثابت نگاشت می شود.
- می توان اثبات کرد که با استفاده از درهم سازی جامع و رفع تصادم توسط زنجیر سازی در یک جدول خالی با m مکان، هر دنباله ای از s عملیات درج، جستجو، و حذف شامل $n = O(m)$ عملیات درج در زمان $\Theta(s)$ انجام می شود.

درهم‌سازی تصادفی

- ویژگی‌های زیادی برای یک خانواده H از توابع درهم‌سازی وجود دارند که در اینجا تعدادی از آنها را اجمالا بیان می‌کنیم.
- فرض کنید H یک خانواده از توابع درهم‌سازی باشد که هر کدام دارای دامنه U و برد $\{0, 1, \dots, m-1\}$ هستند و فرض کنید h یک تابع درهم‌سازی باشد که به صورت یکنواخت و تصادفی از H انتخاب می‌شود.
- خانواده H یکنواخت است اگر به ازای هر کلید k در U و هر مکان q در $\{0, 1, \dots, m-1\}$ احتمال این که $h(k) = q$ باشد $1/m$ باشد.
- خانواده H جامع است اگر برای هر جفت کلید متمایز k_1 و k_2 در U ، احتمال این که $h(k_1) = h(k_2)$ باشد، حداکثر $1/m$ باشد.

- خانواده H از توابع درهم‌سازی جامع ϵ ¹ است اگر به ازای هر جفت کلید k_1 و k_2 در U احتمال این که $h(k_1) = h(k_2)$ باشد حداکثر ϵ باشد. بنابراین یک خانواده جامع از توابع درهم‌سازی، درواقع جامع $1/m$ ² است.
- خانواده H یک خانواده مستقل d است اگر به ازای هر مجموعه از کلیدهای متمایز k_1, k_2, \dots, k_d در U و هر مجموعه مکان‌های q_1, q_2, \dots, q_d در $\{0, 1, \dots, m-1\}$ احتمال این که $h(k_j) = q_i$ باشد به ازای $i = 1, 2, \dots, d$ برابر با $1/m^d$ باشد.
- خانواده‌های توابع درهم‌سازی جامع را به این دلیل مطالعه می‌کنیم که ساده‌ترین نوع پیاده‌سازی توابع درهم‌سازی هستند که اثبات شده است کارایی خوبی دارند.

¹ ϵ -universal

² $1/m$ -universal

- در اینجا دو روش برای طراحی یک خانواده جامع از توابع درهم سازی معرفی می کنیم.
- روش اول : عدد اول p که به اندازه کافی بزرگ است را انتخاب کنید به طوری که هر کلید k در محدوده 0 تا $p - 1$ قرار بگیرد.
- فرض کنید \mathbb{Z}_p مجموعه $\{1, 2, \dots, p - 1\}$ باشد و \mathbb{Z}_p^* مجموعه $\{1, 2, \dots, p - 1\}$.
- از آنجایی که اندازه مجموعه مرجع کلیدها بزرگتر از تعداد مکان ها در جدول درهم سازی است، داریم $p > m$.

- به ازای هر $a \in \mathbb{Z}_p^*$ و $b \in \mathbb{Z}_p$ ، تابع h_{ab} را به صورت تبدیل خطی¹ و کاهشی با استفاده عملگر باقیمانده بر p و سپس باقیمانده بر m تعریف می کنیم.

$$h_{ab}(k) = ((ak + b) \bmod p) \bmod m$$

- برای مثال اگر $p = 17$ و $m = 6$ باشد، داریم

$$\begin{aligned} h_{3,4}(8) &= ((3 \times 8 + 4) \bmod 17) \bmod 6 \\ &= (28 \bmod 17) \bmod 6 \\ &= 11 \bmod 6 \\ &= 5 \end{aligned}$$

¹ linear transformation

- به ازای p و m ، خانواده همه توابع درهم سازی به صورت زیر است.
$$H_{pm} = \{h_{ab} : a \in \mathbb{Z}_p^*, b \in \mathbb{Z}_p\}$$
- هر تابع درهم سازی h_{ab} مجموعه \mathbb{Z}_p را به \mathbb{Z}_m نگاشت می کند.
- یک ویژگی این خانواده از توابع درهم سازی این است که اندازه m می تواند هر مقدار دلخواهی باشد و نیازی نیست اول باشد.
- از آنجایی که می توانید از میان $p - 1$ مقدار برای a و p مقدار برای p انتخاب کنید، خانواده H_{pm} شامل $p(p - 1)$ تابع درهم سازی است.
- می توان اثبات کرد H_{pm} یک خانواده جامع است.

- روش دوم : از این روش در عمل بیشتر استفاده می‌شود زیرا کارایی بالایی دارد.
- فرض کنید H یک خانواده از توابع درهم‌سازی ضرب-انتقال صورت زیر باشد.
 $H = \{h_a : h_a \text{ تابع درهم‌سازی ضرب انتقال است و } 1 \leq a < m \text{ و } a \text{ فرد است}\}$
- تابع درهم‌سازی ضرب-انتقال را قبلاً به صورت زیر تعریف کردیم.
$$h_a(k) = (ka \bmod 2^w) \gg (w - l)$$
- می‌توان اثبات کرد توابع خانواده H یک خانواده جامع $2^w/m$ ¹ هستند. به عبارت دیگر احتمال این که هر دو کلید متمایز برخورد کنند حداکثر $2^w/m$ است.

¹ $2^w/m$ -universal

- درهم سازی ورودی های بزرگ مانند وکتورها یا رشته ها؛
- برخی مواقع ورودی تابع درهم سازی آنقدر بزرگ است که به سادگی نمی توان یک عدد اول p پیدا کرد که در ۶۴ بیت جای بگیرد. برای مثال ممکن است بخواهیم یک رشته را به عنوان کلید در نظر بگیریم.
- یک روش درهم سازی برای چنین ورودی هایی تعمیم روش های شرح داده شده است. روش دوم استفاده از توابع درهم سازی رمزنگاری است.

- یک تابع رمزنگاری یک ورودی با اندازه دلخواه را دریافت می کند و یک خروجی با طول ثابت تولید می کند.
 - برای مثال تابع استاندارد SHA – 256 یک خروجی ۲۵۶ بیتی (۳۲ بایتی) تولید می کند.
 - توابع رمزنگاری به قدری مورد استفاده اند که برخی از پردازنده ها توابع رمزنگاری را به صورت سخت افزاری پیاده سازی می کنند و بنابراین سرعت اجرای بالایی دارند.
 - می توانیم تابع درهم سازی برای یک کلید با طول دلخواه را به صورت زیر تعریف کنیم.
- برای تعریف یک خانواده از توابع درهم سازی می توانیم یک رشته a به ابتدای ورودی الحاق کنیم.

$$h(k) = \text{SHA} - 256(k) \bmod m$$

$$h_a(k) = \text{SHA} - 256(a||k) \bmod m$$

- آدرس دهی باز¹ روشی برای رفع تصادم است که برخلاف زنجیرسازی از فضایی خارج از جدول درهم سازی استفاده نمی کند.
- در آدرس دهی باز همه عناصر در خود جدول درهم سازی ذخیره می شوند و از فضای اضافی استفاده نمی شود. بنابراین جدول درهم سازی ممکن است پر شود. در نتیجه ضریب بار α هیچ گاه نمی تواند بیشتر از ۱ شود.

¹ open addressing

- در این روش تصادم به صورت زیر مدیریت می شود : وقتی می خواهیم یک عنصر در جدول وارد کنیم، این عنصر جدید در صورت امکان در مکانی که به عنوان انتخاب اول آن محاسبه می شود قرار می گیرد. این روند ادامه پیدا می کند تا جایی که یک مکان خالی پیدا شود.
- برای جستجوی یک عنصر مکان های جدول به ترتیب انتخاب ها بررسی می شوند تا جایی که یا عنصر یافته شود و یا به یک خانه خالی برخورد کنیم که در این صورت پیام جستجوی ناموفق صادر می کنیم.
- آدرس دهی باز از هیچ اشاره گر و فضای اضافی استفاده نمی کند، در نتیجه می توان از جدول بزرگتری برای داده های آن استفاده کرد.

- برای درج در جدول توسط آدرس‌دهی باز، به ترتیب همه انتخاب‌ها برای مکان کلید مورد نظر را واری¹ می‌کنیم تا وقتی که یک فضای خالی در جدول پیدا شود.
- تابع درهم‌سازی به ازای یک کلید باید شماره انتخاب یا شماره واری² را به عنوان ورودی دوم دریافت کند. بنابراین دامنه و برد تابع درهم‌سازی به صورت زیر خواهد بود.

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

- در آدرس‌دهی باز نیاز داریم به ازای هر کلید k ، یک دنباله انتخاب‌ها یا دنباله واری³ صورت $\langle h(k, 0), h(k, 1), \dots, h(k, m-1) \rangle$ تولید کنیم تا همه مکان‌های جدول مورد بررسی قرار بگیرند.

¹ probe

² probe number

³ probe sequence

- تابع Hash-Insert جدول T و کلید k را دریافت می کند و مکان ذخیره آن را باز می گرداند و در صورتی که جدول پر باشد پیام خطا صادر می کند.

Algorithm Hash Insert

```
function HASH-INSERT(T,k)
1: i = 0
2: repeat
3:   q = h(k,i)
4:   if T[q] == NIL then
5:     T[q] = k
6:     return q
7:   else i = i + 1
8: until i == m
9: error "hash table overflow"
```

- الگوریتم جستجوی کلید k نیز دنباله ای از مکان های جدول را بررسی می کند تا وقتی که کلید k یافته شود. وقتی به یک مکان خالی می رسیم جستجو با پیامی مبنی بر عدم موفقیت و بازگرداندن NIL به پایان می رسد.

- تابع Hash-Search جدول T و کلید k را دریافت می کند و در صورتی که مکان q شامل کلید k پیدا شد مکان q را باز می گرداند و در صورتی که k در جدول وجود نداشت مقدار NIL را باز می گرداند.

Algorithm Hash Search

```
function HASH-SEARCH(T,k)
1: i = 0
2: repeat
3:   q = h(k,i)
4:   if T[q] == k then
5:     return q
6:   i = i + 1
7: until T[q] == NIL or i == m
8: return NIL
```

- حذف کردن از جدول درهم سازی با آدرس دهی باز کمی پیچیده تر است. وقتی یک کلید از مکان q را حذف می کنیم ممکن است به اشتباه آن مکان به عنوان مکان خالی در نظر گرفته شود. اگر مکان حذف شده را با NIL علامتگذاری کنیم، ممکن است نتوانیم کلید k را که مکان q برای آن واریسی یافته شده بوده در حالی که مکان اصلی آن اشغال بوده است را پیدا کنیم.
- برای حل این مشکل، وقتی یک واریسی کلید را در جدول حذف می کنیم به جای علامتگذاری آن با مقدار NIL، مقدار Deleted را در مکان حذف شده قرار می دهیم.

- در هنگام درج در جدول، تابع Hash-Insert مکانی که با Deleted علامتگذاری شده است را به عنوان مکان خالی در نظر می گیرد. در هنگام جستجو اما وقتی تابع Hash-Search به مکان Deleted برخورد می کند جستجو را ادامه می دهد، زیرا مکانی که با Deleted علامتگذاری شده در هنگام درج کلید مورد جستجو اشغال بوده است.
- زمان جستجو در روش آدرس دهی باز به علت استفاده از این علامتگذاری به ضریب α بستگی ندارد و به همین دلیل وقتی می خواهیم کلیدها را حذف کنیم از روش زنجیرسازی استفاده می کنیم که زمان جستجوی بهتری دارد.

- دو روش برای طراحی توابع درهم سازی توسط آدرس دهی باز وجود دارد که در اینجا در مورد آنها صحبت خواهیم کرد : درهم سازی دوگانه و واریسی خطی.

- درهم سازی دوگانه ¹ : درهم سازی دوگانه یکی از بهترین روش های موجود برای آدرس دهی باز است.
- درهم سازی دوگانه از تابع درهم سازی به صورت زیر استفاده می کند.

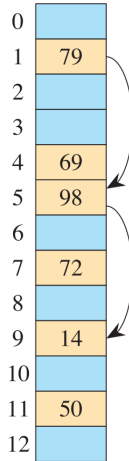
$$h(k, i) = (h_1(k) + ih_2(k)) \bmod m$$

- به طوری که h_1 و h_2 دو تابع درهم سازی هستند.
- اولین مکان انتخابی $T[h_1(k)]$ است و مکان های انتخابی بعدی گام هایی به اندازه ضربی از $h_2(k)$ بعد از مکان انتخابی اول هستند.

¹ double hashing

آدرس دهی باز

- در شکل زیر یک مثال از درهم سازی دوگانه نشان داده شده است.



- اندازه جدول ۱۳ است و $h_1(k) = k \bmod 13$ و $h_2(k) = 1 + (k \bmod 11)$.
- کلید 14 بعد از واریسی مکان‌های 1 و 5 که اشغال شده‌اند وارد مکان 9 می‌شود.

- برای این که همه جدول مورد جستجو قرار بگیرد مقدار $h_2(k)$ باید نسبت به اندازه جدول یعنی m اول باشد. یک روش مناسب برای اطمینان حاصل کردن از برقراری این شرط است که m را توانی از ۲ انتخاب کنیم و h_2 را به گونه ای طراحی کنیم که همیشه عدد فرد تولید کند. یک روش دیگر این است که m را اول انتخاب کنیم و h_2 را به گونه ای طراحی کنیم که عددی مثبت کوچکتر از m تولید کند.
- برای مثال می توانیم m را یک عدد اول انتخاب کنیم و توابع را به صورت زیر طراحی کنیم.

$$h_1(k) = k \bmod m$$

$$h_2(k) = 1 + (k \bmod m')$$

- به طوری که m' عددی است که کمی کوچکتر از m است (مثلاً $m - 1$).

- به عنوان مثال، اگر $k = 123456$ و $m = 701$ باشد و $m' = 700$ باشد آنگاه $h_1(k) = 80$ و $h_2(k) = 257$ خواهد بود، بنابراین اولین واریسی مکان 80 را انتخاب می کند و در واریسی های بعدی ضرایبی از 257 مکان بعدی (با محاسبه باقیمانده بر m) انتخاب می شوند تا وقتی که یک مکان خالی پیدا شود.
- اگرچه مقادیر m که اول یا توانی از ۲ نیستند می توانند در درهم سازی دوگانه استفاده شوند، اما در عمل انتخاب $h_2(k)$ سخت تر می شود، زیرا باید اطمینان حاصل شود که نسبت به m اول است.
- وقتی m عدد اول یا توانی از ۲ باشد، درهم سازی دوگانه تعداد $\Theta(m^2)$ دنباله واریسی تولید می کند زیرا هر جفت $(h_1(k), h_2(k))$ یک دنباله واریسی متفاوت است.

- واریسی خطی¹ : واریسی خطی یک حالت خاص از درهم سازی دوگانه است و یکی از روش های بسیار ساده در آدرس دهی باز برای رفع تصادم است.
- همانند درهم سازی چندگانه ابتدا مکان $T[h_1(k)]$ بررسی می شود. اگر این مکان اشغال بود، مکان بعدی یعنی $T[h_1(k) + 1]$ بررسی می شود. بررسی تا $T[m - 1]$ ادامه پیدا می کند و مجدداً به $T[0]$ باز می گردد تا به $T[h_1(k) - 1]$ برسد.
- بنابراین گام واریسی که در درهم سازی دوگانه برابر با $h_2(k)$ بود در واریسی خطی برابر با ۱ است پس $h_2(k) = 1$ به ازای همه مقادیر k .

¹ linear probing

- پس تابع درهم سازی در واریسی خطی به ازای $i = 0, 1, \dots, m - 1$ برابر است با :

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

- مقدار $h_1(k)$ می تواند هریک از مقادیر $\{0, 1, \dots, m - 1\}$ باشد.

- در اینجا آدرس دهی باز را بر اساس ضریب بار جدول درهم سازی یعنی $\alpha = n/m$ تحلیل می کنیم.
- در آدرس دهی باز، حداکثر یک عنصر در یک مکان ذخیره می شود و بنابراین $n \leq m$ است و نتیجه $\alpha \leq 1$ است.
- از آنجایی که حذف از جدول درهم سازی تأثیری در فرایند جستجو ندارد فرض می کنیم حذف وجود ندارد.

- همانند قبل فرض می‌کنیم درهم‌سازی جایگشت‌ها را یکنواخت و مستقل تولید می‌کند. بنابراین در این شرایط دنباله واریسی $\langle h(k, 0), h(k, 1), \dots, h(k, m-1) \rangle$ که برای درج یا جستجو برای کلید k استفاده می‌شود با احتمال برابر می‌تواند هریک از جایگشت‌های $\langle 0, 1, \dots, m-1 \rangle$ باشد. البته هر کلید یک دنباله واریسی یکتا خواهد داشت.
- حال تعداد مورد انتظار واریسی‌ها برای درهم‌سازی با آدرس دهی باز را با شروع از تعداد مورد انتظار واریسی‌ها در یک جستجوی ناموفق تحلیل می‌کنیم.
- کران به دست آمده برابر است با $1 + \alpha + \alpha^2 + \alpha^3 + \dots$
- اولین واریسی همیشه رخ خواهد داد، اگر مکان اول انتخاب شده با احتمال α اشغال باشد، انتخاب دوم بررسی می‌شود. احتمال این که هر دو انتخاب اول و دوم اشغال باشند برابر است با α^2 و به همین ترتیب الی آخر.

- قضیه : به ازای یک جدول درهم سازی با آدرس دهی باز و ضریب بار $\alpha = n/m < 1$ ، تعداد واریسی ها در یک جستجوی ناموفق حداکثر $1/(1 - \alpha)$ است.

- اگر α ثابت باشد، این قضیه پیش بینی می کند که جستجوی ناموفق در زمان $O(1)$ اجرا می شود.
- برای مثال، اگر جدول درهم سازی نیمه پر لیوان باشد، تعداد متوسط واریسی ها در یک جستجوی ناموفق حداکثر $2 = 1/(1 - 0.5)$ است. اگر ۹۰ درصد اشغال باشد، تعداد متوسط واریسی ها حداکثر $10 = 1/(1 - 0.9)$ است.

- بنابراین درج یک عنصر در جدول درهم سازی با آدرس دهی باز در جدولی با ضریب α حداکثر $1(1 - \alpha)$ واریسی به طور میانگین نیاز دارد.

- قضیه : به ازای یک جدول درهم سازی با آدرس دهی باز و ضریب بار $\alpha > 1$ ، تعداد مورد انتظار واریسی ها در یک جستجوی موفق برابر است با حداکثر

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

- اثبات : در جستجوی کلید k همان دنباله واریسی تولید می شود که در هنگام درج کلید k تولید شده است. اگر $(i + 1)$ امین کلید درج شده در جدول درهم سازی کلید k باشد، ضریب بار در هنگام درج i/m بوده است و بنابراین طبق قضیه قبل تعداد مورد انتظار واریسی ها در جستجوی کلید k برابر با $1(1 - i/m) = m/(m - i)$ خواهد بود.
- اگر بر روی همه n کلید در جدول میانگین بگیریم مقدار زیر را برای مقدار مورد انتظار واریسی ها در یک جستجوی ناموفق به دست می آوریم :

$$\begin{aligned} \frac{1}{n} \sum_{i=0}^{n-1} \frac{m}{m-i} &= \frac{m}{n} \sum_{i=0}^{n-1} \frac{1}{m-i} \\ &= \frac{m}{n} \ln \frac{m}{m-n} \\ &= \frac{1}{\alpha} \ln \frac{1}{1-\alpha} \end{aligned}$$

- در یک جدول درهم سازی نیمه پر ، تعداد مورد انتظار واری‌ها در یک جستجوی موفق برابر با 1.387 است. اگر ۹۰ درصد جدول اشغال باشد، تعداد مورد انتظار واری‌ها کمتر از 2.559 است. اگر $\alpha = 1$ باشد، در یک جستجوی ناموفق، همه m مکان جدول باید واری شوند.