

FIXME: Zatím chybí červeně-černé stromy a trie.

Hashování

Mějme universum U (jeho velikost označíme u), množinu P přihrádek ($p = |P|$) a nějakou funkci $h : U \rightarrow P$, které budeme říkat *hashovací funkce*.

Datová struktura bude fungovat takto: když prvek vkládáme, spočteme hashovací funkci a vložíme prvek do příslušné přihrádky (přihrádky budeme reprezentovat jako pole seznamů). Pokud chceme prvek vyhledat nebo smazat, opět vyhodnotíme hashovací funkci a dozvíme se, ve které jediné přihrádce ho dává smysl hledat.

Budeme-li předpokládat, že výpočet funkce h trvá $\mathcal{O}(1)$, bude vkládání pracovat v konstantním čase a ostatní operace v čase lineárním s počtem prvků v dané přihrádce. Pokud se hashovací funkce bude chovat „rozumně náhodně“, můžeme očekávat, že po vložení n prvků jich bude v každé přihrádce přibližně n/p , takže při volbě $p \approx n$ můžeme získat konstantní časovou složitost operací. (Volit $p \gg n$ nemá smysl, protože pak bychom inicializací pole trávili příliš mnoho času.)

Tento přístup má ale samozřejmě své zadrhele: potřebujeme s prvky universa umět počítat (už si nevystačíme s porovnáváním), ale hlavně potřebujeme sehnat hashovací funkci, která se chová dostatečně rovnoměrně. Často se používají funkce, které se pro obvyklé vstupy chovají „pseudonáhodně“, třeba:

- $x \mapsto ax \bmod p$, pokud je universum číselné (pro nějakou konstantu a ; nejlepší je, když a i p jsou prvočísla);
- $x_1, \dots, x_n \mapsto (\sum_i C^i x_i) \bmod p$, pokud hashujeme řetězce (C a p opět nejlépe prvočíselná, navíc je-li ℓ obvyklá délka řetězce, mělo by být $C^\ell \gg p$).

Nicméně, ať už zvolíme jakoukoliv deterministickou funkci, vždy budou existovat nepříjemné vstupy, pro které skončí všechny prvky v téže přihrádce a operace budou mít lineární složitost namísto konstantní. Pomůžeme si snadno: vybereme hashovací funkci náhodně. Ne ze všech funkcí (ty bychom neuměli reprezentovat), nýbrž z vhodně zvolené třídy funkcí, které umíme snadno popisovat pomocí parametrů.

Definice: Systém funkcí S z U do P nazveme c -universální (pro nějaké $c \geq 1$), pokud pro všechny dvojice x, y navzájem různých prvků z U platí

$$\Pr_{h \in S}[h(x) = h(y)] \leq c/p.$$

(Kdybychom volili náhodně z úplně všech funkcí, vyšla by tato pravděpodobnost právě $1/p$ – c -universální systém je tedy nejvýše c -krát horší.)

Lemma: Buď h funkce náhodně vybraná z nějakého c -universálního systému. Nechť x_1, \dots, x_n jsou navzájem různé prvky universa vložené do struktury a x je nějaký prvek universa. Potom pro očekávaný počet prvků v téže přihrádce jako x platí:

$$\mathbb{E}[\#\{x : h(x) = h(x_i)\}] \leq cn/p.$$

Důkaz: Pro dané x definujeme indikátorové náhodné proměnné:

$$Z_i = \begin{cases} 1 & \text{když } h(x) = h(x_i) \\ 0 & \text{jinak} \end{cases}$$

Jinými slovy, Z_i říká, kolikrát padl prvek x_i do přihrádky $h(x)$, což je buď 0 nebo 1. Proto $Z = \sum_i Z_i$ a díky linearitě střední hodnoty je hledaná hodnota $\mathbb{E}[Z]$ rovna $\sum_i \mathbb{E}[Z_i]$. Přitom $\mathbb{E}[Z_i] = \Pr[Z_i = 1] \leq c/p$ podle definice c -universálního systému. Takže $\mathbb{E}[Z] \leq cn/p$. \heartsuit

FIXME: Doplnit přehashování.

Zbývá dořešit, kde nějaký c -universální systém sehnat. Známých konstrukcí je vícero, zde si předvedeme jednu lineárně algebraickou.

Lemma: Předpokládejme, že p je prvočíslo, přihrádky jsou identifikované prvky konečného tělesa \mathbb{Z}_p a universum U je vektorový prostor dimenze d nad tělesem \mathbb{Z}_p , tedy \mathbb{Z}_p^d . Uvažujme systém funkcí $S = \{h_t \mid t \in \mathbb{Z}_p^d\}$, kde $h_t(x) := t \cdot x$ (skalární součin s vektorem s). Pak tento systém je 1-universální.

Důkaz: Necht $x, y \in \mathbb{Z}_p^d$, $x \neq y$. Potom jistě existuje i , pro nějž $x_i \neq y_i$; bez újmy na obecnosti předpokládáme, že $i = d$. Nyní volíme t náhodně po složkách a počítáme pravděpodobnost kolize (rovnost modulo p značíme \equiv):

$$\begin{aligned} \Pr_{t \in \mathbb{Z}_p^d}[h_t(x) \equiv h_t(y)] &= \Pr[x \cdot t \equiv y \cdot t] = \Pr[(x - y)t \equiv 0] = \\ &= \Pr \left[\sum_{i=1}^d (x_i - y_i)t_i \equiv 0 \right] = \Pr \left[(x_d - y_d)t_d \equiv - \sum_{i=1}^{d-1} (x_i - y_i)t_i \right]. \end{aligned}$$

Pokud už jsme t_1, \dots, t_{d-1} zvolili a nyní náhodně volíme t_d , nastane kolize pro právě jednu volbu (poslední výraz je lineární rovnice tvaru $ax = b$ pro nenulové a a ta má v libovolném tělese právě jedno řešení). Pravděpodobnost kolize je tedy nejvýše $1/p$, jak požaduje 1-universalita. \heartsuit

Věta (Bertandův postulát): Pro libovolné $n \geq 1$ existuje prvočíslo p , které splňuje nerovnost $n < p \leq 2n$.