



# Univerzálne triedy hashovacích funkcií

Veronika Krajčová, Martin Kvočka, Marek Trgiňa

20.5.2010

## 1 Úvod

Rôzne vstupy pre program môžeme vnímať ako prvky z triedy problémov, pričom výstupom programu je správne riešenie daného problému. Ak hovoríme o priemernom výkone programu, spravidla priemerujeme jeho výsledky cez triedu problémov, ktoré program rieši. Páni Gill, Rabin, Strassen a Solovay na niektorých triedach problémov použili iný prístup. Navrhli, aby si program náhodne vybral algoritmus z triedy algoritmov, ktorým bude problém riešiť. Týmto spôsobom boli schopný ohraničiť priemerný výkon triedy algoritmov pre najhorší možný vstup. Tento priemer na najhoršom vstupe môže byť lepší než výkon ľubovoľného konkrétneho algoritmu na tomto vstupe. Tento prístup prekonáva niektoré problémy, ako:

1. Klasická analýza (priemerovanie cez triedu vstupov) musí počítať s distribúciou vstupov. Tieto predpoklady však v určitých prípadoch nemusia platiť.
2. Dôsledkom (1) je že nemôžeme klasicky analyzovať priemerný výkon podprocedúry nezávisle na hlavnej časti programu, pretože hlavná časť môže rozhodnúť distribúciu dát.
3. Ak programu predložíme najhorší možný vstup, tak neexistuje spôsob ako sa vyhnúť jeho slabej výkonnosti. Avšak, pri použití algoritmu z triedy algoritmov by sme mohli detekovať, že je pomalý na danom vstupe, a zvoliť iný algoritmus.

V nasledujúcom texte ukážeme použitie tohto prístupu pri hashovaní. Ukážeme že ak je trieda funkcií zvolená správne, potom priemerný výkon programu na ľubovoľnom vstupe bude aspon taký dobrý ako jednotlivá funkcia, zvolená s vedomím konkrétneho vstupu. Taktiež predvedieme niektoré triedy hashovacích funkcií ktoré zaručia, že každá vzorka z priestoru vstupov bude rovnomerne distribuovaná dostatočným počtom funkcií tak aby kompenzovala slabý výkon algoritmu v prípade nešťastnej volby funkcie.

## 2 Univerzálne triedy hash funkcií

Značenie:

Všetky hashovacie funkcie budú mapovať množinu  $A$  do  $B$ , pričom vždy bude  $|A| > |B|$ .  $A$  budeme nazývať množinou kľúčov,  $B$  zas množinou indexov. Nech  $f$  je hash funkcia a  $\delta_f(x, y)$  je 1 ak  $x \neq y$  a  $f(x) = f(y)$ , inak 0. Teda  $\delta_f(x, y)$  je 1 ak  $x$  a  $y$  sú rôzne prvky  $A$  ktoré sa mapujú na rovnakú hodnotu pomocou funkcie  $f$ . Ak  $f, x, y$  nahradíme množinou, potom sčítavame cez všetky prvky v množine. Teda, ak  $H$  je kolekcia hash funkcií,  $x \in A$  a  $S \subset A$ , potom  $\delta_H(x, S)$  znamená:

$$\sum_{f \in H} \sum_{y \in S} \delta_f(x, y)$$

Vlastnosti univerzálnych tried:

Nech  $H$  je trieda funkcií z  $A$  do  $B$ . O  $H$  vravíme že je  $universal_2$  ak pre všetky  $x, y$  z  $A$ ,  $\delta_H(x, y) \leq \frac{|H|}{|B|}$ . To znamená, že  $H$  je  $universal_2$  ak žiaden pár rôznych kľúčov nie je mapovaný na tie isté indexy viac než jedno  $|B|$ -tinou funkcií.

Tvrdenie 1:

Pre každú kolekciu hash funkcií (nie nutne  $universal_2$ ) existuje  $x, y \in A$  také že

$$\delta_H(x, y) > \frac{|H|}{|B|} - \frac{|H|}{|A|}$$

Tvrdenie 2:

Nech  $x$  je ľubovoľný prvok  $A$  a  $S$  podmnožina  $A$ . Nech  $f$  je funkcia zvolená náhodne z  $universal_2$  triedy funkcií (s rovnakou pravdepodobnosťou). Potom očakávaný počet prvkov z  $S$  s ktorými  $x$  koliduje, teda  $\delta_f(x, S)$ , je  $\leq \frac{|S|}{|B|}$ .

Zaujímá nás použitie týchto funkcií pri operáciach ukladania a získavania dát. Ak máme sekvenciu  $R$  požiadaviek (vloženie alebo výber) do nejakej databáze, a hash funkciu  $f$ , potom definujeme cenu  $R$  vzhľadom na  $f$  ako  $C(f, R)$ , teda ako sumu cien jednotlivých požiadaviek. Cena jednotlivej požiadavky zodpovedajúcej prvku  $x$  je jedna plus počet rôznych predtým vložených  $y$  pre ktoré  $f(x) = f(y)$ .

Táto funkcia ceny odráža najhoršiu cenu vloženia alebo nájdenia prvkov v úložisti a návratovú schému v ktorej každý prvok  $B$  je asociovaný so spojovaným zoznamom, a prvok  $x$  je uložený v tomto zozname asociovaný s  $f(x)$ .

Nasledujúci teorém dáva pekný odhad na očakávanú cenu pomocou  $universal_2$  tried hashovacích funkcií používajúcich spojovaný zoznam na riešenie kolízií.

Tvrdenie 3:

Nech  $R$  je sekvencia  $r$  požiadaviek ktoré zahŕňajú  $k$  vložení. Ďalej nech  $H$  je  $universal_2$  trieda hashovacích funkcií. Potom ak zvolíme  $f$  náhodne z  $H$ , očakávaná cena je:

$$C(f, R) \leq r(1 + \frac{k}{|B|})$$

Špeciálny prípad tohto tvrdenia je ak veľkosť  $k$  je zhruba rovnaká  $B$ , potom je očakávaná cena  $2r$ . Všimnite si, že toto lineárne ohraničenie platí pre ľubovoľnú sekvenciu požiadaviek, nie len pre "priemernú" požiadavku. Vo väčšine aplikácií však existuje horná hranica na počet prvkov ktoré je možno uložiť, teda  $B$  môže byť vhodne zvolené. Ak žiadna horná hranica neexistuje, môžeme jej veľkosť voliť dynamicky a znovu prehashovať ak sa voľba ukáže ako príliš malá. Rehashovanie totiž môžeme urobiť v lineárnom čase, niekedy aj real-time.

Môžeme ukázať, že očakávaná cena (spriemerovaná cez hashovacie funkcie) ľubovoľnej požiadavky je prakticky rovnaká ako očakávaná cena (spriemerovaná cez možné požiadavky) ľubovoľnej jednotlivej hash funkcie aplikovanej na nahodnú požiadavku po náhodných vlozeniach. Dôvod je nasledovný: Nech  $a = |A|$  a  $b = |B|$ . Argument z tvrdenia 1 implikuje že ak je  $f$  ľubovoľná hash funkcia a  $x$  a  $y$  sú zvolené náhodne z  $A$ , tak pravdepodobnosť  $\delta_f(x, y)$  je  $\geq (1/b - 1/a)$ . Potom teda pravdepodobnosť  $\delta_f(x, S)$  je  $\geq |S|(1/b - 1/a)$ , kde  $S$  je náhodná podmnožina  $A$  ktorá bola predtým uložená. Teda cena požiadavky je aspoň  $1 + |S|(1/b - 1/a)$ .

Taktiež je možné ohraničiť pravdepodobnosť že pre danú sekvenciu

požiadaviek  $R$  je výkon náhodne zvolenej funkcie horší než tolerovateľný na  $R$ . Keďže vieme že  $C(f, R)$  musí byť aspoň  $r$ , môžeme predpokladať že ak  $k$  je zhruba rovnaké ako  $|B|$ , pravdepodobnosť že  $C(f, R) > t - r$  je menej ako  $1/(t - 1)$ . Pre niektoré triedy hashovacích funkcií je možné vyvodit' odhad na štandardnú odchýlku alebo vyššie momenty ceny náhodne zvolenej funkcie na konkrétnom  $R$ . Toto nám umožní získať oveľa lepší odhad toho že pravdepodobnosť  $C(f, R)$  bude veľká.

### 3 Niektoré *universal*<sub>2</sub> triedy

Prvá trieda *universal*<sub>2</sub> hashovacích funkcií ktorú uvedieme je vhodná na použitie v situáciách, keď bitové reťazce reprezentujúce kľúče je možné jednoducho násobiť počítačom. Zoberme  $A = 0, 1, \dots, a - 1$  a  $B = 0, 1, \dots, b - 1$ . Nech  $p$  je prvočíslo také že  $p \geq a$ . Nech  $g$  je nejaká funkcia zo  $Z_p$  do  $B$ , ktorá čo najpresnejšie mapuje rovnaký počet prvkov  $Z_p$  na každý prvok z  $B$ . Formálne požadujeme  $|\{y \in Z_p | g(y) = z\}| \leq \lceil \frac{p}{q} \rceil$  pre všetky  $z \in B$ . Prirodzená voľba pre  $g$  je zvyšok modulo  $b$ . Ak  $b = 2^k$  pre nejaké  $k$ , toto predstavuje posledných  $k$  bitov v binárnej reprezentácii  $y$ .

Nech  $m$  a  $n$  sú prvky  $Z_p$  kde  $m \neq 0$ . Definujeme  $h_{m,n} : A \rightarrow Z_p$  ako  $h_{m,n}(x) = (mx + n) \bmod p$ . Potom definujeme  $f_{m,n}(x) = g(h_{m,n}(x))$ . Trieda  $H$  je množina  $\{f_{m,n} | m, n \in Z_p, m \neq 0\}$ . Takáto trieda  $H$  je *universal*<sub>2</sub>.

Algoritmy sú často analyzované za predpokladu že násobenie zaberie jednotku času a počet násobení býva považovaný za cenu algoritmu. Tento model je vhodný v prípadoch, keď neexistujú operácie ktoré by bolo možné vykonať neohraničený počet krát za každé násobenie. Keď je použitá trieda *universal*<sub>2</sub> hashovacích funkcií, počet odkazov do pamäte za každú požiadavku je možné ohraničiť spriemerovaním nad všetkými funkciami v triede ako v tvrdení (3). Takže tento model je vhodný a hashovacie funkcie v triede uvedenej vyššie je v ňom možné vykonať za koštantný čas. Takže v tomto modeli, pre ľubovoľnú sekvenciu požiadaviek má algoritmus lineárnu časovú zložitosť vzhľadom k počtu požiadaviek.

Trieda *universal*<sub>2</sub> popísaná vyššie nie je vhodná v prípadoch keď sú kľúče príliš dlhé na to, aby ich bolo možné násobiť použitím jednej inštrukcie. Nasledujúce tvrdenie dáva spôsob akým rozšíriť triedu funkcií pre dlhé kľúče.

Tvrdenie 4:

Predpokladajme že  $|B|$  je mocninou 2 a  $H$  je trieda funkcií z  $A$  do  $B$  taká, že pre každé  $i \in B$ :  $|\{f \in H \mid f(x) \oplus f(y) = i\}| = \frac{|H|}{|A|}$ , kde  $\oplus$  značí operáciu XOR. Potom môžeme definovať  $universal_2$  triedu hashovacích funkcií z  $A \times A$  do  $B$  takto. Pre  $f, g \in H$  definujeme  $H_{f,g}((x, y)) = f(x) \oplus g(y)$ . Potom táto nová trieda hashovacích funkcií  $J = \{h_{f,g} \mid f, g \in H\}$  je  $universal_2$  a zároveň spĺňa podmienky v tomto tvrdení. Toto tvrdenie sa úplne nevzťahuje na  $universal_2$  triedu funkcií, ktorá bola uvedená vyššie, a to preto že  $H$  nie je mocninou 2 a preto že počet funkcií pre ktoré platí  $f(x) \oplus f(y) = 0$ , i.e.  $\delta_H(x, y)$ , je v skutočnosti menší ako  $|H|/|B|$ .

Nasleduje trieda funkcií, ktoré nevyžadujú násobenie, čo môže byť výhoda v rôznych aplikáciách. Zoberme si  $A$  ako množinu  $i$ -ciferných čísel o základe  $\alpha$  a  $B$  ako množinu binárnych čísel dĺžky  $j$ . Potom  $|A| = \alpha^i |B| = 2^j$ . Nech  $M$  je trieda polí dĺžky  $i\alpha$ , ktorých prvky sú binárne reťazce dĺžky  $j$ . Pre  $m \in M$ , nech  $m(k)$  je binárny reťazec ktorý je  $k$ -ty prvok  $m$ , a pre  $x \in A$ , nech  $x_k$  je  $k$ -ta číslica  $x$  v bázi  $\alpha$ . Definujeme

$$f_m(x) = m(x_1 + 1) \oplus m(x_1 + x_2 + 2) + \dots + (m(\sum_{k=1}^i x_k + k))$$

Trieda  $H$  je množina  $\{f_m \mid m \in M\}$ . Táto trieda spĺňa vlastnosti  $universal_2$ . Pre dané  $B$ , každá funkcia v  $H$  má lineárnu časovú zložitosť vzhľadom k dĺžke kľúča.

## 4 Záver

Programátori často strávia dlhú dobu pri tom, aby odladili hashovacie funkcie pre aplikácie kde je kritické dosiahnutie uniformného rozloženia. Dosiahnutie tohto cieľa môže byť zložité, pretože je potrebné aby očakávaná vstupná množina bola usporiadaná takým spôsobom, že to nespôsobí nedostatočný výkon hashovacej funkcie. Jednou z praktických vlastností triedy  $universal_2$  funkcií je, že sa v triede nachádza mnoho použiteľných funkcií. Jednoduchým náhodným výberom hashovacej funkcie z tejto triedy je možné dosiahnuť uniformné rozloženie. A tým, že sa použitá funkcia mení pri každom vykonaní programu, je možné dosiahnuť dobrý výkon v priemernom prípade. Použitie  $universal_2$  tried dáva možnosť dobrého odhadu medzí, v ktorých sa bude pohybovať priemerný výkon algoritmu používajúceho hashovanie.

## Literatúra

- [1] DILIP V. SARWATE. *A Note on Universal Classes of Hash Functions.*. [cited 2010-5-14]. Available at: journal Information Processing Letters, vol. 10
- [2] GEORGE MARKOWSKY, J.LAWRENCE CARTER, MARK N. WEGMAN. *Analysis of universal class of hash functions* [online]. [cited 2010-5-14]. Available at: <http://laptops.maine.edu/HashFunctions.pdf>.
- [3] MARTIN DIETZFELBINGER AND FRIEDHELM MEYER. *A new universal class of hash functions and dynamic hashing in real time* . [cited 2010-5-14]. Available at: Lecture Notes in Computer Science, vol. 443/1990.