

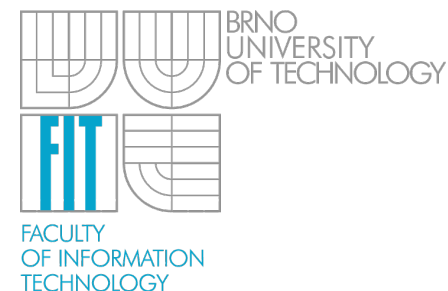
Bloom Filters Basics

Accelerated Network Technologies

Research Group

Martin Žádník

Brno University of Technology, Faculty of Information Technology
Bozetechova 2, 612 00 Brno, CZ
<http://merlin.fit.vutbr.cz/ant/>



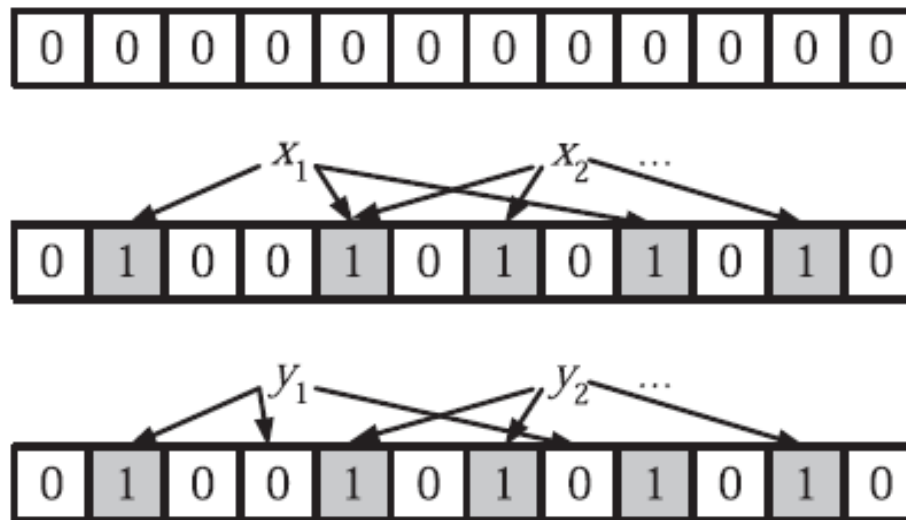
INVESTICE DO ROZVOJE VZDĚLÁVÁNÍ

Motivace

- **Struktura pro uložení informace o náležitosti prvku do množiny**
 - Šetření místem
 - Rychlý a jednoduchý přístup
 - Velikost paměti x přesnost
- **Není možné provést výčet prvků**
- **Obsahuje falešně pozitivní – tedy prvky, které nejsou součástí množiny**

Princip

- Vložení znamená nastavení několika bitů v poli
 - Hash udává, které bity
 - Pokud je bit nastaven na 1, pak 1 zůstává
 - Velikost paměti x přesnost



Vlastnosti

- Notace

- m ... počet bitů filtru
- n ... počet prvků množiny
- k ... počet hash funkcí

- Pravděpodobnost, že nějaký bit po vložení n prvků na k míst je stále nula

Pravděpodobnost nastavení nějakého bitu na 1

$$p' = \left(1 - \frac{1}{m}\right)^{kn} \approx e^{-kn/m}$$

Pravděpodobnost, že bit na 1 nebude nastaven

Vlastnosti

- Praviděpodobnost falešně pozitivních

- s rostoucím počtem nastavených 1 roste falešně pozitivní
- růst je exponenciální

$$f = \left(1 - e^{-kn/m}\right)^k = (1 - p)^k$$

Praviděpodobnost, že bit je stále nula

- m ... počet bitů filtru, n ... počet prvků množiny, k ... počet hash funkcí

Vlastnosti

- Minimalizace falešně pozitivních

- více hash funkcí
- menší počet nastavených jedniček

Sprostá úprava, aby to šlo derivovat

$$f = \left(1 - e^{-kn/m}\right)^k = (1 - p)^k$$

$$f = \exp\left(k \ln\left(1 - e^{-kn/m}\right)\right)$$

Vnitřek derivujeme

$$\frac{\partial g}{\partial k} = \ln\left(1 - e^{-\frac{kn}{m}}\right) + \frac{kn}{m} \frac{e^{-\frac{kn}{m}}}{1 - e^{-\frac{kn}{m}}}$$

- Derivace je nula, když $k = (m / n) * \ln 2 = 0.69 * (m / n)$

- m ... počet bitů filtru, n ... počet prvků množiny, k ... počet hash funkcí

Vlastnosti

- Efektivita, aneb počet bitů m pro dosažení false rate ϵ

Výčet všech filtrů délky m

$$2^m \binom{n + \epsilon(u - n)}{n} \geq \binom{u}{n}$$

Možné kombinace množin nad universem

Kombinace i s chybou

Redundance oproti prosté kombinaci

$$m \geq \log_2 \frac{\binom{u}{n}}{\binom{n + \epsilon(u - n)}{n}} \approx$$

$$n \log_2 e \cdot \log_2(1/\epsilon)$$

- Prostorová složitost $\sim 1.44 \cdot n$

- m ... počet bitů filtru, n ... počet prvků množiny, u ... počet prvků univer

Bloom vs. Hash

- Necht' $m = n * j$, pak falešně pozitivní je $(0.62)^j$
- Hash funkce a různá rozšíření Bloom filtru umožňují dosáhnout lepší hranice $(0.5)^j$
- Jednoduchá hash je vlastně Bloom filtr s jednou hash
- Operace nad Bloom filtry
 - spojení dvou filtrů – OR
 - průnik dvou filtrů – AND

Rozšíření

- Čítající Bloom filtr
 - Místo jednoho bitu, celý čítač
 - Je možné i odebrat
- Velikost čítače by měla postačovat 4 bity pro většinu případů, tj. pokud dimenzujeme správně velikost a hash

$$\mathbb{P}(c(i) = j) = \binom{nk}{j} \left(\frac{1}{m}\right)^j \left(1 - \frac{1}{m}\right)^{nk-j}$$

$$\mathbb{P}(\max_i c(i) \geq 16) \leq 1.37 \times 10^{-15} \times m$$

- Jaká je pravděpodobnost úspěšného odebrání?
 - m ... počet bitů filtru, n ... počet prvků množiny, k ... počet hash funkcí

Aplikace

- Slovníky
- Databáze
- Distribuované cache
- Distribuované uložistě, např. P2P
- Routing, detekce smyček, multicast
- Měření, např. heavy toky