

Hashing of strings (vectors)

• scalar product $h_t(x) = x \cdot t = \sum_{i=0}^{d-1} x_i \cdot t_i \pmod{p}$, $h_t: \mathbb{Z}_p^d \rightarrow \mathbb{Z}_p$
 $(x_0, \dots, x_{d-1}) \in \mathbb{Z}_p^d$, p prime, $d \geq 1$, $t \in \mathbb{Z}_p^d$

• $\mathcal{H} = \{h_t \mid t \in \mathbb{Z}_p^d\}$ is **1-universal** for any $d \geq 1$, any prime p . (lecture 7)

Note: Composing a universal family with mod m may lose universality.

Q: How can \mathcal{H} "gain" 2-independence? ↑ target # buckets
↙ (stronger than universality)

• linear congruence $h_{a,b}(x) = ((ax+b) \pmod{p}) \pmod{m}$
 $h_{a,b}: \mathbb{Z}_p \rightarrow [m]$, $a, b \in \mathbb{Z}_p$, $p \geq m$, p prime

• $\mathcal{L} = \{h_{a,b} \mid a, b \in \mathbb{Z}_p\}$ is $(2,1)$ -independent for any $p \geq m$,
 and **$(2,2)$ -independent** if $p \geq 4m$ where p is prime. (lecture 8)

lem: Let \mathcal{F} be c -universal family of $f: U \rightarrow [r]$, let \mathcal{G} be

$(2,d)$ -independent family of $g: [r] \rightarrow [m]$. Then the family

$\mathcal{H} = \mathcal{F} \circ \mathcal{G} = \{f \circ g \mid f \in \mathcal{F}, g \in \mathcal{G}\}$ is **$(2,c')$ -independent** where
 $(f \circ g)(x) = g(f(x))$ $c' = (\frac{m}{r} + 1)d$

general composition

Pf: $x_1, x_2 \in U$ distinct, $i_1, i_2 \in [m]$ (not a collision)

careful: $(2,d)$ -ind. of \mathcal{G} applies only if $f(x_1) \neq f(x_2)$!

"match event" $M: g(f(x_1)) = i_1 \wedge g(f(x_2)) = i_2$
 "collision event" $C: f(x_1) = f(x_2)$ 0 if $i_1 \neq i_2$

$$\Pr[M] = \underbrace{\Pr[M|C]}_{(2,d)\text{-ind.} \leq 1} \cdot \underbrace{\Pr[\neg C]}_{(1,d)\text{-ind.} \leq \frac{d}{m}} + \underbrace{\Pr[M|C]}_{c\text{-univ.}} \cdot \Pr[C] \leq \frac{d}{m^2} \cdot 1 + \frac{d}{m} \cdot \frac{c}{r} = \frac{c'}{m^2}$$

Note: if $\Pr[C] = 0$, then $\Pr[M] \leq \frac{d}{m^2}$ ($\Pr[C] = 1$ does not happen). \square

$\Rightarrow \mathcal{H} \circ \mathcal{L}$ is $(2,c')$ -ind. for any $p \geq m$, and **$(2,5/2)$ -ind.** if $p \geq 4m$.

Rolling hashing (another option for hashing of strings)

$$h_a(x) = \sum_{i=0}^{d-1} x_i \cdot a^i \pmod{p}, \quad h_a: \mathbb{Z}_p^d \rightarrow \mathbb{Z}_p, \quad a \in \mathbb{Z}_p, \quad d \geq 1, \quad p \text{ prime}$$

(x₀, ..., x_{d-1})

evaluating a polynomial with coefficients x_i at a, or scalar product with (a⁰, a¹, ..., a^{d-1})

Thm: $R = \{h_a \mid a \in \mathbb{Z}_p\}$ is **d-universal** for any $d \geq 1$, p prime.

Pf: $x, y \in \mathbb{Z}_p^d$ distinct

there are at most $d-1$ roots

$$P_a[h_a(x) = h_a(y)] = P_a\left[\sum_{i=0}^{d-1} (x_i - y_i) a^i = 0\right] \leq \frac{d}{p} \quad \square$$

a is a root of a (nonzero) polynomial of degree at most $d-1$

by Lem.

$\Rightarrow R \circ \mathcal{L}$ is **(2, 5/2)**-independent if $p \geq 4dm$.

R versus \mathcal{L} : d-univ. / 1-univ. | single parameter / d parameters, rolling

Rolling

motivation: Robin-Karp substring search - compare hashes in $O(1)$ time instead of entire substrings

... , x₀, x₁, ..., x_{d-1}, x_d, ...

input file (string)

search window x x'

x substring of length d, x' next

$$h_a(x) = x_0 + x_1 a + \dots + x_{d-1} a^{d-1}$$

$$h_a(x') = x_1 + \dots + x_{d-1} a^{d-2} + x_d a^{d-1} \pmod{p}$$

$$h_a(x') = \frac{h_a(x) - x_0}{a} + x_d a^{d-1}$$

if a^{d-1} precomputed in $O(1)$ time

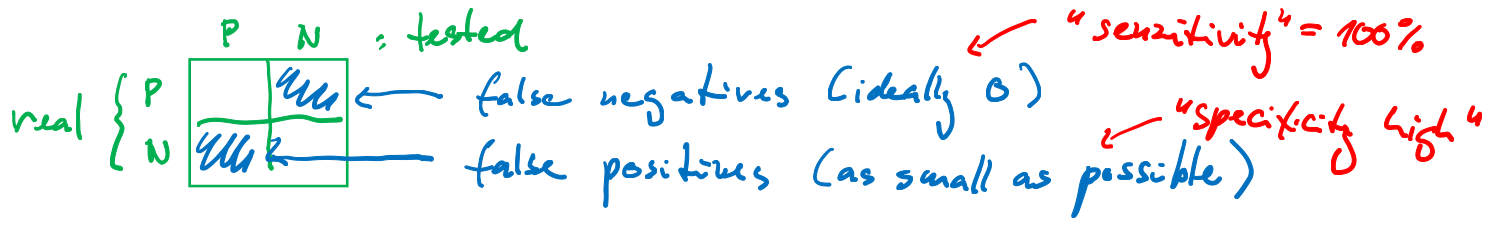
Remarks:

• alphabet size $k \leq p$ (prime), if $k \ll p$ we may encode multiple characters into one: $x_0 | x_1 | x_2$ $x_0 + kx_1 + k^2x_2 < p$ (e.g. for $p \geq 4dm$)

• if strings have variable lengths: use padding with a blank char to a fixed (maximal) length \uparrow encode by 0

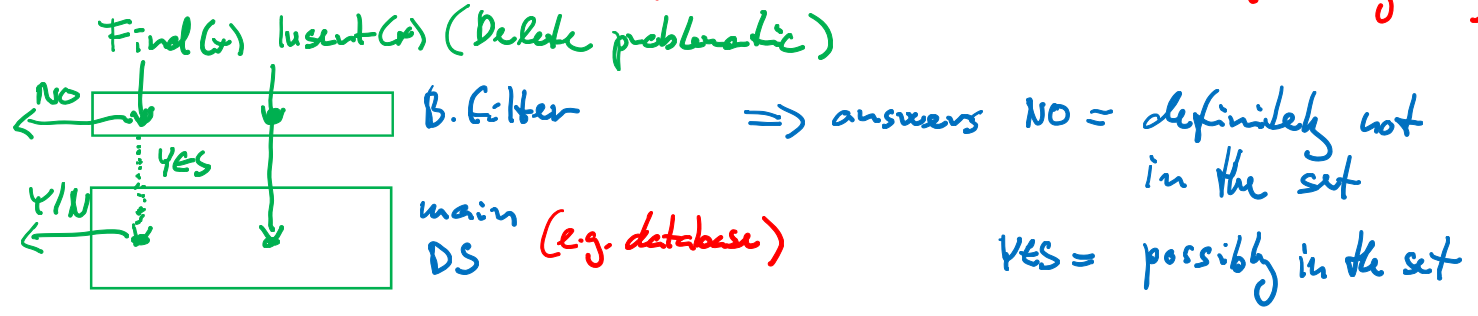
Bloom filters (Application of hashing)

filter: (fast and efficient) test with a one-sided error (ideally)



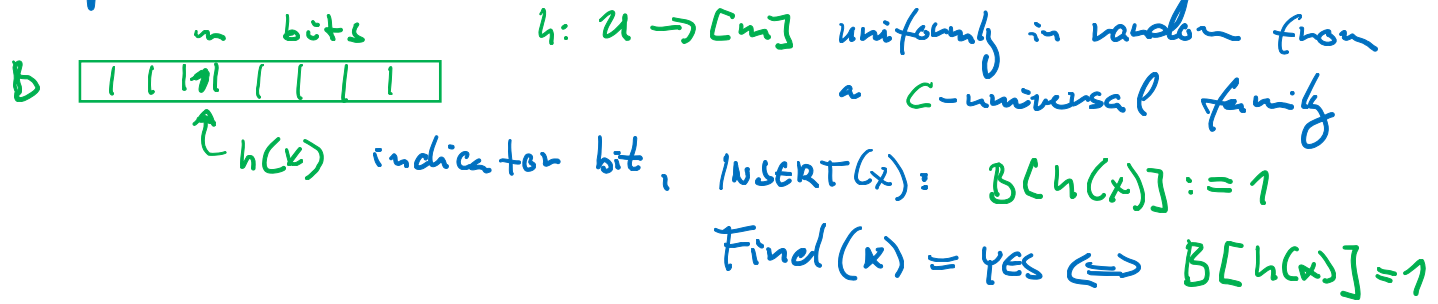
⇒ such test filters out negatives

Bloom filter (approximate representation of sets with 0 false negatives)



aim: memory efficient and small probability of false positives (trade-off)

Simple B. filter



⇒ supports (only) Find, insert, has **no false negatives**

• What is the probability of false positives?

x_1, \dots, x_n already inserted, y distinct (true negative)

$$Pr_h[y \text{ is FP}] = Pr[\exists i: h(y) = h(x_i)] \leq \sum_{i=1}^n Pr[h(y) = h(x_i)] \leq \frac{cn}{m}$$

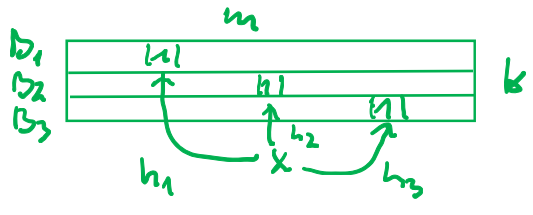
↑ union bound ↑ C-univ.

⇒ for $Pr[FP] \leq \epsilon$ (given target) we need $m = \lceil n/\epsilon \rceil$ bits (say $c=1$)

ex: for $n=10^6$, $\epsilon=0.01$ we have $m=100$ Mb. (linearly dep. on $1/\epsilon$)

multi-band filter (conjunction of simple filters)

k "bands" per n bits, $h_i: U \rightarrow [m]$ chosen **independently** ($i=1..k$) from a c -universal family (say $c=1$)



INSERT(x) = $B_i[h_i(x)] := 1 \quad \forall i=1..k$
FIND(x) = YES $\Leftrightarrow B_i[h_i(x)] = 1 \quad \forall i$

$\Rightarrow O(k)$ time, no false negatives

Probability of false positives: x_1, \dots, x_n inserted, y distinct (true negative)

$$Pr[FP] = Pr[\exists i \exists j h_i(y) = h_i(x_j)] = \prod_{i=1}^k Pr[\exists j h_i(y) = h_i(x_j)] \leq \prod_{i=1}^k \frac{cn}{m} = \frac{1}{2^k}$$

h_1, \dots, h_k independent

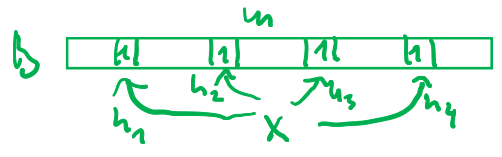
\Rightarrow for $Pr[FP] \leq \epsilon$ we can take $m = 2n$, $k = \lceil \log(1/\epsilon) \rceil$,
uses $M = 2n \lceil \log(1/\epsilon) \rceil$ bits (logarithmically in $1/\epsilon$)

ex: for $n = 10^6$, $\epsilon = 0.01$ we have $k=7 \Rightarrow M = 14 Mb$
 $\epsilon = 0.001$ $k=10$ $20 Mb$

Note: It is known that any DS with no false positives and $Pr[FP] \leq \epsilon$ needs at least $n \log(1/\epsilon)$ bits. (Optimal up to factor 2)

Furthermore, if h_i are **totally random** (and independent) then taking $m = n / \ln 2 \approx 1.44n$, $k = \lceil \log(1/\epsilon) \rceil$ also works.

single-band filter (compressed multiband)



$h_i: U \rightarrow [m]$ chosen independently ($i=1..k$)
(assume they are totally random)

INSERT(x) = $B[h_i(x)] := 1 \quad \forall i=1..k$
FIND(x) = YES $\Leftrightarrow B[h_i(x)] = 1 \quad \forall i$

$\Rightarrow O(k)$ time, no false negatives

x_1, \dots, x_n inserted \Rightarrow at most kn bits set to 1

$p = \Pr[B[i]=0] = (1 - 1/m)^{nk} \approx e^{-nk/m} \Rightarrow k = -m/n \cdot \ln p$
 (fixed bit) \uparrow (totally random independent hashes) $\leftarrow 1+x \approx e^x$ (Taylor's expansion)

δ distinct from x_1, \dots, x_n (true negative)

$\Pr[FP] = \Pr[\exists i: B[h_i(\delta)] = 1] \approx (1-p)^k = (1-p)^{-m/n \ln p} = e^{-m/n \ln p \ln(1-p)}$
 ("independent") \uparrow minimized for $p = 1/2$

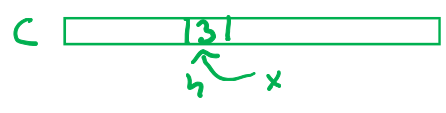
\Rightarrow for $\Pr[FP] \leq \epsilon$ set $k = \lceil \log(1/\epsilon) \rceil$
 and $m = n / \ln 2 \cdot k = 1.44 n \lceil \log(1/\epsilon) \rceil$

(looks like random bits)

Counting filters (to implement Delete)

idea: since shared bits cannot be deleted, replace bits with counters

for simple filter: INSERT/DELETE(x): increment/decrement $C[h(x)]$



FIND(x) = YES $\Leftrightarrow C[h(x)] > 0$

problem: counters are limited, b bits $\Rightarrow 2^b - 1 =: t$ max value!
 if $C[i] = t$, the counter is "stuck" (no inc/dec anymore)
 \Rightarrow new type of false positives!

However, the probability of having some counter stuck is small (assuming a totally random hashing function for simple filter):

$\Pr_h[C[i] \geq t] \leq \binom{n}{t} \left(\frac{1}{m}\right)^t \leq \left(\frac{ne}{mt}\right)^t = \left(\frac{e \ln 2}{t}\right)^t$ for $m = n / \ln 2 \approx 1.44n$.
 if we had unlimited \uparrow here is a t -tuple hashed to i , rest anywhere \leftarrow using $\binom{n}{t} \leq \left(\frac{ne}{t}\right)^t$

ex: $b=4 \Rightarrow t=15, \Pr(C[i] \geq t) \leq 3.06 \cdot 10^{-14}$ "4-bits suffice"
 $m = 10^9 \Rightarrow \Pr[\exists i: \Pr(C[i] \geq t)] \leq 3.06 \cdot 10^{-5}$

Note: if stuck counters accumulated \Rightarrow rebuild