

10. k-independent hash functions

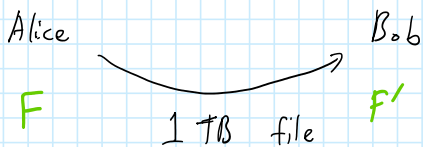
Thursday, September 30, 2021 2:02 AM

Last time: Intro to universal hash functions

Today: k-independent hash families

Application to signatures/checksums

Problem:



How to check for correct receipt?

One solution: hashing with universal hash.

Check if $\text{hash}(F) = \text{hash}(F')$.

Problem: Assign a unique signature $s(x) \forall x \in S$, $|S| = n$.
i.e. want $s(x) \neq s(y) \forall x, y \in S, x \neq y$.

Choose universal hash function $s: U \rightarrow [n^3]$. \leftarrow how big the range?

$$\text{Prob}_s [\exists \{x, y\} \subseteq S : s(x) = s(y)] \leq \sum_{\substack{\{x, y\} \subseteq S \\ \text{union bound}}} \text{Prob}_s [s(x) = s(y)] \leq \binom{n}{2} \frac{1}{n^3} < \frac{1}{2n}.$$

\uparrow birthday problem

Application to sampling from sets

Let $A \subseteq U$. Can we determine some statistic on A without studying all of it?

Recall loglog cardinality estimation

e.g. $|A|$ or polling

Let's use sampling: Let $h: U \rightarrow [m]$ and let $t \ll m$.

Let's sample $x \in U$ if $h(x) < t$. Let $S_{h,t}(A) = \{x \in A \mid h(x) < t\}$.

Then $E[|S_{h,t}(A)|] = |A| \cdot \frac{m}{t}$ Right?

So we can estimate $|A|$.

WRONG Because we didn't actually include uniformity in the definition of universal!

Strong universality (a.k.a. 2-independence)

A random hash function $h: [u] \rightarrow [m]$ is strongly universal if

A random hash function $h: [u] \rightarrow [m]$ is **strongly universal** if

$$\text{Prob}[h(i_1) = j_1 \wedge h(i_2) = j_2] = \frac{1}{m^2}.$$

Lemma: Strong universality implies universality.

Observation 3.1 Strong universality is equivalent to the statement that each key is hashed uniformly into $[m]$ and that every two distinct keys are hashed independently.

proof. Forward case. Let $h: [u] \rightarrow [m]$ be strongly universal. Let $x \neq y \in [u]$. Clearly $\forall q \in [m], \text{Prob}[h(x) = q] = \sum_{r \in [m]} \text{Pr}[h(x) = q \wedge h(y) = r] = \frac{m}{m^2} = \frac{1}{m}$, so uniformity holds.

Furthermore, $\text{Pr}[h(x) = q \mid h(y) = r] = \frac{\text{Pr}[h(x) = q \wedge h(y) = r]}{\text{Pr}[h(y) = r]} = \frac{\frac{1}{m^2}}{\frac{1}{m}} = \frac{1}{m} = \text{Pr}[h(x) = q]$,
↑
 independence.

Backward case: If $h(x)$ and $h(y)$ are independent and uniform,
 $\text{Pr}[h(x) = q \wedge h(y) = r] = \text{Pr}[h(x) = q] \cdot \text{Pr}[h(y) = r] = \frac{1}{m^2}$. □

Also called **2-independence** to focus on the independence of two events, rather than on collision probabilities.

k-independence

Def. \mathcal{H} is a k -wise independent hash family if

$$\forall i_1 \neq i_2 \neq \dots \neq i_k \in [u] \text{ and } \forall j_1, \dots, j_k \in [m],$$

$$\text{Prob}_{h \in \mathcal{H}}(h(i_1) = j_1 \wedge \dots \wedge h(i_k) = j_k) = \frac{1}{m^k}.$$

i.e. independence of k -different events.

Ex. The set \mathcal{H} of ALL functions $[u] \rightarrow [m]$ is k -wise ind. $\forall k$.

$|\mathcal{H}| = m^u$, so $h \in \mathcal{H}$ is representable in $u \lg m$ bits.

↑ This is just picking an ideal hash function which maps everything i.i.d. uniformly.

Ex. Let $u = m = q$, where q is a prime power. ↙ Galois field

Ex. Let $u = m = q$, where q is a prime power ↙ Galois field

Let $\mathcal{P}_{\text{poly-}k}$ be the set of all deg $k-1$ polynomials in $\mathbb{F}_q[x]$.

Claim: $\mathcal{P}_{\text{poly-}k}$ is a k -wise independent family.

proof. Lagrange interpolation. If we know i_1, \dots, i_k distinct and j_1, \dots, j_k , then

$$p(x) = \sum_{r=1}^k \left(\frac{\prod_{y \in [k] \setminus \{r\}} (x - i_y)}{\prod_{y \in [k] \setminus \{r\}} (i_r - i_y)} \right) \cdot j_r \quad \text{satisfies } p(i_r) = j_r \quad \forall r.$$

i.e.

$$p(x) = \frac{(x-i_2)(x-i_3)\dots(x-i_k)}{(i_1-i_2)(i_1-i_3)\dots(i_1-i_k)} j_1 + \frac{(x-i_1)(x-i_3)\dots(x-i_k)}{(i_2-i_1)(i_2-i_3)\dots(i_2-i_k)} j_2 + \dots + \frac{(x-i_1)\dots(x-i_{k-1})}{(i_k-i_1)\dots(i_k-i_{k-1})} j_k$$

When $p(i_2)$ "0 " j_2 "0

Notice the $p(x)$ is still of course a polynomial of degree $k-1$.

Furthermore, $p(x)$ as given above is the unique poly of degree $\leq k-1$ where $p(i_r) = j_r \quad \forall r$.

Why? Suppose some other poly $f(x)$ has $\deg(f) \leq k-1$ and $f(i_r) = j_r \quad \forall r$.

Then $g(x) = p(x) - f(x)$ has deg at most $k-1$ and $g(i_r) = 0 \quad \forall r$.

Further $g(x) \neq 0$, since $p(x) \neq f(x)$.

But also $g(x) = c(x-i_1)(x-i_2)\dots(x-i_k)$, so g is of degree $k > k-1$, a contradiction. Thus, $p(x) = f(x)$.

Thus, $p(x) = \alpha_{k-1} x^{k-1} + \dots + \alpha_1 x + \alpha_0$, so it is determined by k elements of \mathbb{F}_q .

Thus $|\mathcal{P}_{\text{poly-}k}| = q^k$.

But as just shown, exactly one such polynomial goes through all of (i_r, j_r) ,

so $\text{Prob} (h(i_1) = j_1 \wedge \dots \wedge h(i_k) = j_k) = \frac{1}{q^k}$.



Aside: Each $h \in \mathcal{P}_{\text{poly-}k}$ is representable using $k \lg q$ bits.

Practical choices for q

e.g. $p = x^{64} + x^4 + x^3 + x + 1$

Practical choices for q

Could use \mathbb{F}_q , where $q = 2^{64}$?

e.g. $P = x^{64} + x^4 + x^3 + x + 1$

cannot be factored

Recall: $\mathbb{F}_{2^{64}} = \mathbb{F}_2[x] / (P)$, where P is an irreducible poly in $\mathbb{F}_2[x]$ of degree n .

$z \in \mathbb{F}_{2^{64}}$ can be written $a_{n-1}x^{n-1} + \dots + a_1x + a_0 \in \mathbb{F}_2[x]$.

$a_i \in \{0, 1\}$.

\Rightarrow can write as $a_{n-1} a_{n-2} \dots a_0$, a 64-bit binary number.

Addition:

$$\begin{array}{r} a_{n-1}x^{n-1} + \dots + a_1x + a_0 \\ + b_{n-1}x^{n-1} + \dots + b_1x + b_0 \\ \hline (a_{n-1} + b_{n-1})x^{n-1} + \dots + (a_0 + b_0) \end{array}$$

XOR

$$\begin{array}{r} a_{n-1} \dots a_0 \\ b_{n-1} \dots b_0 \\ \hline (a_{n-1} \oplus b_{n-1}) \dots (a_0 \oplus b_0) \end{array}$$

Multiplication requires Euclidean division by P , which can be expensive.

Prime fields: Let p be a large prime.

Then $\mathbb{F}_p \cong \mathbb{Z}/p\mathbb{Z}$, and multiplication just requires mod p .

Alternative choice:

Mersenne Primes are prime numbers of the form $2^n - 1$.

Recall: If n is composite, so is $2^n - 1$

proof.

$$\begin{aligned} 2^{ab} - 1 &= (2^a - 1)(1 + 2^a + 2^{2a} + \dots + 2^{(b-1)a}) \\ &= (2^b - 1)(1 + 2^b + 2^{2b} + \dots + 2^{(a-1)b}) \end{aligned}$$

So n must be prime if $2^n - 1$ is prime.

OEIS A000043: Mersenne exponents $n = 2, 3, 5, 7, 13, 17, 19, 31, 61, 89, 107, 127$

Often, we choose $p = 2^{31} - 1$, $2^{61} - 1$, or $2^{89} - 1$.

Why?

Claim: If $p = 2^q - 1$, p and q prime, then

$$x \equiv x \bmod 2^q + \lfloor x/2^q \rfloor \pmod{p}$$

lemma If p is prime, p and q prime, then

$$x \equiv x \bmod 2^q + \lfloor x/2^q \rfloor \pmod{p}$$

proof. Let $x = a2^q + b$, where $b < 2^q$
 ↑ upper bits of x ↙ lower q bits of x

$$\begin{aligned} x \bmod p &\equiv (a \bmod p) \underbrace{(2^q \bmod p)}_{\substack{2^q \bmod 2^q - 1 \\ = 1}} + (b \bmod p) \\ &\equiv (a + b) \bmod p \end{aligned}$$

But $a = \lfloor x/2^q \rfloor$ is the upper bits of x
 $b = x \bmod 2^q$ is the lower bits of x

⇒ $x \equiv x \bmod 2^q + \lfloor x/2^q \rfloor \pmod{p}$ □

$(x \gg q)$ $\lfloor x/2^q \rfloor$ is a bit-shift to the right by q .

$(x \& p)$ $\lfloor x \bmod 2^q \rfloor$ is a bit-mask operation by p .

Then $y = x \bmod p$ can be computed by $y \leftarrow (x + p) + (x \gg q)$;
 (if $(y \geq p)$, $y \leftarrow y - p$.)

which are all fast bit operations.

Universal hashing of variable-length strings

Consider $x_0 x_1 \dots x_d \in [u]^d$ (each $x_i \in [u]$)

Can we construct an almost universal hash family $[u]^d \rightarrow [q]$?

Let q be a prime power, and work over \mathbb{F}_q .

$$\text{Let } P_{x_0 \dots x_d}(\alpha) = \sum_{i=0}^d x_i \alpha^i$$

Let $h_a(x_0 \dots x_d) = P_{x_0 \dots x_d}(a)$, where $a \in \mathbb{F}_q$ uniformly drawn.

Claim:

If $y_0 \dots y_{d'}$ is some other string with $d' \leq d$, then

$$\text{Prob}_{a \in \mathbb{F}_q} [h_a(x_0 \dots x_d) = h_a(y_0 \dots y_{d'})] \leq \frac{d'}{q}$$

proof. $h_a(x_0 \dots x_d) = h_a(y_0 \dots y_{d'})$

$$\Rightarrow P_{x_0 \dots x_s}(a) - P_{y_0 \dots y_s}(a) = 0$$

degree d poly in $\mathbb{F}_q[x]$.

By fundamental thm of algebra, the poly $P_{x_0 \dots x_s} - P_{y_0 \dots y_s}$ has at most d distinct roots.

So the prob. a random $a \in \mathbb{F}_q$ is a root is at most $\frac{d}{q}$.

