

Hashing

Birthday Problem:

$$n \text{ people } 1 \cdot (1 - \frac{1}{365}) \cdot (1 - \frac{2}{365}) \dots (1 - \frac{n-1}{365})$$

when $n=23$, then prob $\approx 1/2$

Balls into Bins:

linearity of expectation: for random X and Y , $E[X+Y] = E[X] + E[Y]$

n balls into m bins, expected # of empty bins = $m(1 - \frac{1}{m})^n$

$P(\text{ball miss}) = 1 - \frac{1}{m}$ n balls missed.

$P(\text{all } n \text{ ball miss}) = (1 - \frac{1}{m})^n$ by lin. expectation: empty bins $m(1 - \frac{1}{m})^n$

of bins w/ 3 balls: By indep, $P(\text{bin has 3 specific balls}) = (\frac{1}{m})^3 (1 - \frac{1}{m})^{n-3}$

Consider the 1st bin:

$$P(3 \text{ balls go in}) = (\frac{1}{m})^3$$

$$P(n-3 \text{ balls miss}) = (1 - \frac{1}{m})^{n-3}$$

now choose 3 balls into bin: $\binom{n}{3}$ choices of 3 balls.

$$P(1 \text{st bin has 3 balls}) = \binom{n}{3} (\frac{1}{m})^3 (1 - \frac{1}{m})^{n-3}$$

$$w. \text{ lin expect: } m \binom{n}{3} (\frac{1}{m})^3 (1 - \frac{1}{m})^{n-3}$$

Fingerprinting: CHIRP!

① hash sets of $|P|$ consecutive char into a value by taking mod w/ respect to some prime p .

② Compare hashes to hash value of P to see if there's a match.

Ex: Let $|P|=5$ and $p=13$. Find fingerprint for 314152.

$$31415 \bmod 13 = 7$$

$$14152 \bmod 13 = 8$$

fingerprint: 7,8

Successive calc: $O(1)$

Optimize!

> Don't naively take mods

> Every time read char in a file, change at leftmost digit (a) removed and rightmost digit (b) is inserted.

> Can get new hash in $O(1)$:

$$N' = [10(N-10^{p-1}a) + b] \bmod p$$

↑ can be bad!

> collisions: string's hash to same value of p but $s \neq p$

> more accurate by: hash sets of $|P|$ by taking mod p_1, p_2, \dots (multiple)

> tradeoff: space-accuracy

Primality Testing (quickly check if # is prime)

Fermat's Little Thm:

if p is prime and $1 \leq a < p$, $a^{p-1} \equiv 1 \bmod p$

Ex: if $p=7$ and $a=3$, $3^{7-1} = 3^6 = 729 \equiv 1 \bmod 7$

$$a^p \equiv a \bmod p$$

Carichael Num: some composite n that are a-pseudoprime for all a .

Ex: 561, 1105, 1729, inf many

Fermat Test:

① For prime candidate n , pick $1 \leq a < n$

② calculate $a^{n-1} \bmod n$

③ if $a^{n-1} \equiv 1 \bmod n$, then n is a-pseudoprime. else n is composite.

Ex: $n=299$, $a=116$

$$116^{299-1} \equiv 116^{298} \equiv 1 \bmod 299$$

Thus 299 is 116-pseudoprime

Extended-Euler's Example

let $a=51$ and $b=20$

a	b	k	x	y
51	20	2	-9	23
20	11	1	5	-4
11	9	1	-4	5
9	2	4	1	-4
2	1	2	0	1
1	0		1	0

from bottom-up

$$k = \frac{a}{b}$$

$$y = \text{prev } x - \text{current } k - \text{prev } y$$

return $Cd, y, x - Ky$

thus

$$-9(51) + 23(20) = 1$$

always 1 since base case.

now Oliver knows $k_2 = (n, e)$ and he encodes msg m into ciphertext $c = m \bmod n$. Then Lisa calculates $c^d \bmod n$.

Bloom Filters:

trade correctness with Spacetime!

M bits in hash table, k hash functions f_1, \dots, f_k

if $y \in S$, then $f_1(y), \dots, f_k(y) \in S$, $f_1(y), \dots, f_k(y) \in S$

Ex: $m=7, k=2, S = \{x_1, x_2\}$

① insert x_1 :

$$x_1 = 0123456$$

$$010001100$$

② insert x_2 :

$$x_2 = 0123456$$

$$010001110$$

③ find x_3 :

false positive!

So what?

↑ $m \Rightarrow ?$ space \Rightarrow ↓ FP

↑ $k \Rightarrow ?$ # of 1's \Rightarrow ↓ FP

↑ computation time

Still can be good since?

Chances to reject x by hashing to 0.

$P(\text{false positive})$:

occurs when all k indices checked are 1s. $P(\text{bit}=0) = (1 - \frac{1}{m})^k$

$$P(k \text{ indices are 1}) = (1 - (1 - \frac{1}{m})^k)^k$$

Another example:

Partitioned Bloom Filter:

M total bits, k mini tables.

i th hash only hashes to mini table.

Runtime: insert/query $O(1)$, space $O(m)$

Space $O(m/k)$, 0 fpp.

Set Resemblance: are 2 documents similar to each other?

Resemblance: in 2 sets A and B

$$R(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

or Jaccard cost.

$R(A, B) = 0$ when disjoint, $= 1$ when same.

Random Permutations π : some random perm.

need "black box" to output permutation (i.e. $BB(i, x) = \pi(i, x)$)

$\pi_i(A)$ = set of elements from $BB(i, x)$ for every x in A .

Shingling: hash k consec. words/chars into

64-bit hash to get smaller set.

down $5m \Rightarrow$ Set resemb.

To compare sets, compare document sketches,

$$\{\min \pi_i(S_1), \dots, \min \pi_i(S_2)\}$$

Ex: $\pi_1(S_1) = \{7, 14, 12, 3\}$ and $\pi_2(S_2) = \{1, 7, 11, 14\}$

Sketch is min, so $\{3, 13\}$.

$O(n^2)$ naively (compare each A to B)

$O(n)$ w/ hashing

if set resemb. > 0 , then set resemb. algo gives correct estimate of resemb.

$R(A, B) = 0$, then $\frac{|A \cap B|}{|A \cup B|} = 0 \Rightarrow |A \cap B| = 0$

thus 2 disjoint sets, since disjoint, then

$\min(\pi_i(A)) \neq \min(\pi_i(B))$, $(\pi_i(A) \text{ and } \pi_i(B) \text{ don't share elem})$

no min \Rightarrow # match = 0.

perm = 0.

if set resemb. > 0 , then set resemb. algo gives correct estimate of resemb.

$R(A, B) = 0$, then $\frac{|A \cap B|}{|A \cup B|} = 0 \Rightarrow |A \cap B| = 0$

thus 2 disjoint sets, since disjoint, then

$\min(\pi_i(A)) \neq \min(\pi_i(B))$, $(\pi_i(A) \text{ and } \pi_i(B) \text{ don't share elem})$

no min \Rightarrow # match = 0.

perm = 0.

if set resemb. > 0 , then set resemb. algo gives correct estimate of resemb.

$R(A, B) = 0$, then $\frac{|A \cap B|}{|A \cup B|} = 0 \Rightarrow |A \cap B| = 0$

thus 2 disjoint sets, since disjoint, then

$\min(\pi_i(A)) \neq \min(\pi_i(B))$, $(\pi_i(A) \text{ and } \pi_i(B) \text{ don't share elem})$

no min \Rightarrow # match = 0.

perm = 0.

if set resemb. > 0 , then set resemb. algo gives correct estimate of resemb.

$R(A, B) = 0$, then $\frac{|A \cap B|}{|A \cup B|} = 0 \Rightarrow |A \cap B| = 0$

thus 2 disjoint sets, since disjoint, then

$\min(\pi_i(A)) \neq \min(\pi_i(B))$, $(\pi_i(A) \text{ and } \pi_i(B) \text{ don't share elem})$

no min \Rightarrow # match = 0.

perm = 0.

if set resemb. > 0 , then set resemb. algo gives correct estimate of resemb.

$R(A, B) = 0$, then $\frac{|A \cap B|}{|A \cup B|} = 0 \Rightarrow |A \cap B| = 0$

thus 2 disjoint sets, since disjoint, then

$\min(\pi_i(A)) \neq \min(\pi_i(B))$, $(\pi_i(A) \text{ and } \pi_i(B) \text{ don't share elem})$

no min \Rightarrow # match = 0.

perm = 0.

if set resemb. > 0 , then set resemb. algo gives correct estimate of resemb.

$R(A, B) = 0$, then $\frac{|A \cap B|}{|A \cup B|} = 0 \Rightarrow |A \cap B| = 0$

thus 2 disjoint sets, since disjoint, then

$\min(\pi_i(A)) \neq \min(\pi_i(B))$, $(\pi_i(A) \text{ and } \pi_i(B) \text{ don't share elem})$

no min \Rightarrow # match = 0.

perm = 0.

if set resemb. > 0 , then set resemb. algo gives correct estimate of resemb.

$R(A, B) = 0$, then $\frac{|A \cap B|}{|A \cup B|} = 0 \Rightarrow |A \cap B| = 0$

thus 2 disjoint sets, since disjoint, then

$\min(\pi_i(A)) \neq \min(\pi_i(B))$, $(\pi_i(A) \text{ and } \pi_i(B) \text{ don't share elem})$

no min \Rightarrow # match = 0.

perm = 0.

if set resemb. > 0 , then set resemb. algo gives correct estimate of resemb.

$R(A, B) = 0$, then $\frac{|A \cap B|}{|A \cup B|} = 0 \Rightarrow |A \cap B| = 0$

thus 2 disjoint sets, since disjoint, then

$\min(\pi_i(A)) \neq \min(\pi_i(B))$, $(\pi_i(A) \text{ and } \pi_i(B) \text{ don't share elem})$

no min \Rightarrow # match = 0.

perm = 0.

Hashing w/ Chaining

$O(1)$ insertion, $O(1 + \frac{n}{m})$ lookup

Store heads of linked list at every index

however, want $\frac{n}{m} \approx 1$, but don't know n in advance. So use:

Table Doubling

$n < m$ maintained as invariant, when $n=m$,

when $n=m$, reconstruct entire dataset with

Size $2m$ and rehash all n elem. to new table.

Expected hash = $O(1)$, Inverse L.L. $O(\frac{n}{m}) = O(1)$

Amortized inserts = $O(1)$ if: "big budget", "cheap"

(no table-dub = $O(1)$, insertion store time credit w/ $t-1 = O(n)$)

insert takes $2n$, finance with

4. $\frac{n}{2} = 2n$ from prev. $\frac{n}{2}$ insert.

Ex: $A = \{1, 2, 3, 4\}$ $A' = \{1, 2, 3, 4\} \Rightarrow A' = \{102, 33, 50, 7103\}$

$B = \{1, 2, 5, 6\}$ $B' = \{1, 2, 5, 6\} \Rightarrow B' = \{102, 33, 61, 3143\}$

Visually, shared are $\{1, 2\}$

now calc. perm. Hash

each elem w/ some hash

33=33, so return 1. if diff.

hash, wait return 0.

Plotting 1): set resemb.

if resemb. between 2 Doc $R(A, B) > 0$, then set resemb. algo gives correct estimate of resemb.

$R(A, B) = 0$, then $\frac{|A \cap B|}{|A \cup B|} = 0 \Rightarrow |A \cap B| = 0$

thus 2 disjoint sets, since disjoint, then

$\min(\pi_i(A)) \neq \min(\pi_i(B))$, $(\pi_i(A) \text{ and } \pi_i(B) \text{ don't share elem})$

no min \Rightarrow # match = 0.

perm = 0.

if set resemb. > 0 , then set resemb. algo gives correct estimate of resemb.

$R(A, B) = 0$, then $\frac{|A \cap B|}{|A \cup B|} = 0 \Rightarrow |A \cap B| = 0$

thus 2 disjoint sets, since disjoint, then

$\min(\pi_i(A)) \neq \min(\pi_i(B))$, $(\pi_i(A) \text{ and } \pi_i(B) \text{ don't share elem})$

no min \Rightarrow # match = 0.

perm = 0.

if set resemb. > 0 , then set resemb. algo gives correct estimate of resemb.

$R(A, B) = 0$, then $\frac{|A \cap B|}{|A \cup B|} = 0 \Rightarrow |A \cap B| = 0$

thus 2 disjoint sets, since disjoint, then

$\min(\pi_i(A)) \neq \min(\pi_i(B))$, $(\pi_i(A) \text{ and } \pi_i(B) \text{ don't share elem})$

no min \Rightarrow # match = 0.

perm = 0.

if set resemb. > 0 , then set resemb. algo gives correct estimate of resemb.

$R(A, B) = 0$, then $\frac{|A \cap B|}{|A \cup B|} = 0 \Rightarrow |A \cap B| = 0$

thus 2 disjoint sets, since disjoint, then

$\min(\pi_i(A)) \neq \min(\pi_i(B))$, $(\pi_i(A) \text{ and } \pi_i(B) \text{ don't share elem})$

no min \Rightarrow # match = 0.

perm = 0.

if set resemb. > 0 , then set resemb. algo gives correct estimate of resemb.

$R(A, B) = 0$, then $\frac{|A \cap B|}{|A \cup B|} = 0 \Rightarrow |A \cap B| = 0$

thus 2 disjoint sets, since disjoint, then

$\min(\pi_i(A)) \neq \min(\pi_i(B))$, $(\pi_i(A) \text{ and } \pi_i(B) \text{ don't share elem})$

no min \Rightarrow # match = 0.

perm = 0.

if set resemb. > 0 , then set resemb. algo gives correct estimate of resemb.

$R(A, B) = 0$, then $\frac{|A \cap B|}{|A \cup B|} = 0 \Rightarrow |A \cap B| = 0$

thus 2 disjoint sets, since disjoint, then

2-SAT

ex $(\bar{a} \vee b) \wedge (b \vee \bar{c}) \wedge (a \vee c)$

1. implication graph

→ each 2-sat makes 2 implications

ex $\bar{a} \vee b$ if $a = T, b = T (a \rightarrow b)$
if $\bar{b} = T, \bar{a} = T (\bar{b} \rightarrow \bar{a})$

→ construct graph w/ v: each literal and neg

ex $\bar{a} \vee b \quad a \rightarrow b \quad \bar{a} \leftarrow \bar{b}$

→ formula satisfiable if and only if graph doesn't have literal and its negation on s.c.c.

P vs. NP

P: set of all yes/no problems that can be deterministically solved in polynomial time.

NP: set of yes/no problems that can be verified w/ a polynomial length certificate.

P is a subset of NP since can be verified in polynomial by running it.

Problem is NP-complete if $X \in NP$ and ENP -hard
NP-hard if every problem $y \in NP$ reduces to x .

Proving

• in P: come up w/ poly algo and correctness (it and only it)

• in NP: state poly-length witness & how to verify

• NP-hard: reduce known NP-hard problem

1) show that reduction is polynomial

2) prove it's correct

• NP-complete: Both NP and NP-hard.

Heuristic for NP-Hard

Local search:

① Define a cost function f , let's say goal is to minimize f .

② Define neighborhood function $N(x)$ must be:

- easily computable from x
- $N(x)$ is polynomial in input size
- if $y \in N(x)$ then $x \in N(y)$

③ pick starting point x

④ while there is $y \in N(x)$ s.t. $f(y) < f(x)$, $x = y$
return sol

Variations:

hill climbing: generic local search

metropolis: pick random $y \in N(x)$. If $f(y) < f(x)$, $x = y$.
if $f(y) > f(x)$, $x = y$ w/ some prob.

Simu. Annealing: same as metropolis, prob of moving to higher cost neighborhood decreases over time.

Tabu Search: same as metropolis/SA, but w/ memory so won't get stuck.

Parallel Search: do more than 1 search, occasionally replace versions doing poorly w/ ones doing well.

Genetic Algo: keep pop. of searches that changes over time via "breeding" of best performing searches.

For & Fulkerson:

① Let residual graph start as original

② at each step from $s \rightarrow t$ in orig graph, increase flow.

- in resid. graph, take path, reverse flow, and add edge.

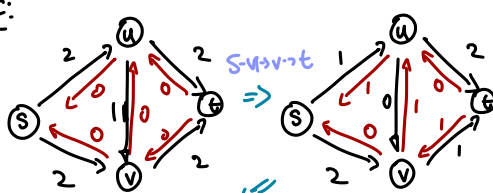
③ repeat until no paths in residual graph

$O(E \cdot f^*)$: every time augment path, increase flow

- augment path, at worst $O(f^*)$, f^* : value of max-fl.

- find aug path in $O(E)$ w/ DFS

F-F:



when flow backward, minus.

• BFS instead, $O(VE^2)$, called Edmonds Karp

Approximation Algor.

k-approx: cost of solution produced by approx algo c is within a factor of k of the cost of an optimal sol.

Examples:

> minimum vertex cover: Given graph $G(V, E)$, find min vertices $S \subseteq V$ s.t. $\forall (u, v) \in E$, have at least 1 of $u \in S$ or $v \in S$.

2-approx: choose edge at random, add both to cover, delete vertices (and edge) from graph and repeat.
2-approx since at worst, add 2 vert. per edge to cover. In opt, only 1 vertex per edge.

$$|OPT| \leq |S| \leq 2|OPT|$$

> max cut: $G(V, E)$, 2 disjoint subsets S and T s.t. $S \cap T = \emptyset$, $S \cup T = V$, and num of edges that crosses cut is maximized. Partition into 2 sets w/ max edges.
1/2 approx algo:

① randomly assign node to S or T s.t. each edge has 1/2 prob. of crossing the cut.

② Starts w/ all vertices on one side of cut, switch vertex if it increases num of edges crossing cut till size can't inc. local search algo.

> Euclidean Traveling Salesman: In set of points (x_i, y_i) in a Euclidean plane, find the tour of the min. length that travels thru all cities.

greedy: run DFS on MST and skip vertices that have already been visited (shortcutting).

2-factor since DFS at most visits edge 2x. Shortcutting reduces travel time due to triangle ineq. Opt. tour is at least size of MST since tour must be spanning tree (w/ last edge back to origin)

$$\text{tour}_{\text{greedy}} \leq 2T (\text{after shortcut}) \leq 2 \text{OPT}$$

can improve to 3/2 approx, or

Christofides-Serdyukov algo.

> max-SAT: most clauses that are satisfiable?

Flip random coin. It has k literals, clause is true w/ prob. 1/2, so prob. not satisf. is $1 - 2^{-k}$, where $k = \#$ literals.

Max Flow min cut:

- cut partitions into S and $V-S$.

- weight of cut = sum of all edges crossing
min cut \geq max flow

If I have cut, all units of flow must pass thru S and $V-S$. If weight of cut is c , then at most c units of flow can pass between s and t .
 $S-T$ = bottleneck.

min cut \leq max flow

At end of F-F, look at all nodes from S and all nodes that reaches t in residual graph. no path from s to t or else algo hasn't term. Some cut of graph and weight is max flow, so min can't be bigger.

intuition: max flow push across the smallest capacity cut.