

CS5234 Algorithms at Scale

$e \approx 2.718$
 $e^{-1} \approx 0.368$
 $e^{-2} \approx 0.135 < \frac{1}{6}$

- Fact: $e^{-2x} \leq 1-x \leq e^{-x}$
- Markov's ineq.: IF X is a nonnegative ran. var then $\forall k > 0: \Pr[X \geq k] \leq \frac{E[X]}{k}$
- Chebyshev ineq.: IF X is a ran. var then $\forall k > 0: \Pr[|X - E[X]| \geq k] \leq \frac{\text{Var}[X]}{k^2}$
- Variance: $\text{Var}[X] = E[(X - E[X])^2] = E[X^2] - E[X]^2$
- Where X, Y are indep.: $E[XY] = E[X] \cdot E[Y]$
 $\text{Var}[X+Y] = \text{Var}[X] + \text{Var}[Y]$

Hoeffding bound: IF X_1, \dots, X_n are indep. ran. vars s.t. $\forall i, X_i \in [a_i, b_i]$,
 then $\forall \delta > 0, \Pr[|Z - E[Z]| \geq \delta] \leq 2 \exp\left(-\frac{2\delta^2}{\sum_{i=1}^n (b_i - a_i)^2}\right)$ where $Z = \sum_{i=1}^n X_i$ and $E[Z] = \sum_{i=1}^n E[X_i]$
 (by linearity of exp.)
 Where $a_i = 0$ and $b_i = 1$, we get: $\forall \delta > 0, \Pr[|Z - E[Z]| \geq \delta] \leq 2 \exp\left(-\frac{2\delta^2}{n}\right)$

Adjlist format graph: get i^{th} neighbour of node u in constant time

- Connectivity problem on sparse graphs (n nodes, m edges, d max. degree) \rightarrow return true if connected
 \rightarrow return false if ϵ -far from connected (i.e. graph cannot be connected even if you can add/remove at most ϵnd entries in adjlist (note: each edge has two entries in the adjlist))
- Lemma: IF G is ϵ -far from connected, then it has $\epsilon nd/4$ connected components (note: when add/removing edges, must ensure max degree is satisfied)
- Lemma: IF G is ϵ -far from connected, then it has $\epsilon nd/8$ connected components of size $\leq 8/\epsilon d$
- Algo: repeat $16/\epsilon d$ times: pick random node and check if component size is less than $8/\epsilon d$.
 Takes $O(d \cdot (8/\epsilon d)) = O(1/\epsilon)$ per check, total time $O(16/\epsilon^2 d)$
 If G is ϵ -far from connected, each iteration has at least $\frac{\epsilon nd}{8}$ prob. of finding small component
 So $P(\text{algo returns true when } \epsilon\text{-far}) \leq \left(1 - \frac{\epsilon d}{8}\right)^{16/\epsilon d} \leq e^{-2} \leq \frac{1}{3}$

Multiplicative approximation: $C^*(1-\epsilon) \leq C \leq C^*(1+\epsilon)$
 correct answer estimate correct answer

Approx. connected components algo:
 $\text{sum} = 0$
 for $j = 1$ to s :
 $u \leftarrow \text{randVertex}()$
 if u has at least $\frac{2}{\epsilon}$ nodes reachable (BFS)
 $\text{sum} = \text{sum} + \frac{\epsilon}{2}$
 else (found $n(u)$ nodes)
 $\text{sum} = \text{sum} + \frac{1}{n(u)}$
 return $n \cdot (\text{sum}/s)$
Result: with probability $> \frac{2}{3}$, output is within $CC(G) \pm \epsilon n$
 takes $O\left(\frac{d}{\epsilon^3}\right)$ time
 $\leftarrow d: \text{max degree}$
Result: with probability $> 1 - \frac{1}{s}$, output is within $CC(G) \pm \epsilon n$
 takes $O\left(\frac{d \ln s}{\epsilon^3}\right)$ time

MST approx. algo: outputs a $(1 \pm \epsilon)$ -multiplicative approximation (for integer edge weights in $\{1, \dots, W\}$, assuming graph is connected)
 $\text{sum} = n - W$
 for $j = 1$ to $W-1$:
 $\text{sum} += \text{ApproxCC}(G_j, d, \epsilon', \delta)$
 return sum
 $\epsilon' := \frac{\epsilon}{W}, \delta := \frac{1}{3W}$
 graph containing only edges with weight $\leq j$
Result: with probability $> \frac{2}{3}$, output is within $MST(G) (1 \pm \epsilon)$
 takes $O\left(\frac{dW^4 \log W}{\epsilon^3}\right)$ time

Maximal matching approx algo:
 $\text{query}(e)$: $\rightarrow \text{query}(e)$ is supposed to return true iff e is part of the greedy maximal matching induced by the hashes
 for all neighbours e' of e :
 if $\text{hash}(e') < \text{hash}(e)$:
 if $\text{query}(e') = \text{true}$ return false
 return true
 Expected time complexity of $\text{query} \leq 2 \sum_{k=1}^{\infty} \frac{d^k}{k!} = O(e^d)$

Chernoff bounds: IF X_1, \dots, X_n are indep ran. vars s.t. $\forall i, X_i \in [0, s]$, and let $X = \sum_{i=1}^n X_i$ and $\mu = E[X]$,
 then: for any $0 \leq \delta \leq 1$: $\Pr[X \geq (1+\delta)\mu] \leq e^{-\frac{\mu \delta^2}{3}}$
 $\Pr[X \leq (1-\delta)\mu] \leq e^{-\frac{\mu \delta^2}{3}}$

Algo to approx. maximal matching size:
 $\text{sum} = 0$
 for $j = 1$ to s :
 choose edge uniformly at random.
 if $\text{query}(e)$ then $\text{sum} = \text{sum} + 1$
 return $m \cdot (\text{sum}/s)$

• Yao's minimax principle: Given a problem, let X be the set of inputs, $\Gamma(X)$ be the set of probability distributions (2) on X , D be the set of deterministic algo, R be the set of randomised algo:

$$\forall A \in R, \forall \gamma \in \Gamma(X) : \max_{x \in X} E_{\text{randomness of } A} [\text{cost}(A, x)] \geq \min_{B \in D} E_{x \sim \gamma} [\text{cost}(B, x)]$$

• To show that the expected cost of any randomised algorithm (on the worst case input) is $\geq T$, it suffices to show that there is an input distribution $\gamma \in \Gamma(X)$ such that

$$\text{for any deterministic algorithm } B \in D, E_{x \sim \gamma} [\text{cost}(B, x)] \geq T$$

• Streaming Algorithms:

• only have a small scratch space, but want to calculate some property of the items in the stream.
E.g.:

• return an approximation of the number of times x appears:

• $\text{count}(x) : N(x) - \epsilon m \leq \text{count}(x) \leq N(x) + \epsilon m$
↑ real count ↑ length of stream

• heavy hitters: return every item appearing $\geq 2\epsilon m$ times but no item appearing $< \epsilon m$ times

• Misra-Gries algorithm:

- set P of $\langle \text{item}, \text{count} \rangle$ pairs
- For each u in stream:
 - if $\langle u, c \rangle$ is in P , increment c .
 - else add $\langle u, 1 \rangle$ to P .
 - if $|P| > k$, decrement c for every $\langle v, c \rangle$ in P
 - remove from P all $\langle v, c \rangle$ where $c = 0$.

• $\text{count}(x)$: return c if $\langle x, c \rangle$ in P otherwise return 0.

• Heavy hitters: return x if $\text{count}(x) \geq \epsilon m$.

• Flajolet-Martin (FM) algorithm:

- $x = 1$
- for each u in the stream:
 - if $h(u) < x$ then $x = h(u)$
- return $\frac{1}{x} - 1$ hash

• FM+ algorithm:

- run a copies of FM, to get X_1, \dots, X_a these are the x from FM
- compute $Z = \frac{1}{a} \sum_{j=1}^a X_j$
- return $\frac{1}{Z} - 1$

• FM++ algorithm:

- run b copies of FM+, to get Y_1, \dots, Y_b $Y_j = \frac{1}{Z_j} - 1$
- return the median of Y_i

• Result: $a = \frac{4}{\epsilon^2}$, $b = 36 \ln \frac{2}{\delta}$, then with probability at least $1 - \delta$, FM++ returns an answer in $\epsilon(1 \pm 4\epsilon)$.

• Streaming a graph: each edge is an element in the stream, and edges arrive in arbitrary order.

• UFDS-based:

- count connected components: $O(n \log n)$ space and $O(\alpha(n, n))$ update cost since $\log n$ bits to store each vertex
- check if graph is bipartite: $O(n \log n)$ space and $O(\alpha(n, n))$ update cost inverse Ackerman function

Shortest path approx.

Find a "spanner": a spanning subgraph $H \subseteq G$ s.t. H is sparse (i.e. not too many edges) and for all $u, v \in V(G)$: $d_G(u, v) \leq d_H(u, v) \leq \alpha d_G(u, v)$.

note: if every edge $uv \in E(G)$ satisfies $\frac{d_H(u, v)}{d_G(u, v)} \leq \alpha$,

then $\forall u, v \in V(G)$ (not necessarily an edge), $\frac{d_H(u, v)}{d_G(u, v)} \leq \alpha$

↑ "stretch"

Algo:

- For each edge uv in stream:
- if $d_H(u, v) > 2k - 1$ then:
- add uv to H
- return H

Thms:

The girth of H is $> 2k$
length of smallest cycle

If $\text{girth}(H) > 2k$ then H has $O(n^{1+\frac{1}{k}})$ edges

- If we pick $k=2$, then $\text{space} = O(n^{\frac{3}{2}} \log n)$
- If we pick $k = \log n$, then $\text{space} = O(n^{1+\frac{1}{\log n}} \log n) = O(n \log n)$

Matching approx:

- Do greedy matching - pick an edge if both vertices are still not matched
- It is a 2-approximation

Weighted matching: Graph edges have weights, want to find max weight matching

Algo:

- M : matching, initially empty.
- for each edge uv in stream:
- let C be the set of edges in M that are incident on u or v
- if $w(uv) > (1+\delta)w(C)$ then:
- remove C from M
- add uv to M

Result: it is a b -approximation of optimal.

Clustering:

- k-centre clustering: choose k points (centres) that minimise the maximum distance to a centre
- k-median clustering: choose k points that minimise the average distance to a centre

$D(P, C) := \sum_{i=1}^n |p_i - c(i)|$ where $P = \langle p_1, \dots, p_n \rangle$ are the points and $C = \langle c_1, \dots, c_k \rangle \subseteq P$ and c is a function mapping a point to a centre

C is an (α, δ) -approx: $|C| \leq \alpha |C^*|$ and $D(P, C) \leq \gamma D(P, C^*)$

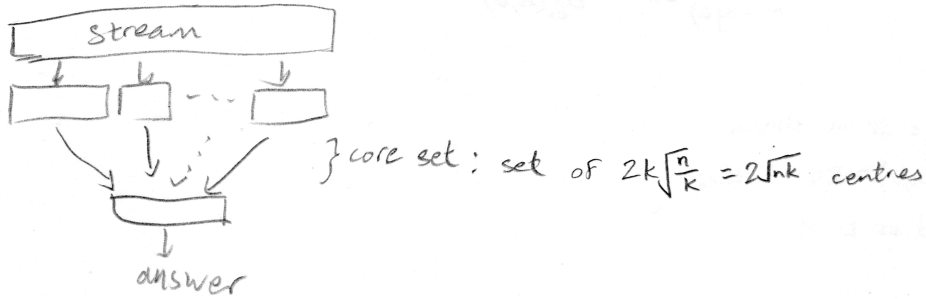
ILP solution:
(NP-hard)

$y_j := \begin{cases} 1 & \text{if } p_j \text{ is a centre} \\ 0 & \text{otherwise} \end{cases}$
 $x_{ij} := \begin{cases} 1 & \text{if } p_i \text{ is assigned to centre } p_j \\ 0 & \text{otherwise} \end{cases}$
 minimise $\sum_{i,j} x_{ij} d(p_i, p_j)$ where $\forall i: \sum_j x_{ij} = 1$ $\forall j: x_{ij}, y_j \in \{0, 1\}$
 $\sum_j y_j \leq k$
 $\forall j: x_{ij} \leq y_j$

Solve the LP version, then do some rounding.

Core-Set algorithm for streaming k-median:

- $C = \emptyset$
- repeat $\sqrt{\frac{n}{k}}$ times:
 - Let $P =$ next \sqrt{nk} points
 - Find (2,4)-approx clustering on P
 - Add $2k$ new cluster centres to C , but weight each cluster centre with the number of points attached to it
- return (2,4)-approx (weighted) clustering on C .

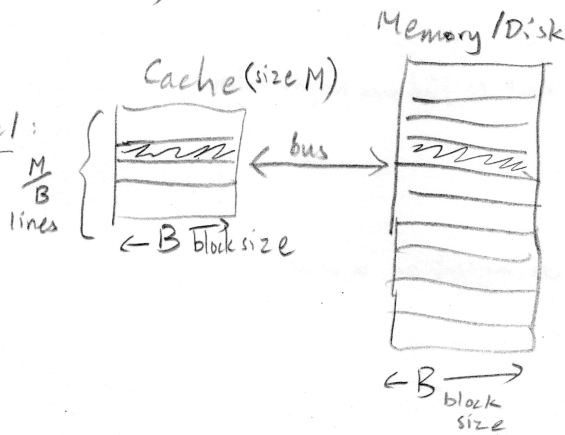


Space: $O(\sqrt{nk})$
 output: (2,80)-approx of k-median
 proof: generally by Δ -ineq.

- can stack more layers if we want a further reduction in space:
 - let $m = n^\epsilon$ (m elements before grouping to the next level)
 - num levels $= \log_m n = \frac{1}{\epsilon}$
 - space $= \frac{2kn^\epsilon}{\epsilon}$
 - approx factor $= O(8^{\frac{1}{\epsilon}})$

Caching

External Memory Model:



- entire cache line gets copied over when we want to access anything on it
- want to minimise the number of times a cache line gets transferred

Examples:

Scanning data:

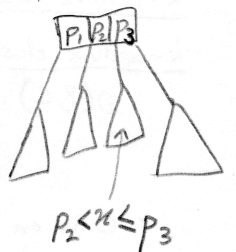
- Linked list: $O(N)$
- Array: $O(N/B)$

Searching data:

- Linked list: $O(N)$
- Red-black tree: $O(\log N)$
- (sorted) Array: $O(\log \frac{N}{B})$
- B-tree: $O(\frac{\log N}{\log B})$

(a,b)-trees:

- tree structure
- satisfies search property
- $b \geq 2a$
- all keys stored in leaves
- internal nodes store pivots to guide search
- root has ≥ 2 children
- non-root nodes have $\geq a$ children
- all nodes have $\leq b$ children
- all leaves have same depth



Properties:

- height of tree $\leq \log_a(\frac{n}{a}) + 1 \in O(\log_B n)$ where $a, b \in O(B)$
- insert: insert into correct leaf, then split from bottom up if $> b$ keys
- delete: remove from correct leaf, then either merge or rebalance siblings from bottom up if $< a$ keys

Sorting data:

- B-tree: $O(N \log_B N)$
- Buffer tree: $O(\frac{N}{B} \log_{\frac{M}{B}} \frac{N}{B})$

amortized cost (w/ parent ptr)

- per node: $O(1)$
- per operation: $O(\log_B n)$

amortized cost:

- per node: $O(\frac{1}{B})$
- per operation: $O(\frac{1}{B} \log_B n)$

Buffer tree: (for fast searching and very fast insertion/deletion)

- Build a (2,4)-tree, but add a buffer of size $2B$ to every node.
- For each leaf: ensure that it has between B and $5B$ keys (inclusive).
- insert:
 - add $ins[key]$ to root buffer
 - Clean buffer: remove any $del[key]$ or duplicate $ins[key]$
 - If $|buffer| \geq B$, flush the buffer
- delete: similar to insert $\rightarrow O(1) + \text{buffer flush} = O(\frac{1}{B} \log n)$
- search: walk from root to leaf, remember to search buffer too: $O(\log n)$
- flush:
 - sort buffer
 - move operations to children's buffers
 - Clean children's buffers
 - recursively flush children's buffers if necessary.

$O(1) + \text{buffer flush} = O(\frac{1}{B} \log n)$

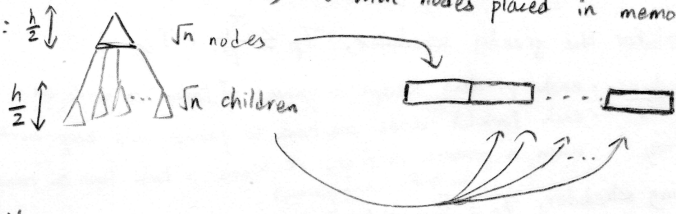
branching factor is 3, not B.
 each item contributes $\Theta(\frac{1}{B})$ to flush operation
 $ins[key]/del[key]$

\sqrt{B} Buffer tree: (each node has \sqrt{B} children)

- insert/delete: $O(\log_B n)$
- search: $O(\frac{\log_B n}{\sqrt{B}})$ (tree depth: $2 \log_B n$)

van Emde Boas search tree: static cache-oblivious search tree

- Like a normal balanced binary search tree, but with nodes placed in memory in a special way.
- recursively layout: $\frac{h}{2}$



search: $O(\log_B n)$

Cache-efficient graph algorithms:

- Breadth-first search: Layer by layer: For each layer: Graph is stored as adjlist format.

- ① $L_{i+1} \leftarrow \text{neighbours of nodes in } L_i$
 $O(|L_i| + \text{edges}(L_i)/B)$
 - ② Sort L_{i+1} $O(\text{sort}(L_{i+1}))$
 - ③ Remove duplicates in L_{i+1}
 $O(\text{edges}(L_i)/B)$
 - ④ $L_{i+1} \leftarrow L_{i+1} \setminus L_i$ $O(|L_i|/B + \text{edges}(L_i)/B)$
 - ⑤ $L_{i+1} \leftarrow L_{i+1} \setminus L_{i-1}$ $O(|L_{i-1}|/B + \text{edges}(L_i)/B)$
- Total cost = $O(|V| + |E|/B + \text{sort}(|E|))$

Count connected components: edges are stored in a single array.

- some recursive algorithm like UFDS
- $O(\text{sort}(E) \log(E))$

MST: $D \frac{E}{M}$ on edge weight:

- Divide E into E_1 and E_2 .
- Recursively find MST T_1 of E_1 .
- Contract E_1 .
- Recursively find MST T_2 of E_2 .
- Expand E_2 .
- Return $T_1 \cup T_2$.
- $O(\text{sort}(E) \log(\frac{E}{M}))$

Parallel Algorithms

- PRAM model:
 - p processors
 - shared memory
 - processors take one step at each clock tick
 - each processor can be programmed separately

Eg. AllZero(A, l, n, p):

```

for i = (n/p) * (j-1) + 1 to (n/p) * j:
    if A[i] != 0 then answer = false
done = done + 1
wait until done = p
return answer
    
```

Fork/Join:

```

[1, ..., n]
E.g. Sum(A, b, e):
    if b = e:
        return A[b]
    mid = (b+e)/2
    fork:
        L ← Sum(A, b, mid)
    join:
        R ← Sum(A, mid+1, e)
    sync
    return L+R
    
```

relies on a good scheduler

- Parallel merge sort:
 - work = $O(n \log n)$
 - span = $O(n)$
 - naive merge
 - span = $O(\log^3 n)$
 - binary search
 - parallel recursive merge is $O(\log^2 n)$ per merge.

- Work: total steps done on all processors: T_1
- Span: longest path in the program: T_∞
- Parallelism: $\frac{T_1}{T_\infty}$ (\approx number of processors that we can use productively)
- On p processors: want: $T_p \approx \frac{T_1}{p} + T_\infty$
 - parallel part
 - sequential part

- Greedy scheduler:
 - If $\leq p$ tasks are ready, execute all of them
 - If $> p$ tasks are ready, execute any p of them
- Brent-Graham thm: For the greedy scheduler, $T_p \leq \frac{T_1}{p} + T_\infty$
- Work-stealing scheduler:
 - each process keeps a queue of tasks to work on
 - each fork() adds one task to queue, and keeps working
 - when a process is free, it steals a task from a random queue.
- Thm: For work-stealing scheduler, $T_p \leq \frac{T_1}{p} + O(T_\infty)$ (no tasks in own queue)

Parallel operations on a balanced binary search tree:

- insert/delete/divide: $T_1 = T_\infty = O(\log n)$
- union/subtraction/intersection/difference: $T_1 = O(n+m)$, $T_\infty = O(\log n + \log m)$
- set symmetric difference

E.g. Union(T_1, T_2):

```

if T1 = null return T2
if T2 = null return T1
key ← root(T1)
(L, R, x) ← split(T2, key)
fork:
    TL ← Union(key.left, L)
    TR ← Union(key.right, R)
sync.
T ← join(TL, TR)
insert(T, key)
return T
    
```

- Basic building blocks operations for BST:
- split(T, k) $\rightarrow (T_1, T_2, x)$ (x could be null if exists, x is the item at k if exists)
 - join(T_1, T_2) $\rightarrow T$ (assumes $\forall x_1 \in T_1, \forall x_2 \in T_2, x_1 < x_2$)
 - root(T) $\rightarrow x$ (get the root item of T , leaving T unchanged)
 - insert(T, x) $\rightarrow T'$

Parallel BFS using BST for storage:

- $F \in \{S\}$
- $D \in \{S\}$
- while $F \neq \emptyset$:
 - $D \leftarrow \text{Union}(D, F)$ Work: $O(m \log n)$, Span: $O(\log^2 m)$
 - $F \leftarrow \text{ProcessFrontier}(F)$ Work: $O(m \log^2 n)$, Span: $O(\log^3 m)$
 - $F \leftarrow \text{SetSubtract}(F, D)$ Work: $O(m \log n)$, Span: $O(\log^2 m)$
- recursive divide/process/union.
- $T_1 = O(m \log^2 n)$, $T_\infty = D \log^3 m$ (diameter of graph)

Map-reduce model:

- separate memory
- loosely synchronized
- data exchanged over fast interconnect
- Data: $\langle \text{key}, \text{value} \rangle$ pairs on distributed shared disk/filesystem
- Round:
 - map(key, value) $\rightarrow (\text{key}, \text{value})$
 - shuffle (group items by key)
 - reduce($\text{key}, [\text{values} \dots]$) $\rightarrow (\text{key}, \text{value})$

Metric: how many rounds needed? (best: $O(1)$)

can optionally produce any number of pairs, but try not to have input/output of more than $O(n^E)$ pairs, or use more than $O(n^E)$ time/memory. each $(\text{key}, \text{value})$ pair should be $O(\text{polylog}(n))$.

- Nice properties to have:
 - Associative reducer: scheduler can perform reduce of the same key on multiple threads
 - Certain (sentinel) pairs come first: e.g. select from multiple datasets (Stream) (assuming small max degree)

E.g. Bellman ford with n iterations