

THE ALGORITHM DESIGN MANUAL

2ⁿ
d

Steven S. Skiena

ch1 - ch18

17 pages

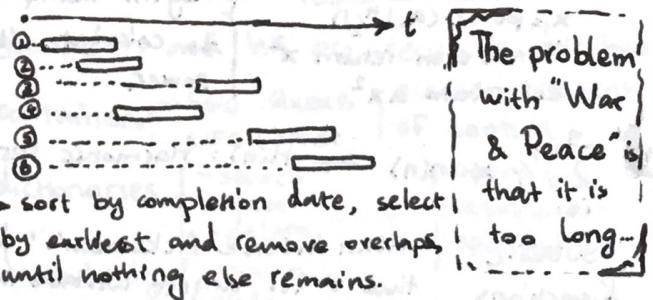
2016,11,10 - 2016,11,13

(ch)

algorithm: a procedure that takes an arbitrary instance of the problem and returns an instance of desired answers.

In the industry, usually a good algorithm is what gets the job done.

*traveling salesman (robot arm) is difficult.



sort by completion date, select that it is by earliest and remove overlaps until nothing else remains.

■ = QED: quod erat demonstrandum.

expressing algorithms

- Human language (e.g. English)
- Programming language (e.g. C/C++/Java)
- Pseudo Code

Proving → think small: target specific shortcomings

Incorrectness → think exhaustively

via Counter → hunt for the weakness

Examples → go for a tie: all inputs are the same

seek extremes

RAM

- simple operations are constant time
- memory access is constant time
- MODEL - loops et al are not simple.

► also talks about Big-O: see [CLRS(4)]

$$n! \gg 2^n \gg n^3 \gg n^2 \gg n \cdot \log n \gg n \gg \sqrt{n} \gg \log n \gg 1$$

POWER (a, n):

if $n=0$ return 1

$x = \text{power}(a, \lfloor n/2 \rfloor)$

if n is even return x^2

else return $a \cdot x^2$

Uses only
 $\log(n)$ multiplications
to calculate the power.

$$\sum_{i=1}^n \frac{1}{i} \approx \ln(n) \rightarrow H(n): \text{Harmonic series}$$

esoteric functions

$$\begin{aligned} & \leftarrow \alpha(n): \text{inverse Ackermann's function} \\ & \text{than } 5 \text{ for largest writable integer} \\ & = 2^{Cn^{1+\epsilon/C}} \approx O(n^{1+\epsilon}) \end{aligned}$$

$\rightarrow f(n)$ dominates $g(n)$ iff. $\lim_{n \rightarrow \infty} \frac{g(n)}{f(n)} = 0$

- abstract data types define the contract
- ↳ replacing implementations does not affect the correctness of the program

Data Structures ① contiguous: single slab of memory.
② linked: need to use pointers.

日本 houses in Japan are numbered in the order in which they were built!

ch2

- arrays:
 - constant access time
 - space efficiency
 - locality
 - fixed size → dynamic array $2N$ effort
- cell phones are pointers to their owners as they move about
- linked lists
 - don't overflow unless memory is full.
 - simpler insertion & deletion.
 - easier to rebase items.

Both arrays and lists are recursive data types

- containers
 - FIFO: Queue
 - LIFO: Stack
- operations are irrelevant of content
- dictionaries
 - search
 - insert
 - delete
- operations require a specific key
- optionally
 - max/min
 - predecessor/successor

► talks about binary search trees: see [CLRS(18)]

► greedy heuristics tend to grab the best possible thing first.

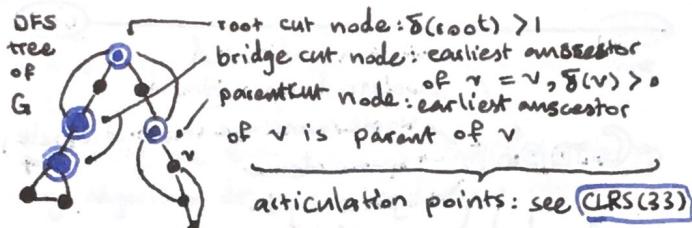
► talks about hashing: see [CLRS(14)]

► talks about Rabin-Karp: see [CLRS(49)]

- uses of hashing
 - duplicate detection
 - anti-plagiarism: hash sentences
 - pre-transmitting the hash (checksum)

- ch4
- A quarter of all mainframe cycles were spent sorting ← Knuth
 - ! sorting the data is one of the first things any algorithm designer should think about.
 - we should abstract comparison into a function
 - SELECTION SORT: see CLRS(4)
 - INSERTION SORT: see CLRS(3)
 - HEAPSORT: see CLRS(9) ↗ doesn't need random access → can use linked lists
 - MERGESORT: see CLRS(8) ↗ if we use linked-lists, we can merge without a buffer by rearranging pointers
 - QUICKSORT: see CLRS(10)
 - BUCKET SORT: see CLRS(12)
 - ↳ great when distribution of data is uniform across some feature of data
 - k-way merge sort with a min-heap (PQ) for managing the k-way merge is great for sorting large files (PQ sits in memory, while sorted blocks go to disk).
 - binary search and problem space halving techniques lie at the heart of divide & conquer solutions

- ch5
- ## Graphs
- directed/undirected
 - weighted/unweighted
 - simple/non-simple (multi-edge, self-loop)
 - sparse/dense
 - cyclic/acyclic
 - embedded/topological
 - implicit/explicit
 - labeled/unlabeled
- representation
- most of the time this is ↗ list of edges the right answer ↗ initialization of the data structure can become a bottleneck!
- talks about BFS: see CLRS(32)
 - ↳ we can use BFS to discover a two-coloring of graph $G \Rightarrow$ if exists, G is bipartite
- talks about DFS: see CLRS(32) → if we store exploration agenda in a stack instead of a queue we have a DFS.
 - we mark the start & finish "time" of exploration for a node, which sandwiches the exploration of its descendants.
 - $\frac{\Delta t}{2} = \# \text{ of descendant nodes, edges, back always point to an ancestor}$
- 5 -



- talks about topological sorting: see CLRS(32)

(ch6)

- Talks about MINIMUM SPANNING TREES: see CLRS(34)
 - { Prim: $O(m + n \lg n)$ using PQ
Kruskal: $O(m \cdot n)$: $O(m \log m)$ using UNION FIND
- Another name for DISJOINT SETS: see CLRS(31)
- Talks about Dijkstra: see CLRS(35)
- Dijkstra for vertex edges: set weight (i,j) to $w(j)$
- Talks about FLOYD-WARSHALL: see CLRS(36)
- Talks about NETWORK FLOW & BIPARTITE MATCHING
see CLRS(37-38)
- Try to model problems as graphs & apply known algorithms.

(ch7)

- 1 Million
 - permutations of 10-11 objects
 - subsets of combinations of 20 items
- employing BACKTRACKING is the same as doing a DFS of problem space, with vertices being states, and edges being transitions.

■ We can construct all subsets of set of size $|k|$ by promoting the states of presence & absence of item "i" to candidacy. We are done when we have a vector of size $|k|$.

■ PRUNING: stopping the search down a path as soon as we have realized this path will not yield an answer.

↳ Proper pruning can speed up algorithms more significantly than any other factor.

Heuristic Approaches

⇒ only consider parts of the solution space



■ All heuristic approaches need to evaluate the "cost/score" of a solution to judge which one is better.

1. Monte Carlo method

- 1 RANDOM SAMPLING: generate possible solutions at random, and return the one that looks the best.
 - ↳ ① requires the ability to generate random numbers from a unified distribution.

→ Hill Climbing

2

LOCAL SEARCH is the process of picking one solution, then slightly tweaking it until no improvements are possible → great at finding the local optima
Not so great in a landscape full of hills!

3

SIMULATED ANNEALING Start at a given temperature. At each step consider "i" neighboring states. Transition if the state is better or a random function tells us to do so. Continue until no further change is found and lower the temperature before the next set of iterations.

We want to avoid getting stuck in local optima by allowing jumps to states that don't improve the solutions. As the temperature decreases, we are less likely to take such risks and will eventually settle on a local optima that is within a larger radius from the starting point than what local searching would allow.



4

GENETICALGORITHMS

model transitions as mutations and allow for evolution and natural selection to kill off harmful mutations.

Often very slow convergence and also extremely complicated to implement.

PARALLELIZATION does speed up computation but adds complexity to code, is hard to debug, and usually can be beaten by some improvements to the sequential solution.

→ Proper load balancing goes a long way in improving parallelization performance

1

DYNAMIC PROGRAMMING is essentially a trade-off of space for time.

→ naturally left-to-right problem spaces

- character strings
- rooted trees
- polygons
- integer sequences

calculating FIBONACCI series

ch8



- ① Come up with the optimizing recursion
DP ② construct bottom-up solution in the right-order

- edit distance

Sample Problems - longest increasing sequence

the cost of each partitioning scheme is the largest sum among its parts. We seek to minimize this.

can S_{ij} be parsed as rule $X \rightarrow YZ$? $M[i:j, X] = \min_{(X \rightarrow YZ) \in G} \left(\sum_{i=k}^{j-1} M[i:k, Y] + M[k:j, Z] \right)$

minimum # of parse errors when parsing S_{ij} to $X \rightarrow YZ$



- triangulating a polygon so that triangles have minimum total circumference.

- When does it work?
- the recurrence has to be correct
 - the number of remembered partial solutions is small
 - each partial solution is not too difficult to calculate

ch9
► By showing that we can solve A by converting it to B which runs in $O(f)$ in $O(g)$ conversion time, we reduce A to B.

⇒ ① We have $A \in \text{P}$

② if $A \in \text{P}(h) \Rightarrow B \in \text{P}(h-g)$, otherwise the reduction would contradict $A \in \text{P}(h)$

→ this is what we use to show hardness.

► If we convert an optimization to a decision of feasibility of k , we can do binary search to answer the optimization.

\hookrightarrow $\text{W(D)} \stackrel{\text{use}}{\approx} \text{TSP}(G, n)$
Hamiltonian cycle $\stackrel{\text{reduce}}{\Rightarrow}$ Traveling salesman

is hard therefore is hard

This is the general formulation for NP proofs.

→ We don't even need to know how the target problem works. We want to know if it is difficult or not.

- Cook's theorem: all problems in NP are the same in terms of difficulty.
- NP-hard: not verifiable in polynomial time

- ① Do I really understand the problem?
 ↗ reiterate the problem
 ↗ articulate all the knowns
- ② Can I find a simple heuristic for my problem?
 → how important is it to solve optimally?
- ③ Are there special cases of the problem that I know how to solve?
 → can I generalize my answer?
- ④ Which algorithmic paradigm is most relevant to my current problem?



The Catalog



Dictionaries

- how many items?
- operation distribution?
 - unsorted linked lists / arrays
 - sorted linked lists / arrays
 - hash tables
 - BSTs
 - B Trees
 - Skip lists

ch10



Priority Queues
 great when we know the range of keys

- additional operations?
- is maximum size known?
- can priorities change?
 - sorted array / list
 - binary heap
 - bounded height PQs
 - BSTs
 - Fibonacci heaps

Suffix Trees



Graph Data Structures

→ a graph that represents all the suffixes of a string

- how big? $|V|$
- how dense? $|E|$
- for which algorithms?
- will be modified?
 - matrices
 - adjacency lists
 - edge lists



Set Data Structures

- bit vectors
- containers / dictionaries
- bloom filters: use k hash functions and set bits $H_i(e)$ for inserting e
- collection of containers
- generalized bit vector: array where $A[i]$ has the name of subset
- dictionary with subset attribute
- the union-find data structure

-13-

Random Number Generation

- Should I have the same random sequence?
- Third-party implementations?
- Implement myself?
- Size of the number?
- Distribution?



Knapsack Problem

- NP-Complete
- Budget/Cost
- Dynamic Programming/Backtracking
- Integer partitioning is a special case of this problem.



SORTING

- range limited? bit vector, counting
- uniformly distributed? buckets
- sortedness is high?
- external? use inorder of a B-Tree, or use k-way merge sort.



Searching

- sequential: good up to 20
- binary: more than 100
- exploit key frequencies: optimal binary search tree using DP by minimizing expected search cost
- self-organizing lists: bump-to-top, splay trees: BST where k gets rotated to root
- one-sided binary search
- external memory? B-Tree, van Emde Boas



Selection → naive: $O(n \lg n)$, best: worst-case $O(n)$
→ if we see each element only once, we can obtain a random sample, & work on that
→ finding the mode is $SL(n \lg n) \leftarrow$ element uniqueness



Generating Permutations

- All: use backtracking
- Random: use the Knuth method
- Next:
 - ① Use rank/unrank where $\text{unrank}(\text{rank}(p), n) = p$
 - ② Use incremental swapping by next() / previous()



Subset generation

- use grey code
- use binary counting



→ undirected: there is a path between every two vertices (BFS / DFS)

Connected Components → directed@weak: for (u, v) there is either $(u \rightarrow v)$ or $(v \rightarrow u)$
② Strong: there is a path between every pair of edges

↳ also useful for finding weakest point of a graph

 → largest $S \subseteq V$ that forms a complete graph
Clique → find maximal by sorting according to \deg .

 → 2-coloring: disjoint covering bipartite components → DFS
Vertex Coloring → finding chromatic number of graph is NP-complete.

 → Vizing's theorem: if $\Delta \text{ max} = k$, we can edge-color with $\Delta + 1$.
Edge Coloring ↳ This offers an $O(\Delta|V||E|)$ algorithm for finding $(\Delta+1)$ -coloring of the graph

 Text Compression → Huffman codes: \approx two scans \approx not adaptive
→ Lempel-Ziv (including LZW): build index as we go and reuse frequent blocks

→ 1977: Lempel-Ziv introduced adaptive search with random restarts and dynamic memory: a form of self-adjusting compression algorithm

Programming Pearls

Jon Bentley

ed 1 - ed 16
2 pages
2nd, 3rd, 4th, 5th, 6th, 7th, 8th