

Lecture 7

Basics of Algorithms

- Aspects of Algorithms
- Important Algorithms

Contents

- Aspects of Algorithms

- Important Algorithms

Computational statistic and algorithms

- ▶ once the statistical problem has been formulated, we have to solve and problem and this relies on **computing**
- ▶ computing: find efficient **algorithms** to solve them

Algorithm

(loosely speaking) a method or a set of instructions for doing something... A program is a set of computer instructions that implement the algorithm.

Top 10 Algorithms of the 20th Century

1. 1946: The **Metropolis Algorithm for Monte Carlo**. Through the use of random processes, this algorithm offers an efficient way to stumble toward answers to problems that are too complicated to solve exactly.
2. 1947: Simplex Method for Linear Programming. An elegant solution to a common problem in planning and decision-making.
3. 1950: Krylov Subspace Iteration Method. A technique for rapidly solving the linear equations that abound in scientific computation.
4. 1951: The **Decompositional Approach to Matrix Computations**. A suite of techniques for numerical linear algebra.
5. 1957: The Fortran Optimizing Compiler. Turns high-level code into efficient computer-readable code.
6. 1959: QR Algorithm for Computing Eigenvalues. Another crucial matrix operation made swift and practical.
7. 1962: **Quicksort Algorithms for Sorting**. For the efficient handling of large databases.
8. 1965: **Fast Fourier Transform**. Perhaps the most ubiquitous algorithm in use today, it breaks down waveforms (like sound) into periodic components.
9. 1977: Integer Relation Detection. A fast method for spotting simple equations satisfied by collections of seemingly unrelated numbers.
10. 1987: Fast Multipole Method. A breakthrough in dealing with the complexity of n-body calculations, applied in problems ranging from celestial mechanics to protein folding.

<https://www.siam.org/pdf/news/637.pdf>

Aspects of computer algorithms

Two most important aspects of algorithms

- ▶ Accuracy: e.g., absolute error $|\hat{\theta} - \theta|$ or relative error $|\hat{\theta} - \theta|/|\theta|$
- ▶ Efficiency: how many computing **resources** are needed to achieve certain **precision**

Other aspects

- ▶ Robustness
- ▶ Stable

Algorithm accuracy

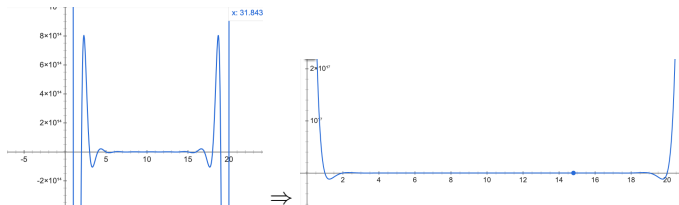
- ▶ Some algorithms are exact (e.g., multiplying two matrices)
- ▶ Other algorithms are approximate because the result to be computed does not have finite closed-form solution
 - ▶ Truncation error: solving $f(x) = 0$ by solving $\tilde{f}(x) = 0$, where $\tilde{f}(x)$ is the Taylor expansion of $f(x)$
 - ▶ Discretization error: solving continuous PDE via discretization $\partial f(x)/\partial x = g(x)$

Algorithm and data

- Performance of algorithm may depend on the data

Example [Wilkinson 59]

$f(x) = (x - 1)(x - 2) \cdots (x - 20)$ If the function is changed to $f(x) = (x - 1 - 2^{-23})(x - 2) \cdots (x - 20)$ the function roots $f(x) = 0$ may change drastically



“Condition of data”

Quantify the condition of a set of data for a particular set of operations

- ▶ finding root: increase in $1/\dot{f}(x)$ near the root (the previous example)
 $f(x) = (x - a_1)(x - a_2) \cdots (x - a_{20})$, the roots varies drastically when we vary (a_1, \dots, a_{20})
- ▶ condition number of a matrix M : $\kappa(M) = \lambda_{\max}(M)/\lambda_{\min}(M)$

Algorithm stability

- ▶ A small perturbation will not result in big difference in solution of a stable algorithm.
- ▶ Perturbation to input data may be due to truncation error
- ▶ *if problem is ill-conditioned, even stable algorithm may produce bad results.*

Example: Matrix completion [Shapiro, Xie, Zhang 2019]

Solving matrix completion problem (Netflix recommender systems)

$$\min \text{rank}(M) \quad s.t. \quad M_{ij} = X_{ij}, (i, j) \in \Omega$$

is unstable. However, an approximation algorithm by least-squares low-rank approximation is stable:

$$\min_{U, V} \sum_{(ij) \in \Omega} (X_{ij} - [UV]_{ij})^2$$

Algorithm robustness

An algorithm is robust if it can be applied reliably to a wide range of data, e.g., random data.

Example: robust optimization

Robust linear programming addresses linear programming problems where the data is uncertain, and a solution which remains feasible despite that uncertainty is sought.

$$(\text{LP}) \quad \min_x c^\top x : \quad a_i^\top x \leq b_i, \quad i = 1, \dots, m.$$

$$(\text{Robust LP}) \quad \min_x c^\top x : \quad a_i \in \mathcal{U}_i, a_i^\top x \leq b_i, \quad i = 1, \dots, m.$$

Robust Optimization, 2009. Aharon Ben-Tal, Laurent El Ghaoui & Arkadi Nemirovski.
Robust statistics: Peter J. Huber, Elvezio M. Ronchetti.

Example failure for linear programming

$$c = \begin{bmatrix} 100 \\ 199.9 \\ -5500 \\ -6100 \end{bmatrix} \quad A = \begin{bmatrix} -.01 & -.02 & .5 & .6 \\ 1 & 1 & 0 & 0 \\ 0 & 0 & 90 & 100 \\ 0 & 0 & 40 & 50 \\ 100 & 199.9 & 700 & 800 \\ & & -I_4 & \end{bmatrix} \quad \text{and} \quad b = \begin{bmatrix} 0 \\ 1000 \\ 2000 \\ 800 \\ 100000 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}.$$

c vector of costs/profits for two drugs, constraints $Ax \preceq b$ on production

- what happens if we vary percentages .01, .02 (chemical composition of raw materials) by .5% and 2%, i.e. $.01 \pm .00005$ and $.02 \pm .0004$?

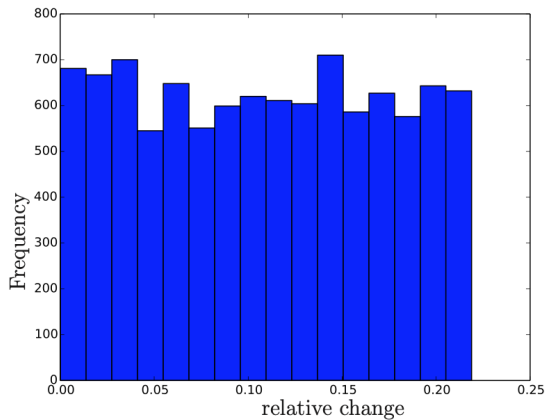


Figure 1: Frequency of fluctuations to a certain level for the non-robust production planning problem with .5% fluctuations in material proportions.

Alternative robust LP

$$\begin{array}{ll}\text{minimize} & c^T x \\ \text{subject to} & (A + \Delta)x \preceq b, \quad \text{all } \Delta \in \mathcal{U}\end{array}$$

where $|\Delta_{11}| \leq .00005$, $|\Delta_{12}| \leq .0004$, $\Delta_{ij} = 0$ otherwise

- solution x_{robust} has degradation *provably* no worse than 6%

Algorithm efficiency

- ▶ usual measure of efficiency is speed, i.e., how long an algorithm takes to produce its result
- ▶ analysis of algorithm: to determine the amount of resources (time and storage) needed to execute an algorithm
- ▶ complexity of algorithm:
 - ▶ time complexity: count the # of operations (flops)
 - ▶ worst-case running time: the longest running time for any input of size n
 - ▶ order of growth of the running time that really interests us
asymptotic analysis and big O notation
e.g. $f(n) = 9 \log n + 5(\log n)^3 + 3n^2 + 2n^3 = \mathcal{O}(n^3)$, as $n \rightarrow \infty$
 - ▶ also care of “memory” complexity

P vs. NP

NP: non-deterministic polynomial

- ▶ polynomial-time algorithms: on inputs of size n , the worst-case running time is $\mathcal{O}(n^k)$ for some constant k : e.g., sorting
- ▶ there are also problems that can be solved but not in time $\mathcal{O}(n^k)$ for any constant k
- ▶ we think of problems that are solvable by polynomial-time algorithms as being tractable, or easy, and problems that require superpolynomial time as being intractable, or hard

Example: Variable selection is NP hard

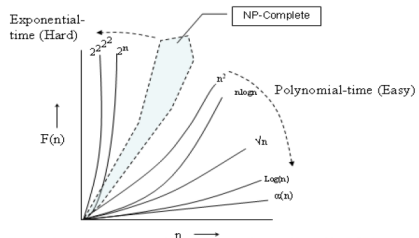
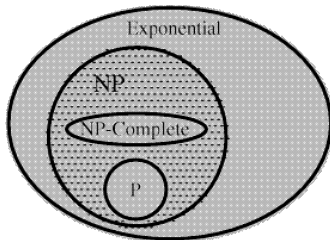
- ▶ Motivation: select the most important k variables in linear regression model
- ▶ Example: predicting housing price
<https://www.kaggle.com/c/house-prices-advanced-regression-techniques>
- ▶ Matrix-vector representation: Solve $y = X\beta$ where $\beta \in \mathbb{R}^p$ is k -sparse.

Exponential complexity

$$\begin{array}{ll} \text{minimize} & \|\beta\|_0 \\ \text{subject to} & y = X\beta \end{array}$$

NP complete

- ▶ an interesting class of problems, called the **NP-complete** problems, whose status is unknown: no polynomial-time algorithm has yet been discovered for an NP-complete problem, nor has anyone yet been able to prove that no polynomial-time algorithm can exist for any one of them
- ▶ NP-complete problems: finding clique, travel salesman



Common approaches to design algorithm

- ▶ Recursion: an algorithm that recursively calls itself (repeat the some form of procedure).

Example: compute mean and variance

Horner's method

$$p_d(x) = c_d x^d + \cdots + c_1 x + c_0$$

evaluated as

$$p_d(x) = x(\cdots x(x(c_d x + c_{d-1}) + \cdots) + c_1) + c_0$$

- ▶ Divide and conquer

A problem is broken into subproblems, each of which is solved, and then the subproblem solutions are combined into a solution for the original problem.

Example: “bubble sort” versus **quick sort**, **FFT**.

- ▶ Greedy algorithm
 - ▶ Each step is as efficient as possible **without** regarding future steps
 - ▶ Greedy algorithm is usually used in the early stages of computation for a problem or when a problem lacks understandable structures.
 - ▶ Example: **gradient descent**, **Newton's method**.
- ▶ Consider “cost-to-go”: Dynamic programming, approximate dynamic programming, Reinforcement learning (RL)
- ▶ Iterative method: **bisection**
Convergence: whether or not it will ends with a (right) fix point if iterate enough steps
- ▶ Approximation: e.g., convex relaxations:
replace non-convex constraints with convex ones.

Example: convex relaxation for variable selection

Solve $y = X\beta$ where $\beta \in \mathbb{R}^p$ is k -sparse.

Exponential complexity (NP hard)

$$\begin{array}{ll} \text{minimize} & \|\beta\|_0 \\ \text{subject to} & y = X\beta \end{array}$$

Polynomial complexity

$$\begin{array}{ll} \text{minimize} & \|\beta\|_1 \\ \text{subject to} & y = X\beta \end{array}$$

Related to **lasso** algorithm and compressed sensing

Contents

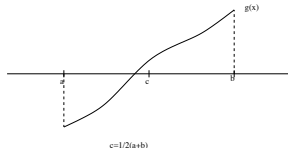
- Aspects of Algorithms
- Important Algorithms

Bisection

- ▶ for **continuous monotone** function $g(x)$, find root such that

$$x^* : g(x^*) = 0$$

- ▶ In statistics: finding the maximum likelihood estimator
- ▶ bisection:
 - ▶ start with $g(a) < 0 < g(b)$
 - ▶ take $c = \frac{1}{2}(a + b)$
 - ▶ if $g(c) < 0$, consider right half interval $[c, b]$
 - ▶ if $g(c) > 0$, consider left half interval $[a, c]$
 - ▶ repeat the above corresponding subinterval



- By doing this, we always have

$$g(x_l) < 0 < g(x_r)$$

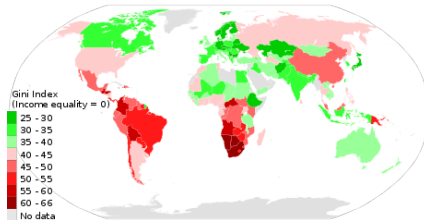
- since g is continuous, there must be a point $x^* \in [x_l, x_r]$ such that $g(x^*) = 0$
- the length of the interval is halved each time
- after n iterations, the final bracketing interval has length $2^{-n}(b - a)$, this means

$$|x^* - x| < 2^{-n}(b - a), \quad \forall x \in [a, b]$$

- the length of the interval converges to 0 as $n \rightarrow \infty$ (convergence rate $e^{-n \log 2}$: exponential convergence rate)

Sorting and statistics

- ▶ Compute median
- ▶ Gini index: A measure of statistical dispersion intended to represent the income or wealth distribution of a nation's residents, most commonly used measurement of inequality.



For population with values y_i , $i = 1, \dots, n$, that are indexed in *non-decreasing order* the Gini index is defined as

$$G = \frac{1}{n} \left(n + 1 - 2 \frac{\sum_{i=1}^n (n + 1 - i) y_i}{\sum_{i=1}^n y_i} \right)$$

Quicksort algorithm

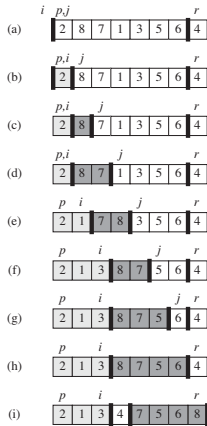
- ▶ sorting: $\{3, 1, 5, 2\} \Rightarrow \{1, 2, 3, 5\}$
- ▶ Native “bubble sort”

Starting from the beginning of the list, compare every adjacent pair, swap their position if they are not in the right order (the latter one is smaller than the former one). After each iteration, one less element (the last one) is needed to be compared until there are no more elements left to be compared.
- ▶ complexity: $\mathcal{O}(n^2)$

- ▶ Quicksort: a recursive algorithm using “divide-and-conquer”
- ▶ a vector of numbers c of length n , start location for sort p , end location for sort q
- ▶ pseudocode
 - quicksort(c , p , q)
 - $r := \text{findpivot}(c, p, q)$
 - quicksort(c , p , $r-1$)
 - quicksort(c , $r+1$, q)

Demo: <http://me.dt.in.th/page/Quicksort/>

example of
PARTITION
(to find pivot)



- ▶ best scenario: [1, 3, 2, 4, 5, 7, 6],
- ▶ worst scenario: [8, 7, 6, 5, 4, 3, 2, 1] or [1, 2, 3, 4, 5, 6, 7]

Complexity of quicksort

Average complexity of quicksort algorithm is $\mathcal{O}(n \log n)$.

Proof:

Homework: show that the worst-case complexity of quicksort is $\mathcal{O}(n^2)$.

Proof sketch

For n numbers. Assume all $n!$ permutations are equally likely.

$$e_0 = 0$$

Complexity: how many elements to be examined.

Recursion:

$$\begin{aligned}e_n &= n - 1 + \frac{1}{n} \sum_{i=1}^n (e_{i-1} + e_{n-i}) \\&= n - 1 + \frac{2}{n} \sum_{i=1}^n e_{i-1} \\ne_n &= n(n-1) + 2 \sum_{i=1}^n e_{i-1}\end{aligned}$$

Take difference

$$\begin{aligned} & ne_n - (n-1)e_{n-1} \\ &= n(n-1) - (n-1)(n-2) + 2 \sum_{i=1}^n e_{i-1} - 2 \sum_{i=1}^{n-1} e_{i-1} \\ &= 2(n-1) + 2e_{n-1} \end{aligned}$$

$$\frac{e_n}{n+1} - \frac{e_{n-1}}{n} = \frac{2(n-1)}{n(n+1)}$$

Recursion leads to

$$\frac{e_n}{n+1} = 2 \sum_{k=1}^n \frac{k-1}{k(k+1)} = \dots = 2 \sum_{k=1}^n \frac{1}{k} + 4 \frac{n}{n+1}$$

e_n on the order of $2(n+1) \log n + 4n$

Convolution

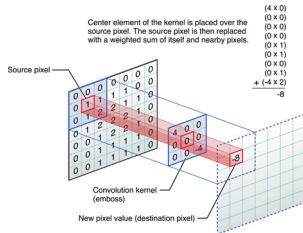
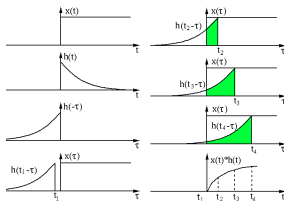
- ▶ convolution of two functions

continuous: $R(t) = x(t) \star h(t) = \int x(t-u)h(u)du$

discrete: $R[m] = x \star h = \sum_j x_{m-j}h_j$

- ▶ very important in statistics, image processing, signal processing, computer science.
 - ▶ distribution of sum of two random variables $U + V$ is the convolution of their PDFs
 - ▶ “filtering” of an image

(Section 3, “Computational statistics” by J. Gentle.)



Convolution and Fourier transform

- Fourier transform, denoted as $\mathcal{F}(x)$ is defined as

$$\mathcal{F}(x) \triangleq X(f) = \int_{-\infty}^{\infty} x(t)e^{-i2\pi ft} dt$$

where $i = \sqrt{-1}$

- Inverse Fourier transform

$$x(t) = \int_{-\infty}^{\infty} X(f)e^{i2\pi ft} df$$

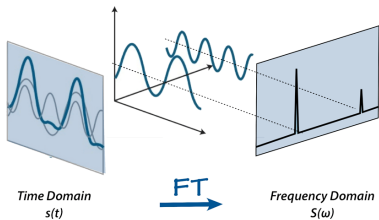
Convolution theorem

$$\mathcal{F}(x \star h) = \mathcal{F}(x) \cdot \mathcal{F}(h)$$

The Fourier Transform .com

$$\mathcal{F}\{g(t)\} = G(f) = \int_{-\infty}^{\infty} g(t)e^{-i2\pi ft} dt$$

$$\mathcal{F}^{-1}\{G(f)\} = g(t) = \int_{-\infty}^{\infty} G(f)e^{i2\pi ft} df$$

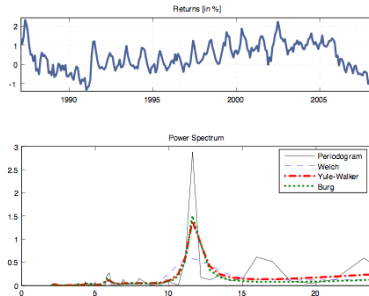


Proof

$$\begin{aligned}\mathcal{F}(x \star h) &= \int \int x(t-u)h(u)e^{-i2\pi ft}dudt \\ &= \int \left\{ \int x(t-u)e^{-i2\pi ft}dt \right\} h(u)du \\ &= X(f) \int e^{-i2\pi fu}h(u)du \\ &= X(f)H(f) = \mathcal{F}(x) \cdot \mathcal{F}(h).\end{aligned}$$

Example: Analysis of Financial Time-Series using Fourier method

- ▶ Case-Shiller home price index for the city of New-York
- ▶ January 1987 to May 2008 and the index is reported on a monthly basis



power spectral $|X(f)|^2$

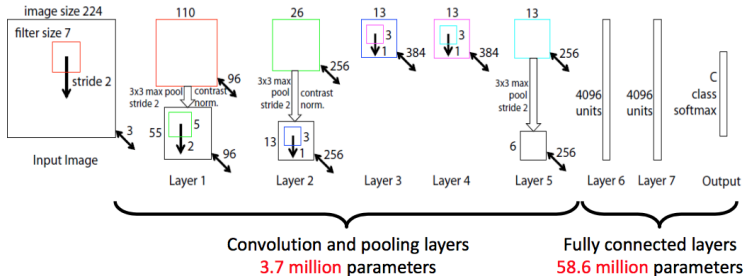
strong seasonalities affect home prices; frequency of recurrence of 12 months.

Example: Convolutional neural network (CNN)

Deep learning using CNN (multi-dimensional Fourier transform)

Convolution + max pooling: highly structured processing pipeline

Fully connected layers: nonlinear classification



Discrete Fourier transform

- ▶ Discrete Fourier transform (DFT), denoted as $\mathcal{F}(x)$ is a vector such that each element is given by

$$\tilde{x}_m = \sum_{j=0}^{N-1} x_j e^{-i \frac{2\pi}{N} jm}$$

where $i = \sqrt{-1}$

- ▶ Inverse DFT

$$x_j = \frac{1}{N} \sum_{m=0}^{N-1} \tilde{x}_m e^{i \frac{2\pi}{N} jm}$$

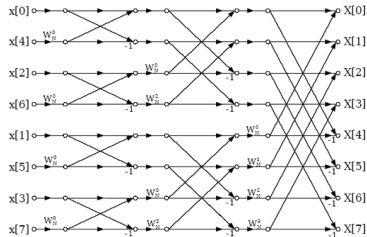
- ▶ equivalently, in matrix-vector representation

$$\tilde{x} = Ax$$

where matrix A is N -by- N and $A_{mj} = e^{-\frac{2\pi i}{N} jm}$

FFT (Fast Fourier Transform)

- ▶ Recall: DFT is equivalent to computing $\tilde{x} = Ax$
- ▶ Normally this is $\mathcal{O}(N^2)$, when the matrix has special form, however, it may be reduced.
- ▶ This is the idea of *fast Fourier Transform (FFT)*.
- ▶ Complexity of FFT is $\mathcal{O}(N \log N)$



butterfly

Derivation of FFT

- ▶ Assume N is even
- ▶ Let $e_n = x_{2n}$ represent the even-indexed samples
- ▶ Let $o_n = x_{2n+1}$ represent the odd-indexed samples
- ▶ One can show that e_n and o_n are zero outside the interval $0 \leq n \leq (N/2) - 1$

- One can show that

$$\tilde{x}_k = \frac{1}{2} \tilde{E}_k + \frac{1}{2} W_N^k \tilde{O}_k, \quad k = 0, 1, \dots, N-1$$

where $W_N = e^{-i\frac{2\pi}{N}}$

the two terms are DFT of the even- and the odd-indexed samples

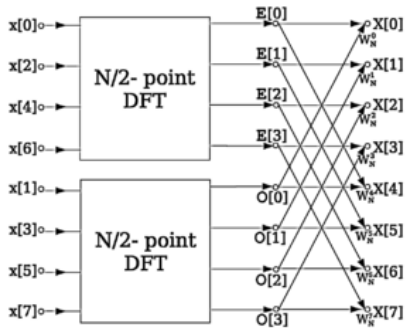
$$\tilde{E}_k = 2 \sum_{n=0}^{N/2-1} e_n W_{N/2}^{nk}, \quad \tilde{O}_k = 2 \sum_{n=0}^{N/2-1} o_n W_{N/2}^{nk}$$

- Moreover, there is symmetry

$$\tilde{E}_{k+N/2} = \tilde{E}_k, \quad \tilde{O}_{k+N/2} = \tilde{O}_k.$$

- **Length- N DFT of x_n can be computed as two DFTs of length $N/2$.**

Divide and conquer



Summary

- ▶ Aspects of algorithms
 - ▶ Accuracy
 - ▶ Efficiency
 - ▶ Robustness and stability
- ▶ Analyzing algorithms
 - ▶ Bisection for finding root
 - ▶ Quicksort algorithm for sorting a sequence of numbers
 - ▶ Convolution and its quick implementation via FFT
- ▶ Some important strategies:
recursion, divide and conquer, convex relaxation etc.