

Lecture 6: Match Data with Algorithms

Xiaotie Deng

AIMS Lab
Department of Computer Science
Shanghai Jiaotong University

October 17, 2017

- 1 Algorithm and Data
- 2 Progressive Data Sequence
- 3 Data of Fixed Parameters
- 4 Data of Fixed Distribution
- 5 Algorithmic Lower Bound for Data of Fixed Distribution

Algorithm and Data

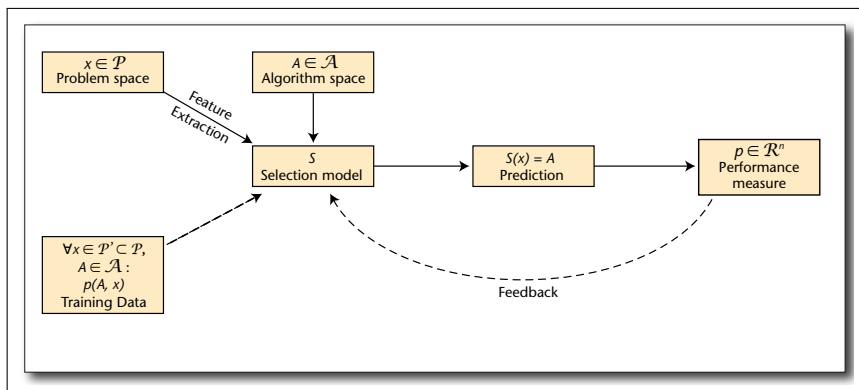
Asymptotic Worst Case Complexity of Algorithm on Data

- Given a problem \mathcal{P} .
 - Each $p \in \mathcal{P}$ is characterized by its input data.
- Design an algorithm for all problem instance in \mathcal{P} .
 - Return a correct output of the problem for each input instance.
- Time(/space/communication) complexity of algorithm A .
 - $t(A, n) = \max\{t(A, x) : \text{time to return output } \forall x : |x| = n\}$
- Complexity of the problem.
 - The best algorithm $\min\{t(A, n) | A \in \mathcal{A}\}$.

Algorithm on Restricted Data

- Data input with fixed parameter.
 - MOOC assignment: <https://ocw.mit.edu/courses/electrical-engineering-and-computer-science/6-046j-design-and-analysis-of-algorithms-spring-2015/lecture-videos/lecture-18-complexity-fixed-parameter-algorithms/>
- Data input following a given distribution.
 - Impagliazzo, Russell. "A personal view of average-case complexity." Structure in Complexity Theory Conference, 1995., Proceedings of Tenth Annual IEEE. IEEE, 1995.
- Fixed data with random noise.

Algorithm Selection



Five Worlds of Impagliazzo

- Algorithmica: $P = NP$, verification equivalent to solving a problem.
- Heuristica: $NP \neq P$ but some tractable on the average for some distribution.
- Pessiland: Hard average problems exist but no one way function.
- Minicrypt: One way function exists but PKC is impossible.
- Cryptomania: The existence of PKC (Public Key Cryptography).

Five Worlds in Big Data

- Propose Some in your assignment

Progressive Data Sequence

An Ensemble of Distribution and Class $AvgP$

- An ensemble of distribution is a sequence of distribution μ_n , $n \in Z$, on the set of positive integers with bit size n .
- A function $T : Z^+ \rightarrow Z^+$ is a polynomial on average with respect to μ_n , $n \in Z^+$, if there is an $\epsilon > 0$ such that the expectation of $T(i)^\epsilon$ when i is chosen according to μ_n is $O(n)$.
- A problem f on μ_n is in $AvgP$ if there is an algorithm to compute f whose running time is polynomial on average with respect to μ_n .

Polynomial Time Benign Algorithm Scheme

- Algorithm $A(x, \delta)$ computes f with benign fault: Output f or '?' and output is the correct function value if not '?';
- runs in polynomial in $|x|$ and $1/\delta$;
- $\forall \delta : 1 > \delta > 0$, and $\forall n \in \mathbb{Z}^+$: $\text{Prob}_{x \in \mu_n \mathbb{Z}^+}(A(x, \delta) = '?') \leq \delta$.

Heuristic Polynomial Time Algorithm \mathcal{HP}

- For x randomly chosen according to $\{\mu_n : n \in \mathbb{Z}^+\}$, and $\forall \delta > 0$, there is a deterministic polynomial time algorithm $A(x, \delta)$ that computes $f(x)$ correctly except an error of upto δ .
- \mathcal{HPP} : probability version.
- $\mathcal{HPP}/poly$: non-uniform algorithm version.

Data of Fixed Parameters

Vertex Cover Set of Fixed Constant Size k

- Idea: There are $\binom{n}{k}$ possible such subsets
- Algorithm: Go over all $\binom{n}{k}$ loops) and check (times m).
- Total time: $O(n^{k+1})$, a polynomial where k is a constant.

FPT: Reduced time to $O(f(k)) * n^{O(1)}$

- Idea: There is a node selected in an edge.
- Algorithm:
 - Go over all edges one by one,
 - binary step: choose u or v if none of nodes u and v is already chosen.
- Total time: $O(2^k + n + m)$
 - no more than depth k , total binary steps bounded by 2^k .

Kernelization in Fixed Parameter Complexity

A preprocessing stage to reduce the input to a smaller input, called a "kernel", that is easier to solve.

- Idea: Remove all vertices of degree $k + 1$.
 - The remaining graph has maximum degree k .
 - The size of kernel, the resulted graph, is no more than $k^2 + k$ since its vertex cover set has no more than k vertices and each connects to no more than k other vertices
- Algorithm:
 - Remove each such vertex one by one
 - Work all choices of the remaining graph
- Total time: $O(f(k) + n + m)$

Data of Fixed Distribution

Order Statistics

- IID random variables: X_1, X_2, \dots, X_n .
- Order Statistics: $X_{i-1,n} < X_{i,n}$, $i = 1, 2, \dots, n$, with $X_{0,n} = 0$.
- For exponential distribution: $Pr[X_{i,n} > t] = ???$

The Case for Exponential Distributions

- CDF (cumulative distribution function)
 - $Pr[X_{1,n} > t] = Pr[X_i > t, i = 1, \dots, n] = e^{-nt}$
 - $Pr[X_{n,n} \leq t] = Pr[X_i < t, i = 1, \dots, n] = (1 - e^{-t})^n$

Joint Distributions of Order Statistics

- Lemma: $f_{X_{1,n}, X_{2,n}, \dots, X_{n,n}}(t_1 < t_2 < \dots < t_n) = n! * \prod_{i=1}^n f_X(t_i)$
- Proof: by symmetry, LHS = $n! * f_{X_1, X_2, \dots, X_n}(t_1, t_2, \dots, t_n)$
 - where $t_1 < t_2 < \dots < t_n$.
- Note: Condition $t_1 < t_2 < \dots < t_n$ is important and useful.
- Corollary:

Generating Order Statistics of Exponential Distributions

- Define $Y_1 = X_{1,n}$, $Y_2 = X_{2,n} - X_{1,n}$, \dots , $Y_n = X_{n,n} - X_{n-1,n}$.
- Then Y is linear in X : $Y = AX$ where A is lower triangle with diagonal terms all 1.
- $f_{X_{1,n}, X_{2,n}, \dots, X_{n,n}}$

Properties of (Negative) Exponential Distribution

- Memoryless: $Pro(X > s + t | X > t) = Pro(X > s)$.
- $Pro((X_{1,n} > t) = \prod_{i=1}^n Pro(X_i > t) = e^{-nt}$.
- $Pro(X_{i,n} - X_{i-1,n} > t) = \prod_{j=i}^n Pro(X_j > t) = e^{-(n-i+1)t}$.
- $X_{n-i+1,n} = Y_i + Y_{i+1} + \dots + Y_n$ where
 - $Y_i = \min\{X_{n-i+1}, X_{n-i}, \dots, X_n\} = \frac{X_i}{i}$

Generating Order Statistics of Exponential Distributions

- Generate $Y_i, i = 1, 2, \dots$.
- Set $X_0 = 0, X_i = X_{i-1} + Y_{n-i+1}$.

Algorithmic Lower Bound for Data of Fixed Distribution

Lower Bound for Sorting

- The Comparison Model.
- The Candidate Permutation Set.
- The Maximum Subset Choosing Adversary.

The Average Complexity of Sorting

- For any deterministic algorithm based on comparison, the average number of comparisons is $\Omega(n \log n)$
- Reference: How Good Is The Information Theory Bound in Sorting? Michael Feldman, Theoretical Computer Science 1 (1976), 355-361.

Proof

- Uniform Distribution
- From the above analysis, we need a binary trees of 2^n leaves.
- We are interested in the average depth of this tree.
- True for balanced tree.
- Proof: Unbalanced tree has no less average depth.

Yao's principle

- The following two measures are equivalent.
 - $Cost(\mathcal{R}, x)$: Expected cost of a randomized algorithm \mathcal{R} on an input x
 - $Cost(A, \mathcal{D})$ Expected cost of a deterministic algorithm A on a random input distribution \mathcal{D} .
- More precisely, the best randomized algorithm on the worst time complexity is equivalent to the worst distribution for the best average algorithmic solution.
 - $\min_{\mathcal{R}} \max_{x \in \mathcal{I}} Cost(\mathcal{R}, x) = \max_{\mathcal{D}} \min_{A \in \mathcal{A}} Cost(A, \mathcal{D})$.
- Proof: Game tree value.

The Deterministic Time Lower Bound of Randomized Sorting

- For any randomized sorting algorithm, there is an input for which the randomized algorithm take time at least $n \log n$.

The Average Time Lower Bound of Randomized Sorting

- Important: The above lower bound is the same for any deterministic algorithm.
- A randomized algorithm is a distribution over a class of deterministic algorithm.
- average $[T(D_i, \mathcal{I})] \geq n \log n$ implies
 - $\sum_i r_i \cdot \text{average}[T(D_i, \mathcal{I})] = \text{average}[T(R, \mathcal{I})] \geq n \log n$
 - where $R = \sum_i r_i T$ and $\sum_i r_i = 1$.

Exercise

- Assume we have proven $T(D_i) \geq c_i f(n)$ for the average time complexity of each deterministic algorithm D_i , prove the some lower bound of $f(n)$ for all randomized algorithm, or give a counter example (based on c_i).

Assignment II (last part)

- 1 For a graph $G = (V, E)$, design a polynomial time algorithm to find a clique, i.e., a subset of vertices which has an edge between each other, of size constant k .