
Algorithms for Big Data (VIII)

Chihao Zhang

Shanghai Jiao Tong University

Nov. 8, 2019

REVIEW

Last week, we learnt a few graph streaming algorithms.

Recall that we have the following simple algorithm for **counting triangles**.

Let $\mathbf{f} = (f_T)_{T \in \binom{[n]}{3}}$ be the vector where for $T = \{x, y, z\}$, $f_T = |\{\{x, y\}, \{x, z\}, \{y, z\}\} \cap E|$.

The algorithm simply returns $F_0 - 1.5F_1 + 0.5F_2$, where $F_i = \|\mathbf{f}\|_i^i$.

We can expand $F_0 - 1.5F_1 + 0.5F_2$ as

$$\sum_{T \in \binom{[n]}{3}} 0.5f_T^2 - 1.5f_T + \mathbf{1}[f_T \neq 0].$$

The “polynomial” $f(x) = 0.5x^2 - 1.5x + \mathbf{1}[x \neq 0]$ satisfies

- ▶ $f(0) = f(1) = f(2) = 0$;
- ▶ $f(3) = 1$.

The **multiplicative error** of the algorithm is unbounded!

I leave the analysis of the algorithm as an **exercise**.

COMMUNICATION COMPLEXITY

Suppose we want to compute some function $f(x, y)$ where $x \in \{0, 1\}^a$ and $y \in \{0, 1\}^b$.

Alice has x and Bob has y , they collaborate to compute f .

The complexity is measured by **bits communicated between the two**.

We consider **one-way communication** model, with possible **public random coins**.

EXAMPLE: EQUALITY

Consider the function $f(x, y) = \text{EQ}(x, y) = \begin{cases} 1 & \text{if } x = y; \\ 0 & \text{otherwise,} \end{cases}$ where $x, y \in \{0, 1\}^n$.

The one-way complexity of EQ is n .

This can be shown by a simple [counting argument](#):

If the number of bits sent by Alice is less than n , then she can send at most $2^1 + 2^2 + \dots + 2^{n-1} = 2^n - 2$ distinct messages.

By the [pigeonhole principle](#), two different strings x and x' share the same message.

Bob can then use $y = x$ to fool the algorithm, a contradiction.

RANDOMNESS IN COMMUNICATION

We can design a more efficient protocol for EQ by tossing coins.

We treat x and y as two integers in $\{0, 2^n - 1\}$.

- ▶ Alice picks a random **prime** $p \in [n^2, 2n^2]$.
- ▶ She sends $(p, x \bmod p)$ to Bob.
- ▶ Bob outputs 1 if $y \bmod p = x \bmod p$, and outputs 0 otherwise.

If $x = y$, the algorithm is always correct.

If $x \neq y$, the algorithm is wrong only if $x = y \bmod p$.

The number of primes between n^2 and $2n^2$ is $\Theta\left(\frac{n^2}{\log n}\right)$ (**prime number theorem**).

At most n primes q satisfy $x = y \bmod q$ since $x, y < 2^n$.

DISJOINTNESS

The function $\text{DISJ}(x, y)$ tests whether two sets represented by x and y respectively share common elements.

Formally, $\text{DISJ}(x, y) = \begin{cases} 1 & \text{if } \langle x, y \rangle > 0 \\ 0 & \text{otherwise.} \end{cases}$, where $x, y \in \{0, 1\}^n$.

Same argument as EQ shows that computing DISJ deterministically requires n bits of one-way communication.

How about randomized protocols?

Unlike EQ, the power of randomness does not help much here...

LOWER BOUND FOR DISJ

Theorem

Randomized protocol for DISJ with correct probability at least $2/3$ needs $\Omega(n)$ bits of one-way communication.

We prove this for a special case of DISJ, the problem of INDEX. So the lower bound is stronger.

INDEX: Alice holds a string $x \in \{0, 1\}^n$, Bob holds an index $i \in [n]$. $\text{INDEX}(x, i) = x_i$.

YAO'S PRINCIPLE

The main tool we will use to derive the lower bound is Yao's principle.

Lemma

If there exists some distribution \mathcal{D} over $\{0, 1\}^a \times \{0, 1\}^b$ such that any **deterministic** one-way communication protocol P with

$$\Pr_{(x,y) \sim \mathcal{D}} [P \text{ is wrong on } (x, y)] \leq \epsilon$$

costs at least k bits, then any **randomized** one-way protocol with error at most ϵ on **any input** also costs at least k bits one-way communication.

LOWER BOUND FOR INDEX

By Yao's principle, we only need to construct a distribution \mathcal{D} over $\{0, 1\}^n \times [n]$ so that for any protocol with costs $o(n)$, it outputs the correct answer with probability less than $7/8$.

We let \mathcal{D} be the uniform distribution over $\{0, 1\}^n \times [n]$.

Assume there exists a protocol P that uses at most $0.1n$ bits of one-way communication.

Namely, Alice holds a function $f : \{0, 1\}^n \rightarrow \{0, 1\}^{0.1n}$. On input x , she sends $f(x)$ to Bob.

Upon receiving $f(x)$, Bob outputs some number $y(f(x))_i$. We collect the outputs (for all possible $i \in [n]$) as a vector $y(f(x)) \in \{0, 1\}^n$.

The algorithm is correct if $x_i = y(f(x))_i$.

Therefore, we only need to upper bound

$$\Pr_{(x,i) \sim \mathcal{D}} [x_i = y(f(x))_i]$$

where both $f : \{0, 1\}^n \rightarrow \{0, 1\}^{0.1n}$ and $y : \{0, 1\}^{0.1n} \rightarrow \{0, 1\}^n$ are **fixed!**

Since i is uniform in $[n]$, for any two strings $x, y \in \{0, 1\}^m$,

$$\Pr_{i \in [n]} [x_i \neq y_i] = \frac{d_H(x, y)}{n}.$$

Let $S = y(x(\{0, 1\}^n)) \subseteq \{0, 1\}^n$ be a set of size at most $\{0, 1\}^{0.1n}$. Since x is uniform in $\{0, 1\}^n$, we only need to show: **there are many $x \in \{0, 1\}^n$ satisfying $d_H(x, S) \geq n/4$.**

This is true since $\text{Ball}(S, \frac{n}{4}) \leq 2^{0.1n} \cdot \sum_{j=0}^{\frac{n}{4}} \binom{n}{j} \leq n2^{0.95n}$.

LOWER BOUND FOR F_∞

Our motivation for introducing communication model is to prove lower bound for streaming problems.

For example, we can use the lower bound for DISJ to derive lower bound for estimating F_∞ .

Theorem

Any randomized algorithm to estimate F_∞ within err $\varepsilon = 0.2$ requires $\Omega(n)$ bits of memory.

Proof.

Treat \mathbf{x} and \mathbf{y} as streams $\{i \in [n] \mid x_i = 1\}$ and $\{i \in [n] \mid y_i = 1\}$ respectively. □

A GENERAL PARADIAM

Previous proof provides a general paradigm for proving streaming lower bound based on communication lower bound:

A streaming algorithm to compute $f(x, y)$ using s bits of memory implies a protocol to compute $f(x, y)$ using at most s bits of one-way communication.

We will see more applications next time!