Aggregating a Data Set: Rankings to Strings

Diptarka Chakraborty (National University of Singapore)

Recent Trends in Algorithms 2022

Finding a Median String (Recall)

• Given a set of strings $S = \{x_1, x_2, ..., x_m\}$ over alphabet Σ , the objective is to find a string $y \in \Sigma^*$ (not necessarily from S) that <u>minimizes</u>

$$Obj(S, y) = \sum_{x_i \in S} ED(x_i, y)$$

- Let y^* be a string that minimizes Obj(S, y)
- y^{*} is referred to as *median*

Questions encountered so far (Recall)

- How to do clustering efficiently? Not in this talk
- How to perform approximate trace reconstruction efficiently?
- What is the connection between approximate trace reconstruction and approximate median?
- How to find an approximate median efficiently?

Rank Aggregation (Recall)

• Given a set of permutations $S = \{x_1, x_2, ..., x_m\}$ over [n], the objective is to find a **permutation** y (not necessarily from S) that minimizes

$$Obj(S, y) = \sum_{x_i \in S} ED(x_i, y)$$

Rank Aggregation - What do we know?

Metric	Upper Bound	Lower Bound
Kendall-tau	PTAS ($(1 + \epsilon)$ -approximation in polytime) [Mathieu, Schudy '07]	NP-hard [Dwork et al. '01] (even for 4 inputs) For 3 inputs, NP-hard or P?
Ulam	$(2 - \epsilon)$ -approximation in polytime [C, Das, Krauthgamer '21] For 3 inputs, in P	NP-hard?

Phase 1: Relax the constraint that the solution need not be a permutation

• Phase 2: Round the relaxed solution into a feasible one

Ulam Median – Phase 1 (Relax)

• The output need not be a permutation (can be an arbitrary string)

• Use dynamic programming for the Edit median, that runs in $O(2^m n^m)$ time (recall, now m = 3)

• Actually, find an *n*-length solution

Exercise 1: Design dynamic programming for the Edit median

Exercise 2: Modify it so that it outputs length restricted solution

Ulam Median – Phase 2 (Round)

- Round the relaxed solution to a feasible one (permutation)
- "Cleverly" delete duplicate symbols and insert missing ones



Ulam Median – High-level Argument

- Consider the following two quantities
 - *OPT_{rank}(S*): Optimum value attained by a permutation (Our final target)
 - $OPT_{n-len}(S)$: Optimum value attained by a *n*-length string (Phase 1 gives)
- Observation: $OPT_{n-len}(S) \le OPT_{rank}(S)$
- Suppose \overline{y} attains $OPT_{n-len}(S)$. Now, $\overline{y} \xrightarrow{Phase 2} Z$



• To show: $Obj(S, z) \le c \cdot Obj(S, \overline{y})$

- We show: $Obj(S, z) \le Obj(S, \overline{y})$ (i.e., c = 1), for output z for three inputs
- Fix an arbitrary optimal *alignment* between \overline{y} and x_1, x_2, x_3
- For each duplicate symbol
 - Only keep the occurrence with <u>maximum matches</u> and delete the rest

Idea: Increase the Obj by 1 (for one input), but decrease by 2 (avoid deletions in two inputs)

- We show: $Obj(S, z) \le Obj(S, \overline{y})$ (i.e., c = 1), for output z for three inputs
- Fix an arbitrary optimal *alignment* between \overline{y} and x_1, x_2, x_3
- For each duplicate symbol
 - Only keep the occurrence with <u>maximum matches</u> and delete the rest
- For each missing symbol
 - Insert while maintaining the match with at least one input

```
Idea: <u>Decrease</u> the Obj by 1 (for one matched input),
but <u>increase</u> by 2 (for the rest two inputs)
```

- We show: $Obj(S, z) \le Obj(S, \overline{y})$ (i.e., c = 1), for output z for three inputs
- Delete each duplicate: <u>Increase</u> the Obj by 1 (for one input), but <u>decrease</u> by 2 (avoid deletions in two inputs)
- Insert each missing: <u>Decrease</u> the Obj by 1 (for one matched input), but <u>increase</u> by 2 (for the rest two inputs)
- No. of duplicate entries = No. of missing (since \overline{y} has length n)
- Hence, *Obj* stays the same

Ulam Median for constant inputs

• Can be extended to $Obj(S, z) \le \frac{3}{2}Obj(S, \overline{y})$ (i.e., $c \approx 3/2$), for O(1) inputs [C, Das, Krauthgamer '21]

• Similar approach works for the center problem [C, Gajjar, Jha '21]

Max. Rank Aggregation / Ulam Center

• Given a set of permutations $S = \{x_1, x_2, ..., x_m\}$ over [n], the objective is to find a **permutation** y (not necessarily from S) that minimizes

 $Obj(S, y) = \max_{x_i \in S} ED(x_i, y)$

• The Ulam center problem is NP-hard

Ulam Center for O(1) inputs [C, Gajjar, Jha '21]

- Phase 1: Find an *n*-length center string (need not be a permutation)
- Phase 2: Round the relaxed solution into a feasible one
- Interestingly, Phase 2 is essentially solving the <u>Hamming center with</u> <u>wildcards</u>

• **Theorem:** Achieve 3/2-approximation for the Ulam center in polytime for O(1) inputs (and exact for 3 inputs)

Ulam Median – What about arbitrary inputs?

the no. of inputs *m*)

time

Phase 1: Relax the constraint that the solution need not be a permutation
 Takes exponential (in

• Phase 2: Round the relaxed solution into a feasible one

- Consider the case when $OPT(S) \ge mn/5$ (Recall, m = no. of inputs, n = length of input)
- **Theorem:** Best input achieves (2ϵ) -approximation, for some $\epsilon > 0$ (independent of n, m)

• Assumption: No input $x \in S$ is at distance $\leq (1 - \epsilon)OPT/m$ from an (unknown) median y^*



• Else, that point achieves $(2 - \epsilon)$ -approximation

• High Regime: $OPT \ge mn/5$



• For each point $x_i \in S$, there is a set I_i of edit (or unaligned) symbols, where $|I_i| \ge (1 - \epsilon)n/10$ (for simplicity, $|I_i| \ge n/15$)

- For each point $x_i \in S$, $|I_i| \ge n/15$
- Claim: There can be at most 30 such sets with pairwise intersection of size $\leq n/450$ (Follows from inclusion-exclusion)



• Observation: For any two x_i, x_j , $ED(x_i, x_j) \le |I_i| + |I_j| - |I_i \cap I_j|$

So the points inside the outer ball forms
 < 30 clusters (with a *buddy point*)

• Where the distance of any point from its buddy is roughly $\leq \left(2 - \frac{1}{30}\right)n/15$



- By averaging, one of the clusters must contributes *OPT*/70 mass to the *OPT*
- Next, consider the buddy point (\overline{y}) of the largest contributing cluster
- It breaks 2-factor for its own cluster, and roughly maintains 2-factor for other points



• So \overline{y} attains (2 - 1/30)-approx. for OPT/70 mass, and roughly 2-approx. for the rest

- Overall, it attains (2ϵ) -approx.
- **Theorem:** Best input achieves (2ϵ) -approx., for some $\epsilon > 0$ (independent of n, m)



Ulam Median (Small Regime)

- Consider the case when OPT(S) is "small" (and distributed over all/many symbols)
- Intuition: Most pairs of symbols (a, b) maintain their relative ordering in "majority" (much more than 50%) of inputs, and thus can be retrieved.
- Of course, if we could identify relative ordering of all the pairs, we can identify the (unknown) median using topological ordering

Ulam Median (Small Regime)

- Issue: Relative ordering of some pair of symbols could be changed in majority of inputs (maybe because of "a few" bad symbols)
- Can be handled by removing (shortest) cycles from the <u>graph with</u> <u>edges representing the relative ordering in "majority" inputs</u>
- Another issue: How to merge high and small regime cases?
- Merging only leads us to (2ϵ) -approx. algorithm

Open Problems

- Can we get a PTAS for the Ulam median?
- Can we show the Ulam median is NP-hard? (For O(1) inputs?)
- Can we extend the result to the Edit median?

Trace Reconstruction (Recall)

- Problem Statement: Reconstructing an unknown string from its noisy observable copies (aka. *traces*)
- There is an unknown string *x* of length *n*
- We observe a set of "noisy" copies (traces) $x_1, x_2, ..., x_m$
- The objective is to recover \boldsymbol{x}

1. Use as few samples as possible

- 2. Minimize the "error"
- 3. Design an efficient algorithm





- Substitution Channel (Each symbol is flipped with probability p)
- **Deletion Channel** (Each symbol is deleted with probability *p*)
- Insertion-Deletion Channel (While scanning, keeps the next symbol as it is w.p. 1 p, deletes it w.p. p/2, and inserts a uniformly randomly chosen symbol before the next symbol w.p. p/2)

We consider this one

Two Cases w.r.t. Unknown Strings (Recall)

- Worst-case: Unknown string is arbitrary
- Average-case: Unknown string is a uniformly randomly chosen string

What about Approximation? (Recall)

- Many applications (including DNA storage system) do not need the exact reconstruction
- It suffices to recover a string z that is "close" to the unknown string x
- Edit distance (ED) is a natural closeness measure



Getting ~*pn* edit distance is trivial (any input trace works)

Result 1 [C, Das, Krauthgamer '21]

Recall, unknown string = x of length n, and noise parameter = p

Only using *three* traces, we can recover a string *z* such that

$$\Pr\left[ED(x,z) \le O\left(p^2\log\frac{1}{p}\right)n\right] \ge 1 - \frac{1}{n}$$

in time $\tilde{O}(n)$ time (using a deterministic algorithm).

In this talk, we only consider Binary alphabet

Questions encountered so far (Recall)

How to do clustering efficiently?

Not in this talk

• How to perform approximate trace reconstruction efficiently?

 What is the connection between approximate trace reconstruction and approximate median?

• How to find an approximate median efficiently?

Result 2: They are the same (at least in average-case) [C, Das, Krauthgamer '21]

- Recall, in the average-case the unknown string x is chosen uniformly at random
- Let \overline{y} be a $(1 + \epsilon)$ -approximate median of x_1, x_2, \dots, x_m (generated from x by the **insertion-deletion** channel), for $\epsilon \in \left[110 \ p \log \frac{1}{p}, \frac{1}{6}\right]$
- Then w.h.p. $ED(x, \overline{y}) \le O(\epsilon) \frac{OPT}{m}$

Result 2 (almost) Result 1

• For three traces $S = \{x_1, x_2, x_3\}, OPT \le 3 p n$ $(ED(x, x_i) \le p n \text{ w.h.p., and thus } Obj(S, x) \le 3 p n)$

• Thus, w.h.p.
$$ED(x, \overline{y}) \le O\left(p^2 \log \frac{1}{p}\right)n$$

- Recall, a standard dynamic programming computes an exact median in time $O(n^3)$ for three input strings
- So, we can reconstruct the unknown string x approximately in time $O(n^3)$ time

Main Idea:

- Divide the input strings into small blocks,
- Then compute medians for those blocks independently, and
- Finally, concatenate those block medians

Achieving Near-linear Time (Intuition)

- *x*₁ 010100010110110110100011110110100111110001
- *x*₂ **101100100101110010001101110100011011110101**
- *x*₃ 01010100011101 11001001001110 1100111100110
- Block-wise x and z are close (by Result 2), and thus as a whole they are close too

Achieving Near-linear Time (Intuition)

- *x*₁ 010100010110110110100011110110100111110001
- *x*₂ **101100100101110010001101110100011011110101**
- *x*₃ 01010100011101110010010011101100110
- *Z* 010100011100111001000011110110001111110001
- **Challenges:** Induced partitions on x_1, x_2, x_3 may not be even. How do we identify them?



• Take **anchors** of size $\log^2 n$ in x_1



• Find the "best match" of the anchors in x₂, x₃



 Find an exact median of "best match" blocks of x₁, x₂, x₃, and concatenate them to form a string z



- For each anchor, look into its corresponding block in x
- Consider their "**true matches**" in x₂, x₃



• Claim: For each anchor, true matches are close to the best matches (It is crucial that x was chosen uniformly at random)



 So, finding median of the best match blocks is "almost same" as the finding median of the true match blocks



 Buffer (of size ω(log n)) helps to keep the best/true matched blocks well-separated



- Now, apply Result 2 block-wise
- Keep the buffer size much smaller compared to the anchors

Achieving Near-linear Time (Running time)



• There are $\frac{n}{\operatorname{poly} \log n}$ anchors, and for each block-median computation takes poly $\log n$ time.

Our Result 1 (Recall)

Recall, unknown string = x of length n, and noise parameter = p

Only using *three* traces, we can recover a string *z* such that

$$\Pr\left[ED(x,z) \le O\left(p^2\log\frac{1}{p}\right)n\right] \ge 1 - \frac{1}{n}$$

in time $\tilde{O}(n)$ time (using a deterministic algorithm).

Our Result 2 (Recall)

- Recall, in the average-case the unknown string x is chosen uniformly at random
- Let \overline{y} be an $(1 + \epsilon)$ -approximate median of x_1, x_2, \dots, x_m (generated from x by the **insertion-deletion** channel), for $\epsilon \in \left[110 \ p \log \frac{1}{p}, \frac{1}{6}\right]$
- Then w.h.p. $ED(x, \overline{y}) \le O(\epsilon) \frac{OPT}{m}$



Can view as x_j generating from x_i with a higher noise rate ($\sim 2p - \Theta(p^2)$)

 Corollary: A near-optimal alignment between x_i and x_j is also "almost" the same as that induced by the insertion-deletion noise channel

Robustness of Approximate Median

- Let \overline{y} be a $(1 + \epsilon)$ -approximate median
- Consider the set of edit operations that transforms $x_i \rightarrow \overline{y} \rightarrow x_j$
- We can prove that the **cost of the above alignment** (i.e., the number of edit operations) is **near-optimal** $((1 + O(\epsilon))$ -approximation)
- Thus, this alignment is almost the same as that induced by the noise channel



Hence, x and y are very similar (close in edit distance)

Concurrent Works

- For deletion channel (with noise rate *p*):
 - Upper bound: Given any $M \leq \Theta(1/p)$ traces, we can reconstruct up to edit distance $n \cdot (p M)^{\Omega(M)}$ in polynomial time
 - Information-theoretic lower bound: Given any $M \leq \Theta(1/p)$ traces, it is not possible to reconstruct up to edit distance better than $n \cdot (p M)^{O(M)}$ irrespective of the running time

[Chen, De, Lee, Servedio, Sinha '21]

• **Upper Bound:** exp $\left(O_p\left(\left(\log \frac{1}{\epsilon}\right)^{\frac{1}{3}}\right)\right)$ traces suffices to reconstruct up to edit distance ϵn [Chase, Peres '21]

Open Problems

- Can we reduce the error by using more samples?
- **Conjecture:** $ED(x, \bar{y})$ could be reduced to $O(\epsilon n)$, for arbitrary small $\epsilon > 0$, using $poly(1/\epsilon)$ traces.

• Can we get a similar result for the worst-case trace reconstruction (i.e., when the unknown string is arbitrary)?

Thank You!