

© Indian Institute of Technology Delhi (IITD), New Delhi, 2023

FPT APPROXIMATIONS FOR CONSTRAINED CLUSTERING PROBLEMS

by

DISHANT GOYAL

Department of Computer Science and Engineering

Submitted

in fulfillment of the requirements of the degree of **Doctor of Philosophy**

to the



Indian Institute of Technology Delhi

March 2023

Certificate

This is to certify that the thesis titled **FPT APPROXIMATIONS FOR CONSTRAINED CLUSTERING PROBLEMS** being submitted by **Mr. DISHANT GOYAL** for the award of **Doctor of Philosophy in Computer Science and Engineering** is a record of bona fide work carried out by him under my guidance and supervision at the **Department of Computer Science and Engineering, Indian Institute of Technology Delhi**. The work presented in this thesis has not been submitted elsewhere, either in part or full, for the award of any other degree or diploma.

Ragesh Jaiswal

Associate Professor

Department of Computer Science and Engineering

Indian Institute of Technology Delhi

New Delhi- 110016

Amit Kumar

Professor

Department of Computer Science and Engineering

Indian Institute of Technology Delhi

New Delhi- 110016

Acknowledgements

I am highly indebted to my advisors Prof. Ragesh Jaiswal and Prof. Amit Kumar, for their constant support and guidance during my PhD and my life in general. I am highly thankful to them for believing in me and encouraging me. A very special thanks to Prof. Ragesh Jaiswal for tolerating my grammar mistakes. My heartfelt gratitude to Prof. Ragesh Jaiswal for keeping a constant check on my well-being all this time.

A very special thanks to Prof. Anup Bhattacharya for his invaluable encouragement and discussion sessions. I am also grateful to my research committee members Prof. Naveen Garg, Prof. Bhawani Sankar Panda, and Prof. Preeti Ranjan Panda, for their valuable and constructive suggestions during the development of this research work. I am deeply fortunate to have been taught and guided by Prof. Venkatesh Raman during my B.Tech. at IIT Jodhpur.

I extend my gratitude to Tata Consultancy Services (TCS) for their support through the TCS Research Scholar Program Fellowship. I am also thankful to the theoretical computer science stack exchange community for the invaluable discussions on the clustering problems.

A special thanks to my SIT 309 friends – Vinayak Gupta, Sandeep Kumar, Arindam Bhattacharya, Dilpreet Kaur, Omais Shafi, and Ovia Seshadri for their constant support via poker nights, cricket games, food hunting trips, and innumerable coffee breaks. My heartfelt thanks to Vinayak Gupta for his realistic life advises. I am also thankful to have immense support from Neetu Jindal during the last phase of my degree. Lastly, I thank my parents and relatives for

standing by me during the course of this degree. I thank my sister for being there with me in all my good and bad times. This thesis is dedicated to my sister – Neha! :)

Dishant Goyal

Abstract

In this work, we study a wide range of *constrained* clustering problems in offline and streaming settings. We study these problems corresponding to three clustering objectives: k -median, k -means, and k -supplier. The (unconstrained) k -median problem is defined as follows. We are given a set of clients C in a metric space \mathcal{X} , with distance function $d(\cdot, \cdot)$. We are also given a set of feasible facility locations $L \subseteq \mathcal{X}$. The goal is to open a set $F \subseteq L$ of k facilities that minimizes the objective function: $\text{cost}(F, C) \equiv \sum_{j \in C} d(F, j)$, where $d(F, j)$ is the distance of client j to the closest facility in F . The k -means problem is defined in similar manner by replacing the distances with squared distances in the cost function, i.e., $\text{cost}(F, C) \equiv \sum_{j \in C} d(F, j)^2$. On the other hand, the k -supplier objective is defined as: $\text{cost}(F, C) \equiv \max_{j \in C} \{d(F, j)\}$. Furthermore, for $L = C$, the k -supplier problem is known as the k -center problem.

In many applications, there are additional constraints imposed on the clusters. For example, to balance the load among the facilities in resource allocation problems, a capacity u is imposed on every cluster. That is, no more than u clients can be assigned to any facility/cluster. This problem is known as the *capacitated* clustering problem. Likewise, various other applications have different constraints, which give rise to different *constrained* versions of the problem. In the past, the constrained versions of clustering problems were studied separately as independent problems. Recently, Ding and Xu [72] gave a unified framework for these problems that they called the *constrained clustering* framework. They proposed this framework in the context

of the k -median and k -means objectives in the continuous Euclidean space where $L = \mathbb{R}^p$ (p -dimensional Euclidean space) and C is a finite subset of \mathbb{R}^p . In this work, we extend this framework to the k -supplier objective and general metric spaces. The unified framework allows us to obtain results simultaneously for the following constrained versions of the problem: r -gather, r -capacity, balanced, chromatic, fault-tolerant, strongly private, ℓ -diversity, and fair clustering problems. We also study the *outlier* versions of these problems. In the outlier version, a clustering is obtained over at least $|C| - m$ clients instead of the entire client set.

For the constrained k -supplier and k -center problems, we obtain the following results:

- (1) We give 3 and 2 approximation algorithms for the constrained k -supplier and k -center problems, respectively, with FPT (fixed-parameter tractable) running time $k^{O(k)} \cdot n^{O(1)}$, where $n = |C \cup L|$. Moreover, we note that the obtained approximation guarantees are tight. That is, for any constant $\varepsilon > 0$, no algorithm can achieve $(3 - \varepsilon)$ and $(2 - \varepsilon)$ approximation guarantees for the constrained k -supplier and k -center problems, respectively, in FPT time parameterized by k , assuming $\text{FPT} \neq \text{W}[2]$.
- (2) For the outlier versions of the constrained k -supplier and k -center problems, we give 3 and 2 approximation guarantees with FPT running time $(k + m)^{O(k)} \cdot n^{O(1)}$, where $n = |C \cup L|$ and m is the number of outliers. Moreover, we note that the obtained approximation guarantees are tight. That is, for any constant $\varepsilon > 0$, no algorithm can achieve $(3 - \varepsilon)$ and $(2 - \varepsilon)$ approximation guarantees for the constrained k -supplier and k -center problems, respectively, in FPT time parameterized by k and m , assuming $\text{FPT} \neq \text{W}[2]$.

For the constrained k -median and k -means problems, we obtain the following results:

- (3) We give $(3 + \varepsilon)$ and $(9 + \varepsilon)$ approximation algorithms for the constrained k -median and k -means problems, respectively, with FPT running time $(k/\varepsilon)^{O(k)} \cdot n^{O(1)}$, where

$n = |C \cup L|$. For the outlier version of the constrained k -median and k -means problems, we give $(3 + \varepsilon)$ and $(9 + \varepsilon)$ approximation algorithms, respectively, with FPT running time $\left(\frac{k+m}{\varepsilon}\right)^{O(k)} \cdot n^{O(1)}$, where $n = |C \cup L|$ and m is the number of outliers.

- (4) We also study the problems when $C \subseteq L$, i.e., a facility can be opened at a client location as well. For this special case, we design $(2 + \varepsilon)$ and $(4 + \varepsilon)$ -approximation algorithms for the constrained k -median and k -means problems, respectively, with FPT running time $(k/\varepsilon)^{O(k)} \cdot n^{O(1)}$, where $n = |L|$. For the outlier version, we obtain the same approximation guarantees with FPT running time $\left(\frac{k+m}{\varepsilon}\right)^{O(k)} \cdot n^{O(1)}$, where $n = |L|$ and m is the number of outliers. Note that the case $C \subseteq L$ subsumes the case $C = L$. Therefore, this result also holds for the case when $C = L$.
- (5) We show that the analysis of our algorithm is tight. That is, there are instances for which our algorithm does not provide better than $(3 - \delta)$ and $(9 - \delta)$ approximation guarantee corresponding to k -median and k -means objectives, respectively, for any arbitrarily small constant $\delta > 0$. Similarly, the analysis of our algorithm is tight for the special case $C \subseteq L$.
- (6) Our algorithms are based on a simple sampling-based approach. This approach allows us to convert these algorithms to constant-pass log-space streaming algorithms.
- (7) We also study the constrained k -median/means problem in continuous Euclidean space where $L = \mathbb{R}^p$ and C is a finite subset of \mathbb{R}^p . We design $(1 + \varepsilon)$ -approximation algorithm for the outlier version of these problems with FPT running time $O\left(np \cdot \left(\frac{k+m}{\varepsilon}\right)^{O(k/\varepsilon^{O(1)})}\right)$, where $n = |C|$ and m is the number of outliers. We also convert these algorithms to constant-pass log-space streaming algorithms.

We also study the *socially fair k -median/ k -means problem*, which is a generalization of the k -supplier and k -median/means problems. The problem is defined as follows. We are given a set of clients C in a metric space \mathcal{X} with a distance function $d(\cdot, \cdot)$. There are ℓ groups:

$C_1, \dots, C_\ell \subseteq C$. We are also given a set L of feasible centers in \mathcal{X} . The goal in the socially fair k -median problem is to find a set $F \subseteq L$ of k centers that minimizes the maximum average cost over all the groups. That is, find F that minimizes the objective function: $\text{fair-cost}(F, C) \equiv \max_j \left\{ \sum_{x \in C_j} d(F, x) / |C_j| \right\}$, where $d(F, x)$ is the distance of x to the closest center in F . The socially fair k -means problem is defined similarly by using squared distances, i.e., $d^2(\cdot, \cdot)$ instead of $d(\cdot, \cdot)$. We obtain the following results for this problem:

- (8) We design $(3+\varepsilon)$ and $(9+\varepsilon)$ approximation algorithms for the socially fair k -median and k -means problems, respectively, in FPT time $f(k, \varepsilon) \cdot n^{O(1)}$, where $f(k, \varepsilon) = (k/\varepsilon)^{O(k)}$ and $n = |C \cup L|$.
- (9) Furthermore, these approximation guarantees are tight; that is, for any constant $\varepsilon > 0$, no algorithm can achieve $(3 - \varepsilon)$ and $(9 - \varepsilon)$ approximation guarantees for the socially fair k -median and k -means problems in FPT time parametrized by k , assuming $\text{FPT} \neq \text{W}[2]$.

Lastly, we give hardness of approximation result for the k -median problem in the continuous Euclidean space where $L = \mathbb{R}^p$ and C is a finite subset of \mathbb{R}^p . This solves an open problem posed explicitly in the work of Awasthi *et al.* [19]. More precisely, we obtain the following result:

- (10) There exists a constant $\varepsilon > 0$ such that the Euclidean k -median problem in $O(\log k)$ dimensional space cannot be approximated to a factor better than $(1 + \varepsilon)$, assuming the Unique Games Conjecture.

Furthermore, we study the hardness of approximation for the Euclidean k -means/ k -median problems in the *bi-criteria setting*. In the bi-criteria setting, algorithms are allowed to output βk centers (for some constant $\beta > 1$), and the approximation ratio is computed with respect to the optimal k -means/ k -median cost. We show the following results:

- (11) For any constant $1 < \beta < 1.015$, there exists a constant $\varepsilon > 0$ such that there is no $(1 + \varepsilon)$ bi-criteria approximation algorithm for the Euclidean k -median problem in $O(\log k)$ dimensional space assuming the Unique Games Conjecture.
- (12) For any constant $1 < \beta < 1.28$, there exists a constant $\varepsilon > 0$ such that there is no $(1 + \varepsilon)$ bi-criteria approximation algorithm for the Euclidean k -means problem in $O(\log k)$ dimensional space assuming the Unique Games Conjecture.

सार

इस काम में, हम ऑफ़लाइन और स्ट्रीमिंग समायोजन में बाध्य क्लस्टरिंग समस्याओं की एक विस्तृत श्रृंखला का अध्ययन करते हैं। हम तीन क्लस्टरिंग उद्देश्यों के अनुरूप इन समस्याओं का अध्ययन करते हैं: k -माध्यिका, k -माध्य, और k -आपूर्तिकर्ता। (अप्रतिबंधित) k -माध्यिका समस्या को निम्नानुसार परिभाषित किया गया है। हमें मेट्रिक जगह X में ग्राहक C का एक सेट दिया गया है, जिसमें डिस्टेंस फंक्शन $d(., .)$ है। हमें सुविधा स्थानों का एक सेट $L \subseteq X$ भी दिया जाता है। हमें लक्ष्य के सुविधाओं का एक सेट $F \subseteq L$ खोलना है जो उद्देश्य लागत को कम करता है: $\text{लागत}(F, C) = \sum_{j \in C} d(F, j)$, जहां $d(F, j)$ दूरी है ग्राहक j कि F में निकटतम सुविधा के लिए। k -माध्य समस्या को कुछ इसी तरह परिभाषित किया गया है- लागत फलन में दूरियों की जगह वर्ग दूरियों का प्रयोग किया जाता है, अर्थात् $d(F, j)$ की जगह $d(F, j)^2$ । दूसरी ओर, k -आपूर्तिकर्ता उद्देश्य निम्नानुसार परिभाषित किया गया है: $\text{लागत}(F, C) = \max_{j \in C} \{d(F, j)\}$ । इसके अलावा, $L = C$ के लिए, k -आपूर्तिकर्ता समस्या को k -केंद्र समस्या के रूप में जाना जाता है।

कई अनुप्रयोगों में, क्लस्टर पर अतिरिक्त प्रतिबंध लगाए गए हैं। उदाहरण के लिए, संमाध्य आवंटन समस्याओं में सुविधाओं के बीच भार को संतुलित करने के लिये, एक क्षमता हर क्लस्टर पर थोपा गया है। इस समस्या को कैपेसिटेड क्लस्टरिंग समस्या के रूप में जाना जाता है। इसी तरह, विभिन्न अन्य अनुप्रयोग, अलग-अलग बाधाएँ लगाते हैं, जो समस्या के विभिन्न बाध्य संस्करणों को जन्म देती हैं। अतीत में, क्लस्टरिंग समस्याओं के बाध्य संस्करणों का स्वतंत्र रूप से अलग से अध्ययन किया गया था। हाल ही में, डिंग और जू [68] ने इन समस्याओं के लिए एक एकीकृत ढांचा दिया जोकि बाध्य क्लस्टरिंग ढांचा कहा जाता है। उन्होंने इस ढांचे को यूक्लिडियन अंतरिक्ष में k -माध्यिका और k -माध्य उद्देश्यों में प्रस्तावित किया जहां $L = \mathbb{R}^p$ (p -डिमेंशनल यूक्लिडियन स्पेस) और C एक परिमित उपसमुच्चय है \mathbb{R}^p का। इस काम में, हम इसका विस्तार करते हैं k -आपूर्तिकर्ता उद्देश्य और सामान्य मीट्रिक स्थान के लिए। एकीकृत ढांचा समस्या के निम्नलिखित बाध्य संस्करणों के लिए एक साथ परिणाम प्राप्त करने की अनुमति देता है: r -इकट्ठा, r -क्षमता, संतुलित, रंगीन, दोष-सहिष्णु, दृढ़ता से निजी, ℓ -विविधता, और निष्पक्ष क्लस्टरिंग समस्याएं। हम इन समस्याओं के बाहरी

संस्करणों का भी अध्ययन करते हैं। बाहरी संस्करण में, एक क्लस्टरिंग में कम से कम $|C| - m$ ग्राहकों की प्राप्त की जाती है बजाय संपूर्ण ग्राहक सेट के। बाध्य k -आपूर्तिकर्ता और k -केंद्र समस्याओं के लिए, हम निम्नलिखित परिणाम प्राप्त करते हैं:

(1) हम बाध्य k -आपूर्तिकर्ता और k -केंद्र समस्याओं के लिए 3 और 2 सन्निकटन एल्गोरिदम देते हैं तथा FPT (फिक्स्ड-पैरामीटर ट्रेक्टबल) चलने के समय $O(f(k)) \cdot n^{O(1)}$ के साथ, जहां $n = |C \cup L|$ । इसके अलावा, ध्यान दें कि प्राप्त सन्निकटन गारंटी तंग हैं। अर्थात्, किसी भी स्थिरांक $\varepsilon > 0$ के लिए, कोई भी एल्गोरिथम $(3-\varepsilon)$ और $(2-\varepsilon)$ अनुकरण गारंटी प्राप्त नहीं कर सकता है- बाध्य k -आपूर्तिकर्ता और k -केंद्र समस्याओं के लिए FPT समय में, $FPT \neq W[2]$ मानकर।

(2) बाध्य k -आपूर्तिकर्ता और k -केंद्र समस्याओं के बाहरी संस्करणों के लिए, हम 3 और 2 सन्निकटन गारंटी देते हैं, FPT चलने के समय $(k + m)^{O(k)} \cdot n^{O(1)}$ के साथ, जहां $n = |C \cup L|$ और m बाहरी ग्राहकों (आउटलेर्स) की संख्या है। इसके अलावा, ध्यान दें कि प्राप्त सन्निकटन गारंटी तंग हैं। अर्थात्, किसी भी स्थिरांक $\varepsilon > 0$ के लिए, कोई एल्गोरिदम $(3 - \varepsilon)$ और $(2 - \varepsilon)$ सन्निकटन गारंटी प्राप्त नहीं कर सकता है सीमित k -आपूर्तिकर्ता और k -केंद्र समस्याओं के लिए, FPT समय में $FPT \neq W[2]$ मानकर।

बाध्य k -माध्यिका और k -माध्य समस्याओं के लिए, हम निम्नलिखित परिणाम प्राप्त करते हैं:

(3) हम बाध्य k -माध्यिका और k -माध्य समस्याएं के लिए $(3 + \varepsilon)$ और $(9 + \varepsilon)$ सन्निकटन एल्गोरिदम देते हैं, $(k/\varepsilon)^{O(k)} \cdot n^{O(1)}$ FPT चलने के समय के साथ, जहां $n = |C \cup L|$ । बाध्य k -माध्यिका और k -माध्य समस्याओं के बाहरी संस्करण के लिए, हम $(3 + \varepsilon)$ और $(9 + \varepsilon)$ सन्निकटन एल्गोरिदम देते हैं, $((k + m)/\varepsilon)^{O(k)} \cdot n^{O(1)}$ FPT चलने के समय के साथ, जहां $n = |C \cup L|$ और m बाहरी ग्राहकों (आउटलेर्स) की संख्या है।

(4) हम उन समस्याओं का भी अध्ययन करते हैं जब $C \subseteq L$, यानी, ग्राहक स्थान पर एक सुविधा भी खोली जा सकती है। इस विशेष मामले के लिए, हम $(2 + \varepsilon)$ और $(4 + \varepsilon)$ - सन्निकटन एल्गोरिदम देते हैं, बाध्य

k -माध्यिका और k -माध्य समस्याओं के लिए, $(k/\epsilon)^{O(k)} \cdot n^{O(1)}$ FPT चलने के समय के साथ, जहां $n = |L|$. बाहरी संस्करण के लिए, हम वही सन्निकटन गारंटी प्राप्त करते हैं, FPT चलने के समय के साथ $((k + m)/\epsilon)^{O(k)} \cdot n^{O(1)}$, जहां $n = |L|$ और m बाहरी ग्राहकों (आउटलेर्स) की संख्या है। ध्यान दें कि स्थिति $C \subseteq L$, स्थिति $C = L$ को समाहित करती है। इसलिए, यह परिणाम उस स्थिति के लिए भी मान्य है जब $C = L$.

(5) हम दिखाते हैं कि हमारे एल्गोरिथ्म का विश्लेषण कड़ा है। अर्थात्, ऐसे उदाहरण हैं जिनके लिए हमारा एल्गोरिथ्म k -माध्यिका और k -माध्य उद्देश्यों के अनुरूप $(3 - \epsilon)$ और $(9 - \epsilon)$ सन्निकटन गारंटी से बेहतर प्रदान नहीं करता है किसी भी मनमाने ढंग से छोटे $\epsilon > 0$ के लिए। इसी तरह, विशेष मामले के लिए हमारे एल्गोरिथ्म का विश्लेषण तंग है।

(6) हमारे एल्गोरिथ्म एक साधारण नमूना-आधारित दृष्टिकोण पर आधारित हैं। यह दृष्टिकोण इन एल्गोरिथ्म को निरंतर-पास लॉग-स्पेस स्ट्रीमिंग एल्गोरिथ्म में बदलने के लिए हमें अनुमति देता है।

(7) हम निरंतर यूक्लिडियन अंतरिक्ष में बाध्य k -माध्यिका/माध्य समस्या का भी अध्ययन करते हैं जहां $L = R^p$ और C एक परिमित उपसमुच्चय है R^p का। हम इन समस्याओं के बाहरी संस्करण के लिए $(1 + \epsilon)$ -सन्निकटन एल्गोरिथ्म $O(np \cdot ((k + m)/\epsilon)^{O(k/\epsilon^{O(1))})}$ FPT चलने के समय के साथ रचना करते हैं, जहां $n = |C|$ और m बाहरी ग्राहकों (आउटलेर्स) की संख्या है। हम इन एल्गोरिथ्म को निरंतर-पास लॉग-स्पेस स्ट्रीमिंग एल्गोरिथ्म में भी परिवर्तित करते हैं।

हम सामाजिक रूप से निष्पक्ष k -माध्यिका/ k -माध्य समस्या का भी अध्ययन करते हैं जिसे निम्नानुसार परिभाषित किया गया है। हमें मेट्रिक स्पेस X में ग्राहक C का एक सेट, डिस्टेंस फंक्शन $d(., .)$ के साथ दिया गया है। समूह हैं: C_1, \dots, C_ℓ . हमें X में व्यवहार्य केंद्रों का एक सेट L भी दिया गया है। सामाजिक रूप से लक्ष्य निष्पक्ष k -माध्यिका समस्या के केंद्रों का एक सेट $F \subseteq L$ ढूँढना है जो अधिकतम समूहों पर औसत लागत को कम करता है। यही है, F खोजें जो उद्देश्य लागत को कम करता है: $\text{लागत}(F, C) = \max\{ \sum_{j=1}^{\ell} d(F, C_j) / |C_j| \}$, जहां $d(F, x)$, F में निकटतम केंद्र से x की दूरी है। सामाजिक रूप से निष्पक्ष k -माध्य समस्या को समान रूप से

वर्ग दूरी का उपयोग करके परिभाषित किया जाता है, अर्थात्, $d(l, l')$ के बजाय $d(l, l')^2$ । हम इस समस्या के लिए निम्नलिखित परिणाम प्राप्त करते हैं:

(8) हम सामाजिक रूप से निष्पक्ष k -माध्यिका और k -माध्य के लिए $(3+\epsilon)$ और $(9+\epsilon)$ सन्निकटन एल्गोरिदम डिजाइन करते हैं, FPT समय में $f(k, \epsilon) \cdot n^{O(1)}$, जहाँ $f(k, \epsilon) = (k/\epsilon)^{O(k)}$ और $n = |C \cup L|$ ।

(9) इसके अलावा, ये सन्निकटन गारंटी तंग हैं; अर्थात्, किसी अचर $\epsilon > 0$ के लिए, एल्गोरिदम सामाजिक रूप से निष्पक्ष k -माध्यिका और k -माध्य के लिए $(3 - \epsilon)$ और $(9 - \epsilon)$ सन्निकटन गारंटी प्राप्त नहीं कर सकता है, FPT = $W[2]$ मानते हुए।

आखिर में हम k -माध्यिका समस्या के लिए सन्निकटन परिणाम की कठोरता देते हैं, यूक्लिडियन अंतरिक्ष में जहाँ $L = R^p$ और $C \subseteq R^p$ । यह एक खुली समस्या का समाधान करता है जो अवस्थी आदि [18] के कार्यों में स्पष्ट रूप से प्रस्तुत किया गया है। हम निम्नलिखित प्राप्त करते हैं:

(10) एक स्थिर $\epsilon > 0$ मौजूद है, जैसे कि $O(\log k)$ यूक्लिडियन k -माध्यिका समस्या को $(1 + \epsilon)$ से बेहतर कारक के रूप में अनुमानित नहीं किया जा सकता है, अद्वितीय खेल अनुमान (UGC) मानते हुए।

इसके अलावा, हम यूक्लिडियन k -माध्य/ k -माध्यिका के लिए सन्निकटन की कठोरता का अध्ययन करते हैं, द्वि-मानदंड सेटिंग में। द्वि-मानदंड सेटिंग में, एल्गोरिदम को βk केंद्र (कुछ स्थिर $\beta > 1$ के लिए) आउटपुट करने की अनुमति है और सन्निकटन अनुपात की गणना इष्टतम k -माध्य/ k -माध्यिका लागत के संबंध में की जाती है। हम निम्नलिखित परिणाम दिखाते हैं:

(11) किसी भी स्थिरांक $1 < \beta < 1.015$ के लिए, एक अचर $\epsilon > 0$ मौजूद है जिस्से कि कोई $(1+\epsilon)$ द्वि-मानदंड सन्निकटन एल्गोरिथम नहीं मौजूद हो सकता है, $O(\log k)$ यूक्लिडियन k -माध्यिका समस्या के लिए, अद्वितीय खेल अनुमान (UGC) मानते हुए।

(12) किसी भी अचर $1 < \beta < 1.28$ के लिए, एक अचर $\epsilon > 0$ मौजूद है जिस्से कि कोई $(1+\epsilon)$ द्वि-मानदंड सन्निकटन एल्गोरिथम नहीं मौजूद हो सकता है, $O(\log k)$ यूक्लिडियन k -माध्यिका समस्या के लिए, अद्वितीय खेल अनुमान (UGC) मानते हुए।

Contents

Certificate	i
Acknowledgements	iii
Abstract	v
1 Introduction	1
1.1 Classical (Unconstrained) Clustering	2
1.1.1 Polynomial time approximation	2
1.1.2 FPT time approximation	4
1.2 Constrained Clustering	5
1.2.1 Constrained clustering framework: <i>k-supplier/center</i>	7
1.2.2 Constrained clustering framework: <i>k-median/means</i>	11
1.3 Socially Fair Clustering Problem	16
1.4 Hardness of Approximation: <i>Euclidean k-Median</i>	19

1.5	Bi-Criteria Hardness of Approximation: <i>Euclidean k-Median and k-Means</i> . . .	21
1.6	Notations	23
1.7	Organization of Thesis	24
2	Tight FPT Approximation for Constrained k-Center/Supplier	25
2.1	Overview	26
2.1.1	Constrained k -supplier framework	29
2.1.2	Constrained k -supplier framework with outliers	36
2.2	Related Work	40
2.3	Notations	46
2.4	Algorithm for List Outlier k -Supplier	47
2.4.1	Bi-criteria approximation	48
2.4.2	Conversion: <i>bi-criteria approximation to list outlier k-supplier algorithm</i>	50
2.5	Partition Algorithms	53
2.5.1	Partition Algorithms: <i>r-Gather, r-Capacity, Balanced, Chromatic, Fault-Tolerant, and Strongly-Private k-Service Problems</i>	53
2.5.2	Partition Algorithms: <i>ℓ-Diversity and Fair Outlier k-Supplier Problems</i>	57
2.6	FPT Hardness: <i>k-Supplier and k-Center</i>	61
2.7	FPT Hardness: <i>Outlier k-Supplier and k-Center</i>	64

3	FPT Approximation for Constrained k-Median/Means	67
3.1	Overview	69
3.1.1	Constrained k -service framework	72
3.1.2	Constrained k -service framework with outliers	77
3.2	Related Work	82
3.3	Notations and Identities	91
3.4	A Simple List k -Service Algorithm	94
3.5	Algorithm for List Outlier k -Service Problem	97
3.5.1	Analysis for low-cost clusters	104
3.5.2	Analysis for high-cost clusters	107
3.6	Analysis of List Outlier k -Service Algorithm: <i>Special Case</i> $C \subseteq L$	114
3.6.1	Analysis for low-cost clusters	116
3.6.2	Analysis for high-cost clusters	119
3.7	A Matching Lower Bound on Approximation	122
3.8	Streaming Algorithms	124
3.9	Partition Algorithms	127
3.9.1	r -gather/ r -capacity/balanced k -service problem	127
3.9.2	Fault-tolerant k -service problem	134
3.9.3	Ordered-weighted-average (OWA) k -service problem	135

3.9.4	Chromatic and strongly private k -service problems	137
3.9.5	Uncertain k -service problem	140
3.9.6	ℓ -diversity and fair k -service problems	141
3.10	Conclusion and Open Problems	145
4	FPT Approximation for Constrained k-Median/Means: Euclidean & Outlier Setting	147
4.1	Overview	148
4.2	Related Work	152
4.3	Notations and Identities	153
4.4	Algorithm for List Outlier k -Service Problem	155
4.4.1	Analysis for low-cost clusters	158
4.4.2	Analysis for high-cost clusters	161
4.5	Streaming Algorithms	167
4.6	Conclusion	169
5	Tight FPT Approximation for Socially Fair k-Median/Means	171
5.1	Overview	172
5.2	Our Results	174
5.3	Related Work	175

5.4	Notations and Identities	177
5.5	FPT Approximation	178
5.5.1	Bi-criteria approximation	179
5.5.2	Conversion: <i>Bi-criteria to FPT approximation</i>	188
5.6	FPT Lower Bounds	191
5.7	Conclusion	195
6	Hardness of Approximation: k-Median	197
6.1	Overview	198
6.2	Related Work	201
6.3	Summary of Our Contributions	203
6.4	Notations and Useful Inequalities	207
6.5	Inapproximability of Euclidean k -Median	212
6.5.1	Completeness	214
6.5.2	Soundness	216
6.6	Vertex Cover of Non-Star Graphs	226
6.6.1	1-median cost of non-star graphs	227
6.6.2	Vertex cover for matching size two	240
6.6.3	Vertex cover for matching size at least three	243

6.7	Bi-criteria Hardness of Approximation	255
6.7.1	Bi-criteria inapproximability: <i>k-Median</i>	256
6.7.2	Bi-criteria inapproximability: <i>k-Means</i>	261
7	Conclusion and Future Work	265
	Bibliography	267
	List of Publications	285
	Biography	287

List of Figures

2.1	The flow network $G = (V, E)$ that is used in the partition algorithm of the hybrid k -supplier problem.	55
2.2	The flow network $G = (V, E)$ that is used by the partition algorithm of the fair outlier k -supplier problem.	60
3.1	An undirected weighted subgraph on $C_i \cup L_i$	123
3.2	The flow network $G = (V, E)$ that is used in the partition algorithm of the hybrid k -service problem.	138
3.3	The flow network $G = (V, E)$ that is used by the partition algorithm of the fair outlier k -service problem.	144
6.1	Adding vertices to V_G by picking both end points of every edge in M'	220
6.2	Adding vertices to V_G by picking both end points of every edge in M_P that is incident on at least two edges of U_P	221
6.3	Adding vertices to V_G on the basis of the edges in U_P that are incident on two edges of M_P	221

6.4	Adding edges to M_G by picking two blue edges each of which is incident on different vertices of a plank edge. Then, adding the remaining non-plank red edges to M_G	222
6.5	Fundamental non-star graphs: $3-P_2$, A_n , and L_n	227
6.6	Decomposition of $3-L_2$	233
6.7	Decomposition of $2-L_n$ for $n \geq 3$	234
6.8	Decomposition of any non-star graph F into the <i>fundamental</i> non-star graphs. .	236
6.9	A Bridge Graph: $L_{p,q}$, for $p, q \geq 1$	237
6.10	Decomposition of a non-star non-bridge graph F into fundamental non-star graphs.	238

List of Tables

1.1	The best-known approximation guarantees for the clustering problems.	3
1.2	The best-known approximation guarantees for the outlier clustering problems. .	4
1.3	The best known FPT time approximation guarantees for the clustering problems.	5
1.4	List of constrained k -supplier problems that we study in this work.	8
1.5	The known approximation guarantees for the Euclidean k -Median and k -Means problems.	20
2.1	List of constrained k -supplier problems with FPT time partition algorithms (see Section 2.5).	31
2.2	The known results for the constrained k -supplier/center problems with and without outliers for $z = 1$	41
3.1	List of constrained k -service problems with FPT time partition algorithms (see Section 3.9).	73
3.2	The known results for the constrained k -median/means problems with and without outliers.	85

6.1 The known approximation guarantees for the Euclidean k -Median and k -Means problems. 203