Clustering with Constraints Practical Theory for Geometric Center Based Clustering

Melanie Schmidt

3rd Workshop on Geometry and Machine Learning, Budapest

June 11th, 2018

Clustering (Workshop@SoCG 2018)

June 11th, 2018 0 / 30

Clustering

Clustering

• Unsupervised learning task:

Find structure in data in form of clusters without labels

Clustering

• Unsupervised learning task:

Find structure in data in form of clusters without labels

Clusters can take many shapes.

Clustering

Unsupervised learning task:

Find structure in data in form of clusters without labels

 Clusters can take many shapes. A quote from Jain's much cited survey [J10]: '... none of the available clustering algorithms can detect all these clusters' [in Fig 2. from [J10], redrawn below]:



[J10] Jain: 50 years beyond k-means, Pattern Recognition Letters, 2010.

Clustering

• Unsupervised learning task:

Find structure in data in form of clusters without labels

 Clusters can take many shapes. A quote from Jain's much cited survey [J10]: '... none of the available clustering algorithms can detect all these clusters' [in Fig 2. from [J10], redrawn below]:



[J10] Jain: 50 years beyond k-means, Pattern Recognition Letters, 2010.

Even worse:

- The correct output may depend on the data set
- The 'natural' clustering may even lie in the eye of the beholder

Even worse:

- The correct output may depend on the data set
- The 'natural' clustering may even lie in the eye of the beholder



Even worse:

- The correct output may depend on the data set
- The 'natural' clustering may even lie in the eye of the beholder



Even worse:

- The correct output may depend on the data set
- The 'natural' clustering may even lie in the eye of the beholder



Consequence: Tons of different objectives and algorithms.

Clustering (Workshop@SoCG 2018)

Even worse:

- The correct output may depend on the data set
- The 'natural' clustering may even lie in the eye of the beholder



This talk: Partitional Clustering, geometric setting

٥ò

- Given a point set P (and $k \in \mathbb{N}$)
- and some pairwise distances,



- Given a point set P (and $k \in \mathbb{N}$)
- and some pairwise distances,
- partition P into C_1, \ldots, C_k



- Given a point set P (and $k \in \mathbb{N}$)
- and some pairwise distances,
- partition *P* into C_1, \ldots, C_k
- to optimize some objective



- Given a point set P (and $k \in \mathbb{N}$)
- and some pairwise distances,
- partition *P* into C_1, \ldots, C_k
- to optimize some objective

- Pick *k* centers \rightsquigarrow induced partitioning
- *k*-center: maximum distance, metric
- *k*-median: sum of distances, metric
- k-means: sum of squared distances, Euclidean space R^d



- Given a point set P (and $k \in \mathbb{N}$)
- and some pairwise distances,
- partition *P* into C_1, \ldots, C_k
- to optimize some objective

- Pick *k* centers \rightsquigarrow induced partitioning
- *k*-center: maximum distance, metric
- *k*-median: sum of distances, metric
- k-means: sum of squared distances, Euclidean space R^d



- Given a point set P (and $k \in \mathbb{N}$)
- and some pairwise distances,
- partition P into C_1, \ldots, C_k
- to optimize some objective

- Pick *k* centers → induced partitioning
- k-center: maximum distance, metric
- *k*-median: sum of distances, metric
- k-means: sum of squared distances, Euclidean space R^d



- Given a point set P (and $k \in \mathbb{N}$)
- and some pairwise distances,
- partition P into C_1, \ldots, C_k
- to optimize some objective

- Pick *k* centers → induced partitioning
- k-center: maximum distance, metric
- *k*-median: sum of distances, metric
- k-means: sum of squared distances, Euclidean space R^d



- Given a point set P (and $k \in \mathbb{N}$)
- and some pairwise distances,
- partition *P* into C_1, \ldots, C_k
- to optimize some objective

- Pick *k* centers \rightsquigarrow induced partitioning
- *k*-center: maximum distance, metric
- *k*-median: sum of distances, metric
- k-means: sum of squared distances, Euclidean space R^d



- Given a point set P (and $k \in \mathbb{N}$)
- and some pairwise distances,
- partition P into C_1, \ldots, C_k
- to optimize some objective

- Pick *k* centers \rightsquigarrow induced partitioning
- *k*-center: maximum distance, metric
- k-median: sum of distances, metric
- k-means: sum of squared distances, Euclidean space R^d



- Given a point set P (and $k \in \mathbb{N}$)
- and some pairwise distances,
- partition P into C_1, \ldots, C_k
- to optimize some objective

The three 'k'-objectives

- Pick *k* centers \rightsquigarrow induced partitioning
- *k*-center: maximum distance, metric
- *k*-median: sum of distances, metric
- k-means: sum of squared distances, Euclidean space R^d

Facility location: like k-median, but opening cost instead of k



- Given a point set P (and $k \in \mathbb{N}$)
- and some pairwise distances,
- partition P into C_1, \ldots, C_k
- to optimize some objective

The three 'k'-objectives

- Pick *k* centers \rightsquigarrow induced partitioning
- *k*-center: maximum distance, metric
- *k*-median: sum of distances, metric
- k-means: sum of squared distances, Euclidean space R^d

Facility location: like k-median, but opening cost instead of k

not center-based: min sum k-clustering, aversion k-clustering

Approximation: State of the art

Approximation

Hardness

- k-center
- facility location
- k-median
- k-means
- min sum k-clustering
- aversion k-clustering

[G85] Gonzalez. Clustering to minimize the maximum intercluster distance, Theoretical Computer Science 1985.

- [GK99] Guha, Khuller. Greedy strikes back: Improved facility location algorithms, J. Algorithms 1999.
- [HN79] Hsu, Nemhauser, Easy and hard bottleneck location problems. Discrete Applied Mathematics 1979.
- [HS85] Hochbaum, Shmoys, A best possible heuristic for the *k*-center problem, Mathematics of Operations Research 1985.

[L11] Li. A 1.488-approximation algorithm for the uncapacitated facility location problem. ICALP 2011.

Approximation: State of the art

	Approximation	Hardness
k-center	2 [G85,HS85]	2 [HN79]
facility location	1.488 [L11]	1.463 [GK99]
<i>k</i> -median	3 +ε [AGK+01]	$1+2/epprox 1.74~[{ m HN79}]$
k-means	9 +ε [KMN+02]	
min sum k-clustering $O(\epsilon^{-1} \log^{1+\epsilon} n)$ [BZR01]		
aversion k-clustering		

[AGK+01] Arya, Garg, Khandekar, Meyerson, Munagala, Pandit. Local search heuristic for k-median [...], STOC 2001.

[BZR01] Bartal, Charikar, Raz. Approximating min-sum k-Clustering in metric spaces. STOC 2001.
 [G85] Gonzalez. Clustering to minimize the maximum intercluster distance, Theoretical Computer Science 1985.
 [GK99] Guha, Khuller. Greedy strikes back: Improved facility location algorithms, J. Algorithms 1999.

[HN79] Hsu, Nemhauser, Easy and hard bottleneck location problems. Discrete Applied Mathematics 1979.
 [HS85] Hochbaum, Shmoys, A best possible heuristic for the *k*-center problem, Mathematics of Operations Research 1985.
 [KMN+02] Kanungo, Mount, Netanyahu, Piatko, Silverman, Wu. A local search approx. alg. for *k*-means clustering, SoCG 2002.
 [L11] Li. A 1.488-approximation algorithm for the uncapacitated facility location problem. ICALP 2011.

Approximation: State of the art

	Approximation	Hardness
k-center	2 [G85,HS85]	2 [HN79]
facility location	1.488 [L11]	1.463 [GK99]
<i>k</i> -median	2.675 + ε [BPRST15]	1+2/ <i>e</i> ≈1.74 [HN79]
<i>k</i> -means	6.357 [ASFW17]	1.0013 [ACKS15], [LSW15]
min sum k-clustering	<i>O</i> (log <i>n</i>) [BFSS15]	
aversion k-clustering	<i>O</i> (1) [GGS16]	

[ACKS15] Awasthi, Charikar, Krishnaswamy, Sinop. The hardness of approximation of euclidean k-means, SoCG 2015. [AHSW17] Ahmadian, Norouzi-Fard, Svensson, Ward: Better Guarantees for k-Means and Euclidean k-Median by Primal-Dual Algori [BFSS15] Behsaz, Frigostad, Salavatipour, Sivakumar, Appr. Alg. for Min-Sum k-Clustering and Balanced k-Median, ICALP 2015. [BPRST15] Byrka, Pensyl, Rybicki, Srinivasan, Trinh, An Improved Approximation for k-median [...], SODA 2015. [G85] Gonzalez. Clustering to minimize the maximum intercluster distance, Theoretical Computer Science 1985. [GK99] Guha, Khuller. Greedy strikes back: Improved facility location algorithms, J. Algorithms 1999. [GGS16] Gupta, Guruganesh, S. Approximation Algorithms for Aversion k-Clustering via Local k-Median, ICALP 2016. Hsu, Nemhauser, Easy and hard bottleneck location problems. Discrete Applied Mathematics 1979. [HN79] Hochbaum, Shmoys, A best possible heuristic for the k-center problem, Mathematics of Operations Research 1985. [HS85] [LSW15] Lee, S. Wright, Improved and Simplified Inapproximability for k-means, IPL 2017. [L11] Li, A 1.488-approximation algorithm for the uncapacitated facility location problem, ICALP 2011.

Many techniques

Center based

- k-center somewhat easier
- Chain: first *k*-center or facility location, then *k*-median, and *k*-means is last

Many techniques

Center based

- k-center somewhat easier
- Chain: first *k*-center or facility location, then *k*-median, and *k*-means is last
- Embedding into tree metrics
- Local Search
- LP-based rounding
 - Filtering
 - Dual fitting
 - Randomized rounding
- Primal-dual framework





capacities















Examples for constraints

Any problems?











Examples for constraints

Any problems?











Examples for constraints

Another motivation
Another motivation

Techniques for non-center based objectives

- Reductions to center-based objectives
- min sum k-clustering: balanced k-median (cluster cost times |C_i|)
- aversion k-clustering: local k-median (centers have a radius)

Capacitated clustering

Every center c has a maximal capacity u(c) for assigned points

Capacitated clustering

Every center c has a maximal capacity u(c) for assigned points

Lower bounded clustering

It is not worthwhile to open a center c unless a lower bound $\ell(c)$ is met

Capacitated clustering

Every center c has a maximal capacity u(c) for assigned points

Lower bounded clustering

It is not worthwhile to open a center c unless a lower bound $\ell(c)$ is met

Clustering with outliers

Up to z outliers may be ignored for the cost

Capacitated clustering

Every center c has a maximal capacity u(c) for assigned points

Lower bounded clustering

It is not worthwhile to open a center c unless a lower bound $\ell(c)$ is met

Clustering with outliers

Up to z outliers may be ignored for the cost

(Exact) Fair clustering

Assume points have a sensitive attribute. Clusters shall have the same composition wrt this attribute as *P*.

Capacitated facility location

Capacitated facility location, all distances zero

Capacitated facility location, all distances zero

Should be simple?

Capacitated facility location, all distances zero

- Should be simple?
- find facilities with ≥ *n* total capacity

Capacitated facility location, all distances zero

Should be simple?

- find facilities with ≥ *n* total capacity
- minimize total opening cost

Capacitated facility location, all distances zero

Should be simple?

- find facilities with ≥ *n* total capacity
- minimize total opening cost
- ~ Minimum Knapsack Problem

Capacitated facility location, all distances zero

Should be simple?

- find facilities with ≥ n total capacity
- minimize total opening cost
- ~ Minimum Knapsack Problem

Standard LP: integrality gap

Capacitated facility location, all distances zero

Should be simple?

- find facilities with ≥ *n* total capacity
- minimize total opening cost
- ~ Minimum Knapsack Problem

Standard LP: integrality gap



Capacitated facility location, all distances zero

Should be simple?

- find facilities with ≥ *n* total capacity
- minimize total opening cost
- ~ Minimum Knapsack Problem

Standard LP: integrality gap

*i*₁ zero distance *i*₂ *n* clients, cap. n - 1, 1 client, cap. n, opening cost 0 opening cost 1 • integral solution: Assign one client to *i*₂ \rightsquigarrow cost 1 • fractional solution: Assign all clients to *i*₂ by $\frac{1}{n} \rightsquigarrow$ cost $\frac{1}{n}$

Capacitated facility location, all distances zero

Should be simple?

- find facilities with ≥ *n* total capacity
- minimize total opening cost
- ~ Minimum Knapsack Problem

Standard LP: integrality gap unbounded

*i*₁ zero distance *i*₂ *n* clients, cap. n - 1, 1 client, cap. n, opening cost 0 opening cost 1 • integral solution: Assign one client to *i*₂ \rightsquigarrow cost 1 • fractional solution: Assign all clients to *i*₂ by $\frac{1}{2} \rightarrow \cos \frac{1}{2}$

Local search based 5-approximation for cap. fac. location [BGG12]

[ABC+] An, Bhaskara, Chekuri, Gupta, Madan, Svensson. Centrality of trees for capacitated k-center, Math. Progr. 2015.
 [ASS14] An, Singh, Svensson. LP-Based Algorithms for Capacitated Facility Location, FOCS 2014.
 [BGG12] Bansal, Garg, Gupta. A 5-Approximation for Capacitated Facility Location, ESA 2012.
 [CHK12] Cvgan, Haijaghavi, Khuller, LP rounding for k-centers with non-uniform hard capacities. FOCS 2012.

- Local search based 5-approximation for cap. fac. location [BGG12]
- O(1)-approximation with knapsack based LP [ASS14]

[ABC+] An, Bhaskara, Chekuri, Gupta, Madan, Svensson. Centrality of trees for capacitated k-center, Math. Progr. 2015.
 [ASS14] An, Singh, Svensson. LP-Based Algorithms for Capacitated Facility Location, FOCS 2014.
 [BGG12] Bansal, Garg, Gupta. A 5-Approximation for Capacitated Facility Location, ESA 2012.
 [CHK12] Cvaan. Haiaahavi, Khuller. LP rounding for k-centers with non-uniform hard capacities. FOCS 2012.

- Local search based 5-approximation for cap. fac. location [BGG12]
- O(1)-approximation with knapsack based LP [ASS14]
- capacitated k-center: ≥ 3 [CHK12] and ≤ 9 [ABC+]

- [ASS14] An, Singh, Svensson. LP-Based Algorithms for Capacitated Facility Location, FOCS 2014.
- [BGG12] Bansal, Garg, Gupta. A 5-Approximation for Capacitated Facility Location, ESA 2012.
- [CHK12] Cygan, Hajiaghayi, Khuller. LP rounding for k-centers with non-uniform hard capacities, FOCS 2012.

- Local search based 5-approximation for cap. fac. location [BGG12]
- O(1)-approximation with knapsack based LP [ASS14]
- capacitated k-center: ≥ 3 [CHK12] and ≤ 9 [ABC+]
- No constant factor known for capacitated k-median or k-means

- [ASS14] An, Singh, Svensson. LP-Based Algorithms for Capacitated Facility Location, FOCS 2014.
- [BGG12] Bansal, Garg, Gupta. A 5-Approximation for Capacitated Facility Location, ESA 2012.
- [CHK12] Cygan, Hajiaghayi, Khuller. LP rounding for k-centers with non-uniform hard capacities, FOCS 2012.

- Local search based 5-approximation for cap. fac. location [BGG12]
- O(1)-approximation with knapsack based LP [ASS14]
- capacitated k-center: ≥ 3 [CHK12] and ≤ 9 [ABC+]
- No constant factor known for capacitated k-median or k-means
- Various bicriteria approximations

- [ASS14] An, Singh, Svensson. LP-Based Algorithms for Capacitated Facility Location, FOCS 2014.
- [BGG12] Bansal, Garg, Gupta. A 5-Approximation for Capacitated Facility Location, ESA 2012.
- [CHK12] Cygan, Hajiaghayi, Khuller. LP rounding for k-centers with non-uniform hard capacities, FOCS 2012.

- Local search based 5-approximation for cap. fac. location [BGG12]
- O(1)-approximation with knapsack based LP [ASS14]
- capacitated k-center: ≥ 3 [CHK12] and ≤ 9 [ABC+]
- No constant factor known for capacitated k-median or k-means
- Various bicriteria approximations
- Research on easier variants: uniform or soft capacities

- [ASS14] An, Singh, Svensson. LP-Based Algorithms for Capacitated Facility Location, FOCS 2014.
- [BGG12] Bansal, Garg, Gupta. A 5-Approximation for Capacitated Facility Location, ESA 2012.
- [CHK12] Cygan, Hajiaghayi, Khuller. LP rounding for k-centers with non-uniform hard capacities, FOCS 2012.

Lower bounds are even worse..

Lower bounds are even worse ..

Standard LP for facility location has unbounded integrality gap

Lower bounds are even worse ...

- Standard LP for facility location has unbounded integrality gap
- Standard local search approach yields arbitrarily bad solutions

Lower bounds are even worse ...

- Standard LP for facility location has unbounded integrality gap
- Standard local search approach yields arbitrarily bad solutions

Lower bounds: State of the art

82.6-approximation for uniform lower bounds and facility loc. [AS12]

[AS12] Ahmadian, Swamy: Improved approximation guarantees for lower-bounded facility location, WAOA 2012.
 [AS16] Ahmadian, Swamy: Approximation Algorithms for Clustering Problems with Lower Bounds and Outliers, ICALP 2016.

Lower bounds are even worse.. (except for *k*-center)

- Standard LP for facility location has unbounded integrality gap
- Standard local search approach yields arbitrarily bad solutions

Lower bounds: State of the art

82.6-approximation for uniform lower bounds and facility loc. [AS12]
O(1) for non-uniform bounds only known for *k*-center [AS16]

[AS12] Ahmadian, Swamy: Improved approximation guarantees for lower-bounded facility location, WAOA 2012.
 [AS16] Ahmadian, Swamy: Approximation Algorithms for Clustering Problems with Lower Bounds and Outliers, ICALP 2016.

State of the art for k-center

just k-center: 2-approx. [G85][HS86]

State of the art for k-center

- just k-center: 2-approx. [G85][HS86]
- capacitated k-center: 9-approx. [ABCGMS15]

State of the art for k-center

- just k-center: 2-approx. [G85][HS86]
- capacitated k-center: 9-approx. [ABCGMS15]
- Iower bounded k-center: 2-approx. [CGK16]

State of the art for k-center

- just k-center: 2-approx. [G85][HS86]
- capacitated k-center: 9-approx. [ABCGMS15]
- Iower bounded k-center: 2-approx. [CGK16]
- k-center with outliers: 2-approx. [CGK16]

State of the art for k-center

- just k-center: 2-approx. [G85][HS86]
- capacitated k-center: 9-approx. [ABCGMS15]
- Iower bounded k-center: 2-approx. [CGK16]
- k-center with outliers: 2-approx. [CGK16]
- (exact and balanced) fair k-center with 2 colors: 3-approx. [CKLV17]

State of the art for k-center

- just k-center: 2-approx. [G85][HS86]
- capacitated k-center: 9-approx. [ABCGMS15]
- Iower bounded k-center: 2-approx. [CGK16]
- k-center with outliers: 2-approx. [CGK16]
- (exact and balanced) fair k-center with 2 colors: 3-approx. [CKLV17]
- (exact) fair k-center with multiple colors: 14-approx.

State of the art for k-center

- just k-center: 2-approx. [G85][HS86]
- capacitated k-center: 9-approx. [ABCGMS15]
- lower bounded k-center: 2-approx. [CGK16]
- k-center with outliers: 2-approx. [CGK16]
- (exact and balanced) fair k-center with 2 colors: 3-approx. [CKLV17]
- (exact) fair k-center with multiple colors: 14-approx.

 [ABCGMS15] An, Bhaskara, Chekuri, Gupta, Madan, Svensson: Centrality of trees for capacitated k-center. Math. Prog. 2015.

 [CGK16]
 Chakrabarty, Goyal, Krishnaswamy. The non-uniform k-center problem, ICALP 2016.

 [CKLV17]
 Chierichetti, Kumar, Lattanzi, Vassilvitskii. Fair clustering through fairlets, NIPS 2017.

Solve new problems on *k*-center first.

Do we want to track all these problems?










Clustering (Workshop@SoCG 2018)

- Approximation ratios have been optimized for centuries
- Constraints require new techniques to transfer results
- What if a better algorithm for the unconstrained problem appears?
- And what if we want combinations of constraints?

- Approximation ratios have been optimized for centuries
- Constraints require new techniques to transfer results
- What if a better algorithm for the unconstrained problem appears?
- And what if we want combinations of constraints?

→ endless chain of improvements and new problem combinations

- Approximation ratios have been optimized for centuries
- Constraints require new techniques to transfer results
- What if a better algorithm for the unconstrained problem appears?
- And what if we want combinations of constraints?

→ endless chain of improvements and new problem combinations

Question

Can we add constraints?

- Approximation ratios have been optimized for centuries
- Constraints require new techniques to transfer results
- What if a better algorithm for the unconstrained problem appears?
- And what if we want combinations of constraints?

→ endless chain of improvements and new problem combinations

Question

Can we add constraints?

Given an approximation algorithm as a subroutine, can we establish an additional constraint?

Goal

Add uniform lower bound *L* to *k*-center with outliers

Goal

Add uniform lower bound *L* to *k*-center with outliers

What makes *k*-center easy?

Goal

Add uniform lower bound *L* to *k*-center with outliers

What makes k-center easy?

Goal

Add uniform lower bound *L* to *k*-center with outliers

What makes *k*-center easy?

Threshold graph

• Value of the optimal solution is a pairwise distance!

Goal

Add uniform lower bound *L* to *k*-center with outliers

What makes *k*-center easy?

Threshold graph

• Value of the optimal solution is a pairwise distance!

• There are $\Theta(n^2)$ pairwise distances \rightsquigarrow can guess optimum value τ

Goal

Add uniform lower bound *L* to *k*-center with outliers

What makes k-center easy?

- Value of the optimal solution is a pairwise distance!
- There are $\Theta(n^2)$ pairwise distances \rightsquigarrow can guess optimum value τ
- Allows to build threshold graph

Goal

Add uniform lower bound *L* to *k*-center with outliers

What makes k-center easy?

- Value of the optimal solution is a pairwise distance!
- There are $\Theta(n^2)$ pairwise distances \rightsquigarrow can guess optimum value τ
- Allows to build threshold graph
- At most two hops between two points in the same optimum cluster









Clustering (Workshop@SoCG 2018)

0







Threshold graph



Approximation in [HS86]

Guess correct threshold

 \circ

Consider two-hop graph



Approximation in [HS86]

- Guess correct threshold
- Consider two-hop graph
- Find maximal independent set

Add uniform lower bound L to k-center with outliers

Add uniform lower bound L to k-center with outliers

Add uniform lower bound *L* to *k*-center with outliers

Main idea

• Use subroutine for approximating k-center with outliers

Add uniform lower bound L to k-center with outliers

- Use subroutine for approximating *k*-center with outliers
- Build a network to move points using the threshold idea

Add uniform lower bound *L* to *k*-center with outliers

- Use subroutine for approximating *k*-center with outliers
- Build a network to move points using the threshold idea
- One of two outcomes:

Add uniform lower bound *L* to *k*-center with outliers

- Use subroutine for approximating k-center with outliers
- Build a network to move points using the threshold idea
- One of two outcomes:
 - We can successfully establish the lower bound L or

Add uniform lower bound *L* to *k*-center with outliers

- Use subroutine for approximating k-center with outliers
- Build a network to move points using the threshold idea
- One of two outcomes:
 - We can successfully establish the lower bound L or
 - We find a set P' that we recluster

Step 1: Use subroutine to get initial solution



Step 1: Use subroutine to get initial solution



Step 1: Use subroutine to get initial solution


Step 2: Compute threshold edges



















Step 4: Compute an integral maximum flow



Step 4: Compute an integral maximum flow



Step 5: Compute reachable points 2 0 • S

k = 4 centers z = 4 outliers L = 5 lower bound

Step 5: Compute reachable points 2 50 • Ċ . Or S

set P"

٥ŧ

k = 4 centers

z = 4 outliers L = 5 lower bound Step 5: Compute reachable points and clusters



Step 5: Compute reachable points and clusters











Step 4 (again): Compute an integral maximum flow



Step 6: Success! Move the points wrt the maximum flow



Points are only moved once, at the very end!

- Points are only moved once, at the very end!
- Recomputations do not increase the factor

- Points are only moved once, at the very end!
- Recomputations do not increase the factor
- After every flow computation, number of clusters (or number of outliers) decreases → ≤ k ⋅ z flow computations

- Points are only moved once, at the very end!
- Recomputations do not increase the factor
- After every flow computation, number of clusters (or number of outliers) decreases → ≤ k · z flow computations

Theorem (Rösner, S., 2018)

Let *A* be an α -approximation algorithm for *k*-center with outliers. Then we can compute an $(\alpha + 2)$ -approximation for *k*-center with outliers and lower bound *L* in polynomial time.

Clustering (Workshop@SoCG 2018)











Reduce increase of approximation factor

- Reduce increase of approximation factor
- Non-uniform lower bounds

- Reduce increase of approximation factor
- Non-uniform lower bounds
- Other constraints. Fairness works, what about outliers?
Follow-up questions

- Reduce increase of approximation factor
- Non-uniform lower bounds
- Other constraints. Fairness works, what about outliers?
- k-median? k-means?



Quick review for k-means



- Quick review for k-means
- What about constraints?



- Quick review for k-means
- What about constraints?



- Quick review for k-means
- What about constraints?



- Quick review for k-means
- What about constraints?



- Quick review for k-means
- What about constraints?



Coresets / Dim. Red

Reduce complexity,

but approximately preserve k-means cost for all sets of k centers

Coresets / Dim. Red

Reduce complexity,

but approximately preserve k-means cost for all sets of k centers

Definition

Let $P \subset \mathbb{R}^d$ be a point set. If $P' \subset \mathbb{R}^d$

sat. $\forall C \subset \mathbb{R}^d, |C| = k$:

 $|\operatorname{cost}(P,C) - \operatorname{cost}(P',C) | \le \varepsilon \cdot \operatorname{cost}(P,C)$

and *P*′...

Coresets / Dim. Red

Reduce complexity,

but approximately preserve k-means cost for all sets of k centers

Definition

Let $P \subset \mathbb{R}^d$ be a point set. If $P' \subset \mathbb{R}^d$

sat. $\forall C \subset \mathbb{R}^d, |C| = k$:

 $|\operatorname{cost}(P,C) - \operatorname{cost}(P',C) | \le \varepsilon \cdot \operatorname{cost}(P,C)$

and *P*′...

• ... satisfies that $|P'| \ll |P|$

Coresets / Dim. Red

Reduce complexity,

but approximately preserve k-means cost for all sets of k centers

Definition

Let $P \subset \mathbb{R}^d$ be a point set. If $P' \subset \mathbb{R}^d$

sat. $\forall C \subset \mathbb{R}^d, |C| = k$:

 $|\operatorname{cost}(P,C) - \operatorname{cost}(P',C) | \le \varepsilon \cdot \operatorname{cost}(P,C)$

and *P*′...

• . . . satisfies that $|P'| \ll |P| \rightsquigarrow$ then P' is a coreset

Coresets / Dim. Red

Reduce complexity,

but approximately preserve k-means cost for all sets of k centers

Definition

Let $P \subset \mathbb{R}^d$ be a point set. If $P' \subset \mathbb{R}^d$

sat. $\forall C \subset \mathbb{R}^d, |C| = k$:

 $|\operatorname{cost}(P,C) - \operatorname{cost}(P',C) | \le \varepsilon \cdot \operatorname{cost}(P,C)$

and *P*′...

- . . . satisfies that $|P'| \ll |P| \rightsquigarrow$ then P' is a coreset
- . . . has intrinsic dimension s for an $s \ll d$

Coresets / Dim. Red

Reduce complexity,

but approximately preserve k-means cost for all sets of k centers

Definition

Let $P \subset \mathbb{R}^d$ be a point set. If $P' \subset \mathbb{R}^d$

sat.
$$\forall C \subset \mathbb{R}^d, |C| = k$$
:

 $|\operatorname{cost}(P,C) - \operatorname{cost}(P',C) | \le \varepsilon \cdot \operatorname{cost}(P,C)$

and *P*′...

- ... satisfies that $|P'| \ll |P| \rightsquigarrow$ then P' is a coreset
- . . . has intrinsic dimension s for an $s \ll d \rightsquigarrow \dim$. reduction

Coresets / Dim. Red

Reduce complexity,

but approximately preserve k-means cost for all sets of k centers

Definition

Let $P \subset \mathbb{R}^d$ be a point set. If $P' \subset \mathbb{R}^d$ and $\Delta \in \mathbb{R}$ sat. $\forall C \subset \mathbb{R}^d, |C| = k$:

 $|\operatorname{cost}(P, C) - \operatorname{cost}(P', C) + \Delta| \leq \varepsilon \cdot \operatorname{cost}(P, C)$

and *P*′...

• ... satisfies that $|P'| \ll |P| \rightsquigarrow$ then P' is a coreset

• ... has intrinsic dimension s for an $s \ll d \rightsquigarrow \dim$ reduction

Coresets / Dim. Red

Reduce complexity,

but approximately preserve k-means cost for all sets of k centers

Definition

Let $P \subset \mathbb{R}^d$ be a point set. If $P' \subset \mathbb{R}^d$ and $\Delta \in \mathbb{R}$ sat. $\forall C \subset \mathbb{R}^d, |C| = k$:

$$|\operatorname{cost}(P, C) - \operatorname{cost}(P', C) + \Delta| \le \varepsilon \cdot \operatorname{cost}(P, C)$$

and *P*′...

- ... satisfies that $|P'| \ll |P| \rightsquigarrow$ then P' is a coreset
- . . . has intrinsic dimension s for an $s \ll d \rightsquigarrow \dim$. reduction

Input and coreset look alike for every possible solution!

Clustering (Workshop@SoCG 2018)

Coresets for *k*-means

Dimensionality reduction for k-means

 [CMEP15]
 Cohen, Elder, Musco, Musco, Persu, Dim. reduction for k-means clustering and low rank approximation, STOC 2015.

 [DFKVV04]
 Drineas, Frieze, Kannan, Vempala, Vinay, Clustering large graphs via the svd, Maschine Learning 2004.

 [FL11]
 Feldman, Langberg, A unified framework for approximating and clustering data, STOC 2011.

 [FSS13]
 Feldman, S., Sohler, Turning Big Data into Tiny Data: Constant-size C. for k-means, PCA and Pr. Cl., SODA 2013.

 [HPM04]
 Har-Peled, Mazumdar, On coresets for k-means and k-median clustering, STOC 2004.

[JL84] Johnson, Lindenstrauss, Extensions of lipschitz mappings into a hilbert space, Contemporary Mathematics 1984.

Clustering (Workshop@SoCG 2018)

Coresets for *k*-means

	Number of points	Remarks
[HPM04]	$\mathcal{O}(k\varepsilon^{-d}\log n)$	First coreset
[FL11]	$\mathcal{O}(\textit{dk}arepsilon^{-4})$	Independent of log n
[FSS13], [CMEP15]	$\mathcal{O}(k^2 \varepsilon^{-5})$	Independent of <i>d</i> , log <i>n</i>

Dimensionality reduction for *k*-means

 [CMEP15]
 Cohen, Elder, Musco, Musco, Persu, Dim. reduction for k-means clustering and low rank approximation, STOC 2015.

 [DFKVV04]
 Drineas, Frieze, Kannan, Vempala, Vinay, Clustering large graphs via the svd, Maschine Learning 2004.

 [FL11]
 Feldman, Langberg, A unified framework for approximating and clustering data, STOC 2011.

 [FSS13]
 Feldman, S., Sohler, Turning Big Data into Tiny Data: Constant-size C. for k-means, PCA and Pr. Cl., SODA 2013.

 [HPM04]
 Har-Peled, Mazumdar, On coresets for k-means and k-median clustering, STOC 2004.

[JL84] Johnson, Lindenstrauss, Extensions of lipschitz mappings into a hilbert space, Contemporary Mathematics 1984.

Coresets for *k*-means

	Number of points	Remarks
[HPM04]	$\mathcal{O}(k\varepsilon^{-d}\log n)$	First coreset
[FL11]	$\mathcal{O}(\textit{dk}arepsilon^{-4})$	Independent of log n
[FSS13], [CMEP15]	$\mathcal{O}(k^2 \varepsilon^{-5})$	Independent of <i>d</i> , log <i>n</i>

Dimensionality reduction for *k*-means

	Target dimension	Quality quarantee
[JL84]	$\Theta(\varepsilon^{-2} \log n)$	any $\varepsilon \in (0, 1) \Rightarrow 1 + \varepsilon$
[DFKVV04]	k	only $\varepsilon = 1 \Rightarrow 2$
[FSS13], [CMEP15]	$\lceil \mathbf{k} / \varepsilon \rceil$	any $\varepsilon \in (0, 1] \Rightarrow 1 + \varepsilon$

[CMEP15] Cohen, Elder, Musco, Musco, Persu, Dim. reduction for k-means clustering and low rank approximation, STOC 2015. [DFKVV04] Drineas, Frieze, Kannan, Vempala, Vinay, Clustering large graphs via the svd, Maschine Learning 2004.

[FL11] Feldman, Langberg, A unified framework for approximating and clustering data, STOC 2011.

[FSS13] Feldman, S., Sohler, Turning Big Data into Tiny Data: Constant-size C. for k-means, PCA and Pr. Cl., SODA 2013. [HPM04] Har-Peled, Mazumdar, On coresets for k-means and k-median clustering, STOC 2004.

[JL84] Johnson, Lindenstrauss, Extensions of lipschitz mappings into a hilbert space, Contemporary Mathematics 1984.

Conversion to a Streaming Algorithm: Merge & Reduce [BS80]



- read data in blocks
- compute a coreset for each block $\rightarrow s$
- merge and reduce coresets
 - in a tree fashion
- \rightsquigarrow space $s \cdot \log n$

Morally: Merge & Reduce causes overhead $O(\log^c n)$

[BS80] Bentley, Saxe, Decomposable searching problems I: Static-to-dynamic transformation, Journal of Algorithms, 1980.

Conversion to a Streaming Algorithm: Merge & Reduce [BS80]



- read data in blocks
- compute a coreset for each block $\rightarrow s$
- merge and reduce coresets
 - in a tree fashion
- \rightsquigarrow space $s \cdot \log n$

Morally: Merge & Reduce causes overhead $\mathcal{O}(\log^c n)$

Necessary Requirement: Composability

 P_1' and P_2' coresets for P_1 and $P_2 \Rightarrow P_1' \cup P_2'$ coreset for $P_1 \cup P_2$

[BS80] Bentley, Saxe, Decomposable searching problems I: Static-to-dynamic transformation, Journal of Algorithms, 1980.

For *k*-means, coresets are composable because the cost is linear:

 $cost(P_1 \cup P_2, C) = cost(P_1, C) + cost(P_2, C) \quad \forall P_1, P_2, C$

For *k*-means, coresets are composable because the cost is linear:

 $cost(P_1 \cup P_2, C) = cost(P_1, C) + cost(P_2, C) \quad \forall P_1, P_2, C$

For problems with constraints, cost can increase or decrease!

For *k*-means, coresets are composable because the cost is linear:

$$\operatorname{cost}(P_1 \cup P_2, C) = \operatorname{cost}(P_1, C) + \operatorname{cost}(P_2, C) \quad \forall P_1, P_2, C$$

For problems with constraints, cost can increase or decrease!



For *k*-means, coresets are composable because the cost is linear:

$$\operatorname{cost}(P_1 \cup P_2, C) = \operatorname{cost}(P_1, C) + \operatorname{cost}(P_2, C) \quad \forall P_1, P_2, C$$

For problems with constraints, cost can increase or decrease!



For *k*-means, coresets are composable because the cost is linear:

$$\operatorname{cost}(P_1 \cup P_2, C) = \operatorname{cost}(P_1, C) + \operatorname{cost}(P_2, C) \quad \forall P_1, P_2, C$$

For problems with constraints, cost can increase or decrease!



For *k*-means, coresets are composable because the cost is linear:

$$\operatorname{cost}(P_1 \cup P_2, C) = \operatorname{cost}(P_1, C) + \operatorname{cost}(P_2, C) \quad \forall P_1, P_2, C$$

For problems with constraints, cost can increase or decrease!





• Composable coresets for fair k-means

- Composable coresets for fair k-means
- Size is $\mathcal{O}(k\varepsilon^{-d}\log n)$

- Composable coresets for fair k-means
- Size is $\mathcal{O}(k\varepsilon^{-d}\log n)$
- Based on geometric coreset construction

- Composable coresets for fair k-means
- Size is $\mathcal{O}(k\varepsilon^{-d}\log n)$
- Based on geometric coreset construction

Final open questions

- Composable coresets for fair k-means
- Size is $\mathcal{O}(k\varepsilon^{-d}\log n)$
- Based on geometric coreset construction

Final open questions

Composable coresets for other constraints?

- Composable coresets for fair k-means
- Size is $\mathcal{O}(k\varepsilon^{-d}\log n)$
- Based on geometric coreset construction

Final open questions

- Composable coresets for other constraints?
- Smaller size by using sampling?

- Composable coresets for fair k-means
- Size is $\mathcal{O}(k\varepsilon^{-d}\log n)$
- Based on geometric coreset construction

Final open questions

- Composable coresets for other constraints?
- Smaller size by using sampling?

Thank you for your attention!