# Hardness of Approximation for Metric Clustering

## Karthik C. S.

(**R**utgers **U**niversity)

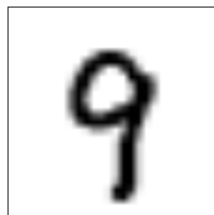March 4th 2022

$28 \times 28$
grayscale image

Task of Classifying Input Data

◎ Reveal internal structure of data

    ○ Clustering gene expression

◎ Reveal internal structure of data
  ○ Clustering gene expression

◎ Partition data
  ○ Market segmentation

◎ Reveal internal structure of data
  ○ Clustering gene expression

◎ Partition data
  ○ Market segmentation

◎ Data Preparation
  ○ Summarize news

◎ Reveal internal structure of data
  ○ Clustering gene expression

◎ Partition data
  ○ Market segmentation

◎ Data Preparation
  ○ Summarize news

◎ Data Exploration
  ○ Underlying rules and Reoccurring patterns

◎ $(\Gamma, \Delta)$ is a metric space

◎ $(\Gamma, \Delta)$ is a metric space

◎ <u>Input:</u> $X \subseteq \Gamma$, $k \in \mathbb{N}$

◎ $(\Gamma, \Delta)$ is a metric space

◎ Input: $X \subseteq \Gamma$, $k \in \mathbb{N}$

◎ Output: A classification $(C, \sigma)$:

◎ $(\Gamma, \Delta)$ is a metric space

◎ <u>Input</u>: $X \subseteq \Gamma$, $k \in \mathbb{N}$

◎ <u>Output</u>: A classification $(C, \sigma)$:

  ○ $C \subseteq \Gamma$ and $|C| = k$

◎ $(\Gamma, \Delta)$ is a metric space

◎ <u>Input:</u> $X \subseteq \Gamma$, $k \in \mathbb{N}$

◎ <u>Output:</u> A classification $(C, \sigma)$:

    ○ $C \subseteq \Gamma$ and $|C| = k$

    ○ $\sigma : X \to C$

◎ $(\Gamma, \Delta)$ is a metric space

◎ Input: $X \subseteq \Gamma$, $k \in \mathbb{N}$

◎ Output: A classification $(C, \sigma)$:

    ○ $C \subseteq \Gamma$ and $|C| = k$

    ○ $\sigma : X \rightarrow C$

    ○ $\sigma$ is *good*

# **Continuous Version**

◎ $(\Gamma, \Delta)$ is a metric space

◎ <u>Input:</u> $X \subseteq \Gamma, k \in \mathbb{N}$

◎ <u>Output:</u> A classification $(C, \sigma)$:

  ○ $C \subseteq \Gamma$ and $|C| = k$

  ○ $\sigma : X \to C$

  ○ $\sigma$ is *good*

4

# Discrete
# ~~Continuous~~ Version

◎ $(\Gamma, \Delta)$ is a metric space

◎ Input: $X \subseteq \Gamma$, $k \in \mathbb{N}$

◎ Output: A classification $(C, \sigma)$:

  ○ $C \subseteq \Gamma$ and $|C| = k$

  ○ $\sigma : X \to C$

  ○ $\sigma$ is *good*

# Discrete ~~Continuous~~ Version

◎ $(\Gamma, \Delta)$ is a metric space

◎ Input: $X \subseteq \Gamma$, $k \in \mathbb{N}$ and $\mathcal{S} \subseteq \Gamma$

◎ Output: A classification $(C, \sigma)$:

    ○ $C \subseteq \overset{\mathcal{S}}{\cancel{X}}$ and $|C| = k$

    ○ $\sigma : X \to C$

    ○ $\sigma$ is *good*

◎ $k$-means, $k$-median, $k$-center, min-sum, correlation clustering . . .

## What is Good Classification?

◎ $k$-means, $k$-median, $k$-center, min-sum, correlation clustering . . .

◎ $k$-center value of $(C, \sigma)$

$$\max_{x \in X} \Delta(x, \sigma(x))$$

## What is Good Classification?

◎ $k$-means, $k$-median, $k$-center, min-sum, correlation clustering ...

◎ $k$-center value of $(C, \sigma)$

$$\max_{x \in X} \Delta(x, \sigma(x))$$

◎ $k$-median value of $(C, \sigma)$

$$\sum_{x \in X} \Delta(x, \sigma(x))$$

## What is Good Classification?

◎ $k$-means, $k$-median, $k$-center, min-sum, correlation clustering ...

◎ $k$-center value of $(C, \sigma)$

$$\max_{x \in X} \Delta(x, \sigma(x))$$

◎ $k$-median value of $(C, \sigma)$

$$\sum_{x \in X} \Delta(x, \sigma(x))$$

◎ $k$-means value of $(C, \sigma)$

$$\sum_{x \in X} \Delta(x, \sigma(x))^2$$

## What is Good Classification?

◎ $k$-means, $k$-median, $k$-center, min-sum, correlation clustering ...

◎ $k$-center value of $(C, \sigma)$

$$\max_{x \in X} \Delta(x, \sigma(x))$$

◎ $k$-median value of $(C, \sigma)$

$$\sum_{x \in X} \Delta(x, \sigma(x))$$

◎ $k$-means value of $(C, \sigma)$

$$\sum_{x \in X} \Delta(x, \sigma(x))^2$$

◎ Don't fit: Facility Location, Hierarchical Clustering ...

Given $(X, S, k)$ as input find a classification $(C, \sigma)$ that minimizes the Clustering objective

Given $(X, S, k)$ as input find a classification $(C, \sigma)$ that
minimizes the Clustering objective

Clustering Problem for objective $\Lambda$

Given $(X, S, k)$ as input find a classification $(C, \sigma)$ that
minimizes the Clustering objective

Clustering Problem for objective $\Lambda$

Yes: There is classification $(C^*, \sigma^*)$, such that $\Lambda(X, \sigma^*) \leq \beta$

Given $(X, S, k)$ as input find a classification $(C, \sigma)$ that minimizes the Clustering objective

Clustering Problem for objective $\Lambda$

Yes: There is classification $(C^*, \sigma^*)$, such that $\Lambda(X, \sigma^*) \leq \beta$

No: For all classification $(C, \sigma)$, we have $\Lambda(X, \sigma) > \beta$

NP-Hard

Efficient Approximation

NP-Hard to Approximate

◎ Many important problems are not <u>tractable</u>

◎ Many important problems are not tractable

◎ Need to cope with the intractability

◎ Many important problems are not tractable

◎ Need to cope with the intractability

◎ Design algorithms that find solutions whose cost or value is close to the optimum

◎ Many important problems are not <u>tractable</u>

◎ Need to cope with the intractability

◎ Design algorithms that find solutions whose cost or value is close to the optimum

◎ For some fundamental problems finding good approximate solutions is as hard as finding optimal solutions

◎ Many important problems are not <u>tractable</u>

◎ Need to cope with the intractability

◎ Design algorithms that find solutions whose cost or value is close to the optimum

◎ For some fundamental problems finding good approximate solutions is as hard as finding optimal solutions

◎ Area studying such results: Hardness of Approximation

Given $(X, S, k)$ as input find a classification $(C, \sigma)$ that approximately minimizes the Clustering objective

Given $(X, S, k)$ as input find a classification $(C, \sigma)$ that
approximately minimizes the Clustering objective

Clustering Problem for objective $\Lambda$

Given $(X, S, k)$ as input find a classification $(C, \sigma)$ that approximately minimizes the Clustering objective

Clustering Problem for objective $\Lambda$

Yes: There is classification $(C^*, \sigma^*)$, such that $\Lambda(X, \sigma^*) \leq \beta$

Given $(X, S, k)$ as input find a classification $(C, \sigma)$ that approximately minimizes the Clustering objective

### Clustering Problem for objective $\Lambda$

Yes: There is classification $(C^*, \sigma^*)$, such that $\Lambda(X, \sigma^*) \leq \beta$

No: For all classification $(C, \sigma)$, we have $\Lambda(X, \sigma) > (1 + \delta) \cdot \beta$

*k*-center

◎ Input: $X, S \subseteq \mathbb{R}^d$, $k \in \mathbb{N}$

# $k$-center modeling

◎ <u>Input:</u> $X, S \subseteq \mathbb{R}^d$, $k \in \mathbb{N}$

◎ <u>Output:</u> A classification $(C, \sigma)$:

  ◦ $C \subseteq S$ and $|C| = k$

  ◦ $\sigma : X \to C$

  ◦ $(C, \sigma)$ minimizes $\max_{x \in X} \|x - \sigma(x)\|_p$

◎ NP-hard [FPT81]

# State-of-the-art: General Metrics

◎ NP-hard [FPT81]

◎ Poly Time 3-approximation (Gonzalez Algorithm)

◎ NP-hard [FPT81]

◎ Poly Time 3-approximation (Gonzalez Algorithm)

◎ NP-Hard to approximate to $3 - o(1)$ factor! [FPT81]

# State-of-the-art: $\ell_p$ Metrics

◎ $\ell_1$ and $\ell_\infty$ metrics

- Poly Time 3-approximation

◎ $\ell_1$ and $\ell_\infty$ metrics

  ○ Poly Time 3-approximation

  ○ NP-Hard to approximate to $3 - o(1)$ factor! [FG88]

◎ Euclidean metric

  ○ Poly Time 2.74-approximation! [NSS13]

- $\ell_1$ and $\ell_\infty$ metrics
  - Poly Time 3-approximation
  - NP-Hard to approximate to $3 - o(1)$ factor! [FG88]
- Euclidean metric
  - Poly Time 2.74-approximation! [NSS13]
  - NP-Hard to approximate to 2.65 factor [FG88]

**Max Coverage**:

◎ <u>Input</u>: Universe and Collection of Subsets $(U, \mathcal{S}, k)$

**Max Coverage**:

◎ Input: Universe and Collection of Subsets $(U, \mathcal{S}, k)$

◎ Objective: Max Fraction of $U$ covered by $k$ subsets in $\mathcal{S}$

**Max Coverage**:

◎ Input: Universe and Collection of Subsets $(U, \mathcal{S}, k)$

◎ Objective: Max Fraction of $U$ covered by $k$ subsets in $\mathcal{S}$

### Theorem (Karp'72)

It is NP-hard to distinguish:

**Max Coverage**:

◎ Input: Universe and Collection of Subsets $(U, \mathcal{S}, k)$

◎ Objective: Max Fraction of $U$ covered by $k$ subsets in $\mathcal{S}$

### Theorem (Karp'72)

It is NP-hard to distinguish:

YES: Max Coverage is 1

**Max Coverage**:

⊚ Input: Universe and Collection of Subsets $(U, \mathcal{S}, k)$

⊚ Objective: Max Fraction of $U$ covered by $k$ subsets in $\mathcal{S}$

### Theorem (Karp'72)

It is NP-hard to distinguish:

YES: Max Coverage is 1

NO: Max Coverage is $< 1$

## Proof Overview: General Metrics

### Theorem (Karp'72)

It is NP-hard to distinguish:

### Theorem (Karp'72)

It is NP-hard to distinguish:

YES: Max Coverage is 1

## Theorem (Karp'72)

It is NP-hard to distinguish:

YES: Max Coverage is 1

NO: Max Coverage is < 1

$$\Downarrow$$

**Theorem (Karp'72)**

It is NP-hard to distinguish:

YES: Max Coverage is 1

NO: Max Coverage is < 1

$$\Downarrow$$

**Theorem (Fowler-Paterson-Tanimoto'81)**

Fix $\varepsilon > 0$. Given input $(X, S, k)$. It is NP-hard to distinguish:

YES: There exists $(C^*, \sigma^*)$ such that $\max_{x \in X} \Delta(x, \sigma^*(x)) \leq 1$

NO: For all $(C, \sigma)$ we have $\max_{x \in X} \Delta(x, \sigma(x)) \geq 3$
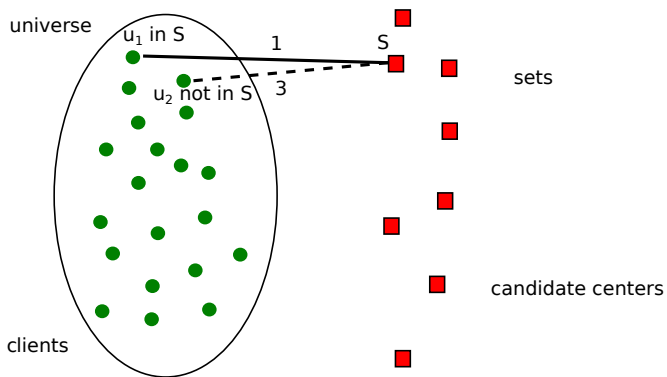
Theorem (Fowler-Paterson-Tanimoto'81)

Given input $(X, S, k)$. It is NP-hard to distinguish:

YES: There exists $(C^*, \sigma^*)$ such that $\max_{x \in X} \Delta(x, \sigma^*(x)) \leq 1$

NO: For all $(C, \sigma)$ we have $\max_{x \in X} \Delta(x, \sigma(x)) \geq 3$

$k$-means & $k$-median

# $k$-means and $k$-median modeling

◎ <u>Input:</u> $X, S \subseteq \mathbb{R}^d$, $k \in \mathbb{N}$

# $k$-means and $k$-median modeling

◎ <u>Input</u>: $X, S \subseteq \mathbb{R}^d$, $k \in \mathbb{N}$

◎ <u>Output</u>: A classification $(C, \sigma)$:

- $C \subseteq S$ and $|C| = k$

- $\sigma : X \to C$

- $k$-means: $(C, \sigma)$ minimizes $\sum_{x \in X} \|x - \sigma(x)\|_p^2$

- $k$-median: $(C, \sigma)$ minimizes $\sum_{x \in X} \|x - \sigma(x)\|_p$

◎ NP-hard when $k = 2$ (Dasgupta'07)

◎ NP-hard when $k = 2$ (Dasgupta'07)

◎ NP-hard in Euclidean plane
   (Megiddo–Supowit'84,
   Mahajan–Nimbhorkar–Varadarajan'12)

## Exact Computation

◎ NP-hard when $k = 2$ (Dasgupta'07)

◎ NP-hard in Euclidean plane
(Megiddo–Supowit'84,
Mahajan–Nimbhorkar–Varadarajan'12)

◎ W[2]-hard in general metric (Guha-Khuller'99)

◎ General metric: $k$-means $\geq 9$
(Ahmadian–Norouzi-Fard–Svensson–Ward'17)

## Approximation Algorithms

◉ General metric:  $k$-means $\geq 9$
 (Ahmadian–Norouzi-Fard–Svensson–Ward'17)

◉ General metric:  $k$-median $\geq 2.67$
 (Byrka–Pensyl–Rybicki–Srinivasan–Trinh'17)

◎ General metric:  $k$-means $\geq 9$
  (Ahmadian–Norouzi-Fard–Svensson–Ward'17)

◎ General metric:  $k$-median $\geq 2.67$
  (Byrka–Pensyl–Rybicki–Srinivasan–Trinh'17)

◎ Euclidean metric $k$-means:

## Approximation Algorithms

- ◎ General metric: $k$-means $\geq 9$
  (Ahmadian–Norouzi-Fard–Svensson–Ward'17)

- ◎ General metric: $k$-median $\geq 2.67$
  (Byrka–Pensyl–Rybicki–Srinivasan–Trinh'17)

- ◎ Euclidean metric $k$-means:

  - ○ Poly time approximation $\approx 6.357$
    (Ahmadian–Norouzi-Fard–Svensson–Ward'17)

## Approximation Algorithms

◎ General metric:   $k$-means $\geq 9$
  (Ahmadian–Norouzi-Fard–Svensson–Ward'17)

◎ General metric:   $k$-median $\geq 2.67$
  (Byrka–Pensyl–Rybicki–Srinivasan–Trinh'17)

◎ Euclidean metric $k$-means:

  ○ Poly time approximation $\approx 6.357$
    (Ahmadian–Norouzi-Fard–Svensson–Ward'17)

  ○ Fixed Dimension: PTAS (Cohen-Addad'18)

## Approximation Algorithms

◎ General metric:  $k$-means $\geq 9$
   (Ahmadian–Norouzi-Fard–Svensson–Ward'17)

◎ General metric:  $k$-median $\geq 2.67$
   (Byrka–Pensyl–Rybicki–Srinivasan–Trinh'17)

◎ Euclidean metric $k$-means:

   ○ Poly time approximation $\approx 6.357$
     (Ahmadian–Norouzi-Fard–Svensson–Ward'17)

   ○ Fixed Dimension: PTAS (Cohen-Addad'18)

   ○ Fixed $k$: PTAS (Kumar–Sabharwal–Sen'10)

# Hardness of Approximation

<u>Discrete Version</u>:

Discrete Version:

◎ General metric: $k$-means $\approx 3.94$, $k$-median $\approx 1.74$
(Guha-Khuller'99)

## Hardness of Approximation

<u>Discrete Version</u>:

- ◎ General metric: $k$-means $\approx 3.94$, $k$-median $\approx 1.74$
  (Guha-Khuller'99)
- ◎ $\ell_2$-metric: $k$-means $\ll 1.01$, $k$-median $\ll 1.01$
  (Trevisan'00)
- ◎ $\ell_1$-metric: $k$-means $\ll 1.01$, $k$-median $\ll 1.01$
  (Trevisan'00)

Discrete Version:

- ◎ General metric: $k$-means $\approx 3.94$, $k$-median $\approx 1.74$
  (Guha-Khuller'99)
- ◎ $\ell_2$-metric: $k$-means $\ll 1.01$, $k$-median $\ll 1.01$
  (Trevisan'00)
- ◎ $\ell_1$-metric: $k$-means $\ll 1.01$, $k$-median $\ll 1.01$
  (Trevisan'00)
- ◎ $\ell_\infty$-metric: $k$-means $\ll 1.01$, $k$-median $\ll 1.01$
  (Guruswami-Indyk'03)

<u>Discrete Version</u>:

- ◎ General metric: $k$-means $\approx 3.94$, $k$-median $\approx 1.74$
  (Guha-Khuller'99)
- ◎ $\ell_2$-metric: $k$-means $\ll 1.01$, $k$-median $\ll 1.01$
  (Trevisan'00)
- ◎ $\ell_1$-metric: $k$-means $\ll 1.01$, $k$-median $\ll 1.01$
  (Trevisan'00)
- ◎ $\ell_\infty$-metric: $k$-means $\ll 1.01$, $k$-median $\ll 1.01$
  (Guruswami-Indyk'03)

<u>Continuous Version</u>:

$k$-means in Euclidean metric $< 1.0013$
(Lee-Schmidt-Wright'17)

Discrete Version:

◎ General metric: $k$-means $\approx 3.94$, $k$-median $\approx 1.74$ (Guha-Khuller'99)

◎ $\ell_2$-metric: $k$-means $\ll$ ~~1.01~~ 1.73, $k$-median $\ll$ ~~1.01~~ 1.27 (Trevisan'00)

◎ $\ell_1$-metric: $k$-means $\ll$ ~~1.01~~ 3.94, $k$-median $\ll$ ~~1.01~~ 1.73 (Trevisan'00)

◎ $\ell_\infty$-metric: $k$-means $\ll$ 1.~~01~~ 3.94, $k$-median $\ll$ ~~1.01~~ 1.73 (Guruswami-Indyk'03)

Continuous Version:

$k$-means in Euclidean metric $<$ ~~1.0013~~ 1.36
(Lee-Schmidt-Wright'17)

Discrete Version:

- ◎ General metric: $k$-means $\approx 3.94$, $k$-median $\approx 1.74$
  (Guha-Khuller'99)
- ◎ $\ell_2$-metric: $k$-means $\ll$ 1.73, 1.17 ~~1.01~~, $k$-median $\ll$ 1.27, 1.06 ~~1.01~~
  (Trevisan'00)
- ◎ $\ell_1$-metric: $k$-means $\ll$ 3.94, 1.56 ~~1.01~~, $k$-median $\ll$ 1.73, 1.14 ~~1.01~~
  (Trevisan'00)
- ◎ $\ell_\infty$-metric: $k$-means $\ll$ 3.94, 3.94 1.~~01~~, $k$-median $\ll$ 1.73, 1.73 ~~1.01~~
  (Guruswami-Indyk'03)

Continuous Version:

$k$-means in Euclidean metric $<$ 1.36, 1.07 ~~1.0013~~
(Lee-Schmidt-Wright'17)

Discrete Version

|  | $k$-means (JCH) | $k$-median (JCH) | $k$-means (UGC) | $k$-median (UGC) |
|---|---|---|---|---|
| $\ell_1$-metric | 3.94 | 1.73 | 1.56 | 1.14 |
| $\ell_2$-metric | 1.73 | 1.27 | 1.17 | 1.06 |
| $\ell_\infty$-metric | 3.94 | 1.73 | 3.94[*] | 1.73[*] |

### Discrete Version

| | $k$-means (JCH) | $k$-median (JCH) | $k$-means (UGC) | $k$-median (UGC) |
|---|---|---|---|---|
| $\ell_1$-metric | 3.94 | 1.73 | 1.56 | 1.14 |
| $\ell_2$-metric | 1.73 | 1.27 | 1.17 | 1.06 |
| $\ell_\infty$-metric | 3.94 | 1.73 | 3.94[*] | 1.73[*] |

### Continuous Version

$k$-means in $\ell_2$-metric $\approx$ 1.36 (JCH), 1.07 (UGC)

$k$-median in $\ell_1$-metric $\approx$ 1.36 (JCH), 1.07 (UGC)

## Discrete Version

|  | $k$-means (JCH) | $k$-median (JCH) | $k$-means (UGC) | $k$-median (UGC) |
|---|---|---|---|---|
| $\ell_1$-metric | 3.94 | 1.73 | 1.56 | 1.14 |
| $\ell_2$-metric | 1.73 | 1.27 | 1.17 | 1.06 |
| $\ell_\infty$-metric | 3.94 | 1.73 | 3.94* | 1.73* |

## Continuous Version

$k$-means in $\ell_2$-metric $\approx$ 1.36 (JCH), 1.07 (UGC)

$k$-median in $\ell_1$-metric $\approx$ 1.36 (JCH), 1.07 (UGC)

A New Embedding Framework to potentially get Strong (tight?) Inapproximability results!

### Theorem (Guha-Khuller'99)

Fix $\varepsilon > 0$. Given input $(X, S, k)$. It is NP-hard to distinguish:

### Theorem (Guha-Khuller'99)

Fix $\varepsilon > 0$. Given input $(X, S, k)$. It is NP-hard to distinguish:

YES: There exists $(C^*, \sigma^*)$ such that $\sum\limits_{x \in X} \Delta(x, \sigma^*(x))^2 \leq |X|$

### Theorem (Guha-Khuller'99)

Fix $\varepsilon > 0$. Given input $(X, S, k)$. It is NP-hard to distinguish:

YES: There exists $(C^*, \sigma^*)$ such that $\sum\limits_{x \in X} \Delta(x, \sigma^*(x))^2 \leq |X|$

NO: For all $(C, \sigma)$ we have $\sum\limits_{x \in X} \Delta(x, \sigma(x))^2 \geq (1 + 8/e - \varepsilon) \cdot |X|$

**Max Coverage**:

◎ <u>Input</u>: Universe and Collection of Subsets $(U, \mathcal{S}, k)$

**Max Coverage**:

◎ Input: Universe and Collection of Subsets $(U, \mathcal{S}, k)$

◎ Objective: Max Fraction of $U$ covered by $k$ subsets in $\mathcal{S}$

**Max Coverage**:

- ◎ Input: Universe and Collection of Subsets $(U, S, k)$

- ◎ Objective: Max Fraction of $U$ covered by $k$ subsets in $S$

### Theorem (Feige'98)

Fix $\varepsilon > 0$. It is NP-hard to distinguish:

**Max Coverage**:

◎ Input: Universe and Collection of Subsets $(U, \mathcal{S}, k)$

◎ Objective: Max Fraction of $U$ covered by $k$ subsets in $\mathcal{S}$

### Theorem (Feige'98)

Fix $\varepsilon > 0$. It is NP-hard to distinguish:

YES: Max Coverage is 1

**Max Coverage**:

- ◎ Input: Universe and Collection of Subsets $(U, \mathcal{S}, k)$

- ◎ Objective: Max Fraction of $U$ covered by $k$ subsets in $\mathcal{S}$

### Theorem (Feige'98)

Fix $\varepsilon > 0$. It is NP-hard to distinguish:

YES: Max Coverage is 1

NO: Max Coverage is at most $1 - 1/e + \varepsilon$

### Theorem (Feige'98)

Fix $\varepsilon > 0$. It is NP-hard to distinguish:

YES: Max Coverage is 1

NO: Max Coverage is at most $1 - 1/e + \varepsilon$

$$\Downarrow$$

## Theorem (Feige'98)

Fix $\varepsilon > 0$. It is NP-hard to distinguish:

YES: Max Coverage is $1$

NO: Max Coverage is at most $1 - 1/e + \varepsilon$

$$\Downarrow$$
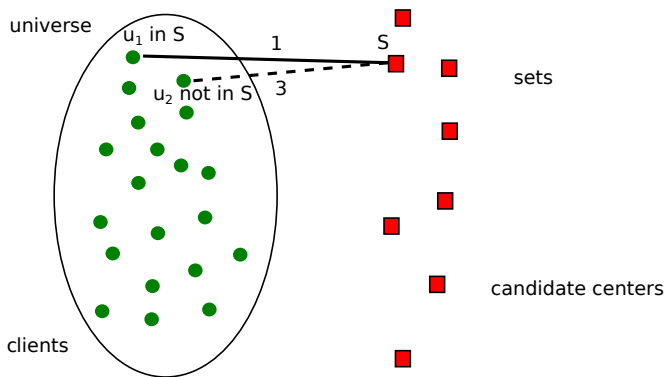
## Theorem (Guha-Khuller'99)

Fix $\varepsilon > 0$. Given input $(X, S, k)$. It is NP-hard to distinguish:

YES: There exists $(C^*, \sigma^*)$ such that $\sum_{x \in X} \Delta(x, \sigma^*(x))^2 \leq |X|$

NO: For all $(C, \sigma)$ we have $\sum_{x \in X} \Delta(x, \sigma(x))^2 \geq (1 + 8/e - \varepsilon) \cdot |X|$

universe

$u_1$ in S

1

S

$u_2$ not in S

3

sets

clients

candidate centers

## Theorem (Guha-Khuller'99)

Fix $\varepsilon > 0$. Given input $(X, S, k)$. It is NP-hard to distinguish:
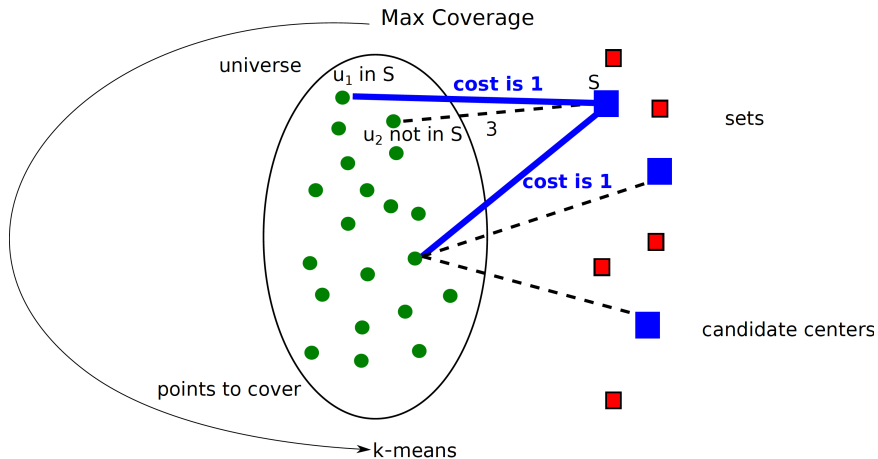
YES: There exists $(C^*, \sigma^*)$ such that $\sum_{x \in X} \Delta(x, \sigma^*(x))^2 \le |X|$

NO: For all $(C, \sigma)$ we have $\sum_{x \in X} \Delta(x, \sigma(x))^2 \ge (1 + 8/e - \varepsilon) \cdot |X|$

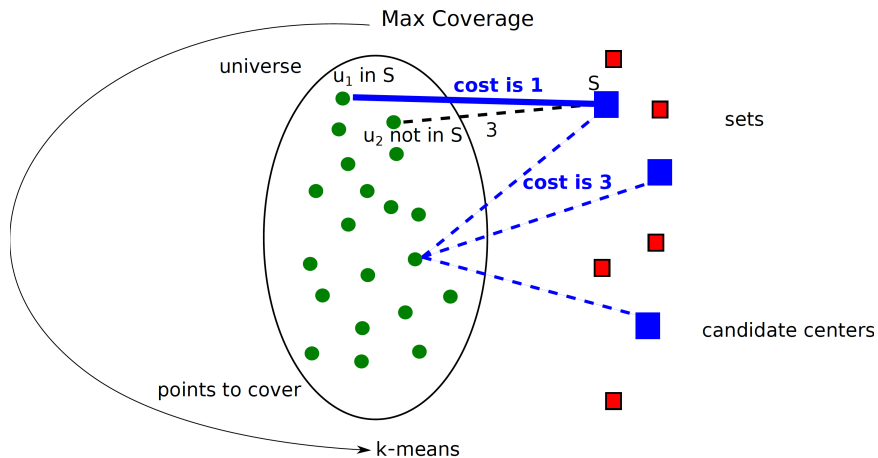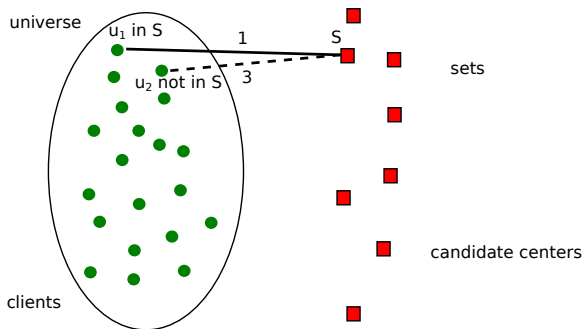# Johnson Coverage Hypothesis

# Johnson Coverage Hypothesis



universe
$u_1$ in S
S
1
$u_2$ not in S
3
sets
candidate centers
clients

## Johnson Coverage Hypothesis (Cohen-Addad–K–Lee)

Fix $\varepsilon > 0$. It is NP-hard to distinguish:

YES: Max Coverage is 1

NO: Max Coverage is at most $1 - 1/e + \varepsilon$

even when set system is induced subgraph of **Johnson graph**.

## $(\alpha, t)$-Johnson Coverage Problem

Given $E \subseteq \binom{[n]}{t}$, and $k$ as input, distinguish between:

**Completeness**: There exists $\mathscr{C} := \{S_1, \ldots, S_k\} \subseteq \binom{[n]}{t-1}$ such that

$$\forall T \in E, \ \exists S_i \in \mathscr{C}, \ S_i \subset T.$$

**Soundness**: For every $\mathscr{C} := \{S_1, \ldots, S_k\} \subseteq \binom{[n]}{t-1}$ we have

$$\Pr_{T \sim E}[\exists S_i, \ S_i \subset T] \leq \alpha.$$

# Johnson Coverage Hypothesis

## $(\alpha, t)$-Johnson Coverage Problem

Given $E \subseteq \binom{[n]}{t}$, and $k$ as input, distinguish between:

**Completeness**: There exists $\mathscr{C} := \{S_1, \ldots, S_k\} \subseteq \binom{[n]}{t-1}$ such that

$$\forall T \in E, \ \exists S_i \in \mathscr{C}, \ S_i \subset T.$$

**Soundness**: For every $\mathscr{C} := \{S_1, \ldots, S_k\} \subseteq \binom{[n]}{t-1}$ we have

$$\Pr_{T \sim E} [\exists S_i, \ S_i \subset T] \leq \alpha.$$

## Johnson Coverage Hypothesis (Cohen-Addad–K–Lee)

$\forall \varepsilon > 0, \exists t_\varepsilon \in \mathbb{N}$ such that $(1 - \frac{1}{e} + \varepsilon, t_\varepsilon)$-Johnson Coverage problem is NP-hard.

◎ $t = 2$: Vertex Coverage problem

◎ $t = 2$: Vertex Coverage problem
  ◦ ≈0.9292 gap is tight!

◎ $t = 2$: Vertex Coverage problem
  ○ ≈0.9292 gap is tight!

◎ 3-Hypergraph Vertex Coverage problem is NP-Hard to approximate to a factor of 7/8

**3 ingredients**

**3 ingredients**

◎ JCH instance

### 3 ingredients

◎ JCH instance

◎ Dimensionality reduction for all $\ell_p$-metrics

  ○ Works only for JCH instances

  ○ Arises from transcript of a communication game

<u>**3 ingredients**</u>

◎ JCH instance

◎ Dimensionality reduction for all $\ell_p$-metrics

  ○ Works only for JCH instances

  ○ Arises from transcript of a communication game

◎ Johnson Graph Embedding into $\ell_p$-metrics

Discrete Version

|  | $k$-means (JCH) | $k$-median (JCH) | $k$-means (UGC) | $k$-median (UGC) |
|---|---|---|---|---|
| $\ell_1$-metric | 3.94 | 1.73 | 1.56 | 1.14 |
| $\ell_2$-metric | 1.73 | 1.27 | 1.17 | 1.06 |
| $\ell_\infty$-metric | 3.94 | 1.73 | 3.94[*] | 1.73[*] |

Discrete Version

|  | $k$-means (JCH) | $k$-median (JCH) | $k$-means (UGC) | $k$-median (UGC) |
|---|---|---|---|---|
| $\ell_1$-metric | 3.94 | 1.73 | 1.56 | 1.14 |
| $\ell_2$-metric | 1.73 | 1.27 | 1.17 | 1.06 |
| $\ell_\infty$-metric | 3.94 | 1.73 | $3.94^*$ | $1.73^*$ |

Continuous Version

$k$-means in $\ell_2$-metric $\approx$ 1.36 (JCH), 1.07 (UGC)

$k$-median in $\ell_1$-metric $\approx$ 1.36 (JCH), 1.07 (UGC)

Discrete Version

|  | $k$-means (JCH) | $k$-median (JCH) | $k$-means (UGC) | $k$-median (UGC) |
|---|---|---|---|---|
| $\ell_1$-metric | 3.94 | 1.73 | 1.56 | 1.14 |
| $\ell_2$-metric | 1.73 | 1.27 | 1.17 | 1.06 |
| $\ell_\infty$-metric | 3.94 | 1.73 | $3.94^*$ | $1.73^*$ |

Continuous Version

$k$-means in $\ell_2$-metric $\approx 1.36$ (JCH), $1.07$ (UGC)

$k$-median in $\ell_1$-metric $\approx 1.36$ (JCH), $1.07$ (UGC)

# Continuous $k$-means and $k$-median

### Theorem (Cohen-Addad–K–Lee'21)

Given input $X \subseteq \mathbb{R}^{O(n)}$, it is NP-hard to distinguish:

# Continuous $k$-means and $k$-median

## Theorem (Cohen-Addad–K–Lee'21)

Given input $X \subseteq \mathbb{R}^{O(n)}$, it is NP-hard to distinguish:

YES: There exists $(C^*, \sigma^*)$ such that $\sum_{x \in X} \|(x - \sigma^*(x)\|_\infty^2 \leq n'$,

NO: For all $(C, \sigma)$ we have $\sum_{x \in X} \|(x - \sigma(x)\|_\infty^2 \geq 4 \cdot n'$.

# Continuous $k$-means and $k$-median

> **Theorem (Cohen-Addad–K–Lee'21)**
>
> Given input $X \subseteq \mathbb{R}^{O(n)}$, it is NP-hard to distinguish:
>
> YES: There exists $(C^*, \sigma^*)$ such that $\sum\limits_{x \in X} \|(x - \sigma^*(x)\|_\infty^2 \leq n'$,
>
> NO: For all $(C, \sigma)$ we have $\sum\limits_{x \in X} \|(x - \sigma(x)\|_\infty^2 \geq 4 \cdot n'$.

◎ $k$-median: 2 inapproximability

# Continuous $k$-means and $k$-median

## Theorem (Cohen-Addad–K–Lee'21)

Given input $X \subseteq \mathbb{R}^{O(n)}$, it is NP-hard to distinguish:

YES: There exists $(C^*, \sigma^*)$ such that $\sum\limits_{x \in X} \|(x - \sigma^*(x)\|_\infty^2 \leq n'$,

NO: For all $(C, \sigma)$ we have $\sum\limits_{x \in X} \|(x - \sigma(x)\|_\infty^2 \geq 4 \cdot n'$.

◎ $k$-median: 2 inapproximability

Continuous is harder than Discrete!

# Continuous $k$-means and $k$-median

## Theorem (Cohen-Addad–K–Lee'21)

Given input $X \subseteq \mathbb{R}^{O(n)}$, it is NP-hard to distinguish:

YES: There exists $(C^*, \sigma^*)$ such that $\sum\limits_{x \in X} \|(x - \sigma^*(x))\|_\infty^2 \le n'$,

NO: For all $(C, \sigma)$ we have $\sum\limits_{x \in X} \|(x - \sigma(x))\|_\infty^2 \ge 4 \cdot n'$.

◎ $k$-median: 2 inapproximability

Continuous is harder than Discrete!

◎ Constant Bicriteria inapproximability

## Theorem (Cohen-Addad–K–Lee'21)

Given input $X \subseteq \mathbb{R}^{O(n)}$, it is NP-hard to distinguish:

YES: There exists $(C^*, \sigma^*)$ such that $\sum\limits_{x \in X} \|(x - \sigma^*(x)\|_\infty^2 \leq n'$,

NO: For all $(C, \sigma)$ we have $\sum\limits_{x \in X} \|(x - \sigma(x)\|_\infty^2 \geq 4 \cdot n'$.

◎ $k$-median: 2 inapproximability

Continuous is harder than Discrete!

◎ Constant Bicriteria inapproximability

◎ Assuming UGC, hardness for $k = 2$!

## Theorem (Cohen-Addad–K–Lee'21)

Given input $X \subseteq \mathbb{R}^{O(n)}$, it is NP-hard to distinguish:

YES: There exists $(C^*, \sigma^*)$ such that $\sum\limits_{x \in X} \|(x - \sigma^*(x)\|_\infty^2 \leq n'$,

NO: For all $(C, \sigma)$ we have $\sum\limits_{x \in X} \|(x - \sigma(x)\|_\infty^2 \geq 4 \cdot n'$.

◎ $k$-median: 2 inapproximability

Continuous is harder than Discrete!

◎ Constant Bicriteria inapproximability

◎ Assuming UGC, hardness for $k = 2$!

◎ Dependency on $d, k$, and $\ell_\infty$ tight

◎ $(\Gamma, \Delta)$ is a metric space

◎ Input: $X \subseteq \Gamma$, $k \in \mathbb{N}$

◎ $(\Gamma, \Delta)$ is a metric space

◎ Input: $X \subseteq \Gamma$, $k \in \mathbb{N}$

◎ Output: A partition $X := X_1 \dot\cup X_2 \dot\cup \cdots \dot\cup X_k$ that minimizes:

◎ $(\Gamma, \Delta)$ is a metric space

◎ <u>Input</u>: $X \subseteq \Gamma$, $k \in \mathbb{N}$

◎ <u>Output</u>: A partition $X := X_1 \dot{\cup} X_2 \dot{\cup} \cdots \dot{\cup} X_k$ that minimizes:

$$\sum_{i \in [k]} \sum_{x,y \in X_i} \Delta(x, y)$$

◎ $(\Gamma, \Delta)$ is a metric space

◎ Input: $X \subseteq \Gamma$, $k \in \mathbb{N}$

◎ Output: A partition $X := X_1 \dot{\cup} X_2 \dot{\cup} \cdots \dot{\cup} X_k$ that minimizes:

$$\sum_{i \in [k]} \sum_{x,y \in X_i} \Delta(x,y)$$

◎ Approximation: $O(\log n)$ [Behsaz et al.'15]

# Minsum (Definition)

◎ $(\Gamma, \Delta)$ is a metric space

◎ Input: $X \subseteq \Gamma$, $k \in \mathbb{N}$

◎ Output: A partition $X := X_1 \dot\cup X_2 \dot\cup \cdots \dot\cup X_k$ that minimizes:

$$\sum_{i \in [k]} \sum_{x,y \in X_i} \Delta(x, y)$$

◎ Approximation: $O(\log n)$ [Behsaz et al.'15]

◎ Hardness: $1 + \varepsilon$ [Guruswami-Indyk'03]

### Theorem (Cohen-Addad–K–Lee'21)

Given input $(X, k)$, it is NP-hard to distinguish:

## Theorem (Cohen-Addad–K–Lee'21)

Given input $(X, k)$, it is NP-hard to distinguish:

YES: There exists a partition $X := X_1 \dot\cup X_2 \dot\cup \cdots \dot\cup X_k$ such that

$$\sum_{i \in [k]} \sum_{x,y \in X_i} \Delta(x, y) \le n',$$

## Theorem (Cohen-Addad–K–Lee'21)

Given input $(X, k)$, it is NP-hard to distinguish:

YES: There exists a partition $X := X_1 \dot\cup X_2 \dot\cup \cdots \dot\cup X_k$ such that

$$\sum_{i \in [k]} \sum_{x,y \in X_i} \Delta(x, y) \leq n',$$

NO: For every partition $X := X_1 \dot\cup X_2 \dot\cup \cdots \dot\cup X_k$ we have

$$\sum_{i \in [k]} \sum_{x,y \in X_i} \Delta(x, y) \geq 1.41 \cdot n'.$$

**Theorem (Cohen-Addad–K–Lee'21)**

Given input $(X, k)$, it is NP-hard to distinguish:

YES: There exists a partition $X := X_1 \dot\cup X_2 \dot\cup \cdots \dot\cup X_k$ such that

$$\sum_{i \in [k]} \sum_{x,y \in X_i} \Delta(x, y) \leq n',$$

NO: For every partition $X := X_1 \dot\cup X_2 \dot\cup \cdots \dot\cup X_k$ we have

$$\sum_{i \in [k]} \sum_{x,y \in X_i} \Delta(x, y) \geq 1.41 \cdot n'.$$

Key Ingredient: Hard Instances of Max-Coverage
with large girth

Improved Inapproximability of

Improved Inapproximability of

⊚ $k$-means and $k$-median in $\ell_p$-metric using JCH framework

Improved Inapproximability of

&#9673; $k$-means and $k$-median in $\ell_p$-metric using JCH framework

&#9673; Continuous versions of $k$-means and $k$-median in General metric

Improved Inapproximability of

◎ $k$-means and $k$-median in $\ell_p$-metric using JCH framework

◎ Continuous versions of $k$-means and $k$-median in General metric

◎ $k$-minsum in General metric

# THANK YOU!