# Replicable Clustering

**Hossein Esfandiari, Amin Karbasi, Vahab Mirrokni, Grigoris Velegkas, Felix Zhou**

Yale, Google Research

# Personnel


Hossein Esfandiari
(Google)


Amin Karbasi
(Yale, Google)


Vahab Mirrokni
(Google)


Grigoris Velegkas
(Yale)


**Felix Zhou**
(Yale)

# Table of Contents

Yale

# Table of Contents

Yale

"1500 Scientists Lift the Lid on Reproducibility."
*Nature* (2016)

# Reproducibility Crisis



Reproducibility Challenge @ NeurIPS 2019

- 2019 NeurIPS/ICLR Reproducibility Challenge
  (`github.com/reproducibility-challenge`)

# Reproducibility Crisis



- 2019 NeurIPS/ICLR Reproducibility Challenge
  (`github.com/reproducibility-challenge`)
- Ongoing ML Reproducibility Challenge
  (`paperswithcode.com/rc2022`)

Trying to develop agreed-upon set of replicable practices is difficult.
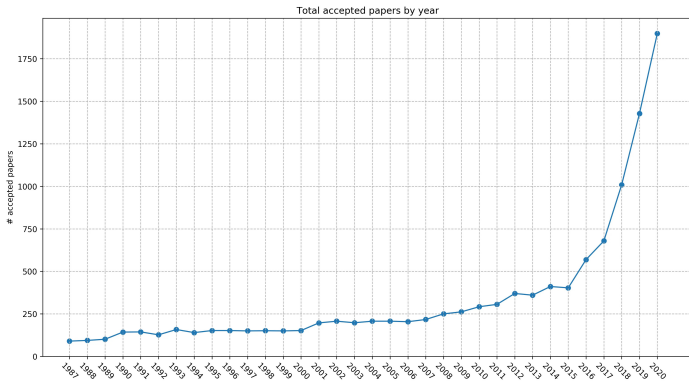


Figure: Number of accepted NeurIPS papers by year.

**Goal:** Design ML algorithms with replicability as theoretical guarantee.

Initiated by [Impagliazzo, Lei, Pitassi, and Sorell '22] (STOC'22).

- $X$ data domain

- $X$ data domain
- $\mathcal{D}$ distribution over $X$

- $X$ data domain
- $\mathcal{D}$ distribution over $X$
- $S_1, S_2 \sim_{i.i.d.} \mathcal{D}^n$ datasets of size $n$

- $X$ data domain
- $\mathcal{D}$ distribution over $X$
- $S_1, S_2 \sim_{i.i.d.} \mathcal{D}^n$ datasets of size $n$
- $\xi$ uniformly random binary string

## Yale

- $X$ data domain
- $\mathcal{D}$ distribution over $X$
- $S_1, S_2 \sim_{i.i.d.} \mathcal{D}^n$ datasets of size $n$
- $\xi$ uniformly random binary string

DEFINITION (REPLICABLE ALGORITHM; [ILPS '22])

A randomized algorithm $\mathcal{A} : X^n \to Y$ is $\rho$-*replicable* if

$$\Pr_{S_1, S_2, \xi} [\mathcal{A}(S_1; \xi) = \mathcal{A}(S_2; \xi)] \geq 1 - \rho.$$

Yale

# Algorithmic Replicability

- $X$ data domain
- $\mathcal{D}$ distribution over $X$
- $S_1, S_2 \sim_{i.i.d.} \mathcal{D}^n$ datasets of size $n$
- $\xi$ uniformly random binary string

## DEFINITION (REPLICABLE ALGORITHM; [ILPS '22])

A randomized algorithm $\mathcal{A} : X^n \to Y$ is *$\rho$-replicable* if

$$\Pr_{S_1, S_2, \xi} [\mathcal{A}(S_1; \xi) = \mathcal{A}(S_2; \xi)] \geq 1 - \rho.$$

**Remark:** Replicability is trivial to obtain by itself!
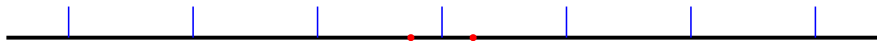
Yale

1. Concentration of Measure



$\mathbb{R}$

1. Concentration of Measure
2. Discretization



$$\mathbb{R}$$

1. Concentration of Measure
2. Discretization
3. Shared Random Offset



$$\mathbb{R}$$

# Table of Contents

## Yale

- **Given:** Sample access to distribution $\mathcal{D}$ over $[0,1]^d$

# Statistical $k$-Means

- **Given:** Sample access to distribution $\mathcal{D}$ over $[0,1]^d$
- **Want:** Choose $k$ centers $y_1, \ldots, y_k \in [0,1]^d$ minimizing expected "cost of travel" to nearest center

# Statistical $k$-Means

- **Given:** Sample access to distribution $\mathcal{D}$ over $[0,1]^d$
- **Want:** Choose $k$ centers $y_1, \ldots, y_k \in [0,1]^d$ minimizing expected "cost of travel" to nearest center
- Solve $\operatorname{argmin}_{y_1, \ldots, y_k \in [0,1]^d} \mathbb{E}_{X \sim \mathcal{D}} \left[ \min_{j \in [k]} \| X - y_j \|_2^2 \right]$

# Statistical $k$-Means

- **Given:** Sample access to distribution $\mathcal{D}$ over $[0,1]^d$
- **Want:** Choose $k$ centers $y_1, \ldots, y_k \in [0,1]^d$ minimizing expected "cost of travel" to nearest center
- Solve $\text{argmin}_{y_1,\ldots,y_k \in [0,1]^d} \, \mathbb{E}_{X \sim \mathcal{D}} \left[ \min_{j \in [k]} \|X - y_j\|_2^2 \right]$
- Sample complexity? Time complexity?

- **Given:** Points $x_1, \ldots, x_n \in [0,1]^d$

- **Given:** Points $x_1, \ldots, x_n \in [0, 1]^d$
- **Want:** Choose $k$ centers $y_1, \ldots, y_k \in [0, 1]^d$ minimizing average "cost of travel" to nearest center
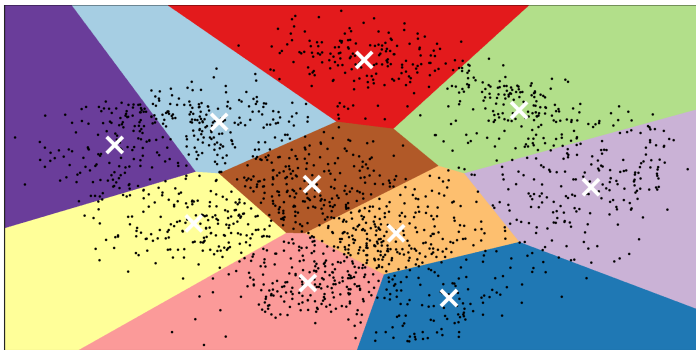
# Sample $k$-Means

- **Given:** Points $x_1, \ldots, x_n \in [0,1]^d$
- **Want:** Choose $k$ centers $y_1, \ldots, y_k \in [0,1]^d$ minimizing average "cost of travel" to nearest center
- Solve $\operatorname{argmin}_{y_1, \ldots, y_k \in [0,1]^d} \frac{1}{n} \sum_{i \in [n]} \left[ \min_{j \in [k]} \| x_i - y_j \|_2^2 \right]$

- **Given:** Points $x_1, \ldots, x_n \in [0,1]^d$
- **Want:** Choose $k$ centers $y_1, \ldots, y_k \in [0,1]^d$ minimizing average "cost of travel" to nearest center
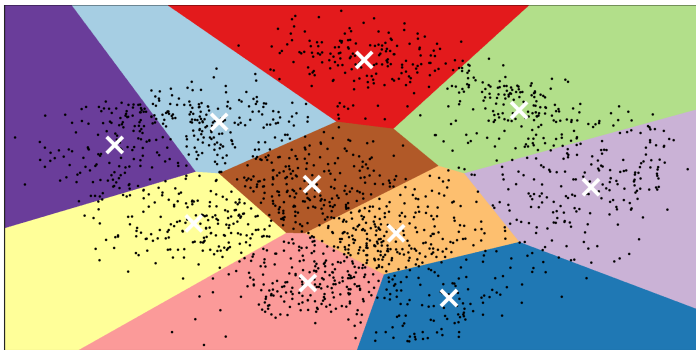- Solve $\operatorname{argmin}_{y_1, \ldots, y_k \in [0,1]^d} \frac{1}{n} \sum_{i \in [n]} \left[ \min_{j \in [k]} \|x_i - y_j\|_2^2 \right]$
- Time complexity?

- Clustering algorithms reveal properties of the underlying population

- Clustering algorithms reveal properties of the underlying population
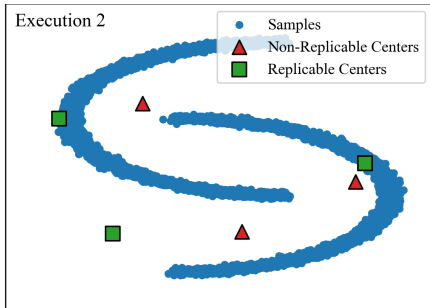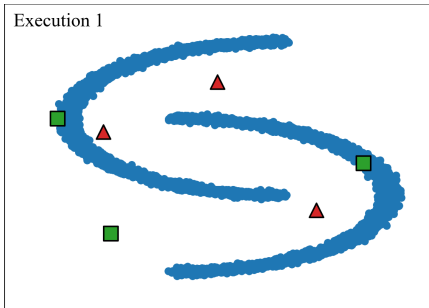
- Replicability is an important property for downstream applications

- **Given:** Sample access to distribution $\mathcal{D}$ over $[0,1]^d$

- **Given:** Sample access to distribution $\mathcal{D}$ over $[0,1]^d$
- **Want:** $\operatorname{argmin}_{y_1,\ldots,y_k \in [0,1]^d} \mathbb{E}_{X \sim \mathcal{D}} \left[ \min_{j \in [k]} \|X - y_j\|_2^2 \right]$

# Replicable Statistical $k$-Means

- **Given:** Sample access to distribution $\mathcal{D}$ over $[0,1]^d$
- **Want:** $\text{argmin}_{y_1,\ldots,y_k \in [0,1]^d} \mathbb{E}_{X \sim \mathcal{D}} \left[ \min_{j \in [k]} \|X - y_j\|_2^2 \right]$
- **And:** $\text{Pr}_{X,X',\xi} \left[ \{y_1, \ldots, y_k\} = \{y_1', \ldots, y_k'\} \right] \geq 1 - \rho$

# Replicable Statistical $k$-Means

- **Given:** Sample access to distribution $\mathcal{D}$ over $[0,1]^d$
- **Want:** $\operatorname{argmin}_{y_1,\ldots,y_k \in [0,1]^d} \mathbb{E}_{X \sim \mathcal{D}} \left[ \min_{j \in [k]} \|X - y_j\|_2^2 \right]$
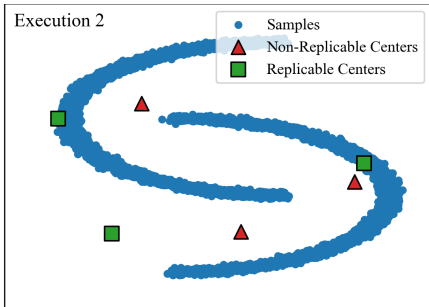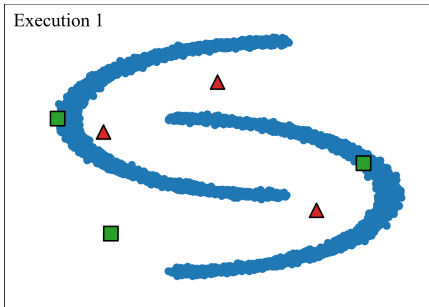- **And:** $\operatorname{Pr}_{X,X',\xi} \left[ \{y_1,\ldots,y_k\} = \{y_1',\ldots,y_k'\} \right] \geq 1 - \rho$
- Sample Complexity? Time Complexity?

# Table of Contents

## Yale

### THEOREM (EKMVZ '23)

▶ Assume black-box polynomial time $\beta$-approximation oracle for sample $k$-means.

### THEOREM (EKMVZ '23)

- Assume black-box polynomial time $\beta$-approximation oracle for sample $k$-means.
- There is a $\rho$-replicable algorithm $\mathcal{A}$ for statistical $k$-means:

THEOREM (EKMVZ '23)

- Assume black-box polynomial time $\beta$-approximation oracle for sample $k$-means.
- There is a $\rho$-replicable algorithm $\mathcal{A}$ for statistical $k$-means:
- with high probability, the cost of the solution is at most $(1 + \varepsilon)\beta \cdot \text{OPT}$.

# Yale

## THEOREM (EKMVZ '23)

- Assume black-box polynomial time $\beta$-approximation oracle for sample $k$-means.
- There is a $\rho$-replicable algorithm $\mathcal{A}$ for statistical $k$-means:
- with high probability, the cost of the solution is at most $(1 + \varepsilon)\beta \cdot \text{OPT}$.
- $\mathcal{A}$ has time and sample complexity

$$\tilde{O}_{\rho,\varepsilon}\left(\text{poly}(k, d)\, k^{\log\log k}\right).$$

Yale

## THEOREM (EKMVZ '23)

- Assume black-box polynomial time $\beta$-approximation oracle for sample $k$-means.
- There is a $\rho$-replicable algorithm $\mathcal{A}$ for statistical $k$-means:
- with probability at least $1 - \delta$, the cost of the solution is at most $(1 + \varepsilon)\beta \cdot \text{OPT}$.
- $\mathcal{A}$ has time and sample complexity

$$\tilde{O}\left(\text{poly}(k, d, 1/\rho)(2^{\sqrt{m}}/\varepsilon)^{O(m)} \log \frac{1}{\delta}\right).$$

where $m = O(\varepsilon^{-2} \log k/\delta\varepsilon)$.

Yale

1. Reduce problem on distribution to problem on samples

1. Reduce problem on distribution to problem on samples
   - Uniform Law of Large Numbers

1. Reduce problem on distribution to problem on samples
   - Uniform Law of Large Numbers

2. Consolidate multiple points into weighted point

1. Reduce problem on distribution to problem on samples
   - Uniform Law of Large Numbers

2. Consolidate multiple points into weighted point
   - Replicable Coreset Construction

0. Data-Oblivious Dimensionality Reduction

1. Reduce problem on distribution to problem on samples
   ▶ Uniform Law of Large Numbers

2. Consolidate multiple points into weighted point
   ▶ Replicable Coreset Construction

0. Data-Oblivious Dimensionality Reduction
   - Johnson-Lindenstrauss transform

1. Reduce problem on distribution to problem on samples
   - Uniform Law of Large Numbers

2. Consolidate multiple points into weighted point
   - Replicable Coreset Construction

# Table of Contents

Yale

THEOREM (BEN-DAVID '07)

- Assume black-box polynomial time $\beta$-approximation oracle for sample *k*-means.
- There is an algorithm $\mathcal{A}$ for statistical *k*-means:
- with high probability, the cost of the solution is at most $(1 + \varepsilon)\beta \cdot (4\,\mathrm{OPT})$.
- $\mathcal{A}$ has time and sample complexity

$$\tilde{O}\left(\mathrm{poly}(d, k, {}^1/_\varepsilon)\right).$$

Yale

▶ **Stable choice of $k$:** Produce clusterings that do not vary much from one sample to another. [BEG '01], [LRBB '04], [VB '05], [B '06], [RC '06], [V '10]

- **Stable choice of $k$:** Produce clusterings that do not vary much from one sample to another. [BEG '01], [LRBB '04], [VB '05], [B '06], [RC '06], [V '10]
- [Ben-David, Pál, Simon; '07]: "for large sample sizes, stability is fully determined by the symmetry within the data."

- **Inspiration:** operations research and statistics [Lloyd '57], [Hakimi '64], [Steinhaus '57]

# Sample *k*-Means

- **Inspiration:** operations research and statistics [Lloyd '57], [Hakimi '64], [Steinhaus '57]
- **Sample Metric *k*-Means:** polynomial time 9-approximation [ANS '19], $(1 + 8/e)$-approximation is NP-hard [CGKLL '19]

# Sample *k*-Means

- **Inspiration:** operations research and statistics [Lloyd '57], [Hakimi '64], [Steinhaus '57]
- **Sample Metric *k*-Means:** polynomial time 9-approximation [ANS '19], $(1 + 8/e)$-approximation is NP-hard [CGKLL '19]
- **Sample Euclidean *k*-Means:** polynomial time 5.912-approximation [CEMN '22], 1.07-approximation is NP-hard [CK '19], [CLK '22]

Yale

# Sample *k*-Means

- **Inspiration:** operations research and statistics [Lloyd '57], [Hakimi '64], [Steinhaus '57]
- **Sample Metric *k*-Means:** polynomial time 9-approximation [ANS '19], $(1 + 8/e)$-approximation is NP-hard [CGKLL '19]
- **Sample Euclidean *k*-Means:** polynomial time 5.912-approximation [CEMN '22], 1.07-approximation is NP-hard [CK '19], [CLK '22]
- **Dimensionality Reduction** Johnson-Lindenstrauss Transform [Johnson '84], [KKM '19]

# Sample $k$-Means

- **Inspiration:** operations research and statistics [Lloyd '57], [Hakimi '64], [Steinhaus '57]
- **Sample Metric $k$-Means:** polynomial time 9-approximation [ANS '19], $(1 + 8/e)$-approximation is NP-hard [CGKLL '19]
- **Sample Euclidean $k$-Means:** polynomial time 5.912-approximation [CEMN '22], 1.07-approximation is NP-hard [CK '19], [CLK '22]
- **Dimensionality Reduction** Johnson-Lindenstrauss Transform [Johnson '84], [KKM '19]
- **Coresets** geometric coresets [FS '05], sampling-based coresets [HSYZ '18], survey [SW '18]

Yale

- **Algorithm Design:** statistical queries, heavy-hitters, medians, learning halfspaces [ILPS '22], stochastic bandits [EKKKMV '23], reinforcement learning [KYGZ '23], [EHKS '23]

- **Algorithm Design:** statistical queries, heavy-hitters, medians, learning halfspaces [ILPS '22], stochastic bandits [EKKKMV '23], reinforcement learning [KYGZ '23], [EHKS '23]
- **Learning Theory:** equivalence with DP [BGHILPSS '23], statistical indistinguishability [KKMV '23], list-replicability [CMY '23], [DPWV '23]

# Table of Contents

**Intuition:** Little hope for replicability with continuous distributions, need to discretize and round similar to mean estimation.

**Idea:** Replicably approximate the input distribution with a finite (discrete) distribution.

Yale

For some $R : \mathcal{X} \to \mathcal{X}$ with small $|R(\mathcal{X})|$, uniformly approximate

$$\text{cost}(y) := \mathbb{E}_X \left[ \min_j \|X - y_j\|_2^2 \right]$$

For some $R : \mathcal{X} \to \mathcal{X}$ with small $|R(\mathcal{X})|$, uniformly approximate

$$\text{cost}(y) := \mathbb{E}_X \left[ \min_j \| X - y_j \|_2^2 \right]$$
$$\approx \mathbb{E}_X \left[ \min_j \| R(X) - y_j \|_2^2 \right]$$
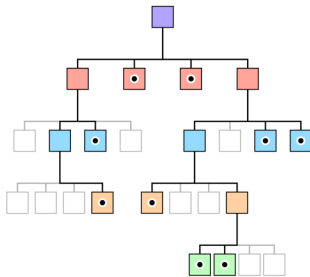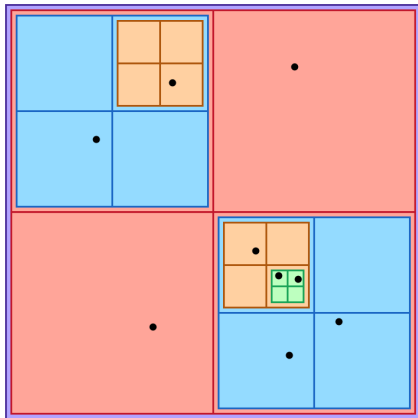
For some $R : \mathcal{X} \to \mathcal{X}$ with small $|R(\mathcal{X})|$, uniformly approximate

$$
\begin{aligned}
\text{cost}(y) &:= \mathbb{E}_X \left[ \min_j \| X - y_j \|_2^2 \right] \\
&\approx \mathbb{E}_X \left[ \min_j \| R(X) - y_j \|_2^2 \right] \\
&= \sum_{z \in R(\mathcal{X})} \min_j \| z - y_j \|_2^2 \cdot \mathbb{P}(R^{-1}(z))
\end{aligned}
$$

Yale

**Idea:** Recursively partition $[0,1]^d$ into subcubes and consolidate mass from entire subcube into single point

- **Given:** Sample access to distribution, $\nu \in [0, 1]$.

- **Given:** Sample access to distribution, $\nu \in [0,1]$.

- **Want:** All elements with mass at least $\nu$.

- **Given:** Sample access to distribution, $\nu \in [0, 1]$.

- **Want:** All elements with mass at least $\nu$.

### THEOREM (ILPS '22)

There is a replicable heavy-hitters algorithm.

**Algorithm** Replicable Quad Tree

---

1: **rQuadTree**(Node $Z$, error $\varepsilon$, depth $i$):
2: **if** $\operatorname{diam}(Z) \geq C_1 \varepsilon$ **then**
3:    **return**
4: **end if**
5: **for** $Z' \in \operatorname{subcubes}(Z)$ **do**
6:    // Heavy Hitters Operation!
7:    **if** $\mathbb{P}(Z') \geq C_2 \cdot {}^{2^i \varepsilon}/k$ **then**
8:       Add $Z'$ as child of $Z$
9:       **rQuadTree**($Z', \varepsilon, i+1$)
10:   **end if**
11: **end for**

---

Yale

0. Data-Oblivious Dimensionality Reduction

# Summary

Yale

0. Data-Oblivious Dimensionality Reduction

1. Uniform Law of Large Numbers

2. Replicable Coreset Estimation

0. Data-Oblivious Dimensionality Reduction

1. Uniform Law of Large Numbers

2. Replicable Coreset Estimation
   - A series of heavy hitter estimations that can be made replicable.

# Future Work

- Polynomial sample/time complexity for replicable $k$-means

# Future Work

- Polynomial sample/time complexity for replicable $k$-means

- Replicable $(k, p)$-clustering

Yale

# Future Work

▶ Polynomial sample/time complexity for replicable $k$-means

▶ Replicable $(k, p)$-clustering

▶ Sample complexity lower bounds for (replicable) clustering

Thank You!