Near-Optimal Dimension Reduction for Facility Location

Lingxiao Huang * Shaofeng Jiang [†] Robert Krauthgamer [‡] Di Yue [†]

June 2025

*Nanjing University [†]Peking University [‡]Weizmann Institute of Science

Dimension Reduction

Challenge: curse of dimensionality

Solution: dimension reduction

- Goal: find a mapping $\pi \colon \mathbb{R}^d \to \mathbb{R}^m$, where $m \ll d$
- Example: Johnson-Lindenstrauss transform

Dimension reduction in algorithm design

- (1) Embed the high-dimensional data to a low-dimensional space: $X \subset \mathbb{R}^d \xrightarrow{\pi} \pi(X) \subset \mathbb{R}^m$
- (2) Compute the problem on the low-dimensional data $\pi(X)$
- (3) Map the solution back to X

Dimension Reduction

Johnson-Lindenstrauss Lemma [JL84]

There exists a random map $\pi \colon \mathbb{R}^d \to \mathbb{R}^m$ for $m = O(\varepsilon^{-2} \log n)$, such that for every *n*-point set $X \subset \mathbb{R}^d$, w.h.p.

$$\forall x, y \in X, \qquad \|\pi(x) - \pi(y)\| \in (1 \pm \varepsilon) \|x - y\|.$$

Dimension Reduction

Johnson-Lindenstrauss Lemma [JL84]

There exists a random map $\pi \colon \mathbb{R}^d \to \mathbb{R}^m$ for $m = O(\varepsilon^{-2} \log n)$, such that for every *n*-point set $X \subset \mathbb{R}^d$, w.h.p.

$$\forall x, y \in X, \qquad \|\pi(x) - \pi(y)\| \in (1 \pm \varepsilon) \|x - y\|.$$

Good properties: linear, data oblivious

- Linear: $\pi: x \mapsto \frac{1}{\sqrt{m}}Gx$ (G is the Gaussian matrix)
- Data oblivious: many applications in streaming

Limitations: target dimension $m = O(\varepsilon^{-2} \log n)$ is tight [LN17]

- ▶ TSP in $\mathbb{R}^{\Theta(\log n)}$ does not admit a PTAS [Tre00]
- Streaming MST in $\mathbb{R}^{\Theta(\log n)}$ requires $\Omega(\sqrt{n})$ bits [CCJ+23]

Doubling Dimension

Going beyond $O(\varepsilon^{-2}\log n):$ seek dependence on the intrinsic dimension of datasets

High ambient dimension vs low intrinsic dimension

- Points in a linear subspace
- Points with sparse vector representation



▶ Doubling dimension ddim(X) [GKL03]: minimum t ≥ 0, such that every ball in X can be covered by at most 2^t balls of half the radius

• If
$$|X| = n$$
, then $\operatorname{ddim}(X) \le \log n$

Dimension Reduction Meets Doubling Dimension

Goal: refine JL lemma, such that $m = m(\varepsilon, \operatorname{ddim}(X))$?

- Remains open
- Not possible for linear maps [IN07]

Dimension Reduction Meets Doubling Dimension

Goal: refine JL lemma, such that $m = m(\varepsilon, \operatorname{ddim}(X))$?

- Remains open
- Not possible for linear maps [IN07]

Goal: JL dimension reduction for specific problems, such that $m = m(\varepsilon, \operatorname{ddim}(X))$ and $\operatorname{opt}(\pi(X)) \in (1 \pm \varepsilon)\operatorname{opt}(X)$

A weaker requirement than preserving pairwise distances

Dimension Reduction for Specific Problems

Problems	Арх.	Target Dimension	Ref.
Max-Cut	$1 + \varepsilon$	$\tilde{O}(\varepsilon^{-2})$	CJK23
k-Median/Means	$1 + \varepsilon$	$O(\varepsilon^{-2}\log k)$	MMR19
k-Subspace Apx.	$1 + \varepsilon$	$\tilde{O}(\varepsilon^{-3}k^2)$	CW25
Nearest Neighbor	$1 + \varepsilon$	$\tilde{O}(\varepsilon^{-2} ddim)$	IN07
k-Center Clustering	$1 + \varepsilon$	$O(\varepsilon^{-2}(\operatorname{ddim} + \log k))$	JKS24
MST	$1 + \varepsilon$	$\tilde{O}(\varepsilon^{-2}(\operatorname{ddim} + \log \log n))$	NSIZ21
UFL	O(1)	O(m ddim)	NSIZ21
UFL	$1 + \varepsilon$	$\tilde{O}(\varepsilon^{-2} ddim)$	This work

Uniform Facility Location (UFL)

Input: *n*-point set $X \subset \mathbb{R}^d$

Objective: find a facility set $F \subset \mathbb{R}^d$, to minimize

$$\operatorname{cost}(X,F) := |F| + \sum_{x \in X} \operatorname{dist}(x,F)$$

Optimal value: $ufl(X) := \min_{F \subset \mathbb{R}^d} cost(X, F)$



Uniform Facility Location (UFL)

Input: *n*-point set $X \subset \mathbb{R}^d$

Objective: find a facility set $F \subset \mathbb{R}^d$, to minimize

$$\operatorname{cost}(X,F) := |F| + \sum_{x \in X} \operatorname{dist}(x,F)$$

Optimal value: $ufl(X) := \min_{F \subset \mathbb{R}^d} cost(X, F)$



Problem (Dimension reduction for UFL) Given $\varepsilon \in (0, 1)$, decide a target dimension $m = m(\varepsilon, \operatorname{ddim}(X))$, such that w.h.p. $\operatorname{ufl}(\pi(X)) \in (1 \pm \varepsilon) \operatorname{ufl}(X)$

Main Result: Dimension Reduction

Theorem (Dimension reduction for UFL) Consider $m = \tilde{O}(\varepsilon^{-2} ddim)$. Then for every finite $X \subset \mathbb{R}^d$ with $ddim(X) \leq ddim$,

 $\Pr[\operatorname{ufl}(\pi(X)) \in (1 \pm \varepsilon) \operatorname{ufl}(X)] \ge 0.99$

Main Result: Dimension Reduction

Theorem (Dimension reduction for UFL) Consider $m = \tilde{O}(\varepsilon^{-2} \text{ddim})$. Then for every finite $X \subset \mathbb{R}^d$ with $\text{ddim}(X) \leq \text{ddim}$,

$$\Pr[\operatorname{ufl}(\pi(X)) \in (1 \pm \varepsilon) \operatorname{ufl}(X)] \ge 0.99$$

Improvements over previous results

- ▶ [NSIZ21]: O(1)-apx, target dimension m = O(ddim)
- [MMR19]: $(1 + \varepsilon)$ -apx, target dimension $m = O(\varepsilon^{-2} \log n)$

Handles a regime between low and high dimension

- Data: low doubling dimension
- Facilities: high dimension

Corollary: Streaming Algorithm

Corollary (Streaming algorithm for UFL)

There is a streaming algorithm that, given as input a set $X \subseteq [\Delta]^d$ presented as a stream, and an upper bound ddim, uses space $\tilde{O}(d \cdot \operatorname{polylog}(\Delta) + (\varepsilon^{-1} \log \Delta)^{\tilde{O}(\operatorname{ddim})})$ and outputs w.h.p a $(1 + \varepsilon)$ -apx to uff(X).

Corollary: Streaming Algorithm

Corollary (Streaming algorithm for UFL)

There is a streaming algorithm that, given as input a set $X \subseteq [\Delta]^d$ presented as a stream, and an upper bound ddim, uses space $\tilde{O}(d \cdot \operatorname{polylog}(\Delta) + (\varepsilon^{-1} \log \Delta)^{\tilde{O}(\operatorname{ddim})})$ and outputs w.h.p a $(1 + \varepsilon)$ -apx to uff(X).

- The first streaming algorithm for UFL that utilize the doubling dimension (generalization of [CLMS13])
- Break the $\Omega(\sqrt{n})$ barrier in [CJK+22]

Main Result: PTAS

Theorem (PTAS for UFL)

There exists a randomized algorithm that computes a $(1+\varepsilon)\text{-apprximation}$ for UFL in time $(2^{m'}d+2^{2^{m'}})\cdot \tilde{O}(n)$, for

$$m' = O\left(\operatorname{ddim}(X) \cdot \log(\operatorname{ddim}(X)/\varepsilon)\right)$$

Main Result: PTAS

Theorem (PTAS for UFL)

There exists a randomized algorithm that computes a $(1+\varepsilon)\text{-apprximation}$ for UFL in time $(2^{m'}d+2^{2^{m'}})\cdot \tilde{O}(n)$, for

$$m' = O\left(\operatorname{ddim}(X) \cdot \log(\operatorname{ddim}(X)/\varepsilon)\right)$$

▶ Facilities are allowed to be picked from the ambient space ℝ^d
 ▶ [CFS21]: PTAS for UFL in time 2^{2^{O(ddim²)}} · d · Õ(n), with facilities restricted to the dataset

Theorem (Dimension reduction for UFL) Consider $m = \tilde{O}(\varepsilon^{-2} \text{ddim})$. Then for every finite $X \subset \mathbb{R}^d$ with $\text{ddim}(X) \leq \text{ddim}$,

 $\Pr[\operatorname{ufl}(\pi(X)) \in (1 \pm \varepsilon) \operatorname{ufl}(X)] \ge 0.99$

Theorem (Dimension reduction for UFL) Consider $m = \tilde{O}(\varepsilon^{-2} \text{ddim})$. Then for every finite $X \subset \mathbb{R}^d$ with $\text{ddim}(X) \leq \text{ddim}$,

 $\Pr[\operatorname{ufl}(\pi(X)) \in (1 \pm \varepsilon) \operatorname{ufl}(X)] \ge 0.99$

• Easy direction: $ufl(\pi(X)) \le (1 + \varepsilon) ufl(X)$

- ► Let F* be the optimal solution for X, then π(F*) is a feasible solution for π(X)
- $\operatorname{ufl}(\pi(X)) \leq \operatorname{cost}(\pi(X), \pi(F^*)) \lesssim \operatorname{cost}(X, F^*) = \operatorname{ufl}(X)$
- Suffices to preserve the cost of one solution

Theorem (Dimension reduction for UFL) Consider $m = \tilde{O}(\varepsilon^{-2} \text{ddim})$. Then for every finite $X \subset \mathbb{R}^d$ with $\text{ddim}(X) \leq \text{ddim}$,

 $\Pr[\operatorname{ufl}(\pi(X)) \in (1 \pm \varepsilon) \operatorname{ufl}(X)] \ge 0.99$

• Easy direction: $ufl(\pi(X)) \le (1 + \varepsilon) ufl(X)$

- Let F^* be the optimal solution for X, then $\pi(F^*)$ is a feasible solution for $\pi(X)$
- $\operatorname{ufl}(\pi(X)) \leq \operatorname{cost}(\pi(X), \pi(F^*)) \lesssim \operatorname{cost}(X, F^*) = \operatorname{ufl}(X)$
- Suffices to preserve the cost of one solution

• Hard direction: $ufl(\pi(X)) \ge (1 - \varepsilon) ufl(X)$

- Optimal solution F_{π}^* for $\pi(X)$ is random
- Need to preserve the cost of all solutions
- Idea: metric decomposition

Technical Overview: Decomposition

Construct a partition ${\bf \Lambda}$ for X, s.t. for parameter $\kappa=\Theta({\rm ddim}/\varepsilon)^{\Theta({\rm ddim})}$

- (a) Every cluster $C \in \mathbf{\Lambda}$ satisfies $\mathrm{ufl}(C) = \Theta(\kappa)$
- (b) $\sum_{C \in \mathbf{\Lambda}} \operatorname{ufl}(C) \in (1 \pm \varepsilon) \cdot \operatorname{ufl}(X)$



Property (a): κ determines the target dimension m = O(ε⁻² log κ) = Õ(ε⁻²ddim)

• Property (b):
$$(1 + \varepsilon)$$
-apx

 $\begin{array}{l} \mbox{Step 1 Construct a partition } {\bf \Lambda} \mbox{ for } X, \mbox{ s.t. for parameter} \\ \kappa = \Theta(\mathrm{ddim}/\varepsilon)^{\Theta(\mathrm{ddim})} \\ \mbox{ (a) Every cluster } C \in {\bf \Lambda} \mbox{ satisfies } \mathrm{ufl}(C) = \Theta(\kappa) \\ \mbox{ (b) } \sum_{C \in {\bf \Lambda}} \mathrm{ufl}(C) \in (1 \pm \varepsilon) \cdot \mathrm{ufl}(X) \end{array}$

 $\begin{array}{l} \text{Step 2 } \operatorname{ufl}(\pi(X)) \geq \sum_{C \in \mathbf{\Lambda}} \operatorname{ufl}(\pi(C)) - \varepsilon \cdot \operatorname{ufl}(X) \\ \quad \bullet \quad \text{Property (b) carries over to the target space} \end{array}$

Step 3
$$\sum_{C \in \mathbf{\Lambda}} \operatorname{ufl}(\pi(C)) \ge (1 - \varepsilon) \sum_{C \in \mathbf{\Lambda}} \operatorname{ufl}(C)$$

Apply *k*-median results in [MMR19] to every cluster $C \in \mathbf{\Lambda}$
Target dimension $m = O(\varepsilon^{-2} \log \kappa)$ suffices

Step 4 Apply property (b) $\sum_{C \in \Lambda} \operatorname{ufl}(C) \ge \operatorname{ufl}(X)$

Our Decomposition Procedure

Hierarchical Decomposition [Tal04]

 A generalization of randomly shifted quadtree to doubling metrics



 $\blacktriangleright \mathsf{Node} \leftrightarrow \mathsf{cluster}, \mathsf{children} \subseteq \mathsf{parent}$

Root: X

- Leaves: singletons
- Level *i*: diameter $\Theta(2^i)$

▶ Each node (cluster) has 2^{O(ddim)} child nodes (clusters)

Cutting Probability



Cutting probability [Tal04]: $\forall x, y \in X$,

 $\Pr_{\mathcal{H}}[x,y \text{ are in different clusters at level } i] \leq O(\operatorname{ddim}) \cdot \operatorname{dist}(x,y)/2^i$

Cutting Probability



Cutting probability [Tal04]: $\forall x, y \in X$,

 $\Pr_{\mathcal{H}}[x, y \text{ are in different clusters at level } i] \leq O(\operatorname{ddim}) \cdot \operatorname{dist}(x, y) / 2^i$

Badly-cut pair [CFS21]: say (x, y) is badly cut, if x, y are in different clusters at level $\log(\varepsilon^{-1} \operatorname{ddim} \cdot \operatorname{dist}(x, y))$

•
$$\Pr[(x, y) \text{ is badly cut}] \leq O(\varepsilon)$$

Property (a): Every cluster $C \in \mathbf{\Lambda}$ satisfies $\mathrm{ufl}(C) = \Theta(\kappa)$



Property (a): Every cluster $C \in \mathbf{\Lambda}$ satisfies $ufl(C) = \Theta(\kappa)$



Property (a): Every cluster $C \in \mathbf{\Lambda}$ satisfies $ufl(C) = \Theta(\kappa)$



Property (a): Every cluster $C \in \Lambda$ satisfies $ufl(C) = \Theta(\kappa)$



Property (a): Every cluster $C \in \mathbf{\Lambda}$ satisfies $ufl(C) = \Theta(\kappa)$



Property (a): Bottom-up Construction Property (a): Every cluster $C \in \Lambda$ satisfies $ufl(C) = \Theta(\kappa)$

Given threshold $\kappa = \Theta(\operatorname{ddim}/\varepsilon)^{O(\operatorname{ddim})}$, find the lowest level "heavy cluster" (ufl $(C) \ge \kappa$) in a bottom-up manner



- Every cluster $C \in \Lambda$ satisfies $\kappa \leq ufl(C) \leq 2^{O(ddim)}\kappa$
- Previous top-down construction [CLMS13]: polylog(n) upper bound

• Extra $\log \log n$ factor in the target dimension

Property (b): Proof Idea



Idea: adding extra facilities to F^* to serve each $C \in \mathbf{\Lambda}$ locally

$$\blacktriangleright F' := F^* \cup (\bigcup_{C \in \Lambda} \underbrace{N_C}_{\text{a net on } C})$$

Connection cost:
$$\forall C \in \mathbf{\Lambda}, \forall x \in C$$
,
 $\operatorname{dist}(x, F' \cap C) \leq (1 + \varepsilon) \operatorname{dist}(x, F^*)$

▶ Opening cost: $|F'| - |F^*| \le \varepsilon \sum_{C \in \Lambda} \operatorname{ufl}(C)$

Property (b): Connection Cost

Goal: dist $(x, F' \cap C) \leq (1 + \varepsilon) \operatorname{dist}(x, F^*)$



Let F^\ast be the optimal solution for X and assume $F^\ast \subseteq X$ for simplicity

Property (b): Connection Cost

Goal: dist $(x, F' \cap C) \leq (1 + \varepsilon) \operatorname{dist}(x, F^*)$



If $F^*(x) \in C$, then \checkmark

Property (b): Connection Cost

Goal: dist $(x, F' \cap C) \leq (1 + \varepsilon) \operatorname{dist}(x, F^*)$



What if $F^*(x) \notin C$?

Property (b): Connection Cost Goal: $dist(x, F' \cap C) \le (1 + \varepsilon) dist(x, F^*)$



Key observation: separation property

- ▶ Badly-cut pair: say (x, y) is badly cut, if x, y are in different clusters at level log(ε⁻¹ddim · dist(x, y))
- Separation property: if $(x, F^*(x))$ is not badly cut, then $x \in C, F^*(x) \notin C \Longrightarrow \operatorname{dist}(x, F^*(x)) \ge \varepsilon \cdot \operatorname{diam}(C)$
 - An $\varepsilon \cdot \operatorname{diam}(C)$ -net N_C suffices

Property (b): Connection Cost Goal: $dist(x, F' \cap C) \le (1 + \varepsilon) dist(x, F^*)$



Key observation: separation property

- ▶ Badly-cut pair: say (x, y) is badly cut, if x, y are in different clusters at level log(ε⁻¹ddim · dist(x, y))
- Separation property: if $(x, F^*(x))$ is not badly cut, then $x \in C, F^*(x) \notin C \Longrightarrow \operatorname{dist}(x, F^*(x)) \ge \varepsilon \cdot \operatorname{diam}(C)$
 - An $\varepsilon \cdot \operatorname{diam}(C)$ -net N_C suffices

Property (b): Opening Cost





▶ $|N_C|$ is small: $|N_C| \le (1/\varepsilon)^{O(\text{ddim})} \le \varepsilon \kappa \le \varepsilon \cdot \text{ufl}(C)$ ▶ $|F'| - |F^*| = \sum_{C \in \mathbf{\Lambda}} |N_C| \le \varepsilon \sum_{C \in \mathbf{\Lambda}} \text{ufl}(C)$

Eliminate Badly-Cut Pairs

Remaining issue: what if $(x, F^*(x))$ is badly-cut?

Prior solution [CFS21]

• Move x to $F^*(x)$, creating a new instance X'

- Algorithmically, only feasible to eliminate badly-cut pairs $(x, F_0(x))$ for some approximation solution F_0
- X' is random (randomness comes from \mathcal{H})
- ▶ Hard to deal with remaining badly-cut pairs $(x, F^*(x))$

Eliminate Badly-Cut Pairs

Remaining issue: what if $(x, F^*(x))$ is badly-cut?

Prior solution [CFS21]

• Move x to $F^*(x)$, creating a new instance X'

- Algorithmically, only feasible to eliminate badly-cut pairs $(x, F_0(x))$ for some approximation solution F_0
- X' is random (randomness comes from \mathcal{H})
- ▶ Hard to deal with remaining badly-cut pairs $(x, F^*(x))$

Our approach: modify \mathcal{H} instead of X

- At each level, force close points to be clustered together
- Construct Λ accordingly
- Can still apply cutting probability to handle remaining badly-cut pairs

Thank you!