



**HAL**  
open science

## Near-neighbor preserving dimension reduction via coverings for doubling subsets of $\ell_1$

Ioannis Z. Emiris, Vasilis Margonis, Ioannis Psarros

► **To cite this version:**

Ioannis Z. Emiris, Vasilis Margonis, Ioannis Psarros. Near-neighbor preserving dimension reduction via coverings for doubling subsets of  $\ell_1$ . *Theoretical Computer Science*, 2022, 942, pp.169-179. 10.1016/j.tcs.2022.11.031 . hal-04294296

**HAL Id: hal-04294296**

**<https://hal.science/hal-04294296>**

Submitted on 28 Nov 2023

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

# Near-Neighbor Preserving Dimension Reduction via Coverings for Doubling Subsets of $\ell_1$

Ioannis Z. Emiris<sup>a,b,1</sup>, Vasilis Margonis<sup>a,2,\*</sup>, Ioannis Psarros<sup>a,3</sup>

<sup>a</sup>*Department of Informatics & Telecommunications,  
National & Kapodistrian University of Athens, Athens 15784, Greece*

<sup>b</sup>*ATHENA Research & Innovation Center, Maroussi 15125, Greece*

---

## Abstract

Nearest neighbor-preserving embeddings exist for  $\ell_2$  (Euclidean) and  $\ell_1$  (Manhattan) metrics, as well as doubling subsets of  $\ell_2$ , where doubling dimension is today the most effective way of capturing input structure. These randomized embeddings bound the distortion only for distances between the query point and a point-set. Motivated by the study of fast Approximate Nearest Neighbor search in  $\ell_1$ , this paper settles the missing case of doubling subsets of  $\ell_1$ . In particular, we introduce a randomized dimensionality reduction by means of a *near* neighbor-preserving embedding; the latter is related to the decision-with-witness problem. The input set gets represented via appropriate covering point-sets. For this, we leverage either approximate  $r$ -nets or randomly shifted grids, with different tradeoffs between preprocessing time and target dimension. We exploit Cauchy random variables, and derive a concentration bound of independent interest.

---

\*Corresponding author

*Email addresses:* [emiris@di.uoa.gr](mailto:emiris@di.uoa.gr) (Ioannis Z. Emiris), [basilis.math@gmail.com](mailto:basilis.math@gmail.com) (Vasilis Margonis), [ipsarros@uni-bonn.de](mailto:ipsarros@uni-bonn.de) (Ioannis Psarros)

<sup>1</sup>Funding: Partially supported by the European Union's H2020 research and innovation programme under grant agreement No. 734242 (LAMBDA).

<sup>2</sup>Present address: Institute of Informatics and Telecommunications, National Center of Scientific Research "Demokritos", Agia Paraskevi 15341, Greece.

<sup>3</sup>Present address: Institute of Computer Science, University of Bonn, Endenicher Allee 19a, 53115 Bonn, Germany. Funding: this research has been co-financed by Greece and the European Union (European Social Fund) through the Operational Programme "Human Resources Development, Education and Lifelong Learning" in the context of the project "Strengthening Human Resources Research Potential via Doctorate Research" (MIS-5000432), implemented by the State Scholarships Foundation (IKY) of Greece.

*Keywords:* Approximate Nearest Neighbor, Manhattan metric, doubling dimension, randomized embedding, dimensionality reduction, Cauchy distribution, covering point set

---

## 1. Introduction

Proximity search is a key computational question with a wide variety of applications. The corresponding problems in metric spaces of low dimension have been usually handled by space tessellation. Such solutions are affected by the curse of dimensionality, which makes them too costly in high dimensions. In the past two decades, the increasing need for manipulating high-dimensional data has led to randomized and approximation algorithms with polynomial dependence on the dimension, since the latter cannot be assumed as fixed but is part of the input parameters.

Approximate Nearest Neighbor search, which is an optimization problem, constitutes a cornerstone question in this area. Existing reductions (see, e.g. [1]), which incur a polylogarithmic time overhead, reduce this search to the following decision problem with witness, namely the  $(c, R)$ -Approximate *Near* Neighbor question:

**Definition 1** (Approximate Near Neighbor). *Let  $(X, d_X)$  be a metric space. Given  $P \subseteq X$  and reals  $R > 0$ ,  $c \geq 1$ , construct a data structure  $\mathcal{S}$  which, given a query point  $q \in X$ , performs as follows:*

- *If the nearest neighbor of  $q$  lies within distance at most  $R$ , then  $\mathcal{S}$  is allowed to report any point  $p^* \in P$ , such that  $d_X(q, p^*) \leq cR$ .*
- *If all points lie at distance more than  $cR$  from  $q$ , then  $\mathcal{S}$  returns  $\perp$ .*

*The data structure  $\mathcal{S}$  always returns either a point at distance  $\leq cR$  or  $\perp$ , even when none of the above two cases occurs.*

We shall now suppose that  $R = 1$ , because this can always be achieved by re-scaling the input point set, and we refer to this problem as  $c$ -ANN, or simply

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

ANN. We focus on subsets of  $\ell_1^d$ , in other words, the input dataset consists of  $n$  vectors in  $\mathbb{R}^d$  endowed with the standard  $\ell_1$  distance function, also known as Manhattan metric, whose norm is denoted by  $\|\cdot\|_1$ . Note that all logarithms are base 2.

**Previous work.** Major landmarks in the study of data structures for high-dimensional normed spaces are the different variants, proofs, and applications of the celebrated Johnson-Lindenstrauss Lemma (e.g. [2, 3, 4]), sketches based on  $p$ -stable distributions [5], and Locality Sensitive Hashing (e.g. [6, 7, 8]). In the core of most high-dimensional solutions, lies the fact that for certain metric spaces, e.g.  $\ell_p$  for  $p \in [1, 2]$ , the distance can be efficiently sketched. Spaces which are considered to be harder in this context, such as  $\ell_\infty$ , can also be treated [9], and are quite interesting since they can be used as host spaces for various norms [10].

Significant amount of work has been undertaken for point sets of low *doubling* dimension, since it is today one of the primary paradigms for capturing input structure (formal definitions in the next section). For any finite metric space  $X$  of doubling dimension  $\dim(X)$ , there exists a data structure [11, 12] with expected preprocessing time  $O(2^{\dim(X)}n \log n)$ , space usage  $O(2^{\dim(X)}n)$  (or even  $O(n)$ ) and query time  $O(2^{\dim(X)} \log n + \varepsilon^{-O(\dim(X))})$ .

Indyk and Naor [13] introduced the notion of nearest-neighbor preserving embeddings, which are randomized embeddings between two metric spaces with guaranteed bounded distortion only for the distances between a query point and a point set. They achieved a dimension reduction for doubling subsets of  $\ell_2$ , with the target dimension depending only on the input dataset's doubling dimension. Even before, Indyk [5] had introduced a randomized embedding for dimension reduction in  $\ell_1$ , which is suitable for proximity search purposes, and it achieves target dimension polylogarithmic in the size of the point set. Naturally, such approaches can be easily combined with any known data structure to be used in the projection space. Randomized embeddings have been recently used in the ANN context [14], for doubling subsets of  $\ell_p$ ,  $2 < p < \infty$ .

Dimensionality reduction in  $\ell_1$  cannot be achieved in the same generality as in  $\ell_2$ , even by assuming that the point set is of low doubling dimension [15]: There are arbitrarily large  $n$ -point subsets  $P \subseteq \ell_1$  which are doubling with constant 6, such that every embedding with distortion  $D$  of  $P$  into  $\ell_1^k$  requires  $k = n^{\Omega(1/D^2)}$ . Aiming at more restrictive guarantees, e.g. preserving distances within some pre-defined range, is a relevant workaround. Then, dimension reduction techniques for doubling subsets of  $\ell_p$ ,  $p \in [1, 2]$ , have been proposed in [16], but they rely on partition algorithms which require the whole point set to be known in advance. Hence, applicability of such techniques is very limited and, specifically, they do not seem amenable to an online setting where query points are not known in advance.

**Our results.** A new dimensionality reduction scheme is proposed, by means of a *near* neighbor-preserving embedding for doubling subsets of  $\ell_1$ , thus settling a case that remained open. Our definition is essentially a modified version of the nearest neighbor preserving embedding of [13]: the required guarantees are weaker, since we consider the decision version of the problem, therefore the embedding depends on some range parameter  $R > 0$ .

**Definition 2** (Near-neighbor preserving embedding). *Let  $(Y, d_Y)$ ,  $(Z, d_Z)$  be metric spaces and  $X \subseteq Y$ . A distribution over mappings  $f : Y \rightarrow Z$  is a near-neighbor preserving embedding with range  $R > 0$ , distortion  $D \geq 1$  and probability of correctness  $\mathcal{P} \in [0, 1]$  if for every  $\alpha \geq D$  and any  $q \in Y$ , if  $x \in X$  is such that  $d_Y(x, q) \leq R$ , then with probability at least  $\mathcal{P}$ ,*

- $d_Z(f(x), f(q)) \leq D \cdot R$ ,
- $\forall p \in X : d_Y(p, q) > D \cdot \alpha \cdot R \implies d_Z(f(p), f(q)) > \alpha \cdot R$ .

Considering a point set  $P \subset \ell_1^d$  of cardinality  $n$ , our results concern  $\ell_1^k$  as the target space, where  $k$  depends on the doubling dimension of  $P$ . We assume that  $R = 1$ , since we can re-scale the given point set. More specifically:

1. In Theorem 10, we prove that for every  $\varepsilon \in (0, 1/2)$  and  $c \geq 1$ , there is a randomized mapping  $h : \ell_1^d \rightarrow \ell_1^k$  that can be computed in time

1  
2  
3  
4  
5  
6  
7  
8  
9  $\tilde{O}(dn^{1+1/\Omega(c)})$  and is *near* neighbor-preserving for  $P$  with distortion  $1+6\varepsilon$   
10 and probability of correctness  $\Omega(\varepsilon)$ , where

$$11 \quad k = (\log \lambda_P \cdot \log(c/\varepsilon))^{\Theta(1/\varepsilon)} / \zeta(\varepsilon),$$

12  
13  
14  
15 for a function  $\zeta(\varepsilon) > 0$  depending only on  $\varepsilon$ . Although the mapping  $h$   
16 depends on the point set, the parameter  $c$  is user-defined and therefore  
17 provides a trade-off between preprocessing time and target dimension.  
18  
19

- 20  
21 2. In Theorem 13, we show that for every  $\varepsilon \in (0, 1/2)$ , there is a randomized  
22 mapping  $h' : \ell_1^d \rightarrow \ell_1^k$  that can be computed in time  $O(dkn)$ , and is  
23 *near* neighbor-preserving for  $P$  with distortion  $1+6\varepsilon$  and probability of  
24 correctness  $\Omega(\varepsilon)$ , where

$$25 \quad k = (\log \lambda_P \cdot \log(d/\varepsilon))^{\Theta(1/\varepsilon)} / \zeta(\varepsilon),$$

26  
27  
28 for a function  $\zeta(\varepsilon) > 0$  depending only on  $\varepsilon$ . In this case, the function  
29  $h'$  is oblivious to  $P$  and well-defined over the whole space, but the target  
30 dimension depends on  $d$ .  
31  
32  
33  
34  
35

36 On the low-preprocessing-time extreme, one can embed the dataset in near-  
37 linear time, but the target dimension is polynomial in  $\log \log n$ . This is to  
38 be juxtaposed to the analogous result in [5], which achieves target dimension  
39 polynomial in  $\log n$ , without any assumption on the doubling dimension of the  
40 dataset. On the other hand, it is possible to obtain a preprocessing time of  
41  $dn^{1+\delta}$  for any constant  $\delta > 0$ , and target dimension which depends solely on  
42 the doubling dimension.  
43  
44  
45  
46  
47

48 **Methodology.** Both of the aforementioned embeddings  $h, h'$  consist of two  
49 basic components. First, we represent the point set  $P$  with an  $\varepsilon$ -covering set,  
50 and then we apply a random linear projection à la Indyk [5] to that set, using  
51 Cauchy variables.  
52  
53  
54

55 The role of the covering set is to exploit the doubling dimension of  $P$ . In the  
56 analogous result for  $\ell_2$  in [13], no representative sets were used; the mapping was  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

just a random linear projection of  $P$ . In the case of  $\ell_1$  however, a similar analysis of a linear projection with Cauchy variables without these representative sets seems to be impossible, since the Cauchy distribution is heavy tailed.

In Theorem 10, we consider  $c$ -approximate  $r$ -nets as a covering set. Inspired by the algorithm of [17] for  $\ell_2$ , we design an algorithm that computes a  $c$ -approximate  $r$ -net in  $\ell_1$  in subquadratic, but superlinear, time. On the other hand, Theorem 13 relies on randomly shifted grids, which can be computed in linear time, but are inferior to  $r$ -nets in terms of capturing the doubling dimension of the point set, hence inferior in terms of the target dimension.

To bound the distortion incurred by the randomized projection, we exploit the 1-stability property of the Cauchy distribution. To this end, we establish a new concentration bound for sums of independent Cauchy variables, which should be of interest beyond the scope of this paper. To overcome the technical difficulties associated with the heavy tails of the Cauchy distribution, we study sums of *square roots* of Cauchy variables, whereas in [5], Indyk considers sums of *truncated* Cauchy variables instead. Although our concentration bound is rather weak, it is sufficient for our purposes and its analysis is much simpler compared to that in [5].

**Algorithmic implications.** Our results show that efficient dimension reduction for doubling subsets of  $\ell_1$  is possible, in the context of ANN. In particular, these results imply efficient sketches, meaning that one can solve ANN with minimal storage per point. Dimensionality reduction also serves as a problem reduction from a high-dimensional hard instance to a low-dimensional easy instance. Since the algorithms presented in this paper are quite simple, they should also be of practical interest: they easily extend the scope of any implementation which has been optimized to solve the problem in low dimension, so that it may handle high-dimensional data.

Our embedding can be combined with the bucketing method in [1] for the  $(1+\varepsilon)$ -ANN problem in  $\ell_1^d$ . For instance, setting  $c = \log n$  in Theorem 10, yields preprocessing time  $dn^{1+o(1)}$ , space  $n^{1+o(1)}$  and query time  $O(d) \cdot (\log \lambda_P \cdot$

1  
 2  
 3  
 4  
 5  
 6  
 7  
 8  
 9  $\log \log n)^{O(1/\varepsilon)}$  assuming that the doubling dimension is a fixed constant. This  
 10 improves upon existing results: the query time of [18] depends on the aspect ra-  
 11 tio of the dataset, while the data structures of [11, 12] support queries with time  
 12 complexity which depends exponentially on the doubling dimension. However,  
 13 it is worth noting that one could potentially improve the results of [11, 12, 18] in  
 14 the special case of  $\ell_1$ , by employing ANN data structures with fast query time,  
 15 in order to accelerate the traversal of the net-tree. Hence, while our result gives  
 16 a simple framework for exploiting the intrinsic dimension of doubling subsets of  
 17  $\ell_1$ , it is unlikely that it shall improve upon simple variants of previous results  
 18 in terms of complexity bounds.  
 19  
 20  
 21  
 22  
 23  
 24  
 25  
 26  
 27  
 28  
 29  
 30  
 31  
 32  
 33  
 34  
 35  
 36  
 37  
 38  
 39

This paper is the complete and final version of our results that appeared in preliminary form in [19].

**Paper structure.** The next section introduces basic concepts and some relevant existing results. Section 3 establishes a new concentration bound on sums of independent Cauchy variables. Section 4 achieves dimensionality reduction by means of representing the point set by a carefully chosen net, while Section 5 employs randomly shifted grids for the same task. We conclude with a discussion of our results, and of potential improvements.

## 2. Preliminaries

In this section, we define basic notions about doubling metrics and present some relevant existing results.

**Definition 3.** Consider any metric space  $(X, d_X)$  and let  $B(p, r) = \{x \in X \mid d_X(x, p) \leq r\}$ . The doubling constant of  $X$ , denoted  $\lambda_X$ , is the smallest integer  $\lambda_X$  such that for any  $p \in X$  and  $r > 0$ , the ball  $B(p, r)$  can be covered by  $\lambda_X$  balls of radius  $r/2$  centered at points in  $X$ .

The doubling dimension of  $(X, d_X)$  is defined as  $\log \lambda_X$ . Nets play an important role in the study of embeddings, as well as in designing efficient data structures for doubling metrics.

**Definition 4.** For  $c \geq 1$ ,  $r > 0$  and metric space  $(V, d_V)$ , a  $c$ -approximate  $r$ -net of  $V$  is a subset  $\mathcal{N} \subseteq V$  such that no two points of  $\mathcal{N}$  are within distance  $r$  of each other, and every point of  $V$  lies within distance at most  $c \cdot r$  from some point of  $\mathcal{N}$ .

**Theorem 5.** Let  $P \subset \ell_1^d$  such that  $|P| = n$ . Then, for any  $c > 0$ ,  $r > 0$ , one can compute a  $c$ -approximate  $r$ -net of  $P$  in time  $\tilde{O}(dn^{1+1/c'})$ , where  $c' = \Omega(c)$ . The result is correct with high probability. The algorithm also returns the assignment of each point of  $P$  to the point of the net which covers it.

*Proof.* We employ some standard ideas from [1]. An analogous result for  $\ell_2$  can be found in [17]. Let us first suppose that  $r = 1$ , since we are able to re-scale the point set. Let us now consider a randomly shifted grid with side-length 2. The probability that two points  $p, q \in P$  fall into the same grid cell is at least  $1 - \|p - q\|_1/2$ . For each non-empty grid cell, we snap points to a grid: each coordinate is rounded to the nearest multiple of  $\delta = 1/10dc$ . Then, coordinates are multiplied by  $1/\delta$  and each point  $x = (x_1, \dots, x_d) \in [2\delta]^d$  is mapped to  $\{0, 1\}^{2d/\delta}$  by a function  $G$  as follows:  $G(x) = (g(x_1), \dots, g(x_d))$ , where  $g(z)$  is a binary string of  $z$  ones followed by  $2/\delta - z$  zeros. For any two points  $p, q$  in the same grid cell, let  $f(p), f(q)$  be the two binary strings obtained by the above mapping. Notice that,

$$\|f(p) - f(q)\|_1 \in (2/\delta) \cdot \|p - q\|_1 \pm 1.$$

Hence,

$$\|p - q\|_1 \leq 1 \implies \|f(p) - f(q)\|_1 \leq (2/\delta) + 1,$$

$$\|p - q\|_1 \geq c \implies \|f(p) - f(q)\|_1 \geq (2/\delta) \cdot c - 1.$$

Let us now employ the LSH family of [1], for the Hamming space. After standard concatenation, we may assume that the family is  $(\rho, c'\rho, n^{-1/c'}, n^{-1})$ -sensitive, where  $\rho = (2/\delta) + 1$  and  $c' = \Omega(c)$ . Let  $\alpha = n^{-1/c'}$  and  $\beta = n^{-1}$ .

Notice that for the above two-level hashing table, we obtain the following guarantees: Any two points  $p, q \in P$ , such that  $\|p - q\|_1 \leq 1$ , fall into the same

1  
2  
3  
4  
5  
6  
7  
8  
9 bucket with probability  $\geq \alpha/2$ . Any two points  $p, q \in P$ , such that  $\|p - q\|_1 \geq c$ ,  
10 fall into the same bucket with probability  $\leq \beta$ .

11 Finally, we independently build  $k = \Theta(n^{1/c'} \log n)$  hashtables as above,  
12 where the random hash function is defined as a concatenation of the function  
13 which maps points to their grid cell id and one LSH function. We pick an arbi-  
14 trary ordering  $p_1, \dots, p_n \in P$ . We follow a greedy strategy in order to compute  
15 the approximate net. We start with point  $p_1$ , and we add it to the net. We  
16 mark all (unmarked) points which fall at the same bucket with  $p_1$ , in one of  
17 the  $k$  hashtables, and are at distance  $\leq cr$ . Then, we proceed with point  $p_2$ .  
18 If  $p_2$  is unmarked, then we repeat the above. Otherwise, we proceed with  $p_3$ .  
19 The above iteration stops when all points have been marked. Throughout the  
20 procedure, we are able to store one pointer for each point, indicating the center  
21 which covered it.  
22  
23  
24  
25  
26  
27  
28  
29

30  
31 **Correctness.** The probability that a good pair  $p, q$  does not fall into the  
32 same bucket for any of the  $k$  hashtables is  $\leq (1 - \alpha/2)^k \leq n^{-10}$ . Hence, with  
33 high probability, the packing property holds, and the covering property holds  
34 because the above algorithm stops when all points are marked.  
35  
36  
37

38 **Running time.** The time to build the  $k$  hashtables is  $k \cdot n = \tilde{O}(n^{1+1/c'})$ .  
39 Then, at most  $n$  queries are performed: for each query, we investigate  $k$  buck-  
40 ets and the expected number of false positives is  $\leq k \cdot n^2 \cdot \beta = \tilde{O}(n^{1+1/c'})$ .  
41 Hence, if we stop after having seen a sufficient amount of false positives, we ob-  
42 tain time complexity  $\tilde{O}(n^{1+1/c'})$  and the covering property holds with constant  
43 probability. We can repeat the above procedure  $O(\log n)$  times to obtain high  
44 probability of success.  $\square$   
45  
46  
47  
48  
49

50 The main result in the context of randomized embeddings for dimension  
51 reduction in  $\ell_1^d$  is the following theorem, which exploits the 1-stability property  
52 of Cauchy random variables and provides with an asymmetric guarantee: The  
53 probability of non-contraction is high, but the probability of non-expansion  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

is constant. Nevertheless, this asymmetric property is sufficient for proximity search.

**Theorem 6** (Thm 5, [5]). *For any  $\varepsilon \leq 1/2$ ,  $\delta > 0$ ,  $\varepsilon > \gamma > 0$  there is a probability space over linear mappings  $f : \ell_1^d \rightarrow \ell_1^k$ , where  $k = (\ln(1/\delta))^{1/(\varepsilon-\gamma)}/\zeta(\gamma)$ , for a function  $\zeta(\gamma) > 0$  depending only on  $\gamma$ , such that for any pair of points  $p, q \in \ell_1^d$ :*

$$\begin{aligned} \Pr [\|f(p) - f(q)\|_1 \leq (1 - \varepsilon) \|p - q\|_1] &\leq \delta, \\ \Pr [\|f(p) - f(q)\|_1 \geq (1 + \varepsilon) \|p - q\|_1] &\leq \frac{1 + \gamma}{1 + \varepsilon}. \end{aligned}$$

Note that the embedding is defined as  $f(u) = Au/T$ , where  $A$  is a  $k \times d$  matrix with each element being an i.i.d. Cauchy random variable. In addition,  $T$  is a scaling factor defined as the expectation of a sum of truncated Cauchy variables, such that  $T = \Theta(k \log(k/\varepsilon))$ , see [5, Lem. 5].

One key observation here is that, given a point set  $P$  in a space of bounded aspect ratio  $\Phi$ , one can directly employ Theorem 6. The number of points can be upper bounded by a function of  $\lambda_P$  and  $\Phi$ , and hence the new dimension  $k$  depends only on these parameters. Here, we establish better bounds than those in Theorem 6 for doubling subsets of  $\ell_1^d$ , without any assumption on the aspect ratio.

### 3. Concentration bounds for Cauchy variables

In this section, we prove some basic properties of the Cauchy distribution, which serves as our main embedding tool.

Let  $C_{\mathcal{D}}$  denote the Cauchy distribution with density  $c(x) = (1/\pi)/(1 + x^2)$ . One key property of the Cauchy distribution is the so-called 1-stability property: Let  $v = (v_1, \dots, v_k) \in \mathbb{R}^k$  and  $X_1, \dots, X_k$  be i.i.d. random variables following  $C_{\mathcal{D}}$ , then  $\sum_{i=1}^k X_i v_i$  is distributed as  $X \cdot \|v\|_1$ , where  $X \sim C_{\mathcal{D}}$ .

The Cauchy distribution has undefined mean. However, for  $0 < q < 1$ , the mean of the  $q$ -th power of a Cauchy random variable can be defined. More

specifically, for some  $X \sim C_{\mathcal{D}}$  we have

$$\mathbb{E}\left[|X|^{1/2}\right] = \frac{2}{\pi} \int_0^{\infty} \frac{\sqrt{x}}{1+x^2} dx = \frac{2}{\pi} \frac{\pi}{\sqrt{2}} = \sqrt{2}.$$

The following lemma provides a bound for the moment-generating function of  $|X|^{1/2}$ .

**Lemma 7.** *Let  $X \sim C_{\mathcal{D}}$ . Then for any  $\beta > 1$ :*

$$\mathbb{E}\left[\exp(-\beta|X|^{1/2})\right] \leq \frac{2}{\beta}.$$

*Proof.* For any constant  $\beta$ ,

$$\int_0^1 e^{-\beta x^{1/2}} dx = \frac{2}{\beta^2} \left(1 - \frac{\beta+1}{e^{\beta}}\right).$$

Then, for any  $\beta > 1$ ,

$$\begin{aligned} \mathbb{E}\left[\exp(-\beta|X|^{1/2})\right] &= \int_{-\infty}^{\infty} e^{-\beta|x|^{1/2}} \cdot c(x) dx = \frac{2}{\pi} \int_0^{\infty} e^{-\beta x^{1/2}} \cdot \frac{1}{1+x^2} dx \\ &= \frac{2}{\pi} \int_0^1 e^{-\beta x^{1/2}} \cdot \frac{1}{1+x^2} dx + \frac{2}{\pi} \int_1^{\infty} e^{-\beta x^{1/2}} \cdot \frac{1}{1+x^2} dx \\ &\leq \frac{2}{\pi} \int_0^1 e^{-\beta x^{1/2}} dx + \frac{2}{\pi} \int_1^{\infty} e^{-\beta} \cdot \frac{1}{1+x^2} dx \\ &= \frac{2}{\pi} \cdot \frac{2}{\beta^2} \left(1 - \frac{\beta+1}{e^{\beta}}\right) + \frac{1}{2e^{\beta}} \\ &\leq \frac{4}{\pi\beta^2} + \frac{1}{2e^{\beta}} \\ &\leq \frac{2}{\beta}. \quad \square \end{aligned}$$

Let  $S := \sum_{j=1}^k |X_j|$  where each  $X_j$  is an i.i.d. Cauchy variable. To prove concentration bounds for  $S$ , we study the sum  $\tilde{S} := \sum_{j=1}^k |X_j|^{1/2}$ . By Hölder's Inequality, for any  $x \in \mathbb{R}^d$  and  $p > q > 0$ ,

$$\|x\|_p \leq \|x\|_q \leq d^{1/q-1/p} \|x\|_p.$$

Consequently, for  $x = (X_1, \dots, X_k) \in \mathbb{R}^k$ ,  $p = 1$  and  $q = 1/2$  we have that  $S \leq \tilde{S}^2 \leq k \cdot S$ , hence for any  $t > 0$ ,

$$\Pr[S \leq t] \leq \Pr[\tilde{S} \leq \sqrt{tk}]. \quad (1)$$

We use the bound on the moment-generating function, to prove a Chernoff-type concentration bound for  $\tilde{S}$ , which by Eq. (1) translates into a concentration bound for  $S$ .

**Lemma 8.** *For every  $D > 1$ ,*

$$\Pr \left[ \tilde{S} \leq \frac{\mathbb{E}[\tilde{S}]}{D} \right] \leq \left( \frac{10}{D} \right)^k.$$

*Proof.* Since  $X_j$ 's are independent,  $\mathbb{E}[\tilde{S}] = \sqrt{2}k$ . Then, by Lemma 7 and Markov's inequality, for any  $\beta > 1$ , it follows that

$$\begin{aligned} \Pr \left[ \tilde{S} \leq \frac{\mathbb{E}[\tilde{S}]}{D} \right] &= \Pr \left[ \exp(-\beta\tilde{S}) \geq \exp \left( -\beta \cdot \frac{\mathbb{E}[\tilde{S}]}{D} \right) \right] \\ &\leq \frac{\mathbb{E}[\exp(-\beta\tilde{S})]}{\exp(-\beta \mathbb{E}[\tilde{S}]/D)} \\ &= \frac{\mathbb{E}[\exp(-\beta|X_j|^{1/2})]^k}{\exp(-\beta\sqrt{2}k/D)} \\ &\leq \left( \frac{2}{\beta} \right)^k \cdot e^{\sqrt{2}\beta k/D}. \end{aligned}$$

Setting  $\beta = D$  completes the proof.  $\square$

#### 4. Net-based dimension reduction

This section describes the dimension reduction mapping for  $\ell_1$  via  $r$ -nets. Let  $P \subset \ell_1^d$  be a set of  $n$  points with doubling constant  $\lambda_P$ . For some point  $x \in \mathbb{R}^d$  and  $r > 0$ , we denote by  $B_1(x, r)$  the  $\ell_1$ -ball of radius  $r$  around  $x$ . The embedding is non-linear and is carried out in two steps.

First, let us compute a  $c$ -approximate  $(\varepsilon/c)$ -net  $\mathcal{N}$  of  $P$  with the algorithm of Theorem 5. Moreover, the algorithm assigns each point of  $P$  to the point of  $\mathcal{N}$  which covered it. Let  $g : P \rightarrow \mathcal{N}$  be this assignment. In the second step, for every  $s \in \mathcal{N}$  and any query point  $q \in \ell_1^d$ , we apply the linear map of Theorem 6. That is,  $f(s) = As/T$ , where  $A$  is a  $k \times d$  matrix with each element being an i.i.d. Cauchy random variable. Recall that value  $T$  satisfies

$$T = \Theta(k \log(k/\varepsilon)).$$

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

By the 1-stability property of the Cauchy distribution,  $f(s)$  is distributed as  $\|s\|_1 \cdot (Y_1, \dots, Y_k)$ , where each  $Y_j$  is i.i.d. and  $Y_j \sim C_{\mathcal{D}}$ . Hence,  $\|f(s)\|_1 = \|s\|_1 \cdot S$  where  $S := \sum_j |Y_j|$ .

We now define the embedding to be  $h = f \circ g$ . We apply  $h$  to every point in  $P$ , and  $f$  to any query point  $q$ . It is clear from the properties of the net that  $g$  incurs an additive error of  $\pm\varepsilon$  on the distance between  $q$  and any point in  $P$ , so it is sufficient to consider the distortion of  $f$ .

Our analysis consists of studying separately the following disjoint subsets of  $\mathcal{N}$ : Points that lie at distance at most  $D_0$  from the query and points that lie at distance at least  $D_0$ , for some  $D_0 > 1$  chosen appropriately. For the former set, we directly apply Theorem 6, since it has bounded diameter.

The next lemma guarantees the low distortion for points of the latter set, namely those that are sufficiently far from the query. We consider the sum of the square roots of each  $|Y_j|$ , i.e.,  $\tilde{S} = \sum_j |Y_j|^{1/2}$ , in order to employ the tools of Section 3.

**Lemma 9.** *Fix a query point  $q \in \ell_1^d$ . For any  $\varepsilon \leq 1/2$ ,  $c \geq 1$ ,  $\delta \in (0, 1)$ , there exists  $D_0 = O(\log(k/\varepsilon))$  such that for  $k = \Theta(\log^2 \lambda_P \cdot \log(c/\varepsilon) + \log(1/\delta))$ , with probability at least  $1 - \delta$ ,*

$$\forall s \in \mathcal{N} : \|s - q\|_1 \geq D_0 \implies \|f(s) - f(q)\|_1 \geq 4.$$

*Proof.* Assume wlog that the query point is the origin  $(0, \dots, 0)$ . For some  $D_0 > 1$ , we define the following subsets of  $\mathcal{N}$ :

$$N_i := \{s \in \mathcal{N} \mid D_i \leq \|s\|_1 < D_{i+1}\}, \quad D_i = 2^{2^i} D_0, \quad i = 0, 1, 2, \dots$$

By the definition of doubling constant and the fact that two points of  $\mathcal{N}$  lie at distance at least  $\varepsilon$ ,

$$|N_i| \leq \lambda_P^{\lceil \log(4cD_{i+1}/\varepsilon) \rceil} \leq \lambda_P^{4 \log(cD_{i+1}/\varepsilon)}.$$

Therefore, by the union bound, and Eq. (1):

$$\begin{aligned}
\Pr \left[ \exists i \exists s \in N_i : \|f(s)\|_1 \leq \frac{4\|s\|_1}{D_i} \right] &= \Pr \left[ \exists i \exists s \in N_i : S \leq \frac{4T}{D_i} \right] \\
&\leq \sum_{i=0}^{\infty} |N_i| \Pr \left[ \tilde{S} \leq \frac{\sqrt{4kT}}{\sqrt{D_i}} \right] \\
&= \sum_{i=0}^{\infty} |N_i| \Pr \left[ \tilde{S} \leq \mathbb{E}[\tilde{S}] \cdot \sqrt{\frac{2T}{k2^{2i}D_0}} \right].
\end{aligned}$$

By Lemma 8, for  $D_0 = \lceil 800T/k \rceil = \Theta(\log(k/\varepsilon))$  and  $k > 4 \cdot \log \lambda_P \cdot \log(cD_0/\varepsilon) + 2 \log(2\lambda_P/\delta)$ :

$$\begin{aligned}
\sum_{i=0}^{\infty} |N_i| \Pr \left[ \tilde{S} \leq \frac{\mathbb{E}[\tilde{S}]}{10 \cdot 2^{i+1}} \right] &\leq \sum_{i=0}^{\infty} \lambda_P^{4 \log(cD_{i+1}/\varepsilon)} \left( \frac{1}{2^{i+1}} \right)^k \\
&= \sum_{i=0}^{\infty} \frac{2^{\log(\lambda_P)(4 \log(cD_0/\varepsilon) + 2i + 2)}}{2^{k(i+1)}} \\
&\leq \sum_{i=0}^{\infty} \frac{2^{\log(\lambda_P) \cdot 4 \log(cD_0/\varepsilon)} \cdot 2^{2 \log(\lambda_P)(i+1)}}{2^{(4 \cdot \log \lambda_P \cdot \log(cD_0/\varepsilon))(i+1)} \cdot 2^{2 \log(2\lambda_P/\delta)(i+1)}} \\
&\leq \sum_{i=0}^{\infty} 2^{-2 \log(2/\delta)(i+1)} \\
&= \sum_{i=0}^{\infty} \left( \frac{\delta^2}{4} \right)^i - 1 \\
&= \frac{\delta^2}{4 - \delta^2} \\
&\leq \delta.
\end{aligned}$$

Finally, for some large enough constant  $C$ , we demand that

$$k > C (\log \lambda_P \cdot \log(c \log k/\varepsilon) + \log(1/\delta)) > 4 \cdot \log \lambda_P \cdot \log(cD_0/\varepsilon) + 2 \log(2\lambda_P/\delta)$$

which is satisfied for  $k = \Theta(\log^2 \lambda_P \cdot \log(c/\varepsilon) + \log(1/\delta))$ .  $\square$

**Theorem 10.** *Let  $P \subset \ell_1^d$  such that  $|P| = n$ . For any  $\varepsilon \in (0, 1/2)$  and  $c \geq 1$ , there is a non-linear randomized embedding  $h = f \circ g : \ell_1^d \rightarrow \ell_1^k$ , where*

$$k = (\log \lambda_P \cdot \log(c/\varepsilon))^{\Theta(1/\varepsilon)} / \zeta(\varepsilon),$$

for a function  $\zeta(\varepsilon) > 0$  depending only on  $\varepsilon$ , such that, for any  $q \in \ell_1^d$ , if there exists  $p^* \in P$  such that  $\|p^* - q\|_1 \leq 1$ , then, with probability  $\Omega(\varepsilon)$ :

$$\|h(p^*) - f(q)\|_1 \leq 1 + 3\varepsilon,$$

$$\forall p \in P : \|p - q\|_1 > 1 + 9\varepsilon \implies \|h(p) - f(q)\|_1 > 1 + 3\varepsilon.$$

Set  $P$  can be embedded in time  $\tilde{O}(dn^{1+1/\Omega(c)})$ , and any query  $q \in \ell_1^d$  can be embedded in time  $O(dk)$ .

*Proof.* Let  $f, g$  be the mappings defined in the beginning of the section and  $D_0 = \Theta(\log(k/\varepsilon))$ . Assume wlog for simplicity that  $q = 0^d$ . Then, by Lemma 9 for  $k = \Theta(\log^2 \lambda_P \cdot \log(c/\varepsilon))$ , with probability at least  $1 - \varepsilon/5$ , we have:

$$\forall p \in P : \|p - q\|_1 \geq D_0 + \varepsilon \implies \|h(p) - f(q)\|_1 \geq 4.$$

By Theorem 6, for  $\gamma = \varepsilon/10$  and  $\delta = \varepsilon/(5\lambda_P^{8 \log(cD_0/\varepsilon)})$ , with probability at least  $1 - \varepsilon/5$ , we get:

$$\forall p \in P : \|p - q\|_1 \in (1+9\varepsilon, D_0 + \varepsilon) \implies \|h(p) - f(q)\|_1 > (1+8\varepsilon)(1-\varepsilon) \geq 1+3\varepsilon.$$

Moreover,

$$\Pr [\|h(p^*) - f(q)\|_1 \leq 1 + 3\varepsilon] \geq 1 - \frac{1 + \varepsilon/10}{1 + \varepsilon} \geq 1 - (1 - \varepsilon/2).$$

Then, the target dimension needs to satisfy the following inequality:

$$k \geq \frac{(\ln(5\lambda_P^{8 \log(cD_0/\varepsilon)}/\varepsilon))^{2/\varepsilon}}{\zeta(\varepsilon)} = \frac{(\Theta(\log \log k \cdot \log \lambda_P + \log \lambda_P \cdot \ln(c/\varepsilon)))^{2/\varepsilon}}{\zeta(\varepsilon)}.$$

Hence, for  $k = (\log \lambda_P \cdot \log(c/\varepsilon))^{\Theta(1/\varepsilon)} / \zeta(\varepsilon)$ , we achieve a total probability of success in  $\Omega(\varepsilon)$ , which completes the proof.  $\square$

## 5. Dimension reduction based on randomly shifted grids

In this section, we explore some properties of randomly shifted grids, and we present a simplified embedding which consists of a first step of snapping points to a grid, and a second step of randomly projecting grid points.

Let  $w > 0$  and  $t$  be chosen uniformly at random from the interval  $[0, w]$ . The function

$$h_{w,t}(x) = w \cdot \left\lfloor \frac{x-t}{w} \right\rfloor$$

induces a random partition of the real line into segments of length  $w$ . Hence, the function

$$g_w(x) = (h_{w,t_1}(x_1), \dots, h_{w,t_d}(x_d)),$$

for  $t_1, \dots, t_d$  independent uniform random variables in the interval  $[0, w]$ , induces a randomly shifted grid in  $\mathbb{R}^d$ . For a set  $X \subseteq \mathbb{R}^d$ , we denote by  $g_w(X)$ , the image of  $X$  on the randomly shifted grid points defined by  $g_w$ . For some  $x \in \mathbb{R}^d$  and  $r > 0$ , the number of grid cells of  $g_w(\ell_1^d)$  that  $B_1(x, r)$  intersects per axis is independent, and in expectation is  $1+2r/w$ . Then, the expected total number of grid cells that  $B_1(x, r)$  intersects is at most  $(1+2r/w)^d$ .

Now let  $P \subset \ell_1^d$  be a set of  $n$  points with doubling constant  $\lambda_P$  and  $q \in \ell_1^d$  a query point. For  $w = \varepsilon/d$ , the  $\ell_1$ -diameter of each cell is  $\varepsilon$  and therefore  $g_w(P)$  is an  $\varepsilon$ -covering set of  $P$ .

**Lemma 11.** *Let  $\mathcal{R} > 1$  and  $P' := B_1(q, \mathcal{R}) \cap P$ . Then, for  $w = \varepsilon/d$*

$$\mathbb{E} [|g_w(P')|] \leq 8\lambda_P^{2\log(d\mathcal{R}/\varepsilon)}.$$

*Proof.* By the doubling constant definition, there exists a set of balls of radius  $\varepsilon/d^2$  centered at points in  $P'$ , of cardinality at most  $\lambda_P^{2\log(d\mathcal{R}/\varepsilon)}$  which covers  $P'$ . For each ball of radius  $\varepsilon/d^2$ , the expected number of intersecting grid cells is  $(1+2/d)^d \leq e^2$ . The lemma follows by linearity of expectation.  $\square$

The next lemma shows that, with constant probability, the growth on the number of representatives, as we move away from  $q$ , is bounded.

**Lemma 12.** *Let  $\{D_i\}_{i \in \mathbb{N}}$  be a sequence of radii such that, for any  $i$ ,  $D_{i+1} = 4D_i$ . Let  $A_i$  be the points of  $g_w(P)$  within distance  $D_{i+1} = 2^{2(i+1)}D_0$  from  $q$ . Then, with probability at least  $1/3$ ,*

$$\forall i \in \{-1, 0, \dots\} : |A_i| \leq 4^{i+3} \lambda_P^{2\log(dD_{i+1}/\varepsilon)}.$$

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

*Proof.* By Lemma 11,  $\mathbb{E}[|A_i|] \leq 8\lambda_P^{2\log(dD_{i+1}/\varepsilon)}$  for every  $i \in \{-1, 0, \dots\}$ . Then, a union bound followed by Markov's inequality yields

$$\Pr [\exists i \in \{0, 1, \dots\} : |A_i| \geq 4^{i+1} \mathbb{E}[|A_i|]] \leq 1/3.$$

In addition,

$$\Pr [|A_{-1}| \geq 4 \mathbb{E}[|A_i|]] \leq 1/4. \quad \square$$

**Theorem 13.** *Let  $P \subset \ell_1^d$  such that  $|P| = n$ . For any  $\varepsilon \in (0, 1/2)$ , there is a non-linear randomized embedding  $h' : \ell_1^d \rightarrow \ell_1^k$ , where*

$$k = (\log \lambda_P \cdot \log(d/\varepsilon))^{\Theta(1/\varepsilon)} / \zeta(\varepsilon),$$

for a function  $\zeta(\varepsilon) > 0$  depending only on  $\varepsilon$ , such that for any  $q \in \ell_1^d$ , if there exists  $p^* \in P$  such that  $\|p^* - q\|_1 \leq 1$ , then with probability  $\Omega(\varepsilon)$ ,

$$\|h'(p^*) - f(q)\|_1 \leq 1 + 3\varepsilon,$$

$$\forall p \in P : \|p - q\|_1 > 1 + 9\varepsilon \implies \|h'(p) - f(q)\|_1 > 1 + 3\varepsilon.$$

Any point can be embedded in time  $O(dk)$ .

*Proof.* We follow the same reasoning as in the proof of Theorem 10. The embedding is  $h' = f \circ g_{\varepsilon/d}$ , where  $f$  is the randomized linear map defined in Section 4. As before, we apply  $h'$  to every point in  $P$ , and only  $f$  to queries. The randomly shifted grid incurs an additive error of  $\varepsilon$  in the distances between  $q$  and  $P$ .

Assume wlog that  $q = 0^d$  and let  $A_i$  be the points of  $g_{\varepsilon/d}(P)$  within distance  $D_{i+1} = 2^{2(i+1)}D_0$  from  $q$ . Hence, by Lemma 12,

$$\begin{aligned} \Pr \left[ \exists i \exists s \in A_i : \|f(s)\|_1 \leq \frac{4\|s\|_1}{D_i} \right] &\leq \sum_{i=0}^{\infty} |A_i| \Pr \left[ S \leq \frac{4T}{D_i} \right] \\ &\leq \sum_{i=0}^{\infty} 4^{i+3} \lambda_P^{2\log(dD_{i+1}/\varepsilon)} \Pr \left[ \tilde{S} \leq \frac{\sqrt{4kT}}{\sqrt{D_i}} \right]. \end{aligned}$$

As in Lemma 9, for  $D_0 = \lceil 800T/k \rceil = \Theta(\log(k/\varepsilon))$ ,  $k \geq 20 \log \lambda_P \cdot \log(\frac{dD_0}{\varepsilon\delta})$  and  $\delta = \varepsilon/5$ ,

$$\begin{aligned} \sum_{i=0}^{\infty} 4^{i+3} \lambda_P^{2\log(dD_{i+1}/\varepsilon)} \Pr \left[ \tilde{S} \leq \frac{\sqrt{4kT}}{\sqrt{D_i}} \right] &\leq \sum_{i=0}^{\infty} \frac{2^{2i+6+2\log \lambda_P [\log(dD_0/\varepsilon)+2(i+1)]}}{2^{k(i+1)}} \\ &\leq \varepsilon/5. \end{aligned}$$

Hence, for  $k = \Omega((\log^2 \lambda_P \cdot \log(d/\varepsilon)))$ , with probability at least  $1 - \varepsilon/5$ , we have:

$$\forall p \in P : \|p - q\|_1 \geq D_0 + \varepsilon \implies \|h'(p) - f(q)\|_1 \geq 4.$$

Now, we are able to use Theorem 6 for points which are at distance at most  $D_0 + \varepsilon$  from  $q$ , and the near neighbor. By Lemma 12, with constant probability, the number of grid points at distance  $\leq D_0 + \varepsilon$ , is at most  $32 \cdot \lambda_P^{4 \log(dD_0/\varepsilon)}$ . Hence, by Theorem 6, for  $\gamma = \varepsilon/10$  and  $\delta = \varepsilon/(160 \lambda_P^{4 \log(dD_0/\varepsilon)})$ , with probability at least  $1 - \varepsilon/5$ , it holds:

$$\forall p \in P : \|p - q\|_1 \in (1 + 9\varepsilon, D_0 + \varepsilon) \implies \|h'(p) - f(q)\|_1 > 1 + 3\varepsilon.$$

Moreover, with probability at least  $\varepsilon/2$ , we obtain:

$$\|h'(p^*) - f(q)\|_1 \leq 1 + 3\varepsilon.$$

As in Theorem 10, the target dimension needs to satisfy the following:

$$k \geq \frac{(\ln(160 \lambda_P^{4 \log(dD_0/\varepsilon)}/\varepsilon))^{2/\varepsilon}}{\zeta(\varepsilon)}.$$

Hence, for  $k = (\log \lambda_P \cdot \log(d/\varepsilon))^{\Theta(1/\varepsilon)} / \zeta(\varepsilon)$  we achieve total probability of success  $\Omega(\varepsilon)$ .  $\square$

## 6. Conclusion

We have filled in a gap in the spectrum of randomized embeddings with bounded distortion only for distances between the query and a point set: such embeddings existed for  $\ell_2$  and  $\ell_1$  and for doubling subsets of  $\ell_2$ . Here we settle the case of doubling subsets of  $\ell_1$  with a *near* neighbor-preserving embedding. In the meantime, we obtain a concentration bound on sums of independent Cauchy variables. Our algorithms are quite simple, therefore they should also be of practical interest.

We rely on approximate  $r$ -nets or randomly shifted grids. For the former, Theorem 10 provides with a trade-off between the preprocessing time required and the target dimension. On the other hand, Theorem 13 has the advantage

1  
2  
3  
4  
5  
6  
7  
8  
9 of fast preprocessing: any point is embedded in  $O(dk)$  time, and the embedding  
10 is oblivious to the point set. In regards to the near-linear preprocessing time,  
11 the two results are comparable, since the dimension in Theorem 13 can be  
12 substituted by the target dimension of Theorem 6.  
13  
14

15 Let us underline that any potential improvements to Theorem 6 should lead  
16 to improvements to Theorems 10 and 13. The target dimension in these theo-  
17 rems follows from a direct application of Theorem 6 to the representative data  
18 points which lie inside a bounding ball centered at the query.  
19  
20  
21

### 22 *Acknowledgements*

23  
24 IZE is member of team AROMATH, joint between INRIA Sophia-Antipolis  
25 (France) and NKUA. IP thanks Robert Krauthgamer for useful discussions on  
26 the topic.  
27  
28  
29

### 30 **References**

- 31  
32  
33 [1] S. Har-Peled, P. Indyk, R. Motwani, Approximate nearest neighbor: To-  
34 wards removing the curse of dimensionality, *Theory of Computing* 8 (1)  
35 (2012) 321–350.  
36  
37  
38 [2] D. Achlioptas, Database-friendly random projections: Johnson-  
39 Lindenstrauss with binary coins, *J. Comput. Syst. Sci.* 66 (4) (2003)  
40 671–687.  
41  
42  
43 [3] N. Ailon, B. Chazelle, The fast Johnson-Lindenstrauss transform and ap-  
44 proximate nearest neighbors, *SIAM J. Comput.* 39 (1) (2009) 302–322.  
45  
46  
47 [4] E. Anagnostopoulos, I. Z. Emiris, I. Psarros, Randomized embeddings with  
48 slack and high-dimensional approximate nearest neighbor, *ACM Trans. Al-*  
49 *gorithms* 14 (2) (2018) 18:1–18:21.  
50  
51  
52 [5] P. Indyk, Stable distributions, pseudorandom generators, embeddings, and  
53 data stream computation, *J. ACM* 53 (3) (2006) 307–323.  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

- 1  
2  
3  
4  
5  
6  
7  
8  
9 [6] A. Andoni, P. Indyk, Near-optimal hashing algorithms for approximate  
10 nearest neighbor in high dimensions, *Commun. ACM* 51 (1) (2008) 117–  
11 122.  
12  
13  
14 [7] A. Andoni, T. Laarhoven, I. P. Razenshteyn, E. Waingarten, Optimal  
15 hashing-based time-space trade-offs for approximate near neighbors, in:  
16 Proc. ACM-SIAM Symposium on Discrete Algorithms, SODA, Barcelona,  
17 Spain, 2017, pp. 47–66.  
18  
19 [8] P. Indyk, R. Motwani, Approximate nearest neighbors: Towards removing  
20 the curse of dimensionality, in: Proc. ACM Symp. Theory of Computing,  
21 1998, pp. 604–613.  
22  
23 [9] P. Indyk, On approximate nearest neighbors under  $L_\infty$  norm, *J. Comput.*  
24 *Syst. Sci.* 63 (4) (2001) 627–638.  
25  
26 [10] A. Andoni, H. L. Nguyen, A. Nikolov, I. P. Razenshteyn, E. Waingarten,  
27 Approximate near neighbors for general symmetric norms, in: Proc. ACM  
28 Symposium on Theory of Computing, STOC, Montreal, Canada, 2017, pp.  
29 902–913.  
30  
31 [11] R. Cole, L. A. Gottlieb, Searching dynamic point sets in spaces with  
32 bounded doubling dimension, in: Proc. ACM Symp. Theory of Computing,  
33 ACM, New York, USA, 2006, pp. 574–583.  
34  
35 [12] S. Har-Peled, M. Mendel, Fast construction of nets in low dimensional  
36 metrics, and their applications, in: Proc. Symp. Computational Geometry,  
37 2005, pp. 150–158.  
38  
39 [13] P. Indyk, A. Naor, Nearest-neighbor-preserving embeddings, *ACM*  
40 *Trans. Algorithms* 3 (3) (2007) 31.  
41  
42 [14] Y. Bartal, L. A. Gottlieb, Approximate nearest neighbor search for  $\ell_p$ -  
43 spaces ( $2 < p < \infty$ ) via embeddings, in: Proc. LATIN: Theoretical Informatics -  
44 13th Latin American Symp., Buenos Aires, Argentina, 2018, pp.  
45 120–133.  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

- 1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65
- [15] J. Lee, M. Mendel, A. Naor, Metric structures in  $L_1$ : dimension, snowflakes, and average distortion, *Europ. J. Combin.* 26 (8) (2005) 1180–1190.
  - [16] Y. Bartal, L. A. Gottlieb, Dimension reduction techniques for  $\ell_p$ , ( $1 < p < 2$ ), with applications, in: 32nd International Symposium on Computational Geometry, SoCG 2016, June 14-18, 2016, Boston, MA, USA, 2016, pp. 16:1–16:15.
  - [17] D. Eppstein, S. Har-Peled, A. Sidiropoulos, Approximate greedy clustering and distance selection for graph metrics, *CoRR* abs/1507.01555 (2015).
  - [18] R. Krauthgamer, J. R. Lee, Navigating nets: Simple algorithms for proximity search, in: *Proc. 15th Annual ACM-SIAM Symp. Discrete Algorithms, SODA'04, 2004*, pp. 798–807.
  - [19] I. Emiris, V. Margonis, I. Psarros, Near-neighbor preserving dimension reduction for doubling subsets of  $\ell_1$ , in: *Approximation, Randomization, & Combinatorial Optimization: Algorithms & Techniques, APPROX/RANDOM 2019, September 2019, MIT, USA, 2019*, pp. 47:1–47:13. doi:10.4230/LIPIcs.APPROX-RANDOM.2019.47.