# Embedding Dimension of Contrastive Learning and k-Nearest Neighbors

## Dmitrii Avdiukhin<sup>1</sup>, Vaggos Chatziafratis, Orr Fischer<sup>3</sup>, Grigory Yaroslavtsev<sup>4</sup>



<sup>1</sup>Northwestern University, <sup>2</sup>UC Santa Cruz, <sup>3</sup>Weizmann Institute of Science, <sup>4</sup>George Mason University

### **OUR RESULT**

### **CONTRASTIVE LEARNING**

Consider a set of *m* non-contradictory constraints. There exists embedding into  $\ell_p$  space satisfying all constraints if d is chosen as follows.  $d \in \Theta(\sqrt{m})$  for p = 2•  $O(\sqrt{m})$  always suffices, and some set of constraints require  $\Omega(\sqrt{m})$ Neural Networks ★  $d \in O(m) \cap \Omega(\sqrt{m})$  for positive integer *p k*-NN For a positive integer p, there exists embedding into  $\ell_p$  space of dimension  $d = poly(k) \cdot polylog(n)$  preserving the k-nearest neighbors of each point. No polynomial dependence on n • Very surprising since k-NN information encodes  $\Theta(n^2)$  constraints  $\clubsuit$  No dependence on pMAIN TECHNIQUES Intuition For  $m = \Theta(n^2)$ , have  $d = \Theta(n)$  [Chatziafratis, Indyk`23] • Gives the  $\Omega(\sqrt{m})$  lower bound If there are O(n) constraints: • If the constraints are spread over all n points,  $d = \Theta(1)$ • If these constraints are concentrated on  $\sqrt{n}$  points,  $d = \Theta(\sqrt{n})$ . The density seems to matter **Constraint graph** Create edges for all constraints (contrastive learning) or points (k-NN) x with its nearest neighbors  $y_1, \dots, y_k$  $(x, y^+, z^-)$  $y_{1}$  $y_1$  $y_k$ ν Z... Arboricity **Definition**: the minimum number of forests to cover the graph Measure of graph density K Æ High arboricity Low arboricity **Bounds on arboricity** • Clique has arboricity  $\sqrt{m/2}$ • The  $\sqrt{m/2}$  bound is tight









$$||f(x) - f(y)||^2 =$$

	C
	(