

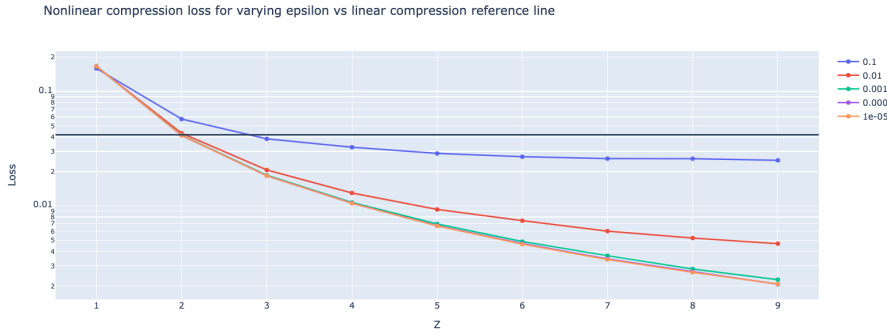
Regardless of whether large models end up using nonlinear compression, it should be possible to view directions being used with nonlinear compression as linear feature directions and reverse engineer the computation being used for compression like any other circuit. If this kind of encoding is pervasive throughout the network then it may merit some kind of automated decoding. It shouldn't pose a fundamental challenge to interpretability unless the model learns a scheme for doing complex computation while staying in a complicated nonlinear representation, which we suspect is unlikely.

To help provide intuition, the simplest example of what a nonlinear compression scheme might look like is compressing two  $[0,1]$  dimensions  $\mathbf{x}$  and  $\mathbf{y}$  into a single  $[0,1]$  dimension  $\mathbf{t}$ :

$$\mathbf{t} = \frac{\lfloor Z\mathbf{x} \rfloor + \mathbf{y}}{Z}$$

This works by quantizing the  $\mathbf{x}$  dimension using some integer  $Z$  such that the floating point precision of  $\mathbf{t}$  is split between  $\mathbf{x}$  and  $\mathbf{y}$ . This particular function needs the discontinuous floor function to compute, and the discontinuous fmod function to invert, but models can't compute discontinuous functions. However it's possible to replace the discontinuities with steep linear segments that are only some epsilon value wide.

We can compare the mean squared error loss on random uniform dense values of  $\mathbf{x}$  and  $\mathbf{y}$  and see that even with epsilons as large as 0.1 and  $Z$  values as small as 3 the nonlinear compression outperforms linear compression such as picking one of the dimensions or using the average:



## Connection between compressed sensing lower bounds and the toy model

Here, we formalize the relationship between a compressed sensing lower bound and the toy model.

Let  $T(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}^n$  be the complete toy model autoencoder defined by  $T(\mathbf{x}) = \text{ReLU}(W_2 W_1 \mathbf{x} - \mathbf{b})$  for an  $m \times n$  matrix  $W_1$  and an  $n \times m$  matrix  $W_2$ .

We derive the following theorem:

**Theorem 1.** Suppose that the toy model recovers all  $\mathbf{x}$  with  $T(\mathbf{x})$  such that  $\|T(\mathbf{x}) - \mathbf{x}\|_2 \leq \epsilon$  for sufficiently small  $\epsilon$  and  $W_1$  has the  $(\delta, k)$  restricted isometry property. The inner dimension of the projection matrix  $W$  is  $m = \Omega(k \log(n/k))$ .

We prove this result by framing our toy model as a compressed sensing algorithm. The primary barrier to doing so is that our optimization only searches for vectors that are close in  $\ell_2$  distance to the original vector and may not itself be exactly  $k$ -sparse. The following lemma resolves this concern through a denoising step:

**Lemma 1.** Suppose that we have a toy model  $T(\mathbf{x})$  with the properties in Theorem 1. Then there exists a compressed sensing algorithm  $\mathbf{f}(\mathbf{y}) : \mathbb{R}^m \rightarrow \mathbb{R}^n$  for the measurement matrix  $W_1$ .

*Proof.* We construct  $\mathbf{f}(\mathbf{y})$  as follows. First, compute  $\tilde{\mathbf{x}} = \text{ReLU}(W_2 \mathbf{y} - \mathbf{b})$ , as in  $T(\mathbf{x})$ . This produces the vector  $\tilde{\mathbf{x}} = T(\mathbf{x})$  and so by supposition  $\|T(\mathbf{x}) - \mathbf{x}\|_2 \leq \epsilon$ . Next, we threshold  $\tilde{\mathbf{x}}$  to obtain  $\tilde{\mathbf{x}}'$  by dropping all but its  $k$  largest entries. Lastly, we solve the optimization problem:  $\min_{\mathbf{x}'} \|\mathbf{x}' - \tilde{\mathbf{x}}'\|$  subject to  $W_1 \mathbf{x}' = \mathbf{y}$ , which is convex because  $\mathbf{x}'$  and  $\tilde{\mathbf{x}}'$  have the same support. For sufficiently small  $\epsilon$  (specifically,  $\epsilon$  smaller than the  $(k+1)$ th largest entry in  $\mathbf{x}$ ), both  $\tilde{\mathbf{x}}$  and the nearest  $k$ -sparse vector to  $\mathbf{x}$  have the same support, and so the the convex optimization problem has a unique solution: the nearest  $k$  sparse vector to  $\mathbf{x}$ . Therefore,  $\mathbf{f}$  is a compressed sensing algorithm for  $W_1$  with approximation factor 1. .

Lastly, we use the deterministic compressed sensing lower bound of Do Ba, Indyk, Price, and Woodruff [49]:

**Theorem 2 (Corollary 3.1 in [49]).** Given a  $k \times n$  matrix  $A$  with the restricted isometry property, a sparse recovery algorithm find a  $k$ -sparse approximation  $\hat{x}$  of  $x \in \mathbb{R}^n$  from  $Ax$  such that

$$\|x - \hat{x}\|_1 \leq C(k) \min_{x', \|x'\|_0 \leq k} \|x - x'\|_1$$

for an approximation factor  $C(k)$ . If  $C(k) = O(1)$ , then a sparse recovery algorithm exists only if  $m = \Omega(k \log(n/k))$ .

Theorem 1 follows directly from Lemma 1 and Theorem 2.