

---

# Embedding game: dimensionality reduction as a two player zero-sum game

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1            Dimensionality reduction is often formulated as a minimization containing a sparse  
2            sum of attractive interactions and a dense sum of repulsive interactions  $\sum_{ij} f(\|\mathbf{y}_i -$   
3             $\mathbf{y}_j\|)$  between embedding vectors. This dense sum is usually subsampled to avoid  
4            computing all  $N^2$  terms. In this paper we provide a novel approximation to the  
5            repulsive sum by deriving a landmark-based lower bound and then maximizing this  
6            lower bound with respect to the landmarks. After inserting this approximation into  
7            the original objective we are left with a minimax problem where the embedding  
8            vectors minimize the objective by pulling on their neighbors and running away  
9            from the landmarks while the landmarks maximize the objective by pulling on  
10           the embedding vectors and running away from other nearby landmarks. We use  
11           gradient descent ascent to find saddle points and show that our method can produce  
12           high quality visualizations without ever explicitly computing any pairwise repulsion  
13           between embedding vectors.

## 14    1 Introduction

15           Dimensionality reduction algorithms can be useful in a wide variety of contexts. Reducing the dimen-  
16           sionality of vectors can reduce the computational burden on downstream tasks such as recognition  
17           or neighborhood searches. Reducing the dimensionality of inputs can also be a method to filter out  
18           unwanted variability in the original inputs. In the extreme case of reduction to 2 or 3 dimensions, it  
19           can be used to produce visualization of the input [14]. This is the case we will be concerned with in  
20           this paper.

21           Recent algorithms have yielded very impressive looking visualizations of complicated datasets  
22           [13, 8, 12, 11, 1]. Common to each of these algorithms is an objective function containing non-  
23           linear interactions between all (or nearly all) pairs of embedding vectors. There are a variety of  
24           approximations that have been proposed to approximate these all-pairs interactions. t-distributed  
25           Stochastic Neighbor Embedding (T-SNE) has taken inspiration from physical simulation and used  
26           the Barnes-Hut algorithm to cleverly discretize embedding space in a manner that allows for efficient  
27           approximation of all-pairs nonlinear interactions [13]. LargeVis and UMAP both use weighted edge  
28           sampling. At each iteration (or after several iterations), a random subset of interactions are chosen  
29           (with higher weighted interactions more likely to be chosen) and the subsampled objective is instead  
30           optimized [8, 12]. PyMDE uses a similar idea, but instead samples negative edges uniformly at  
31           random and these edges are fixed throughout training [1]. The tractable Latent Variable Model used  
32           landmarks to approximate the repulsion and relied on a sophisticated coarse graining scheme to  
33           reduce the number of pairwise interactions [11].

34           In this paper we will use a landmark approach for approximating nonlinear all-pairs interactions. We  
35           will derive a lower bound to a sum of all-pairs interactions which we then maximize with respect to  
36           the landmarks. We will show that this requires much fewer landmarks than if we simply randomly  
37           sampled some fixed set of landmarks. Our algorithm can be interpreted as a two player game where

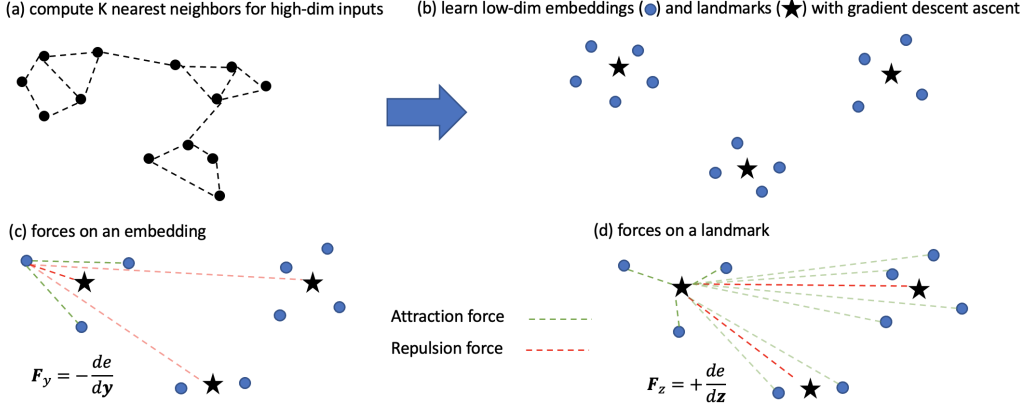


Figure 1: Diagram of our method. Each embedding is attracted to the neighbors we compute in step a and repelled by all landmarks. There is no direct embedding-embedding repulsion. Landmarks are attracted to all embedding vectors and repelled from all other embedding vectors. These forces are derived in Eqs. 4, 5 as the gradients of the objective in Eq. 3.

38 the embedding vectors minimize the objective by pulling on their neighbors and running away from  
 39 the landmarks while the landmarks maximize the objective by pulling on the embedding vectors  
 40 and running away from other nearby landmarks. We use gradient descent ascent and show that our  
 41 method can produce high quality visualizations.

## 42 2 Embedding Game

43 Our starting point closely follows [1] on Minimum Distortion embedding vectors (MDE). We  
 44 assume we have a set of high-dimensional inputs  $\{\mathbf{x}_i\}_{i=1}^N$  and we wish to produce a corre-  
 45 sponding set of low-dimensional *embedding vectors*  $\{\mathbf{y}_i\}_{i=1}^N$  which reveal interesting structure  
 46 in the original inputs. To accomplish this we compute a sparse set of neighborhood edges:  
 47  $\mathcal{N} = \{(i, j) : j \text{ is a } k\text{-nearest neighbor of } i\}$ . We define two functions:  $f(d)$  will penalize large  
 48 distances between embedding vectors and  $g(d)$  will penalize small distances between embedding  
 49 vectors. We wish to find embedding vectors  $\mathbf{y}$  which minimize:

$$\min_{\mathbf{y}} \sum_{i,j \sim \mathcal{N}} f(\|\mathbf{y}_i - \mathbf{y}_j\|) + \sum_{i,j} g(\|\mathbf{y}_i - \mathbf{y}_j\|) \quad (1)$$

50 Like many dimensionality reduction objectives, this contains a sum over  $N^2$  terms. One method  
 51 to avoid this unwieldy sum is simply to sample edges [8, 12, 1]. We proceed in a different fashion  
 52 by defining a set of *landmarks*  $\{\mathbf{z}_a\}_{a=1}^L$  and making a landmark-based approximation to the sum.  
 53 Perhaps the simplest landmark-based approximation is simply to set the landmarks to be randomly  
 54 chosen samples  $\mathbf{z}_1 = \mathbf{y}_{i_1}, \mathbf{z}_2 = \mathbf{y}_{i_2}, \dots$  and approximate the all-pairs sum with  $\frac{N}{L} \sum_{ia} g(\|\mathbf{y}_i - \mathbf{z}_a\|)$ .  
 55 This can actually be interpreted as another method of sampling edges. However we show in the  
 56 Experiments section that this idea has several problems when using a small number ( $< 300$ ) of  
 57 landmarks.

58 Fortunately we can make a more powerful landmark-based approximation if  $g(\|\mathbf{y}_i - \mathbf{y}_j\|)$  defines a  
 59 positive semi-definite kernel function, in other words if the  $N \times N$  matrix of pairwise interactions  
 60 defined by  $g$  is positive semi-definite for all  $\{\mathbf{y}_i\}$ . Some examples of  $g$  satisfying this property are  
 61  $g(d) = \exp(-d^2)$  and  $g(d) = 1/(1 + d^2)$ . Unfortunately our bound will not work with  $g(d)$  that go  
 62 to infinity as  $d \rightarrow 0$  which are sometimes used in the literature. However we'll show that we can  
 63 learn high quality embedding vectors without this "infinite as  $d$  goes to zero" property.

64 **Key inequality:** if  $g(\|\mathbf{y}_i - \mathbf{y}_j\|)$  is a positive semi-definite kernel function, then for any  $\{\mathbf{y}_i\}_{i=1}^N$  and  
 65  $\{\mathbf{z}_a\}_{a=1}^L$ :

$$\left[ \sum_{i,a} g(\|\mathbf{y}_i - \mathbf{z}_a\|) \right]^2 \leq \left[ \sum_{ij} g(\|\mathbf{y}_i - \mathbf{y}_j\|) \right] \left[ \sum_{a,b} g(\|\mathbf{z}_a - \mathbf{z}_b\|) \right] \quad (2)$$

66 We are unaware of this inequality being presented in the literature but we prove it in the Appendix.  
 67 The key to our proof is to replace  $g(\|\mathbf{y}_i - \mathbf{y}_j\|)$  with the inner product between high dimensional  
 68 vectors  $\phi_i \cdot \phi_j$ . This is allowed because of our assumption that  $g$  defines a kernel function. A useful  
 69 property of our inequality is that with at most  $N$  landmarks, there exist  $\mathbf{z}$  which yield *equality*, rather  
 70 than inequality. This can be seen by simply setting the  $\mathbf{z}_a = \mathbf{y}_i$  when  $L = N$ . So long as we  
 71 have enough landmarks, our approximation should yield the same result as the original “all-pairs”  
 72 repulsion objective. Of course we hope that our approximation is useful for  $L \ll N$ .

73 To generate our landmark-based approximation to the all-pairs sum in Eq. 1 we divide both sides of  
 74 Eq. 2 by  $\sum_{ab} g_{ab}$  and then maximize  $\sum_{ia}^2 / \sum_{ab}$  with respect to  $\mathbf{z}$ . This yields the tightest lower  
 75 bound to the sum  $\sum_{ij} g(\|\mathbf{y}_i - \mathbf{y}_j\|)$ . We then replace the pairwise sum in Eq. 1 with our tightest  
 76 lower bound to yield the minimax problem we ultimately try to solve:

$$\min_{\mathbf{y}} \max_{\mathbf{z}} \sum_{i,j \sim \mathcal{N}} f(\|\mathbf{y}_i - \mathbf{y}_j\|) + \frac{\left[ \sum_{i,a} g(\|\mathbf{y}_i - \mathbf{z}_a\|) \right]^2}{\sum_{a,b} g(\|\mathbf{z}_a - \mathbf{z}_b\|)} \quad (3)$$

77 The denominator in the right-hand term contains  $LN$  interactions (between all landmark-embedding  
 78 vector pairs) and the numerator contains  $L^2$  interactions (between all landmark-landmark pairs).

## 79 2.1 Gradient descent ascent (GDA)-based optimization

```
# x: input vectors (shape: (n,m))
# k,l: num neighbors, num landmarks
# f,g: attractive, repulsive penalty functions

edges = knn_edges(x,k) # k nearest neighbors for each input
y = laplacian_eigenmap(edges) # init embedding vectors
z = sample_landmarks(y,l) # init landmarks as randomly chosen embedding
  vectors

for i in range(n_iter):
  # pairwise distances
  yy = norm((y[edges[0]] - y[edges[1]]),dim=1) # shape: (n*k)
  yz = cdist(y,z) # shape: (n,l)
  zz = cdist(z,z) # shape: (l,l)

  # energy
  e = f(yy).sum() + (g(yz)**2).sum() / g(zz).sum() # Eq. 3

  # gradients
  e.backward()

  # updates
  y -= eta_y / (y.grad**2).mean().sqrt() * y.grad
  z += eta_z / (z.grad**2).mean().sqrt() * z.grad
```

**Algorithm 1:** PyTorch-style pseudocode for embedding game

80 We provide PyTorch-style pseudocode in Alg. 1. We use a rescaled gradient descent-ascent algorithm  
 81 to find saddle points of Eq. 3. This rescaling is helpful as the embedding gradients are much smaller  
 82 than the landmark gradients and the rescaling lets us set  $\eta_y, \eta_z$  to be similar magnitudes. In principle  
 83 this rescaling could lead to convergence issues, however this was not problem in practice.

84 As is standard practice [8, 1, 11], we initialize the embedding vectors using Laplacian eigenmaps [3].  
 85 This method sets  $\mathbf{y}$  to be the 2nd and 3rd smallest eigenvectors of the normalized graph Laplacian  
 86 defined by the k-nearest neighbor graph. The  $\mathbf{y}$  are rescaled so each dimension is unit variance.  
 87 In practice we don’t find exact eigenvectors but rather 100 power iterations to approximate these  
 88 eigenvectors. We initialize the landmarks by randomly sampling embedding vectors after laplacian  
 89 initialization.

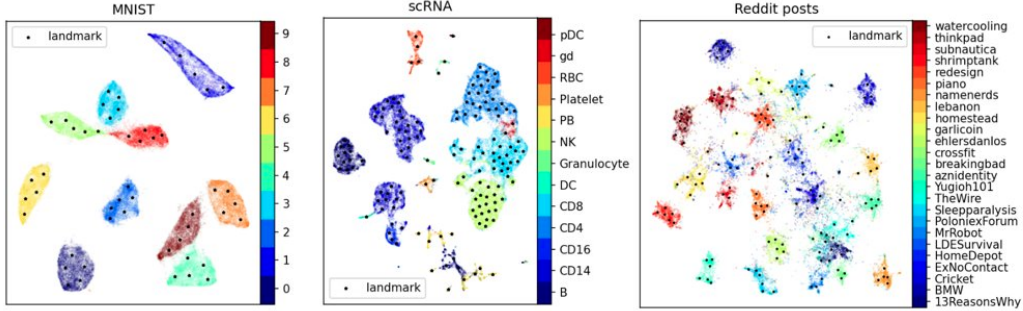


Figure 2: embedding vectors (color) and landmarks (black) generated by our algorithm on three datasets.

90 Both in theory and in practice, choosing the learning rates for GDA problems can be much more  
 91 complicated than for simple gradient descent problems. In the experiments section we explore how  
 92 various learning rate choices impact convergence.

## 93 2.2 Force-based interpretation of the game

94 It is useful to consider the forces on each player of the game. The force on each embedding takes the  
 95 form:

$$\mathbf{f}_i = -\frac{dE}{d\mathbf{y}_i} = \alpha_1 \sum_{j \sim \mathcal{N}_i} f'_{ij} \frac{\mathbf{y}_j - \mathbf{y}_i}{\|\mathbf{y}_j - \mathbf{y}_i\|} + \alpha_2 \sum_a g'_{ia} \frac{\mathbf{y}_i - \mathbf{z}_a}{\|\mathbf{y}_i - \mathbf{z}_a\|} \quad (4)$$

96 where  $\alpha_1, \alpha_2$  are non-negative constants. The force has two sources: each embedding feels attraction  
 97 to its neighbors and repulsion from all landmarks. In practice,  $g'$  will decay with distance so the  
 98 repulsion is strongest from nearby landmarks. In simple terms *the embedding vectors run towards*  
 99 *their neighbors and away from the landmarks.*

100 Because the landmarks are maximizing the objective, the force is now the positive gradient. This takes  
 101 the form:

$$\mathbf{f}_a = +\frac{dE}{d\mathbf{z}_a} = \beta_1 \sum_i g'_{ia} \frac{\mathbf{y}_i - \mathbf{z}_a}{\|\mathbf{y}_i - \mathbf{z}_a\|} + \beta_2 \sum_b \frac{\mathbf{z}_a - \mathbf{z}_b}{\|\mathbf{z}_a - \mathbf{z}_b\|} \quad (5)$$

102 where  $\beta, \beta_2$  are non-negative constants. This force also has two source. Landmarks are attracted to  
 103 the embedding vectors (but prefer close embedding vectors) and repelled by other landmarks. In  
 104 simple terms *the landmarks run towards the closest embedding vectors and away from the other*  
 105 *nearby landmarks.*

## 106 3 Visualization results

107 We show visualizations produced by our algorithm on three different datasets (i) the classic MNIST  
 108 dataset (ii) a single cell RNA dataset (iii) 25k reddit posts from 25 different subreddits. The resulting  
 109 embedding vectors and landmarks are shown in Fig. 2. The penalty functions we use are:

$$f(d) = \log(1 + d^2) \quad g(d) = \frac{\lambda}{1 + d^2} \quad (6)$$

110  $\lambda$  is a hyperparameter which we tune for each dataset. The “log 1-plus” attractive penalty has seen  
 111 successes in previous works so we stick with it [8, 1]. Intuitively this function penalizes large  
 112 distances less than a more intuitive  $d^2$  penalty, which may be important in the presence of noisy  
 113 neighborhood graphs. Our “cauchy” repulsive penalty  $g$  is a more unusual choice. Unlike other  
 114 works,  $g$  does not approach infinity as the distance goes to zero. This is important as our bound  
 115 does not work when  $g(0) \rightarrow \infty$ . However, this penalty still decays to zero as  $d \rightarrow \infty$ , meaning that  
 116 repulsion is strongest between nearby vectors. We show empirically that we can generate high quality  
 117 visualizations with this class of distortion function.

118 For every dataset we adopt the same learning rate schedule. We perform 3k GDA iterations with  
 119  $\eta_y = 0.03, \eta_z = 0.3$ . We divide both learning rates by 10 and perform 3k more iterations. There is

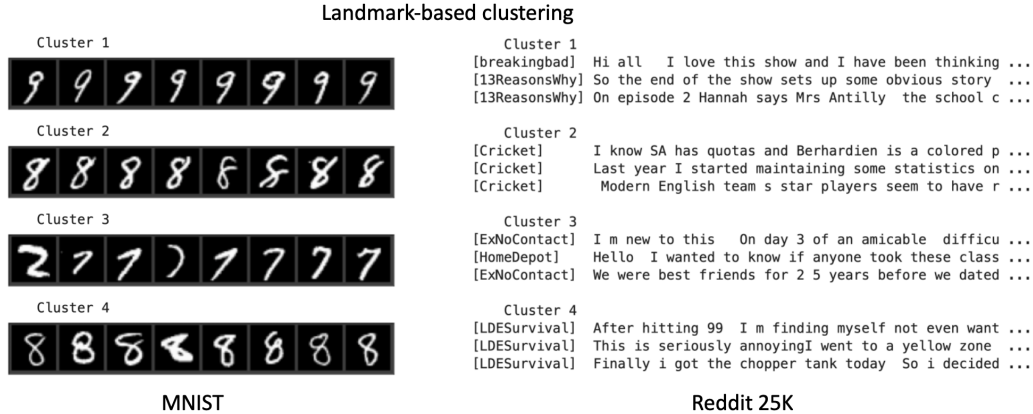


Figure 3: Clustering derived by assigning each embedding to the nearest landmark. We show random subsets of each cluster for the first four clusters in the MNIST dataset and the Reddit25K dataset

120 no stochasticity in our system as we are doing full batch updates, however the learning rate decay  
 121 still seems helpful for improving the resulting visualization.

### 122 3.1 Datasets

123 **MNIST** The classic MNIST dataset contains 70,000 grayscale images of size 28x28, each containing  
 124 a single handwritten digit. The generated 2D visualization is shown for  $K = 15$  neighbors,  $L = 50$   
 125 landmarks and  $\lambda = 0.001$ .

126 **scRNA** The scRNA dataset contains 40,000 PCA embedding vectors of single cell mRNA transcrip-  
 127 tomes from COVID-19 patients. The input vectors are 30 dimensional. This dataset originates from  
 128 [15] and can be conveniently be download from [1]. We use  $K = 15$ ,  $L = 150$  and  $\lambda = 0.01$ .

129 **25K Reddit Posts** This dataset contains 25,000 reddit posts from 25 different subreddits (1000 posts  
 130 per subreddit). This is generated by sampling 25 subreddits from the larger 1 million post dataset  
 131 which can be downloaded from Kaggle [6]. To generate feature vector for each post, we use the  
 132 strategy detailed in [2] and create a weighted average of GloVe vectors used in the post. These 300  
 133 dimensional “post vectors” are then fed into our algorithm. We use  $K = 15$ ,  $L = 150$  and  $\lambda = 0.01$ .

### 134 3.2 Landmark-based clustering

135 In each case the landmarks appear to tile the high density regions of embedding space. This suggests  
 136 a method to cluster the data. Assign each embedding to the closest landmark in embedding space:

$$c_i = \underset{a}{\operatorname{argmin}} \|y_i - z_a\| \quad (7)$$

137 One might expect similar results by running KMeans on the embedding vectors, but a difference is  
 138 that we already have the landmarks (cluster centers) as a result of running our algorithm. In Fig. 3  
 139 we show randomly chosen inputs that are assigned to selected landmarks for the MNIST and Reddit  
 140 post datasets. This a can provide an interesting way to quickly visualize data. Instead of directly  
 141 visualizing the embedding vectors, one can examine clusters derived from the learned landmarks.

## 142 4 Duality and learning rates

143 Choosing the learning rates for this problem is non-trivial because we have a non-convex non-concave  
 144 minimax optimization. However we provide a rule of thumb which is to choose the landmark learning  
 145 rate sufficiently large relative to the embedding learning rate. This is motivated by experiment and  
 146 extrapolation from theoretical results on nonconvex-concave optimization.

147 **Duality** In general there is a duality gap for our optimization in Eq. 3. In other words the order of the  
 148 optimization (min-max vs. max-min) is extremely important for this problem. In fact the reversed  
 149 “max-min” problem admits a completely degenerate set of landmarks. For any fixed set of landmarks

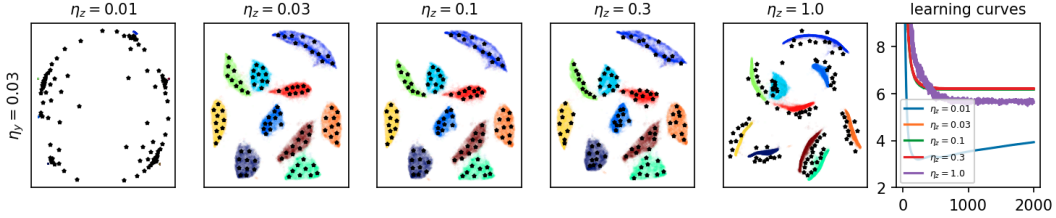


Figure 4: Varying the landmark learning rate  $\eta_z$  with the embedding learning rate fixed at  $\eta_y = 0.03$  (for the MNIST dataset). Too small, and the embedding vectors collapse to a few points ( $\eta_z = 0.01$ ). Too large and the embedding vectors again appear to begin to collapse, although not as extreme here. There appears to be a window of learning rates where we observe nice clustering.

150 ( $\mathbf{z}$ ), the embedding vectors ( $\mathbf{y}$ ) can achieve 0 error by: one, setting the first sum in Eq. 3 to zero  
 151 by collapsing to a point  $\mathbf{y}_i = \mathbf{y}_j$  for all  $i, j$ , and two: setting the numerator of the second term in  
 152 Eq. 3 to zero by running off to infinity. So for any fixed  $\mathbf{z}$ , we have  $\min_{\mathbf{y}} e(\cdot, \mathbf{z}) = 0$  and therefore  
 153  $\max_{\mathbf{z}} \min_{\mathbf{y}} e = 0$ . Ultimately we care about generating useful visualizations of data, so having all  
 154 our embedding vectors collapsed at a single point, infinitely far away from the landmarks, is not good.

155 **Learning rates** Intuitively we should set the landmark learning rate to be relatively fast compared  
 156 to the embedding learning rate. This way, the landmarks can approximately perform  $\max_{\mathbf{z}} e(\mathbf{y}, \cdot)$   
 157 before the embedding vectors have a chance to update appreciably. Then the embedding vectors can  
 158 perform approximate gradient descent on the objective  $\max_{\mathbf{z}} e(\cdot, \mathbf{z})$  and the algorithm is more likely  
 159 to find solutions to the original min max problem.

160 Rigorously justifying this intuition, that setting  $\eta_y \ll \eta_z$  will find a solution of  $\min_{\mathbf{y}} \max_{\mathbf{z}}$ , is  
 161 rather challenging in the case of nonconvex-nonconcave objective like ours. In the simpler case of a  
 162 nonconvex-concave problem, this intuition (choose a fast learning rate for the inside-maximization)  
 163 can be rigorously shown to be correct [7]. We’ll show via experiment that this intuition seems to be  
 164 useful.

165 **Experiment** We run our algorithm on the MNIST dataset with  $k = 15$ ,  $l = 100$ , and  $\lambda = 0.001$ .  
 166 We’ll fix  $\eta_y = 0.03$  and vary  $\eta_z$ . Learning curves and visualizations after 2k iterations of Alg. 1 are  
 167 shown in Fig. 4. Too small, and the embedding vectors collapse to a few points ( $\eta_z = 0.01$ ). Too  
 168 large and the embedding vectors again appear to begin to collapse, although not as extreme here.  
 169 There appears to be a window of learning rates where we observe nice clustering. This seems to be  
 170 explained by our intuitions on the duality problem.

171 When  $\eta_z$  is too small, the embedding vectors can begin to minimize first and perform  $\min_{\mathbf{y}} e(\cdot, \mathbf{z})$ ,  
 172 which we argued can be 0 when the embedding vectors collapse to a point and move away from  
 173 the landmarks. When  $\eta_z$  is too large, then again the landmarks are not performing the inner loop  
 174 maximization, and the embedding vectors again can perform  $\min_{\mathbf{y}} e(\cdot, \mathbf{z})$ , by collapsing to a point.  
 175 In Fig. 4 we see the embedding vectors begin to collapse, although its not as extreme as when  $\eta_z$  was  
 176 too large.

## 177 5 Comparison with sample-based techniques

178 We’ll compare to two fixed-sample-based methods for approximating the all-pairs sum in Eq. 1. The  
 179 first method is used by [1] and simply samples  $L$  edges for each node. The second method randomly  
 180 designates  $L$  embedding vectors at the start of training to be the landmarks and replaces the sum using  
 181 the landmarks. This method is not widely used in practice, and as we’ll see it leads to surprisingly  
 182 low quality visualizations. We call these methods “fixed-sample” because these edges/landmarks are  
 183 sampled once at the start of training and then fixed for all subsequent iterations.

184 We evaluate each method a) qualitatively by looking at the resulting visualizations and b) quantitatively  
 185 by comparing the value of the “all-pairs” objective (Eq. 1) using the generated embedding vectors.  
 186 We observe that the sampled-edge-based algorithm seems to outperform our algorithm for fixed  $L$ ,  
 187 while the sample-landmark-based algorithm dramatically underperforms.

188 **Sampled edges** At the start of training, we sample  $L$  edges  $(i, j)$  uniformly at random for each node  
 189  $i$ . We call this set of edges  $\mathcal{N}^-$ . The all-pairs sum  $\sum_{ij}$  is then replaced with a sum over these edges  
 190 and reweighted by a factor  $N/L$ . So we optimize the objective:

$$\min_{\mathbf{y}} \sum_{i,j \sim \mathcal{N}} f(\|\mathbf{y}_i - \mathbf{y}_j\|) + \frac{N}{L} \sum_{i,j \sim \mathcal{N}^-} g(\|\mathbf{y}_i - \mathbf{y}_j\|) \quad (8)$$

191 **Sampled landmarks** At the start of training, we designate  $L$  embedding vectors uniformly at random  
 192 to be landmarks. We call this set of embedding vectors  $\mathcal{E}$ . We again replace the all-pairs sum in Eq. 1  
 193 with a reweighted sum over landmarks:

$$\min_{\mathbf{y}} \sum_{i,j \sim \mathcal{N}} f(\|\mathbf{y}_i - \mathbf{y}_j\|) + \frac{N}{L} \sum_i \sum_{j \sim \mathcal{E}} g(\|\mathbf{y}_i - \mathbf{y}_j\|) \quad (9)$$

194 By defining the set of negative edges  $\mathcal{N}^- = \{(i, j) : i = 1, 2, \dots, N \text{ and } j \in \mathcal{E}\}$ , we can rewrite Eq.  
 195 9 to look identical to Eq. 8. In other words the sampled landmark scheme can be interpreted as just  
 196 another method to sample negative edges.

197 Averaging over edge or landmark choices, both schemes yield unbiased estimators  $\langle \frac{N}{L} \sum_{ij \sim \mathcal{N}^-} \rangle =$   
 198  $\sum_{ij}$ . But because these are chosen once at the start of training and then fixed, these actually yield  
 199 zero variance but biased estimates. In other words both methods generates a biased estimate of  $\sum_{ij}$   
 200 and its gradients. If our method exactly computed the global maxima with respect to the landmarks  
 201 for a fixed set of embedding vectors, it could as well be regarded as a zero-variance biased estimate  
 202 of the all-pairs sum. However, it does not find the global maxima, so there is some variance in our  
 203 estimates (due to for instance random initialization of the landmarks).

204 **Experiment training details** In all three experiments (sampled edges, sampled landmarks, and  
 205 our method) we use the same hyperparameters. We train on the MNIST dataset. We use the  $f, g$   
 206 described in Eq. 6. We set  $\lambda = 0.001$ . We use  $K = 15$  for the neighborhood graph. We compare  
 207  $L \in \{1, 3, 10, 30, 100, 300\}$ . For the sampled edges and sampled landmarks experiment, we only  
 208 have a learning rate for the embedding vectors. We use the same learning rate schedule for  $\eta_y$  in  
 209 both settings: we perform 2k iterations with  $\eta_y = 0.03$ , then 2k more with  $\eta_y = 0.003$ . For the  
 210 experiment using our method, we have the additional parameter  $\eta_z$  and we set it at  $\eta_z = 10\eta_y$ .

## 211 5.1 Visualization (qualitative comparison)

212 The resulting visualizations are shown in Fig. 5. As the number of landmarks  $L$  increase beyond  
 213 30, both the optimized-landmark and sampled-edge algorithms yield nearly identical visualizations.  
 214 This seems reasonable, as  $L$  increases each algorithm yields a better and better estimate of the  
 215 all-pairs objective so beyond a certain threshold of  $L$ , all algorithms should ultimately yield the  
 216 same visualizations. For  $L = 1, 10, 30$  the optimized landmark algorithm suffers from degeneracies  
 217 not seen in the sampled-edge algorithm. In particular we observe a number of "clusters" appear to  
 218 collapse to single points.

219 Intuitively when there are more clusters than landmarks, there is no mechanism for the optimized  
 220 landmark algorithm to prevent some of the clusters from collapsing. In the landmark algorithms  
 221 there is no direct embedding-embedding repulsive forces. If there is no landmark inside a cluster  
 222 of embedding vectors which are attracting each other, there is no outward repulsion preventing this  
 223 cluster from collapsing to a point as we see in the cases  $L = 1, 3, 10$  where there are not very many  
 224 landmarks. The sampled-edge method seems to provide a more robust mechanism for avoiding this  
 225 embedding collapse.

226 The sampled-landmark algorithm gave nearly complete embedding collapse for all  $L$  we tested  
 227 (note that at some point, for sufficiently large  $L$  all three methods should yield the same result). All  
 228 embedding vectors which were not designated as landmarks simply collapsed towards a single point  
 229 while the landmark vectors were repelled to an exterior ring around the origin. It may seem surprising  
 230 how different this result is from the sampled edges experiment, given that it can be interpreted as  
 231 another way to generate negative edges. This result shows the importance of the exact method used  
 232 to approximate the all-pairs repulsive sum.

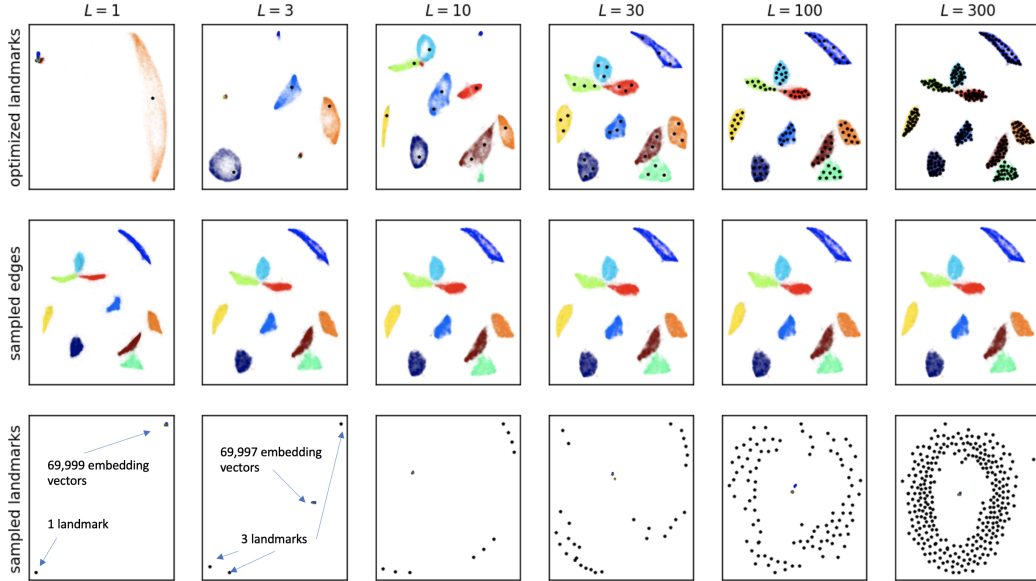


Figure 5: Comparing our optimized landmark method for approximating the all-pairs repulsion in Eq. 1 to simple edge-sampling and landmark-sampling methods. Black dots indicated landmarks (not relevant/present in sampled edges figures). The edge sampling method seems to outperform our method for a fixed  $L$ . When we don't have enough landmarks ( $L=1,3,10$ ) we see clusters of embedding collapse to points. However our method is much better than a simple random sampling of landmarks. When we randomly sample and fix landmarks, nearly all the embedding vectors which were not designated as landmarks collapse to a single point, while the landmark vectors repel from each other and other embedding vectors.

233 **5.2 Quantitative Results**

234 We also compare the sampled-edge and the optimized-landmark algorithms quantitatively in Fig. 6.  
 235 We don't show the sampled-landmark method as it is far worse than the sampled-edge or optimized-  
 236 landmark methods. All curves in this Figure are for  $L = 10$ .

237 In (a) we show the energy we are actually optimizing (Eq. 3 for optimized-landmarks and Eq. 8 for  
 238 sampled-edges). This indicates that both algorithms are at least optimizing the approximations, and  
 239 the degeneracies we observed in Fig. 5 are not a failure of the optimization routine. In (b) we show  
 240 an unbiased approximation to the "true" energy at each iteration (Eq. 1). This is done by randomly  
 241 sampling a large number of edges from the all pairs sum, and these edges are chosen i.i.d. at each  
 242 iteration

243 For both orange and blue, the energy we optimize is less than the all-pairs energy, and this difference  
 244 can be regarded as analogous to a generalization gap. Each algorithm "overfits" to the energy we  
 245 optimize, but this "overfitting" appears worse for our method than the sampled edge method.

246 We observe that the sampled-edge algorithms yields a lower all-pairs energy (the energy we truly  
 247 wish to optimize). This is in agreement with the fact that the sampled edge method yields the highest  
 248 quality visualizations with a small number of sampled edges (Fig. 5). In (c) we plot the root-mean-  
 249 squared error between the gradients of the approximate energy and gradients of the estimated true  
 250 energy. Our method produces a more faithful gradient estimate at nearly all times in training. This is  
 251 extremely surprising as the final true energy, after many gradient updates is lower for the sampled  
 252 edge method which gives a worse gradient estimate in terms of mean squared error. This result shows  
 253 how the exact details used to approximate the sum can be very important.



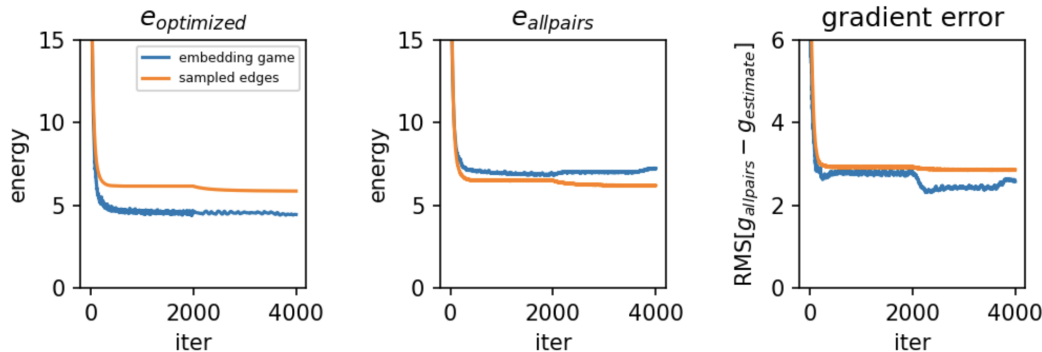


Figure 6: Quantitative comparison between our method (optimized landmarks) and sampled edge method. The left figure show the energy each method actually optimizes. The middle figure shows the “all-pairs” energy (Eq. 1).

## 254 6 Related Work

255 **Landmark methods** Instead of dealing with  $N^2$  nonlinear pairwise interactions, landmark based  
 256 approaches instead designate a small set of  $n$  landmarks, and instead work a smaller set with  $nN$   
 257 iterations. Unlike our method where individual landmarks do not correspond to any particular  
 258 sample, most landmark approaches use sampling to designate certain samples as landmarks. For  
 259 dimensionality reduction, landmark approaches have been applied ISOMAP by [4], and to stochastic  
 260 neighbor embedding vectors by [10]. Finally more sophisticated landmark sampling schemes have  
 261 been used by [11].

262 Perhaps the most well-known class of landmark methods is the Nyström method [16], which is used  
 263 to approximate  $N \times N$  kernel similarity matrices  $g(\mathbf{x}_i, \mathbf{x}_j)$  with two smaller  $N \times n$  and  $n \times n$   
 264 matrices. Our method actually can be understood from the kernel perspective. The heart of our  
 265 method is approximating the sum over all pairs of interactions  $\sum_{i,j} g(\mathbf{x}_i, \mathbf{x}_j)$  which we can interpret  
 266 as the inner product of the vector of all ones and the kernel similarity matrix  $\mathbf{1}^\top \mathbf{G} \mathbf{1}$ .

267 **Game formulations of learning** This work falls into a category of works formulating well-established  
 268 algorithms like multidimensional scaling, principle components analysis, and whitening as a game  
 269 [9, 5]. This work is formulating a certain class of nonlinear embedding problems as a game.

## 270 7 Discussion

271 This paper presents a novel method for approximately the all-pairs repulsive term present in many  
 272 manifold learning algorithms. When the nonlinear repulsion terms are described by a kernel function,  
 273 we can derive a lower bound for the all-pairs sum which we then maximize to find the tightest lower  
 274 bound. We show that this optimization requires much fewer landmarks than would be required  
 275 if we instead just randomly designated embedding vectors to be landmarks. However, compared  
 276 to sampling edges randomly this scheme still requires more computation to achieve comparable  
 277 visualizations.

278 To make this method more useful, future work should find automated schemes for performing the  
 279 minimax optimization, so a user does not have to specify learning rates. This might be much more  
 280 challenging here than for a minimization problem because in general there is no guarantee that  
 281 an increase or decrease in the objective means we are getting closer to a saddle point. In practice  
 282 however, we observed success by setting larger learning rates for  $\eta_z$  and smaller for  $\eta_y$ . This might  
 283 suggest that finding a quick and robust inner loop maximization, with outer loop gradient steps could  
 284 be a promising direction.

285 It would be interesting to apply our method to other regimes. In particular finding high dimensional  
 286 embedding vectors. Additionally, exploring the behavior of this method in the online setting is  
 287 promising, as our method for approximating a sum of all-pairs pairwise interactions does not actually  
 288 require any pairwise distances between embedding vectors to be computed.

## 289 References

- 290 [1] Akshay Agrawal, Alnur Ali, and Stephen P. Boyd. Minimum-distortion embedding. *CoRR*,  
291 abs/2103.02559, 2021.
- 292 [2] Sanjeev Arora, Yingyu Liang, and Tengyu Ma. A simple but tough-to-beat baseline for sentence  
293 embeddings. In *International conference on learning representations*, 2017.
- 294 [3] Mikhail Belkin and Partha Niyogi. Laplacian eigenmaps for dimensionality reduction and data  
295 representation. *Neural computation*, 15(6):1373–1396, 2003.
- 296 [4] Vin De Silva and Joshua B Tenenbaum. Global versus local methods in nonlinear dimensionality  
297 reduction. In *NIPS*, volume 15, pages 705–712, 2002.
- 298 [5] Ian Gemp, Brian McWilliams, Claire Vernade, and Thore Graepel. Eigengame: Pca as a nash  
299 equilibrium. *arXiv preprint arXiv:2010.00554*, 2020.
- 300 [6] Mike Swarbrick Jones. Reddit self posts dataset.  
301 <https://www.kaggle.com/mswarbrickjones/reddit-selfposts>.
- 302 [7] Tianyi Lin, Chi Jin, and Michael Jordan. On gradient descent ascent for nonconvex-concave  
303 minimax problems. In *International Conference on Machine Learning*, pages 6083–6093.  
304 PMLR, 2020.
- 305 [8] Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation  
306 and projection for dimension reduction, 2020.
- 307 [9] Cengiz Pehlevan, Anirvan M Sengupta, and Dmitri B Chklovskii. Why do similarity matching  
308 objectives lead to hebbian/anti-hebbian networks? *Neural computation*, 30(1):84–124, 2017.
- 309 [10] Nicola Pezzotti, Thomas Höllt, B Lelieveldt, Elmar Eisemann, and Anna Vilanova. Hierarchical  
310 stochastic neighbor embedding. In *Computer Graphics Forum*, volume 35, pages 21–30. Wiley  
311 Online Library, 2016.
- 312 [11] Lawrence K. Saul. A tractable latent variable model for nonlinear dimensionality reduction.  
313 *Proceedings of the National Academy of Sciences*, 117(27):15403–15408, 2020.
- 314 [12] Jian Tang, Jingzhou Liu, Ming Zhang, and Qiaozhu Mei. Visualizing large-scale and high-  
315 dimensional data. Republic and Canton of Geneva, CHE, 2016. International World Wide Web  
316 Conferences Steering Committee.
- 317 [13] Laurens van der Maaten. Accelerating t-sne using tree-based algorithms. *Journal of Machine  
318 Learning Research*, 15(93):3221–3245, 2014.
- 319 [14] Laurens Van Der Maaten, Eric Postma, Jaap Van den Herik, et al. Dimensionality reduction: a  
320 comparative. *J Mach Learn Res*, 10(66-71):13, 2009.
- 321 [15] Aaron J Wilk, Arjun Rustagi, Nancy Q Zhao, Jonasel Roque, Giovanni J Martínez-Colón,  
322 Julia L McKechnie, Geoffrey T Ivison, Thanmayi Ranganath, Rosemary Vergara, Taylor Hollis,  
323 et al. A single-cell atlas of the peripheral immune response in patients with severe covid-19.  
324 *Nature medicine*, 26(7):1070–1076, 2020.
- 325 [16] Christopher Williams and Matthias Seeger. Using the nyström method to speed up kernel  
326 machines. In T. Leen, T. Dietterich, and V. Tresp, editors, *Advances in Neural Information  
327 Processing Systems*, volume 13. MIT Press, 2001.

## 328 Checklist

- 329 1. For all authors...
- 330 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
331 contributions and scope? [Yes] We provide a game theoretic objective in Eq. 3 and  
332 show visualizations in Figure 2
- 333 (b) Did you describe the limitations of your work? [Yes] See discussion

- 334 (c) Did you discuss any potential negative societal impacts of your work? [No]  
 335 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
 336 them? [Yes]
- 337 2. If you are including theoretical results...
- 338 (a) Did you state the full set of assumptions of all theoretical results? [Yes] We state the  
 339 assumptions in section 2  
 340 (b) Did you include complete proofs of all theoretical results? [Yes] In the appendix, we  
 341 prove our main theoretical result, that we can bound the sum of pairwise interactions
- 342 3. If you ran experiments...
- 343 (a) Did you include the code, data, and instructions needed to reproduce the main exper-  
 344 imental results (either in the supplemental material or as a URL)? [No] We provide  
 345 PyTorch style pseudo-code in the main text to run our algorithm.  
 346 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
 347 were chosen)? [Yes] In section 3 and 5 we provide these details  
 348 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
 349 ments multiple times)? [No]  
 350 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
 351 of GPUs, internal cluster, or cloud provider)? [No] These are small scale experiments,  
 352 taking no more than a minute on a single GPU to complete
- 353 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 354 (a) If your work uses existing assets, did you cite the creators? [Yes] We cite mnist, scRNA,  
 355 and Reddit posts dataset  
 356 (b) Did you mention the license of the assets? [No]  
 357 (c) Did you include any new assets either in the supplemental material or as a URL? [No]  
 358 (d) Did you discuss whether and how consent was obtained from people whose data you're  
 359 using/curating? [No] We are using publicly available datasets  
 360 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
 361 information or offensive content? [No] We are using publicly available datasets
- 362 5. If you used crowdsourcing or conducted research with human subjects...
- 363 (a) Did you include the full text of instructions given to participants and screenshots, if  
 364 applicable? [N/A]  
 365 (b) Did you describe any potential participant risks, with links to Institutional Review  
 366 Board (IRB) approvals, if applicable? [N/A]  
 367 (c) Did you include the estimated hourly wage paid to participants and the total amount  
 368 spent on participant compensation? [N/A]

## 369 Proof of Inequality in Eq. 2

370 We restate the claim first. If  $g(\|\mathbf{y}_i - \mathbf{y}_j\|)$  is a positive semi-definite kernel function, then for any  
 371  $\{\mathbf{y}_i\}_{i=1}^N$  and  $\{\mathbf{z}_a\}_{a=1}^L$ :

$$\left[ \sum_{i,a} g(\|\mathbf{y}_i - \mathbf{z}_a\|) \right]^2 \leq \left[ \sum_{ij} g(\|\mathbf{y}_i - \mathbf{y}_j\|) \right] \left[ \sum_{a,b} g(\|\mathbf{z}_a - \mathbf{z}_b\|) \right] \quad (10)$$

372 *Proof:* Because we have assumed that  $g$  is a kernel, Mercer's theorem allows us to replace all the  $g(\cdot)$   
 373 with inner products between high dimensional vectors. Specifically for any  $\{\mathbf{y}_i\}_{i=1}^N$  and  $\{\mathbf{z}_a\}_{a=1}^L$   
 374 there exist  $\{\phi_i\}_{i=1}^N$  and  $\{\psi_a\}_{a=1}^L$  such that

$$g(\|\mathbf{y}_i - \mathbf{z}_a\|) = \phi_i \cdot \psi_a \quad g(\|\mathbf{y}_i - \mathbf{y}_j\|) = \phi_i \cdot \phi_j \quad g(\|\mathbf{z}_a - \mathbf{z}_b\|) = \psi_a \cdot \psi_b \quad (11)$$

375 Now proof of Eq. 10 is a simple matter of proving vector inequalities. First we define the sums  
 376  $\bar{\phi} := \sum_i \phi_i$  and  $\bar{\psi} := \sum_a \psi_a$ . We can replace the sums of  $g$  with our vectors. So the left hand side  
 377 is

$$\left[ \sum_{ia} g(\|\mathbf{y}_i - \mathbf{z}_a\|) \right]^2 = (\bar{\phi} \cdot \bar{\psi})^2 = \|\bar{\phi}\|^2 \|\bar{\psi}\|^2 \text{Cos}[\bar{\phi}, \bar{\psi}]^2 \quad (12)$$

378 And the right hand side of Eq. 10 is:

$$\left[ \sum_{i,j} g(\|\mathbf{y}_i - \mathbf{y}_j\|) \right] \left[ \sum_{a,b} g(\|\mathbf{z}_a - \mathbf{z}_b\|) \right] = \|\bar{\phi}\|^2 \|\bar{\psi}\|^2 \quad (13)$$

379 Because  $\text{Cos}^2 \leq 1$ , the left hand side is always less than the right hand side so we have therefore  
380 proved Eq. 10.