

# t-SNE Notes

Iakovidis Ioannis  
email iakoviid@auth.gr

April 4, 2020

## Abstract

t-distributed Stochastic Neighborhood Embedding (t-SNE) is a widely used dimensionality reduction technique, that is particularly well suited for visualization of high-dimensional datasets. On this diploma thesis we introduce a high performance GPU-accelerated implementation of the t-SNE method in CUDA. We base our approach on the work “Spaceland Embedding of Sparse Stochastic Graphs, HPEC 2019”. Obtaining an embedding essentially requires two steps, namely a dense and a sparse computation. One of the main bottlenecks is that usually sparse computation does not scale well in GPU’s because of the irregularity of the memory accesses. To overcome this problem, we use a fast data translocation technique, that leads to locally dense data, which are better suited for GPU processing. The dense computation is performed with an interpolation-based fast Fourier Transform accelerated method. Finally, for datasets that cannot be loaded into the memory, we use randomized and distributed principal component analysis variants, so that the top principal components are used without fully loading the dataset, hence allowing our implementation to be executed by personal computers with superior efficiency.

## Contents

|          |  |           |
|----------|--|-----------|
| <b>1</b> | <b>Introduction</b>                                    | <b>3</b>  |
| 1.1      | Deriving the Gradient . . . . .                        | 4         |
| 1.2      | Gradient Interpretation . . . . .                      | 7         |
| 1.3      | Hyperparameters and Embedding interpretation . . . . . | 8         |
| 1.3.1    | Perplexity . . . . .                                   | 8         |
| 1.3.2    | Early and Late Exaggeration . . . . .                  | 8         |
| 1.4      | First Complexity Discussion . . . . .                  | 8         |
| 1.5      | Use of Momentum . . . . .                              | 9         |
| <b>2</b> | <b>Attractive Term approximation</b>                   | <b>10</b> |
| <b>3</b> | <b>Repulsive Term approximation</b>                    | <b>11</b> |
| 3.1      | Fast Multipole Methods . . . . .                       | 12        |
| 3.2      | Separation of Variables . . . . .                      | 12        |
| 3.3      | Interpolative Separation of Variables . . . . .        | 13        |
| 3.4      | Use of the FFT . . . . .                               | 15        |
| 3.5      | Avoid Padding . . . . .                                | 16        |
| 3.6      | Improving Accuracy . . . . .                           | 17        |

---

|           |  |           |
|-----------|--|-----------|
| 3.7       | Generalizing to Higher Dimensions . . . . .      | 19        |
| 3.8       | Accuracy Comments . . . . .                      | 22        |
| 3.9       | Barnes–Hut . . . . .                             | 22        |
| 3.10      | Barnes–Hut an algebraic interpretation . . . . . | 23        |
| <b>4</b>  | <b>Starting Matlab Implementation</b>            | <b>24</b> |
| 4.1       | Parameter Selection . . . . .                    | 24        |
| 4.2       | Accuracy Analysis . . . . .                      | 24        |
| 4.2.1     | 1 Dimension . . . . .                            | 24        |
| 4.2.2     | 2 Dimensions . . . . .                           | 27        |
| 4.2.3     | 3 Dimensions . . . . .                           | 30        |
| 4.2.4     | 4 Dimensions . . . . .                           | 34        |
| 4.3       | GPU-Implementation thoughts . . . . .            | 36        |
| <b>5</b>  | <b>FIt-SNE</b>                                   | <b>36</b> |
| <b>6</b>  | <b>SG-SNE</b>                                    | <b>36</b> |
| <b>7</b>  | <b>Out of core PCA</b>                           | <b>36</b> |
| <b>8</b>  | <b>Method Comparisons</b>                        | <b>36</b> |
| <b>9</b>  | <b>GPU implementation</b>                        | <b>36</b> |
| <b>10</b> | <b>Data translocations and Memory Hierarchy</b>  | <b>36</b> |

# 1 Introduction

t-Distributed Stochastic Neighbor Embedding (t-SNE) is a technique for dimensionality reduction that is particularly well suited for the visualization of high-dimensional datasets.

Let  $x_i \in \mathbf{R}^d, i = 1, \dots, n$  be the data samples in high dimensional space and  $y_i \in \mathbf{R}^s, i = 1, \dots, n$  be the data samples in low dimensional space,  $s$  is usually 2 or 3. The goal of t-SNE is to preserve local structure, points that are similar (close) in the high dimensional space are mapped to similar (close) points in the low dimensional space. This is achieved by minimizing the divergence between similarity weights that are defined in the high and low dimensional space.

More specifically from the distances  $d_{ij}$  of  $(x_i)_{i=1}^n$  we define

$$p_{i|j} = \frac{e^{-d_{ij}^2/2\sigma_i^2}}{\sum_{l \neq i} e^{-d_{il}^2/2\sigma_i^2}} \quad \text{and} \quad p_{ij} = \frac{p_{i|j} + p_{j|i}}{2n}$$

as the similarity weights in the high dimensional space. The value  $p_{j|i}$  is interpreted as the distribution of the other points given  $x_i$  or the probability that point  $j$  is a neighbor of point  $i$ . The parameters  $\sigma_i$  are chosen (more detail in the next paragraphs). For the low dimensional space we define

$$q_{ij} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_{k \neq l} (1 + \|y_k - y_l\|^2)^{-1}}$$

as similarity weights but this time we use the t-distribution. We use the t-distribution instead of Gaussian in embedding space to avoid the problem of overcrowding. For notation purposes let  $P = \{p_{ij}\}$ ,  $Q = \{q_{ij}\}$   $n \times n$  matrices and  $X = [x_1, \dots, x_n]$ ,  $Y = [y_1, \dots, y_n]$   $d \times n$  and  $s \times n$  matrices respectively.

And then we take KL divergence

$$C(Y) = KL(P \mid Q) = \sum_{i \neq j} p_{ij} \log \frac{p_{ij}}{q_{ij}}.$$

Minimization can be done with gradient descent. Note that this function is non convex and hence we will halt at a local minimum. But the method still gives great results and even has theoretical guaranties about the output embedding despite the fact that is based on a non convex optimization problem. In figure 1 we see the result of embedding the MNIST Dataset in a two dimensional space with the t-SNE method.

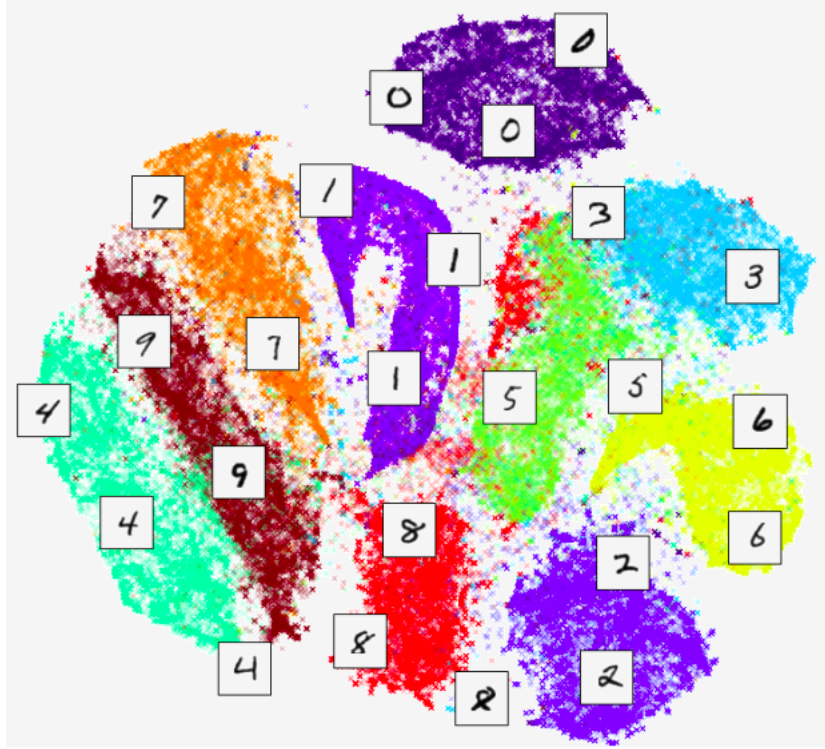


Figure 1: Embedding the MNIST Dataset.

### 1.1 Deriving the Gradient

Here with straight calculations we compute the gradient. First of all:

$$\begin{aligned}
 C(Y) &= \sum_{i \neq j} p_{ij} \log \frac{p_{ij}}{q_{ij}} \\
 &= \sum_{i \neq j} p_{ij} \log p_{ij} - \sum_{i \neq j} p_{ij} \log q_{ij} \\
 &= \sum_{i \neq j} p_{ij} \log p_{ij} + \sum_{i \neq j} p_{ij} \log f_{ij} + \sum_{i \neq j} p_{ij} \log Z \\
 &= \sum_{i \neq j} p_{ij} \log p_{ij} + \sum_{i \neq j} p_{ij} \log f_{ij} + \log Z.
 \end{aligned}$$

Where  $f_{ij} = 1 + \|y_i - y_j\|^2$  and  $Z = \sum_{k \neq l} (1 + \|y_k - y_l\|^2)^{-1}$ .

With some computation:

$$\begin{aligned}
 \frac{\partial \sum_{i \neq j} p_{ij} \log f_{ij}}{\partial y_m} &= \sum_{i \neq m} p_{im} \frac{\partial \log f_{im}}{\partial y_m} + \sum_{i \neq m} p_{mi} \frac{\partial \log f_{mi}}{\partial y_m} \\
 &= 2 \sum_{i \neq m} p_{mi} \frac{\partial \log f_{mi}}{\partial y_m} \\
 &= 2 \sum_{i \neq m} p_{mi} \frac{1}{f_{mi}} \frac{\partial f_{mi}}{\partial y_m} \\
 &= 2 \sum_{i \neq m} p_{mi} \frac{2(y_m - y_i)}{1 + \|y_m - y_i\|^2} \\
 &= 4 \sum_{i \neq m} p_{mi} \frac{Z(y_m - y_i)}{Z(1 + \|y_m - y_i\|^2)} \\
 &= 4 \sum_{i \neq m} p_{mi} q_{mi} Z(y_m - y_i).
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial \log Z}{\partial y_m} &= \frac{\partial Z}{\partial y_m} \frac{1}{Z} \\
 &= \frac{2}{Z} \sum_{i \neq m} \frac{\partial (1 + \|y_m - y_i\|^2)^{-1}}{\partial y_m} \\
 &= -\frac{2}{Z} \sum_{i \neq m} \frac{(y_m - y_i)}{(1 + \|y_m - y_i\|^2)^2} \\
 &= -\frac{2}{Z} \sum_{i \neq m} \frac{2(y_m - y_i)}{1 + \|y_m - y_i\|^2} \\
 &= -4 \sum_{i \neq m} q_{mi}^2 Z(y_m - y_i).
 \end{aligned}$$

So

$$\frac{\partial C(Y)}{\partial y_i} = 4 \sum_{j \neq i} p_{ij} q_{ij} Z(y_i - y_j) - 4 \sum_{j \neq i} q_{ij}^2 Z(y_i - y_j)$$

Also we can perform same the calculation so that it generalizes more easily to changes in the objective function or the similarity weights. Observe how the objective function  $C$  depends on the embedded points. Let more general  $q_{ij} = \frac{w_{ij}}{\sum_{kl} w_{kl}}$  where  $w_{ij}$  is the similarity weight of  $y_i$  and  $y_j$  depended on their distance  $d_{ij}$ .

So when computing the gradient we compute:

1. The distances  $d_{ij}$  are generated from the coordinates,  $y_i, y_j$ .
2. The similarity weights  $w_{ij}$  are generated as a function of the distances,  $d_{ij}$ .
3. The output probabilities,  $q_{ij}$ , are normalized versions of the similarity weights,  $w_{ij}$ .
4. The cost function,  $C$  is normally a divergence of some kind, and hence expressed in terms of the output probabilities,  $q_{ij}$ .

We're going to chain those individual bits together via the chain rule for partial derivatives. The chain of variable dependencies is  $C \rightarrow q \rightarrow w \rightarrow d \rightarrow y$ . So

$$\frac{\partial C}{\partial \mathbf{y}_m} = \sum_{ij} \frac{\partial C}{\partial q_{ij}} \sum_{kl} \frac{\partial q_{ij}}{\partial w_{kl}} \sum_{pq} \frac{\partial w_{kl}}{\partial d_{pq}} \frac{\partial d_{pq}}{\partial \mathbf{y}_m}$$

In the relationship between  $w$ , and  $d$  is such that any cross terms are 0, i.e. unless  $k = p$  and  $l = q$  those derivatives evaluate to 0. Also, either  $k = m$  or  $l = m$ , otherwise  $\partial d_{kl} / \partial \mathbf{y}_m = 0$ . So

$$\begin{aligned} \frac{\partial C}{\partial \mathbf{y}_m} &= \sum_{ij} \frac{\partial C}{\partial q_{ij}} \sum_{kl} \frac{\partial q_{ij}}{\partial w_{kl}} \sum_{pq} \frac{\partial w_{kl}}{\partial d_{pq}} \frac{\partial d_{pq}}{\partial \mathbf{y}_m} \\ &= \sum_{ij} \frac{\partial C}{\partial q_{ij}} \sum_{kl} \frac{\partial q_{ij}}{\partial w_{kl}} \frac{\partial w_{kl}}{\partial d_{kl}} \frac{\partial d_{kl}}{\partial \mathbf{y}_m} \\ &= \sum_{ij} \frac{\partial C}{\partial q_{ij}} \sum_l \frac{\partial q_{ij}}{\partial w_{ml}} \frac{\partial w_{ml}}{\partial d_{ml}} \frac{\partial d_{ml}}{\partial \mathbf{y}_m} + \sum_{ij} \frac{\partial C}{\partial q_{ij}} \sum_k \frac{\partial q_{ij}}{\partial w_{km}} \frac{\partial w_{km}}{\partial d_{km}} \frac{\partial d_{km}}{\partial \mathbf{y}_m} \end{aligned}$$

Exchanging the summations:

$$\begin{aligned} \frac{\partial C}{\partial \mathbf{y}_m} &= \sum_l \left( \sum_{ij} \frac{\partial C}{\partial q_{ij}} \frac{\partial q_{ij}}{\partial w_{ml}} \right) \frac{\partial w_{ml}}{\partial d_{ml}} \frac{\partial d_{ml}}{\partial \mathbf{y}_m} + \sum_k \left( \sum_{ij} \frac{\partial C}{\partial q_{ij}} \frac{\partial q_{ij}}{\partial w_{km}} \right) \frac{\partial w_{km}}{\partial d_{km}} \frac{\partial d_{km}}{\partial \mathbf{y}_m} \\ &= \sum_l f_{ml} \frac{\partial w_{ml}}{\partial d_{ml}} \frac{\partial d_{ml}}{\partial \mathbf{y}_m} + \sum_k f_{km} \frac{\partial w_{km}}{\partial d_{km}} \frac{\partial d_{km}}{\partial \mathbf{y}_m} \\ &= \sum_k f_{mk} \frac{\partial w_{km}}{\partial d_{km}} \frac{\partial d_{km}}{\partial \mathbf{y}_m} + \sum_k f_{km} \frac{\partial w_{km}}{\partial d_{km}} \frac{\partial d_{km}}{\partial \mathbf{y}_m}, \text{ Assume } w \text{ and } d \text{ are symmetric} \\ &= \sum_k \left( f_{mk} + f_{km} \right) \frac{\partial w_{km}}{\partial d_{km}} \frac{\partial d_{km}}{\partial \mathbf{y}_m} \end{aligned}$$

Now

$$\frac{\partial d_{ij}}{\partial \mathbf{y}_i} = \frac{1}{d_{ij}} (\mathbf{y}_i - \mathbf{y}_j) \text{ for } d_{ij} = \left[ \sum_l (y_{il} - y_{jl})^2 \right]^{1/2}$$

And

$$\frac{\partial w_{ij}}{\partial d_{ij}} = -\frac{2d_{ij}}{(1 + d_{ij}^2)} \text{ for } w_{ij} = \frac{1}{(1 + d_{ij}^2)}$$

So

$$\frac{\partial C}{\partial \mathbf{y}_m} = -2 \sum_k (f_{mk} + f_{km}) \frac{(y_m - y_k)}{(1 + d_{km}^2)^2}$$

And for specifically the KL divergence:

$$\frac{\partial C}{\partial q_{ij}} = -\frac{p_{ij}}{q_{ij}}$$

$$\frac{\partial q_{ij}}{\partial w_{km}} = -\frac{w_{ij}}{(\sum_{lt} w_{lt})^2} = -\frac{q_{ij}}{Z}$$

While

$$\frac{\partial q_{ij}}{\partial w_{ij}} = -\frac{w_{ij}}{(\sum_{lt} w_{lt})^2} + \frac{1}{\sum_{lt} w_{lt}} = -\frac{q_{ij}}{Z} + \frac{1}{Z}$$

Ending

$$f_{km} = -\frac{p_{km}}{Z q_{km}} + \sum_{ij} \frac{p_{ij}}{Z}$$

And as a result we get the same gradient expression.

$$\frac{\partial C(Y)}{\partial y_i} = 4 \sum_{j \neq i} p_{ij} q_{ij} Z(y_i - y_j) - 4 \sum_{j \neq i} q_{ij}^2 Z(y_i - y_j)$$

## 1.2 Gradient Interpretation

The t-SNE gradient computation can be reformulated as an N-body simulation problem:

$$F_{attr,i} = \sum_{j \neq i} p_{ij} q_{ij} Z(y_i - y_j)$$

$$F_{rep,i} = \sum_{j \neq i} q_{ij}^2 Z(y_i - y_j)$$

$$\frac{1}{4} \frac{\partial C(Y)}{\partial y_i} = F_{attr,i} - F_{rep,i}$$

The forces can be also viewed with matrix-vector computations in the following ways:

•

$$F_{attr} = (P \odot Q) O \odot Y - (P \odot Q) Y$$

$$F_{rep} = (Q \odot Q) O \odot Y - (Q \odot Q) Y$$

Where  $Y$  is the  $N \times 2$  matrix of embedded points,  $O$  is an  $N \times 2$  matrix of ones.

•

$$F_{attr} = (\text{diag}(\text{sum}((P \odot Q), 1)) - (P \odot Q)) * Y$$

$$F_{rep} = (\text{diag}(\text{sum}((Q \odot Q), 1)) - (Q \odot Q)) * Y$$

in matlab notation (`diag` creates a diagonal matrix with an input vector and `sum(,1)` sums across the rows).

Notice that only the attractive forces term depend on the input similarities, the repulsive forces term depends on the embedded points relative position. This leads to the interpretations that the first term moves  $y_i$  to a weighted average of the other  $y_j$ , with weights bigger of points closer together in the high dimensional space. And the second term pushes points apart when they are too close on the embedding space.

### 1.3 Hyperparameters and Embedding interpretation

t-SNE plots are a useful way for visualizing high-dimensional data. But sometimes can be misleading because some observations do not translate directly from t-SNE plots. Also the parameter choice can have a great effect on the output embedding. By exploring how simple cases behave, we can learn to use t-SNE plots more effectively.

#### 1.3.1 Perplexity

A key hyperparameter for constructing the t-SNE embedding is perplexity. Perplexity is used to compute the  $\sigma_i$  in the similarity weights of the high-dimensional space. Perplexity loosely says how to balance attention between local and global aspects of your data. The parameter is, in a sense, a guess about the number of close neighbors each point has. Choice of perplexity can greatly affect the resulting embedding. Getting the most from t-SNE may mean analyzing multiple plots with different perplexities. Usual values of perplexity range from 5 to 50. Below we see the embedding of 3 Gaussians of 1000 points each, one pair being 10 times as far apart as another pair for different perplexities.

From the above figures we see that some observations we may have from t-SNE plots are not always true. For example distances between clusters might not mean anything. Also cluster sizes in the embedding may be misleading. For this observation we see the next figure that pictures the embedding of 2 Gaussians one with 5 times more points than the other.

#### 1.3.2 Early and Late Exaggeration

As we described above the gradient can be naturally interpreted as a sum of attractive and repulsive forces that move the points in the embedding space. One difficulty when computing the algorithm is that convergence slows down as the number of points increases. A way to improve the optimization process is to multiply the attractive term by an exaggeration parameter  $a$ , this is often performed for the first few hundred iterations (early exaggeration). This trick enables t-SNE to find a better global structure creating tight clusters of points that can more easily move around in the embedding space. Common values for  $a$  are 4 and 12. As we increase the exaggeration parameter we find out that our embedding look smaller and circular. This trick can be used also for the last iterations to shrink the size of the clusters and make them more distinguishable. Ending as it turns out the use of the exaggeration parameter can be rigorously analysed and there is a canonical choice of  $a$  and  $h$  learning rate that leads to exponential convergence and perfect separation of clustered data.

### 1.4 First Complexity Discussion

What is the complexity of this algorithm. First of all we perform gradient descent on a non-convex function so we cannot have an accurate analysis about convergence and condi-



tioning of the problem. But there are theoretical guaranties and t-SNE works surprisingly well. What we will be concerned in this paragraph is the complexity of one step of gradient descent.

To perform gradient descent we need to compute the derivative. By the above computation we have that computing the attractive and repulsive term takes  $O(n^2)$ . But this is infeasible in the following paragraphs we will try to reduce this time by approximating the derivative taking advantage of the structure of  $P \odot Q$  and  $Q \odot Q$ .

## 1.5 Use of Momentum

## 2 Attractive Term approximation

In practice computing the weights for all pairs of points in the high dimensional space ( $P$ ) is too expensive. Instead will approximate this computation by using for every point only the weights of it's  $k$  nearest neighbors. This is justified because the actual weights decrease exponentially with the square of the distance. So after a threshold distance we will have weights very close to zero and dropping them will non change significantly the output embedding. In other words the nearest neighbors capture the local structure of the dataset.

This leads to our new  $P$  and hence  $P \odot Q$  matrix being sparse. And so we can accelerate the computation of the attractive forces by using only the non-zero elements. The computation of the attractive forces resembles a sparse matrix vector product. Computing the  $k$  nearest neighbours can be implemented by the use of kd or vantage point trees also we can compute approximate nearest neighbours for a faster implementation. The resulting complexity for computation of the attractive term is  $O(k \cdot n)$

Experimental we find out that  $k \geq u$  where  $u$  the specified perplexity is enough for a faithful approximation. More specifically we redefine the pairwise similarities between the input objects as:

$$p_{j|i} = \begin{cases} \frac{\exp(-d(x_i, x_j)^2/2\sigma_i^2)}{\sum_{k \in N_i} \exp(-d(x_i, x_k)^2/2\sigma_i^2)}, j \in N_i & \text{and } p_{ij} = \frac{p_{i|j} + p_{j|i}}{2n} \\ 0, & \text{otherwise} \end{cases}$$

Here  $N_i$  are the set of the  $k$  nearest neighbours from  $X$  to  $x_i$ .

### 3 Repulsive Term approximation

Computing the repulsive term is vastly different than the attractive term. The matrix  $Q \odot Q$  can't be approximated to a sparse matrix as it's elements are the weights of the embedded points and change in every iteration. Also the t-distribution has a longer tail (thing that gives us greater results solving the Crowding problem) than the normal distribution we couldn't do this approximation. So we need to approach the computation for  $Q \odot Q$  being dense likely we have methods from numerical analysis that can handle this problem.

Now we will express the Repulsive forces term in the way that we will use in the next paragraphs. Remember that we have:

$$F_{rep,i} = \sum_{j=1, j \neq i}^n \frac{y_i - y_j}{(1 + \|y_j - y_i\|^2)^2} \cdot \frac{1}{\sum_{k \neq l} (1 + \|y_k - y_l\|^2)^{-1}}$$

We can compute this expression by computing sums of the form:

$$u_i = \sum_{j=1}^N G(y_i, y_j) v_j$$

where  $G$  is either  $G_1(y_i, y_j) = \frac{1}{1 + \|y_i - y_j\|^2}$  or  $G_2(y_i, y_j) = \frac{1}{(1 + \|y_i - y_j\|^2)^2}$ .

More precisely for a 2-dimensional embedding we can compute  $F_{rep}$  from:

$$\begin{aligned} h_{1,i} &= \sum_{\substack{j=1 \\ j \neq i}}^N G_1(y_i, y_j) \\ h_{2,i} &= \sum_{\substack{j=1 \\ j \neq i}}^N G_2(y_i, y_j) y_j(1) \\ h_{3,i} &= \sum_{\substack{j=1 \\ j \neq i}}^N G_2(y_i, y_j) y_j(2) \\ h_{4,i} &= \sum_{\substack{j=1 \\ j \neq i}}^N G_2(y_i, y_j) \end{aligned}$$

Then at each step of the gradient descent, the repulsive forces can then be expressed in terms of these 4 sums as follows:

$$F_{rep,i}(1) = (h_{2,i} - y_i(1)h_{4,i})/Z$$

$$F_{rep,i}(2) = (h_{3,i} - y_i(2)h_{4,i})/Z$$

and

$$Z = \sum_{j=1}^N h_{1,i}.$$

### 3.1 Fast Multipole Methods

Essentially a subproblem of computing the gradient is to compute sums of the form:

$$u_i = \sum_{j=1}^N G(y_i, y_j) v_j$$

in our case for computing the repulsive forces  $G(y_i, y_j) = \frac{1}{1 + \|y_i - y_j\|^2}$ . More generally the fast multipole methods are constructed to compute sums of the form:

$$u_i = \sum_{j=1}^N G(t_i, y_j) v_j$$

the points  $\{t_i\}_1^M$  are called targets and  $\{y_i\}_1^N$  are called sources, both of dimension  $s$ . While the kernel  $G(t_i, y_j)$  depends on the distance of  $t_i$  and  $y_j$ . Notice that computing these sums naively takes  $O(NM)$  (quadratic) time.

Putting all evaluations of the kernel in a matrix  $A : A_{ij} = G(t_i, y_j)$  we can see that  $u = Av$ . It is known that if the domain of  $T$  is different than the domain of  $Y$  then we can separate the variables and have  $A \approx B(X)C(Y)$ , let  $B$  be  $M \times P$  and  $C$  be  $P \times N$  matrices. This is true because if the kernel have separate domains (does not blow up) then  $A$  can be efficiently approximated by a low rank matrix.

$$\begin{aligned} u_i &= \sum_{j=1}^N G(t_i, y_j) v_j \\ &= \sum_{j=1}^N \sum_{k=1}^P B_{ik} C_{kj} v_j \\ &= \sum_{k=1}^P B_{ik} \left( \sum_{j=1}^N C_{kj} v_j \right) \\ &= \sum_{k=1}^P B_{ik} m_k \end{aligned}$$

Note that  $m_k$  is only computed once for all  $i$ , meaning that the sums can be computed in  $P \cdot (N + M)$  computations—a dramatic improvement over  $M \cdot N$ .

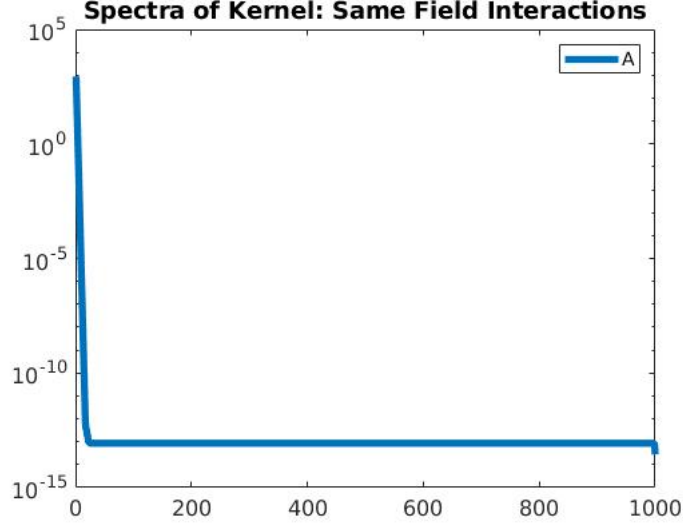
### 3.2 Separation of Variables

But how can we find such a separation the variables. First we will discuss the optimal separation (optimal in terms of accuracy). Assume the Singular Value Decomposition (SVD) of  $A$ ,  $A \approx U \Sigma V^T$  taking the  $p$  highest singular values. Let  $u_i, v_i$  be the columns of  $U$  and  $V$  while  $\sigma_i$  be the singular values. Notice that:

$$\begin{aligned} A &\approx U \Sigma V^T = \sum_{k=1}^p \sigma_k u_k v_k^T \\ \implies A_{ij} &= \sum_{k=1}^p \sigma_k u_k(i) v_k^T(j) \end{aligned}$$

This is a separation of variables like what described above. But how is it's properties in terms of accuracy? The SVD for the  $p$  highest singular values is the optimal approximation of the matrix in terms of the  $L_2$  norm (Eckart-Young-Minsky theorem). The error is given in terms of the singular that we didn't consider. In our case because of our kernel such low rank approximation is very accurate.

Figure 2: Singular values of the Kernel for 1000 random points.



Also notice that because in our case the set of targets and sources is the same we have  $U = V$  and so:

$$A = V \Sigma V^T$$

### 3.3 Interpolative Separation of Variables

But finding the SVD of a matrix even approximately is slow (for  $p$  largest singular values:  $O(N \cdot M \cdot p)$ ). So we need another way to separate variables. We will use an interpolation based technique. Starting let  $s = 1$  (1-dimensional embedding space) and then we will generalize. Instead of the Kernel we use the Lagrange interpolation polynomial for the Kernel defined by equidistant points in the target and source fields.

Let  $R_y$  and  $R_z$  the intervals denoting the ranges of the source and target points,  $y_i \in R_y, \forall i = 1, \dots, N$  and  $z_i \in R_z, \forall i = 1, \dots, M$ . Now suppose that  $\tilde{z}_1, \dots, \tilde{z}_p$  are equidistant points in  $R_z$  and  $\tilde{y}_1, \dots, \tilde{y}_p$  are equidistant points in  $R_y$ , the number of points  $p$  will control the error of our approximation method. The importance of the points being equidistant will be explored in the next paragraph. Let  $L_{l,\tilde{y}}$  and  $L_{l,\tilde{z}}$  be the Lagrange polynomials :

$$L_{l,\tilde{y}}(y) = \frac{\prod_{j \neq l}^p (y - \tilde{y}_j)}{\prod_{j \neq l}^p (\tilde{y}_l - \tilde{y}_j)} \text{ and } L_{l,\tilde{z}}(z) = \frac{\prod_{j \neq l}^p (z - \tilde{z}_j)}{\prod_{j \neq l}^p (\tilde{z}_l - \tilde{z}_j)}$$

With these polynomials we can make an interpolation polynomial  $G_p$  for the Kernel.

$$G_p(z, y) = \sum_{l=1}^p \sum_{j=1}^p G(\tilde{z}_l, \tilde{y}_j) L_{l, \tilde{z}}(z) L_{j, \tilde{y}}(y)$$

Now with the use of the approximation of  $G$  with  $G_p$  we can approximate the

$$u_i = \sum_{j=1}^N G(t_i, y_j) v_j$$

by

$$\tilde{u}_i = \sum_{j=1}^N G_p(t_i, y_j) v_j.$$

So

$$\begin{aligned} \tilde{u}_i &= \sum_{j=1}^N G_p(t_i, y_j) v_j \\ &= \sum_{j=1}^N \sum_{l=1}^p \sum_{m=1}^p G(\tilde{z}_l, \tilde{y}_m) L_{l, \tilde{z}}(z_i) L_{m, \tilde{y}}(y_j) v_j \\ &= \sum_{l=1}^p L_{l, \tilde{z}}(z_i) \left( \sum_{m=1}^p G(\tilde{z}_l, \tilde{y}_m) \left( \sum_{j=1}^N L_{m, \tilde{y}}(y_j) v_j \right) \right). \end{aligned}$$

By this separation we can compute the values using three convolutions (by the nested parentheses) of  $\{\tilde{u}_i\}_1^M$  in  $O((M + N)p + p^2)$  time.

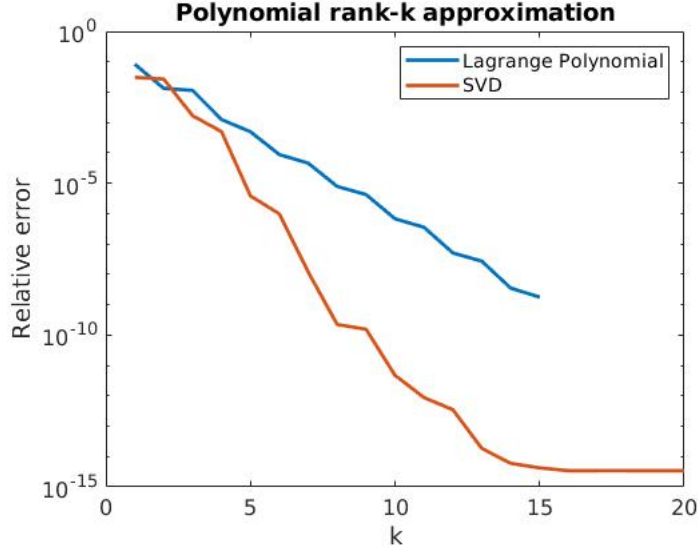
But what are the accuracy properties of this method. The error

$$\begin{aligned} |u_i - \tilde{u}_i| &= \left| \sum_{j=1}^N (G(z_i, y_j) - G_p(z_i, y_j)) \cdot v_j \right| \\ &\leq \sum_{j=1}^N |G(z_i, y_j) - G_p(z_i, y_j)| \cdot |v_j| \\ &\leq \epsilon \sum_{j=1}^N |v_j|. \end{aligned}$$

Provided that the polynomial uniformly approximates the kernel in the ranges. That is:

$$\exists \epsilon : \sup |G(z, y) - G_p(z, y)| \leq \epsilon$$

Figure 3: Relative Error of SVD and Lagrange approximations.



### 3.4 Use of the FFT

For the following text we will consider the case where the targets and the sources are the same set. In computing the convolution  $\sum_{m=1}^p G(\tilde{y}_l, \tilde{y}_m) w_m$  notice that can be expressed as a matrix vector product  $K \cdot w$  where  $K$  the matrix of evaluations of the kernel and  $w$  the vector of the values  $\{w_m\}_1^p$ . But because the points  $\tilde{y}_l$  are equidistributed the matrix  $K$  is Toeplitz and hence the product can be computed efficiently with the FFT in time  $O(p \log p)$ . Note that a matrix vector product takes  $O(p^2)$  but here we can use the FFT to exploiting the structure of a Toeplitz matrix.

The FFT can be used in the following manner. First of all let

$$K = \begin{bmatrix} t_0 & t_1 & \dots & t_{p-1} \\ t_1 & t_0 & \dots & \dots \\ \dots & \dots & \dots & t_1 \\ t_{p-1} & \dots & t_1 & t_0 \end{bmatrix}$$

we construct the matrix

$$C = \begin{bmatrix} K & B \\ B & K \end{bmatrix}$$

where

$$B = \begin{bmatrix} 0 & t_{p-1} & \dots & t_1 \\ t_{p-1} & 0 & \dots & \dots \\ \dots & \dots & \dots & t_{p-1} \\ t_1 & \dots & t_{p-1} & 0 \end{bmatrix}.$$

The matrix  $C$  is circulant and so we can compute it's matrix vector products with the FFT. Now zero padding the vector  $w$  and performing the product should give us the answer.

$$C \cdot \begin{bmatrix} w \\ 0 \end{bmatrix} = \begin{bmatrix} Kw \\ Bw \end{bmatrix}$$

```

1  %FFT matrix-vector product of a Toeplitz matrix
2  p=1024;
3  y = randn(p,1) ;
4  w = randn(p,1) ;
5  K=toeplitz(y);
6
7  z=K*w;
8
9  a=[0;y(p:-1:2)];
10 B=toeplitz(a);
11 C2=[K B;B K];
12
13 w2=[w;zeros(p,1)];
14
15 b=C2*w2;
16 b_fft=ifft(fft(w2).*fft(C2(1,:))');
17 disp(norm(b_fft(1:p)-z));

```

### 3.5 Avoid Padding

As we will see later the fact that we pad the vectors using twice the size memory for the FFT computation is sub-optimal. This will be more apparent when we will generalize to higher dimensions as it will increase our memory by a factor of  $2^d$ ,  $d$  the dimension of the embedding space. For that reason we will see how to remedy this problem here for one dimension, because the visualization much is simpler.

Recapping we want to perform a Toeplitz matrix-vector product. In the previous paragraph we saw that this can be computed by:

$$F^{-1}\left(F([t_0 \ \dots \ t_{p-1} \ 0 \ t_{p-1} \ \dots \ t_1]) \odot F([w_0 \ \dots \ w_{p-1} \ 0 \ \dots \ 0])\right)$$

Where  $t_0, \dots, t_{p-1}$  define the Toeplitz matrix and  $F, F^{-1}$  perform the DFT and inverse DFT. We will advocate that this can be calculated as  $b$  where:

$$b = \frac{b_1 + b_2}{2}$$

$$b_1 = F^{-1}\left(F([t_0 \ \dots \ t_{p-1}] + [0 \ t_{p-1} \ \dots \ t_1]) \odot F([w_0 \ \dots \ w_{p-1}])\right)$$

$$b_2 = \bar{c} \odot F^{-1}\left(c \odot (F([t_0 \ \dots \ t_{p-1}] - [0 \ t_{p-1} \ \dots \ t_1]) \odot F([w_0 \ \dots \ w_{p-1}])\right)$$

$$c = \left[ e^{\frac{2\pi i \cdot 0}{2n}} \ \dots \ e^{\frac{2\pi i \cdot (p-1)}{2n}} \right] \text{ and } \bar{c} \text{ it's conjugate.}$$

Let  $T_1 = [t_0 \ \dots \ t_{p-1}]$  and  $T_2 = [0 \ t_{p-1} \ \dots \ t_1]$ .

Now taking the DFT of the padded Toeplitz vector  $T = [t_0 \ \dots \ t_{p-1} \ 0 \ t_{p-1} \ \dots \ t_1]$  we see that:

$$F(T)[k] = \sum_{j=0}^{2N-1} T[j]w^{jk}, w^{jk} = e^{\frac{-2\pi i}{2n}jk}$$



$$\begin{aligned}
&= \sum_{j=0}^{N-1} T_1[j] w^{jk} + \left( \sum_{j=0}^{N-1} T_2[j] w^{jk} \right) \cdot w^{kn} \\
&= \begin{cases} F(T_1)[k] + F(T_2)[k], & \text{for } k \text{ even} \\ F(c \odot T_1)[k] - F(c \odot T_2)[k], & \text{for } k \text{ odd.} \end{cases}
\end{aligned}$$

Because  $w^{2kj} = 1$  and  $w^{(2k+1)j} = w^j$ . So performing the Hadamard product with the input vector  $V$  gives us:

$$(F(T) \cdot F(V))[k] = \begin{cases} (F(T_1 + T_2)[k]) \cdot F(v)[k], & \text{for } k \text{ even} \\ F(c \odot (T_1 - T_2))[k] \cdot F(c \odot v)[k], & \text{for } k \text{ odd.} \end{cases}$$

where  $v$  the not zero padded input vector. Now we have everything that we need to compute the matrix-vector product, but because we want to use the  $n$  point IDFT instead of the  $2n$  point IDFT we need to compute the above formulas for  $k$  in every parity. As it turns out the result from the cases on opposite parities cancel and we get the above expression.

### 3.6 Improving Accuracy

The procedure above is not numerically stable, and given a large number of the equispaced interpolation points, it will suffer from the Runge phenomenon. Runge phenomenon states that increasing the number of interpolation points does not always increase the accuracy of a function approximation, as a high degree polynomial can have large variations between interpolation points. So, instead of using  $p$  interpolation points for the whole interval, we split the interval into  $N_{int}$  sub-intervals, each with  $p$  interpolation points.

Specifically we subdivide the interval in  $N_{int}$  intervals of equal length  $I_j, j = 1, \dots, N_{int}$ . Let  $y_{j,l}$  denote the  $j$ 'th point in the interval  $I_l$ .

$$y_{j,l} = \frac{h}{2} + ((j-1) + (l-1) \cdot p) \cdot h$$

where  $h = 1/(N_{int} \cdot p)$  the length of an interval.

The algorithm for computing the sum is the same as described in the previous paragraphs but for approximating the kernel we use piecewise polynomial interpolants from the corresponding intervals. We interpolate each point using the interpolation points within its interval. Summarizing:

**Step 1** For each interval  $I_l, l = 1, \dots, N_{int}$  we compute the coefficients  $w_{m,l}$  defined by the formula:

$$w_{m,l} = \sum_{y_j \in I_l} L_{m,l}(y_j) v_j$$

**Step 2** Compute the values  $u_{m,n}$  at the equispaced points  $y_{m,n}$  defined by the formula:

$$u_{m,n} = \sum_{j=1}^{N_{int}} \sum_{j=1}^p G(y_{m,n}, y_{l,j}) w_{l,j}$$

**Step 3** For each point  $y_i$  we compute the output by the interval that it belongs if  $y_i \in I_l$  so:

$$f(y_i) = \sum_{j=1}^p L_{j,l}(y_i) u_{j,l}$$

Notice because all the interpolation points are still equispaced step 2 can be computed with the FFT.

```

1  %% FLT SNE with improved accuracy through subintervals
2  clear all;
3
4  rng(1)
5  a = 0; b=1; n = 1000;
6  [locs,~] = sort(rand(n,1)*(b-a));
7  distmatrix = squareform(pdist(locs));
8  kernel = 1./(1+distmatrix.^2);
9
10 v = sin(10*locs) + cos(2000*locs); %Anything could be used here
11
12 f = kernel*v;%Result
13
14 %% Number of Nint intervals k interpolation points per interval h interval
15 %length
16 k = 2;
17 Nint = 5;
18 h = 1/(Nint *k);
19
20
21 %k interpolation points in each interval
22 interp_points = zeros(k,Nint);
23 for j=1:k
24     for int=1:Nint
25         interp_points(j,int) = h/2 + ((j-1)+(int-1)*k)*h;
26     end
27 end
28
29
30 %% We need to be able to look up which interval each point belongs to
31 int_lookup = zeros(n,1);
32 current_int = 0;
33 for i=1:n
34     if (k*h*(current_int) < locs(i))
35         current_int = current_int +1;
36     end
37     int_lookup(i) = current_int;
38 end
39
40
41 %% Make V, which is now n rows by Nint*k columns
42 V = zeros(n,Nint*k);
43 for ti=1:k
44     for yj=1:n
45         current_int = int_lookup(yj);
46         num = 1;
47         denom = 1;
48         for tii=1:k
49             if (tii ≠ ti)
50                 denom = denom*(interp_points(ti,current_int) ...
51                     -interp_points(tii,current_int));
52                 num= num*(locs(yj) - interp_points(tii,current_int));
53             end
54         end
55         V(yj,(current_int-1)*k+ti) = num/denom;

```

```

56     end
57 end
58
59 %% Make S, which is k*Nint by k*Nint
60 S = ones(k*Nint,k*Nint);
61 for int1=1:Nint
62     for i=1:k
63         for int2=1:Nint
64             for j=1:k
65                 S((int1-1)*k+i, (int2-1)*k+j) = ...
                    1/(1+norm(interp.points(i,int1)-interp.points(j,int2))^2);
66             end
67         end
68     end
69 end
70 %% FLT SNE
71 f.poly_approx=V'*v;
72
73 a=[0;S(k*Nint:-1:2,1)];
74 B=toeplitz(a);
75 C2=[S B;B S];
76
77 v2=[f.poly_approx;zeros(k*Nint,1)];
78
79 b_fft=ifft(fft(v2).*fft(C2(1,:))');
80 b_fft=b_fft(1:k*Nint);
81 f.poly_approx = V*b_fft;
82
83 relativeError=norm(f.poly_approx-f)/norm(f);

```

### 3.7 Generalizing to Higher Dimensions

Here we will describe how we can use the same methods as described above when the embedding space is 2 dimensional further generalization is obvious. To use this interpolation process over a embedding space we will must have an equispaced grid of interpolation points. Lets say we have  $p$  points per row of the grid and  $p^2$  total. We take the approximation.

$$\begin{aligned}
 \tilde{u}_i &= \sum_{j=1}^N G_p(y_i, y_j) v_j \\
 &= \sum_{j=1}^N \sum_{l=1}^{p^2} \sum_{m=1}^{p^2} G(\tilde{y}_l, \tilde{y}_m) L_{l,\tilde{y}}(y_i) L_{m,\tilde{y}}(y_j) v_j \\
 &= \sum_{l=1}^{p^2} L_{l,\tilde{y}}(y_i) \left( \sum_{m=1}^{p^2} G(\tilde{y}_l, \tilde{y}_m) \left( \sum_{j=1}^N L_{m,\tilde{y}}(y_j) v_j \right) \right).
 \end{aligned}$$

Mind that  $u_i$  and  $v_j$  are 2 dimensional vectors.  $L_m$  is the Lagrange polynomial in the  $m$ 'th interpolation point.

$$L_m = L_{xm} L_{ym}$$

Where  $L_{xm}, L_{ym}$  the interpolation polynomials over the  $x$  and  $y$  direction.

We will perform the same algorithm described in the above paragraph computing the sum with the order specified by the parentheses. But in these case the matrix  $S$ ,  $S_{ij} = G(\tilde{y}_i, \tilde{y}_j)$

is not Toeplitz as in the one dimensional case and can not use the same procedure to perform the product with the use of the FFT.

Instead, we observe that the matrix of evaluations of the kernel is a block-Toeplitz matrix with Toeplitz blocks (BTBT). Example for a grid of four points:

$$S = \begin{bmatrix} s_0 & s_1 & s_1 & s_2 \\ s_1 & s_0 & s_2 & s_1 \\ s_1 & s_2 & s_0 & s_1 \\ s_2 & s_1 & s_1 & s_0 \end{bmatrix} = \begin{bmatrix} S_0 & S_1 \\ S_1 & S_0 \end{bmatrix}$$

So what we need is a fast matrix vector product for matrices of this structure. We can compute this product with the use of the FFT as the output elements can be expressed as a convolution of the non-redundant elements of  $S$  and the vector zero padded.

But writing and processing the full BTBT matrix is terribly inefficient (especially in memory). An other way to use the FFT in order to make this product is to construct a matrix of the form.

$$S' = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & s_2 & s_1 & s_2 \\ 0 & s_1 & s_0 & s_1 \\ 0 & s_2 & s_1 & s_2 \end{bmatrix}$$

or in general:

$$S' = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & \dots & 0 \\ 0 & s_{0,(p^2-1)} & \dots & s_{0,(p \cdot (p-1)+1)} & s_{0,p \cdot (p-1)} & s_{0,(p \cdot (p-1)+1)} & \dots & s_{0,(p^2-1)} \\ 0 & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & s_{0,(p-1)} & \dots & s_{0,1} & s_{0,0} & s_{0,1} & \dots & s_{0,(p-1)} \\ 0 & s_{0,(2p-1)} & \dots & s_{0,(p+1)} & s_{0,p} & s_{0,(p+1)} & \dots & s_{0,(2p-1)} \\ 0 & s_{0,(p-1)} & \dots & s_{0,1} & s_{0,0} & s_{0,1} & \dots & s_{0,(p-1)} \\ 0 & s_{0,(2p-1)} & \dots & s_{0,(p+1)} & s_{0,p} & s_{0,(p+1)} & \dots & s_{0,(2p-1)} \\ 0 & s_{0,(p^2-1)} & \dots & s_{0,(p \cdot (p-1)+1)} & s_{0,p \cdot (p-1)} & s_{0,(p \cdot (p-1)+1)} & \dots & s_{0,(p^2-1)} \end{bmatrix}$$

where  $s_{0,i}$  is the evaluation of the kernel between the 0'th and the  $i$ 'th point in the interpolation grid (these are all the non redundant values of  $S$ ). To perform the multiplication  $S \cdot v$  with the use of  $S'$  and the FFT we compute:

$$R = \text{ifft2}\left(\text{fft2}(S') \odot \text{fft2}(V)\right)$$

where

$$V = \left[ \begin{array}{c|c} O_{p \times p} & O_{p \times p} \\ \hline O_{p \times p} & \begin{bmatrix} v(1:p)^T \\ v(p+1:2p)^T \\ \dots \\ v(p(p-1)+1:p^2)^T \end{bmatrix} \end{array} \right]$$

Our end result will be the  $R(1:p, 1:p)$  row major order. Why does this work?

Let's take the same one dimensional procedure for the BTBT matrix  $S$ . Let

$$C = \begin{bmatrix} B & S \\ S & B \end{bmatrix} \text{ where } B = \begin{bmatrix} 0 & S_1 \\ S_1 & 0 \end{bmatrix}$$

The matrix  $C$  exhibits a Block circulant structure and hence we could use a block FFT to accelerate the computation of the product.

$$C \cdot \begin{bmatrix} 0_{2 \times 2} \\ w \end{bmatrix}$$

Essentially this block FFT is the column transform that is part of the 2 dimensional FFT. Note that the 2 dimensional dft can be written as a column and a row dft:

$$\text{fft2}(X) = \text{fft}(\text{fft}(X, [], 2), [], 1) = W^*(xW^*)$$

This explains the column structure of  $S'$ . And the row transform computes the matrix vector products of the individual blocks in the block matrix multiplication.

This procedure doesn't change when we partition our domain into blocks because all the points across the blocks form an equispaced grid. Also it easily generalizes into higher dimensions.

**Theorem 1** (FFT-formula). *Let  $S$  a  $k$ -level BTTB matrix and  $v$  a vector. Then we can perform matrix vector product with the following computation:*

$$S \cdot v = i\text{FFTk}(\text{FFTk}(S') \odot \text{FFTk}(V))$$

Where  $S'$  and  $V$  are computed from  $S$  and  $v$  as described above.

*Proof.* We will use induction on  $k$ .

First remember the base case let  $S$   $n \times n$  Toeplitz and  $v$   $n \times 1$  vector.

We can take the matrix

$$C = \begin{bmatrix} B & S \\ S & B \end{bmatrix} \text{ where } B = \text{Toeplitz}(S(1, n : -1 : 2))$$

Then

$$C \cdot \begin{bmatrix} 0 \\ v \end{bmatrix} = \begin{bmatrix} S \cdot v \\ B \cdot v \end{bmatrix}$$

But  $C$  is a circulant matrix and so it is diagonalizable by the DFT matrices  $F, F^1$ . So

$$C \cdot \begin{bmatrix} 0 \\ v \end{bmatrix} = \begin{bmatrix} S \cdot v \\ B \cdot v \end{bmatrix} = F^{-1} \cdot \text{diag}(\text{DFT}(C(1, :))) \cdot F \cdot V = \text{ifft}(\text{fft}(S') \odot \text{fft}(V))$$

Now let the statement for  $k+1$  let  $l \times l$  be the  $k+1$ -level block size. Then take

$$C = \begin{bmatrix} B & S \\ S & B \end{bmatrix} \text{ where } B = \text{Toeplitz}(S(1 : l, n : -1 : 2))$$

$C$  is block circulant so if we use the block DFT matrices  $F_b, F_b^{-1}$  we get.

$$C \cdot \begin{bmatrix} 0 \\ v \end{bmatrix} = \begin{bmatrix} S \cdot v \\ B \cdot v \end{bmatrix} = F_b^{-1} \text{Diag}(\text{DFT}(C(1 : l, :))) F_b \cdot V$$

Now we can use the inductive step because  $\text{blockDiag}(\text{DFT}(C(1 : l, :))) \cdot V$  is a product of  $k$ -level BTTB matrices and the column-vector of  $V$ . So:

$$\begin{aligned} C \cdot \begin{bmatrix} 0 \\ v \end{bmatrix} &= F_b^{-1} \cdot i\text{FFTk}(\text{FFTk}(\text{Diag}(\text{DFT}(C(1 : l, :))))' \odot \text{FFTk}((F_b \cdot V)')) \\ &= i\text{FFT}(k+1)(\text{FFT}(k+1)(S') \odot \text{FFT}(k+1)(V')) \end{aligned}$$

□

Snippet below performs the 3 dimensional computation.

```

1  kernel_tilde=zeros(2*N1d,2*N1d,2*N1d);
2  for i = 0:N1d-1
3      for j =0:N1d-1
4          for z=0:N1d-1
5              tmp=kernel([y_tilde(1,1) y_tilde(1,2) y_tilde(1,3) ...
6                  ],[y_tilde(i+1,1) y_tilde(j+1,2) y_tilde(z+1,3)],squared);
7              for signi=-1:2:1
8                  for signj=-1:2:1
9                      for signz=-1:2:1
10
11                          kernel_tilde((N1d +signi*i)+1 , (N1d + ...
12                              signj*j)+1,(N1d + signz*z)+1 ) = tmp;
13                      end
14                  end
15              end
16          end
17      end
18
19  fft_kernel=fftn(kernel_tilde);
20
21
22  b=zeros(N1d^3,nsums);
23  for nterms=1:nsums
24      fa=zeros(2*N1d,2*N1d,2*N1d);
25      for(i=1:N1d)
26          for(j=1:N1d)
27              for(z=1:N1d)
28                  fa(i+N1d,j+N1d,z+N1d)=w((i-1)*N1d+j+(z-1)*N1d^2,nterms);
29              end
30          end
31      end
32      result=ifftn(fftn(fa).*fft_kernel);
33
34      result= result(1:N1d,1:N1d,1:N1d);
35      for(i=1:N1d)
36          for(j=1:N1d)
37              for(z=1:N1d)
38                  b((i-1)*N1d+(z-1)*N1d^2+j,nterms)=result(i,j,z);
39              end
40          end
41      end
42
43  end

```

### 3.8 Accuracy Comments

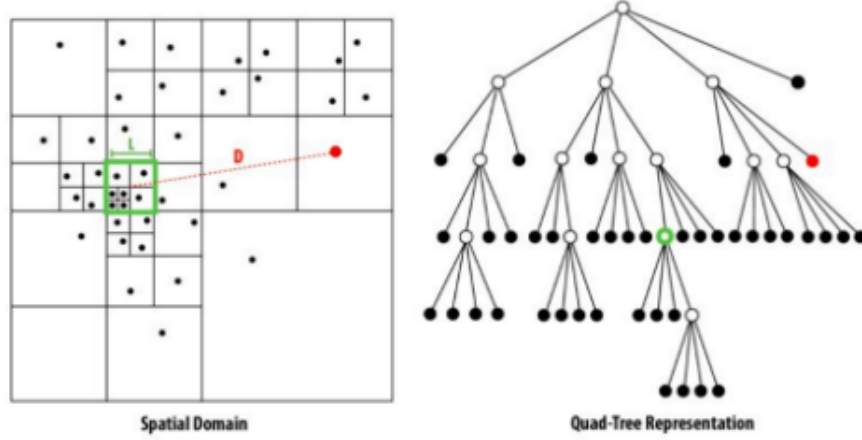
### 3.9 Barnes–Hut

The Barnes–Hut simulation is an approximation algorithm for performing an n-body simulation. It is notable for having order  $O(n \log n)$  compared to a direct-sum algorithm which would be  $O(n^2)$ . The space is hierarchically divided into cells (rectangular or cubic), so that only points from nearby cells need to be treated individually, and points in distant

cells can be treated as a single large point centered at the cell's center of mass. This can dramatically reduce the number of particle pair interactions that must be computed.

More specifically the Barnes–Hut algorithm constructs a octtree or quadtree and for each point it does a depth first search on that tree at each node testing a condition that says if the approximation for this cell is valid.

Figure 4: Barnes-Hut quad-tree representation.



### 3.10 Barnes–Hut an algebraic interpretation

It is clear that given a constructed space tree we can reorder the matrix  $A = \{q_{ij}^2\}$  to a matrix  $A'$  that has the recursive partition:

$$A' = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$$

for a binary tree  $y \in \mathbf{R}$  or

$$A' = \begin{bmatrix} A_{11} & A_{12} & A_{13} & A_{14} \\ A_{21} & A_{22} & A_{23} & A_{24} \\ A_{31} & A_{32} & A_{33} & A_{34} \\ A_{41} & A_{42} & A_{43} & A_{44} \end{bmatrix}$$

for a quadtree  $y \in \mathbf{R}^2$  similarly for an octtree.

Then for the computation for

$$u = A'v = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \begin{bmatrix} v_{near} \\ v_{far} \end{bmatrix} = \begin{bmatrix} A_{11}v_{near} + A_{12}v_{far} \\ A_{21}v_{near} + A_{22}v_{far} \end{bmatrix} \approx \begin{bmatrix} A_{11}v_{near} + a_1 \\ a_2 + A_{22}v_{far} \end{bmatrix}$$

$a_1$  and  $a_2$  represent the approximation of the the terms of the cells with the combined force from the center of mass. This approximation can be done recursively at any level.

## 4 Starting Matlab Implementation

This section is devoted in illustrating the properties of the method and it's approximations. For this task we made a simple MATLAB implementation that computes the exact and FFT-interpolation based embeddings. For all the experiments below we use the MNIST dataset. And we describe the quality of approximation embeddings up to 4 dimensions.

### 4.1 Parameter Selection

Choosing the right parameters is crucial in getting an interpretable embedding. We use perplexity 30, natural choice for the MNIST dataset. This leads us to use 90 nearest neighbours ( $k$ ) for the attractive forces approximation. Generally we use the rule  $k = 3 \cdot \text{perplexity}$ . Also as we will illustrate later the repulsive approximation error increases if for the same interpolation parameters we increase the width of the dataset. So to output generally smaller embeddings we use an medium learning rate of 200 and an early exaggeration of 12. Getting smaller embeddings make our implementation also use way less memory for the same accuracy. This extremely useful when we have a high number of dimensions (3 and 4) that demand a lot of memory.

### 4.2 Accuracy Analysis

The overall error of the implementation can be measured as an error form the attractive and repulsive term approximations. We will look at both of them to determine the the effectiveness of the implementation as well as the total error. We use the measure Relative Square Root Error (RSRE using  $L_2$  norm):

$$error = \frac{\|exact - aprox\|}{\|exact\|}.$$

#### 4.2.1 1 Dimension

First we consider the Repulsive error for embedding in 1 dimension 500 randomly sampled points of the MNIST dataset. Here we measure the RSRE between the exact repulsive term and the interpolation-based approximation It is expected the error will increase for a wider dataset with the same interpolation parameters. So we will change our interpolation parameters according to the difference  $width = \max(y_i) - \min(y_i)$ . More specifically we set the number of interpolation points in each interval to 3 or 5 and the number of interpolation intervals to:

$$\max(20, \lceil width \rceil).$$

In the graph below we see the Repulsive error graph for 5 interpolation points per interval.



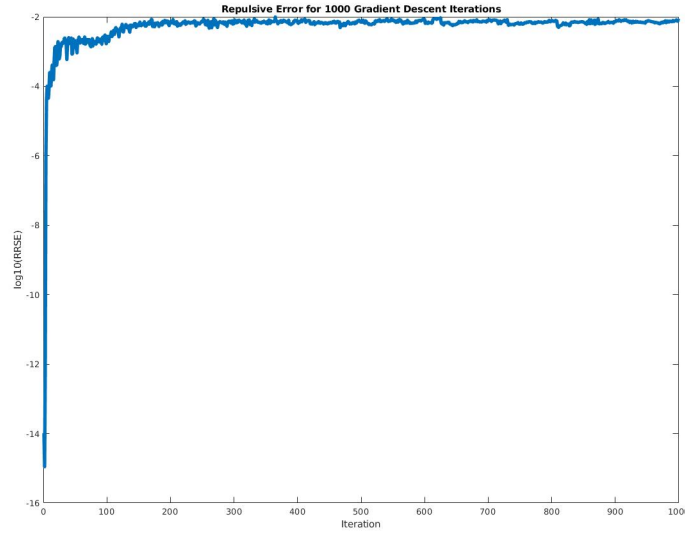


Figure 5: RSRE for 1000 gradient descent iterations 5 points per interval.

Next we graph the error for different interpolation parameters at the 800 iteration, an iteration that the width of the embedding has essentially stop growing. We see that we can reach great precision for 3 and 5 interpolation points per interval and almost exact approximation for 10 points, even though choosing very large numbers can result to time consuming execution.

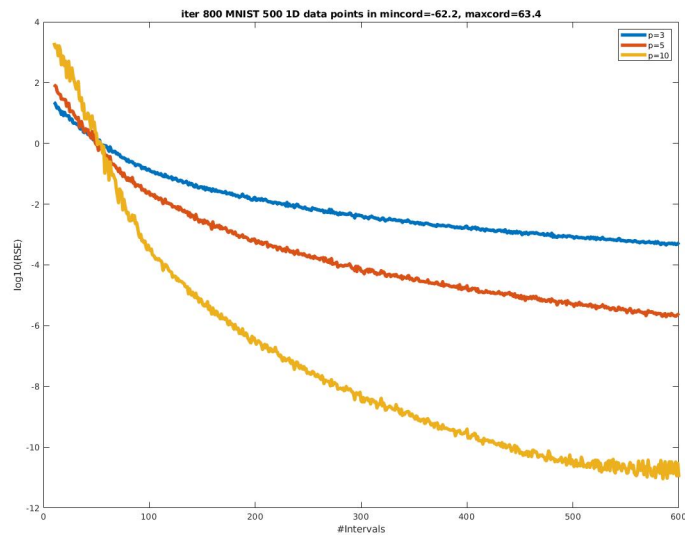


Figure 6: RSRE for the 800 iteration of 1D embedding MNIST using different parameters.

Next comparing the 2 embeddings. We see that there is no essential difference as they capture the same characteristics of the dataset. Clustering nicely same labeled points and having almost the same relative position between clusters.

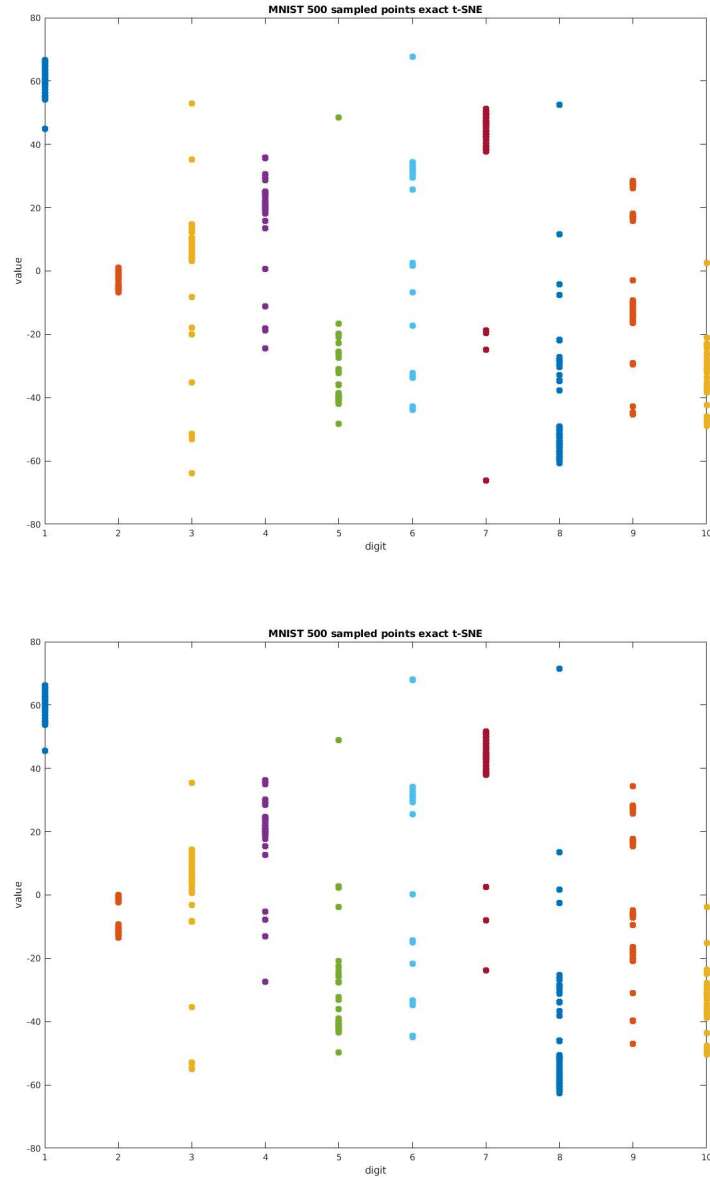


Figure 7: Exact (top) and approximate (bellow) embeddings.

### 4.2.2 2 Dimensions

Now we do the same experiments for 2 dimensional embeddings for 2000 points of the MNIST dataset. Here we use a different rule for updating the number of intervals. Specifically:

$$N_{int} = \max(20, 1.4 \cdot \lceil width \rceil) \text{ where } width = \max_{i,k}(y_i(k)) - \min_{i,k}(y_i(k))$$

Keeping in mind that in two dimension we would have lower accuracy for the same interpolation parameters and width. Below we see the same graph for RSRE using 5 interpolation points per interval ( 25 per box) and the above rule for computing  $N_{int}$ .

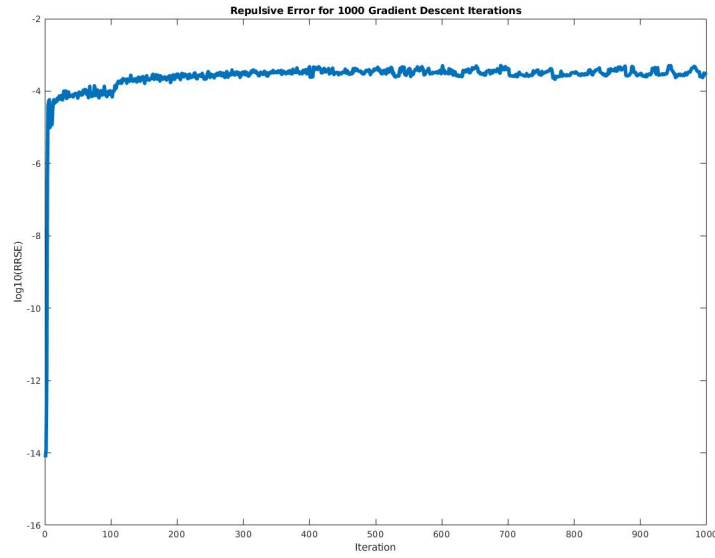


Figure 8: RSRE of Repulsive Forces for 1000 gradient descent iterations 5 points per interval.

We see that we have great accuracy properties. Keep in mind that the 1.4 constant makes for quite an intensive computation, lowering it decreases the accuracy and the time you need to compute the embedding. Next we graph the error of computing the repulsive forces for an iteration for different interpolation parameters.

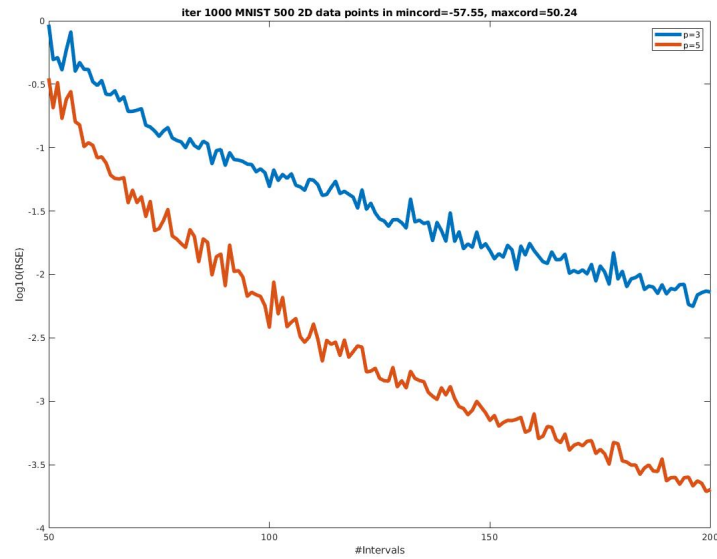


Figure 9: RSRE of Repulsive Forces for the 1000 iteration of 2D embedding MNIST using different parameters.

Also we see the attractive and total gradient RSRE for using 60 nearest neighbours and 30 perplexity.

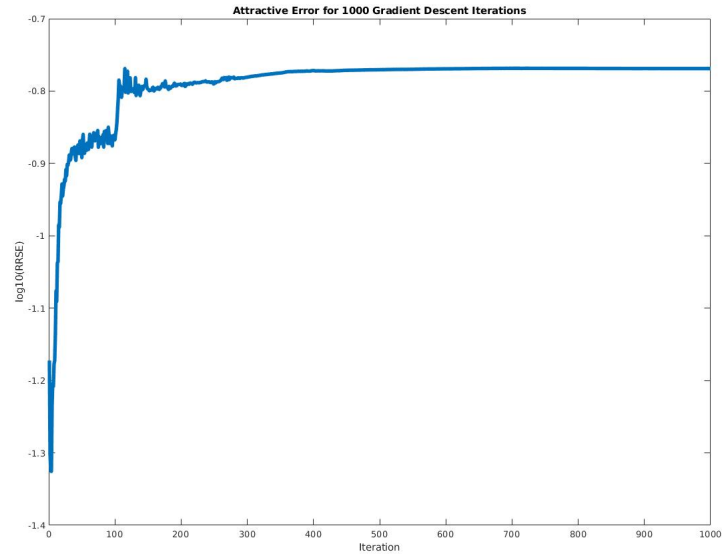


Figure 10: RSRE of the Attractive Forces term for 1000 iterations of 2D embedding MNIST.

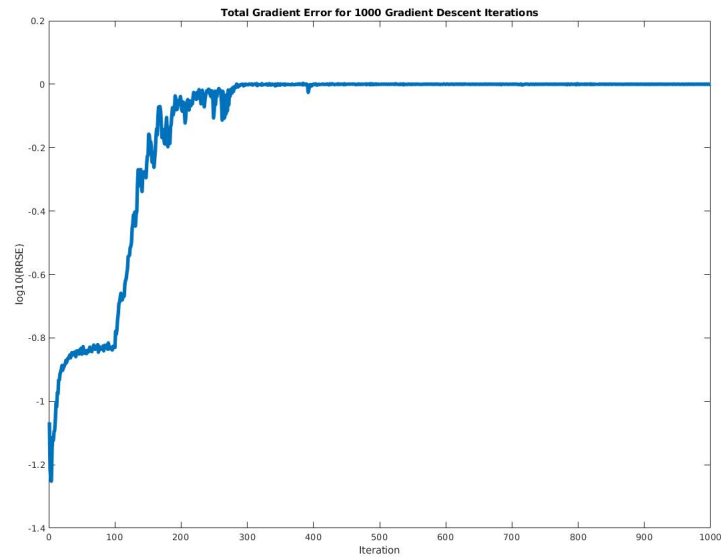


Figure 11: RSRE of Gradient approximation for the 1000 iteration of 2D embedding MNIST.

Finally we need to compare the output embeddings using the exact and approximate methods for the same initial conditions. From the images bellow we see that the embeddings are extremely visually similar. Except the 5 and 6 labeled points.

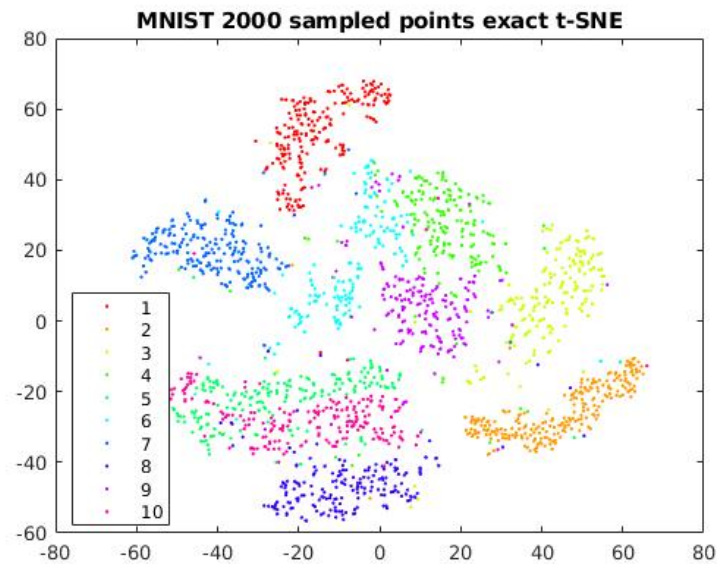


Figure 12: Exact embedding for 2000 MNIST points.

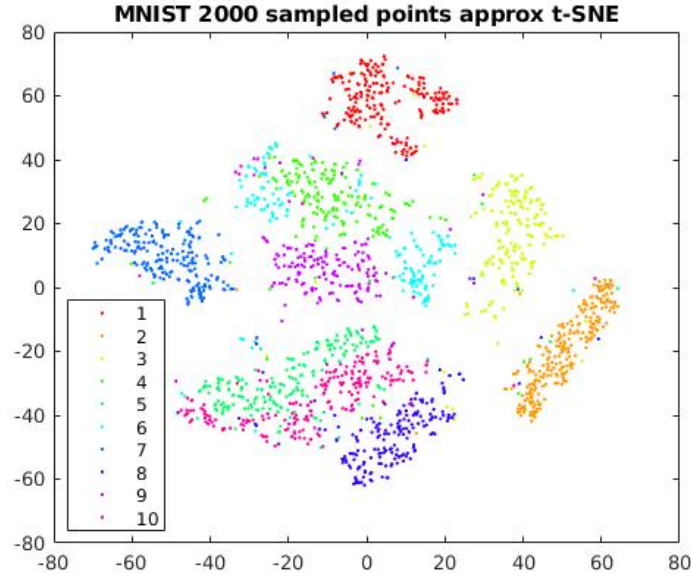


Figure 13: Approximate embedding for 2000 MINST points.

Ending we note that we as the width of the dataset gets larger we stop growing  $N_{int}$  to speed up the computation and use less memory. Below we see the embedding off 60 thousand MNIST points in our implementation and in Flt-SNE.

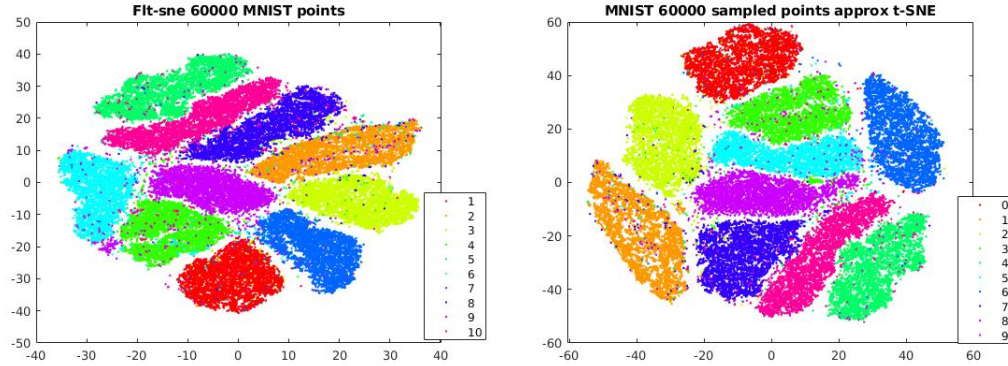


Figure 14: Embedding 60000 MNIST points left Fl-tsne right ours

#### 4.2.3 3 Dimensions

Here we have the challenge that the kernel matrix we need to do the FFT accelerated method is of size

$$2^3 \cdot (p \cdot N_{int})^3.$$

And so we can run out of memory. Using the matlab limit for defining array we can have an array of maximum size about  $2^{24}$ . We need to incorporate this limit to the update rule for the number of intervals. So as an update rule for 3 points per interval we use:

$$N_{int} = \min(\text{limit}, \max(14, \lceil \text{width} \rceil)) \text{ where limit} = 35.$$

Below we present the RSRE graphs. Notice that the Repulsive term error is increasing because of our memory limitations. Also it can saturate as the width of the dataset increases very slowly after a high number of iterations. This problem leads us to either output the embedding when we pass a specified accuracy threshold either continue use our approximate gradient checking if the value of the objective function decreases or normalize the data. In all cases we don't get a accurate t-SNE embedding but for many datasets we have a good approximation of the final embedding when we stop iteration a few hundred iterations after the exaggeration period. An other solution is to use the non-padded convolution with memory  $(p \cdot N_{int})^3$  this leads to an update rule

$$N_{int} = \min(\text{limit}, \max(14, \lceil \text{width} \rceil)) \text{ where limit} = 60$$

and thus allows safely approximate twice as wide embeddings.

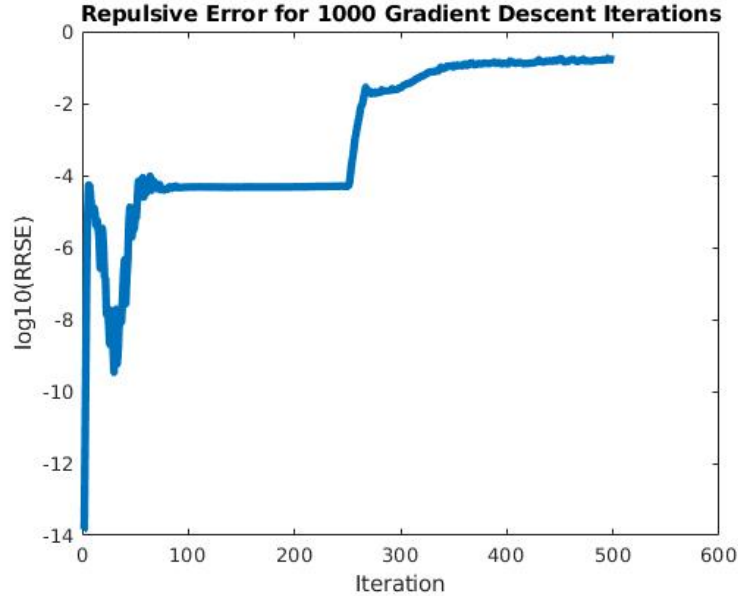


Figure 15: RSRE of Repulsive Forces for 300 gradient descent iterations using padded convolution.

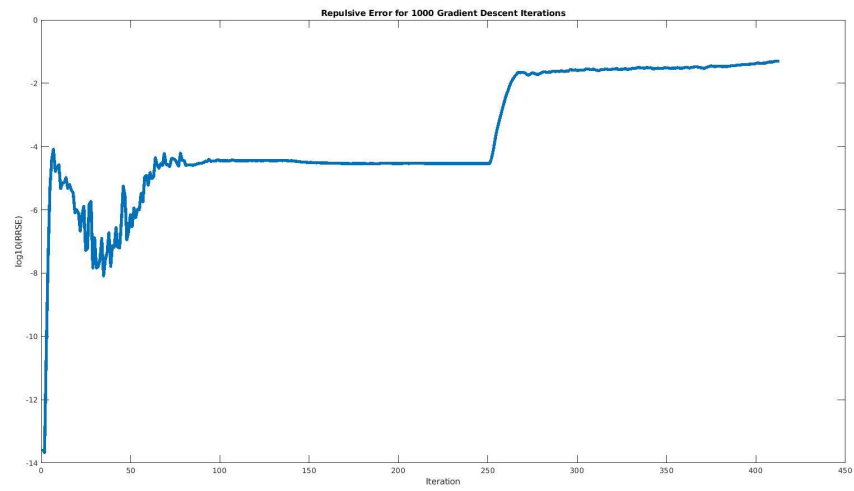


Figure 16: RSRE of Repulsive Forces for 400 iterations using non-padded convolution.

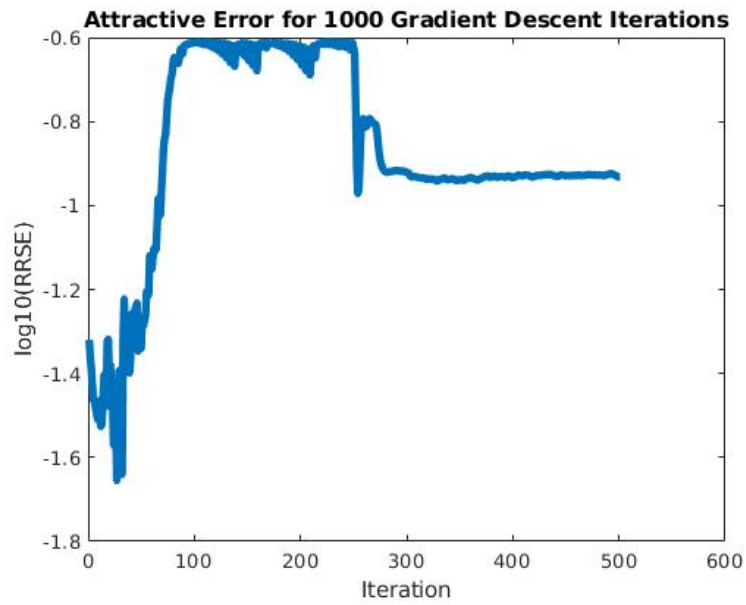


Figure 17: RSRE of Attractive Forces for 300 gradient descent iterations 60 nearest neighbours.



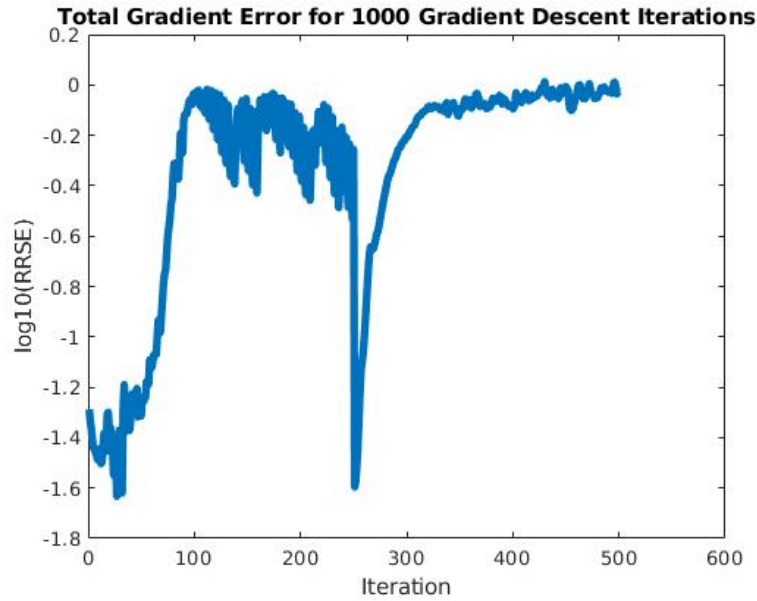


Figure 18: RSRE of the Gradient approx for 300 gradient descent iterations.

Finally we compare the resulting exact and approximate embeddings for the same initial conditions. We observe that the embeddings are very similar.

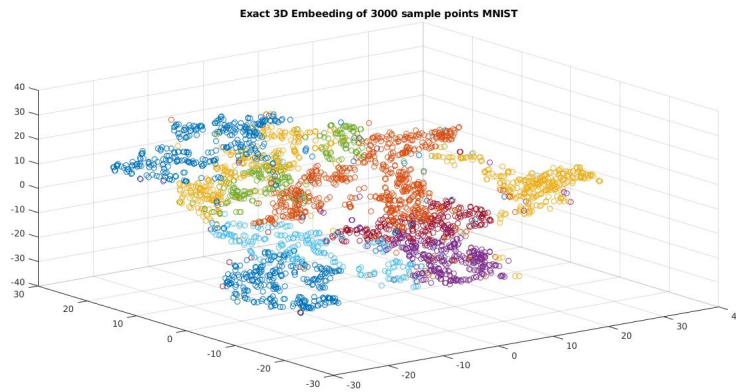


Figure 19: Exact 3D embedding.

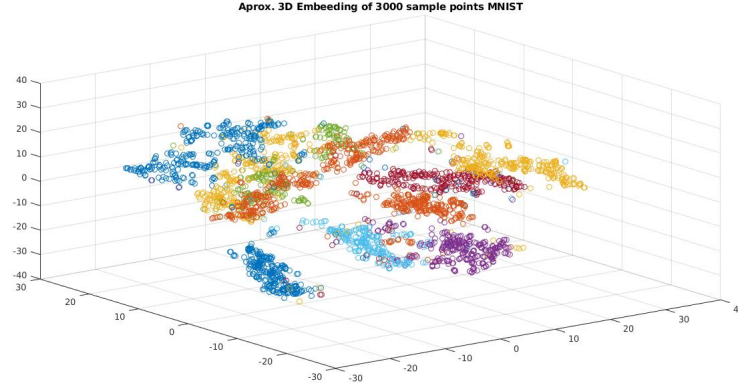


Figure 20: Approximate 3D embedding.

#### 4.2.4 4 Dimensions

In this case we must define an array of size  $2^4 \cdot (p \cdot Nint)^4$  in order to do the Repulsive approximation. And same issues as for the 3 dimensional case apply. So we choose the update rule for 3 points per interval:

$$Nint = \min(\text{limit}, \max(14, \lceil \text{width} \rceil)) \text{ where limit} = 9.$$

Here we see the error graphs for 400 gradient descent iterations on 2000 MNIST points. Again using non-padded convolution we can increase the limit variable to 15 allowing us to a feasible approximation.

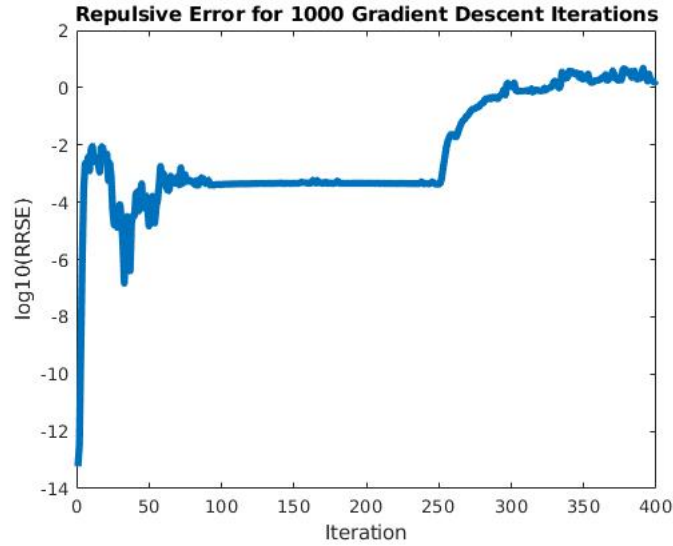


Figure 21: Repulsive Forces RSRE graph.

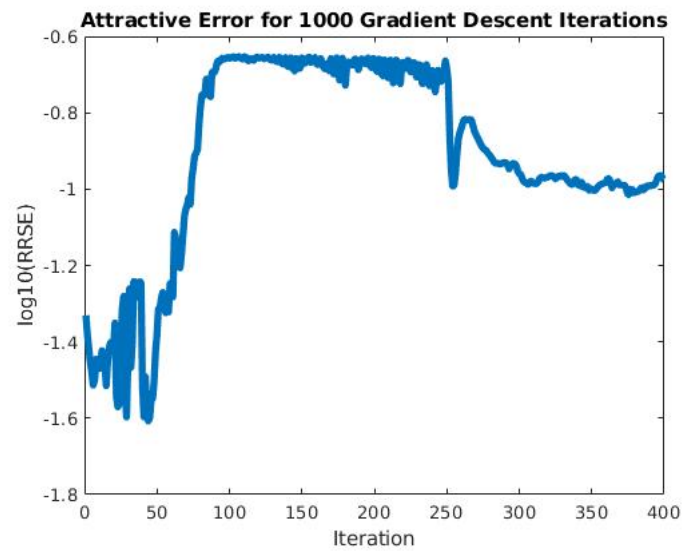


Figure 22: Attractive Forces RSRE graph.

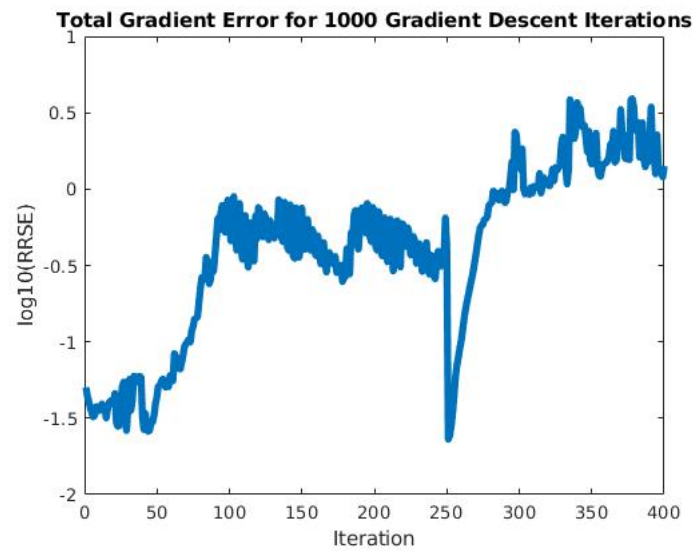


Figure 23: Gradient Error graph.

### 4.3 GPU-Implementation thoughts

## 5 FIt-SNE

## 6 SG-SNE

## 7 Out of core PCA

## 8 Method Comparisons

## 9 GPU implementation

## 10 Data translocations and Memory Hierarchy