

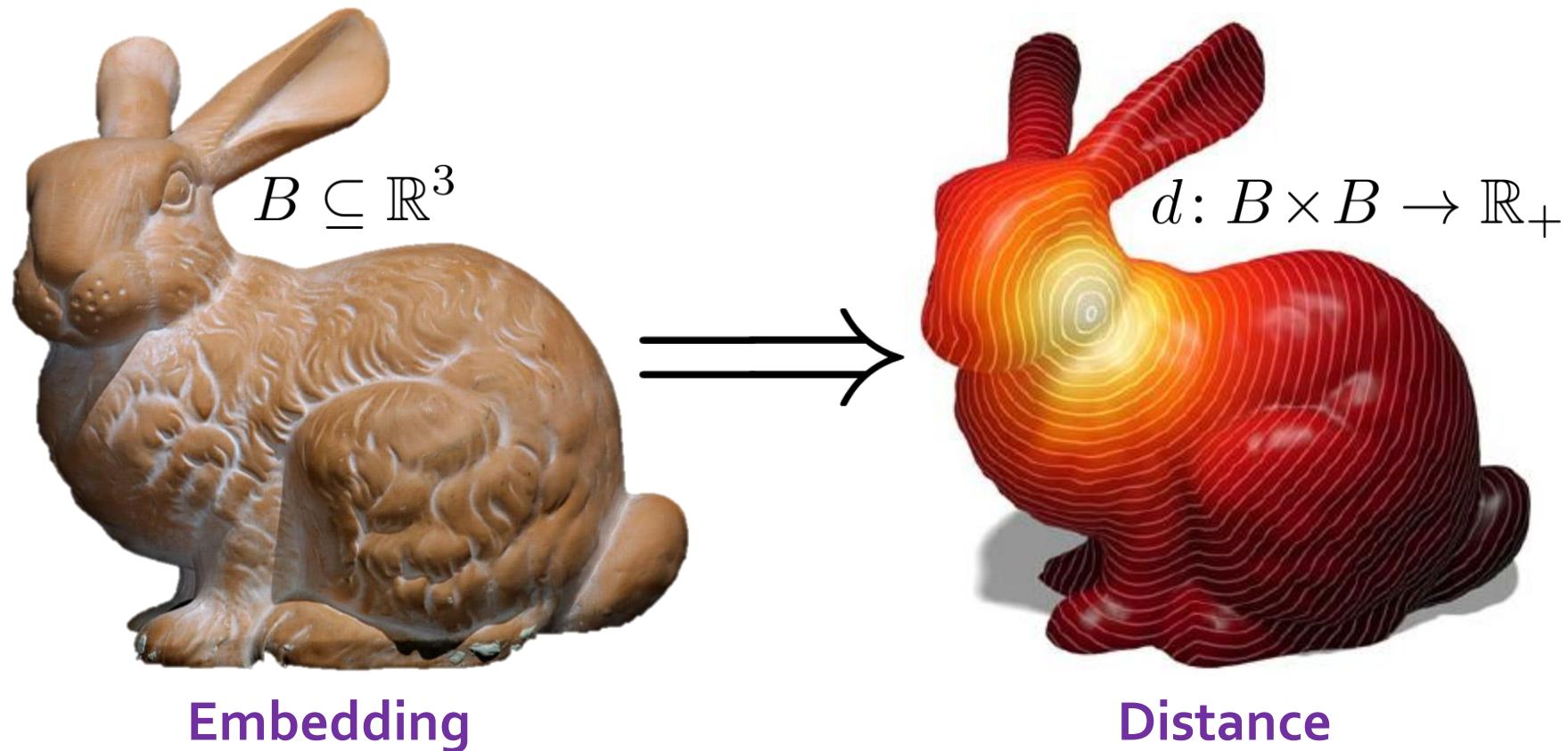


Embedding

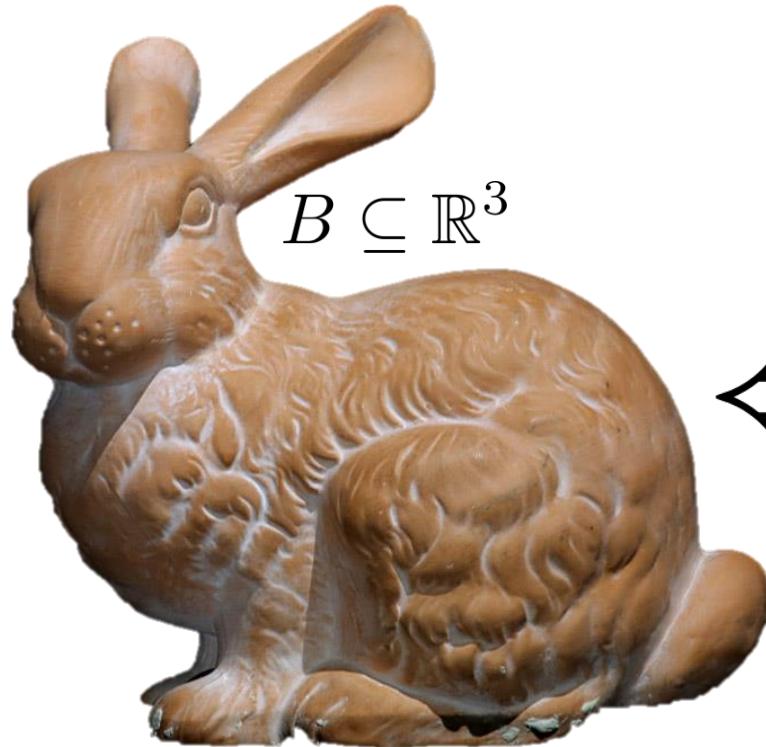
Justin Solomon
MLSS 2019



Forward Problem

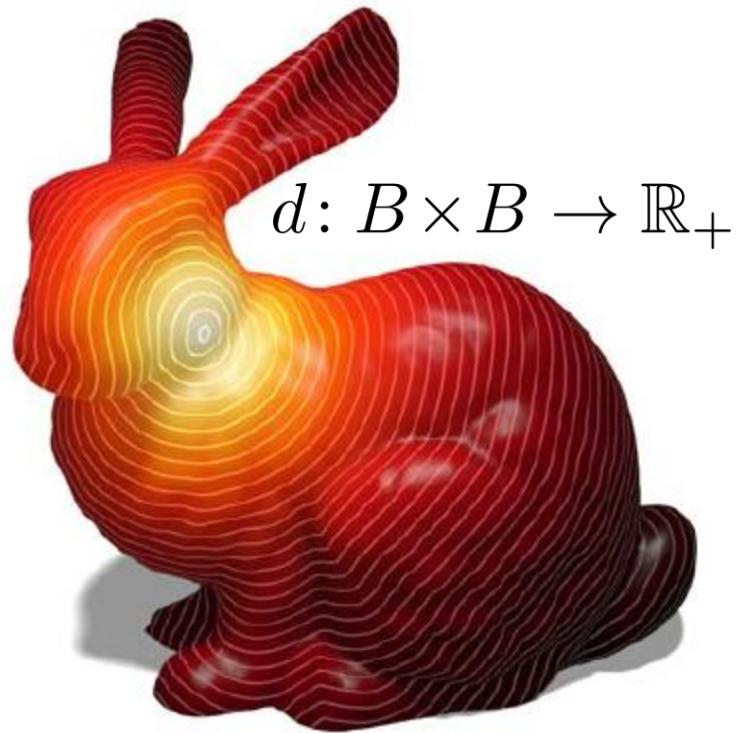
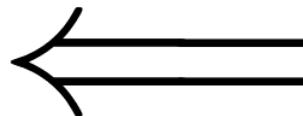


Inverse Problem



$$B \subseteq \mathbb{R}^3$$

Embedding



$$d: B \times B \rightarrow \mathbb{R}_+$$

Distance

Extrinsic vs. Intrinsic Embedding

“Swiss roll”



Extrinsic: 3D



Intrinsic: 2D

Many Names for Similar Tasks

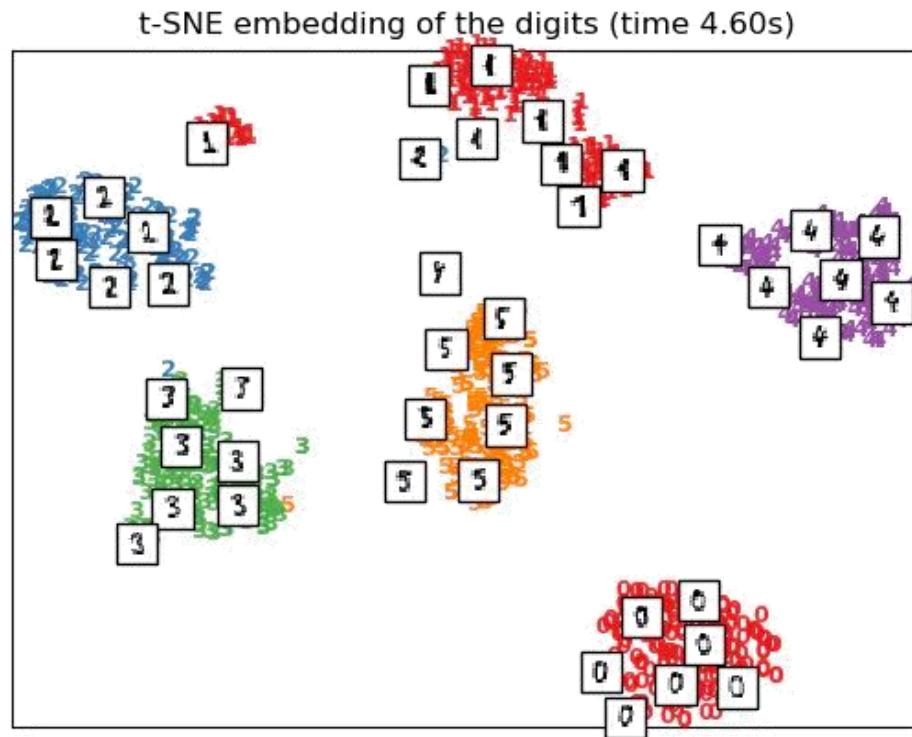
- Dimensionality reduction
 - Embedding
- Parameterization
- Manifold learning

...



Why bother
embedding?

Why Embed?



https://scikit-learn.org/stable/auto_examples/manifold/plot_lle_digits.html

Visualize high-dimensional data

Why Embed?

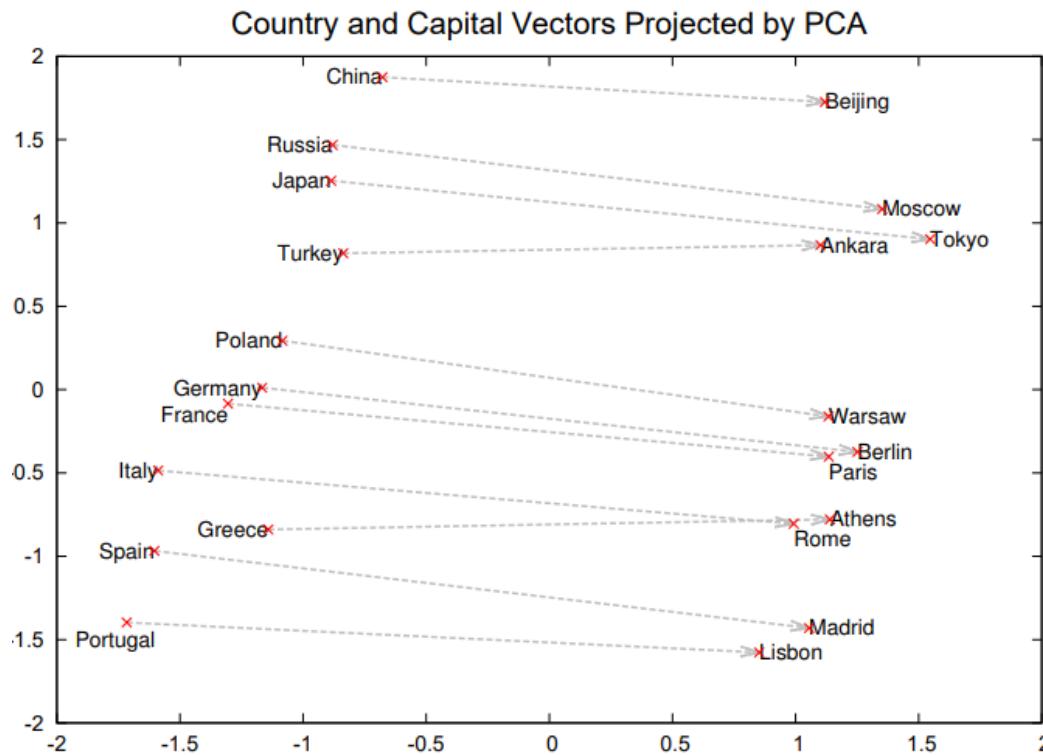


Image from "Distributed Representations of Words and Phrases and their Compositionality" (Mikolov et al.)

Expose latent structure

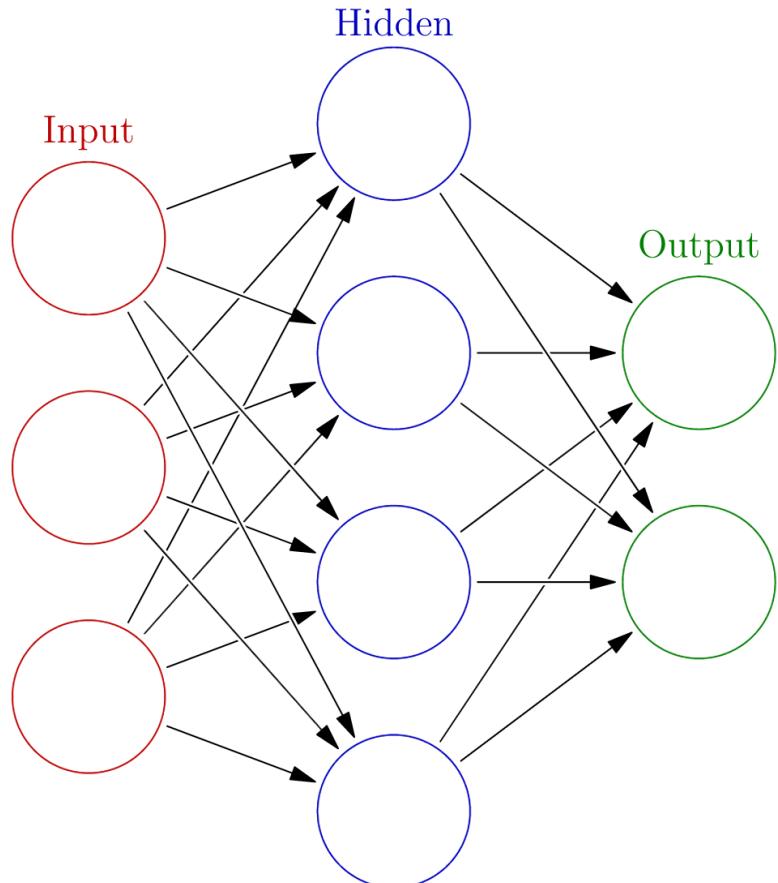
Why Embed?

Naïve nearest neighbor search:

$$O(nd)$$

Accelerates algorithms

Why Embed?



Neural network
layers **embed** your
data into \mathbb{R}^n !

Image from Wikipedia, "Neural Network"

We're doing it anyway!

Basic Task

Given pairwise distances
extract an embedding.

Is it always possible?
What dimensionality?

Metric Space

Ordered pair (M, d) where M is a set and $d: M \times M \rightarrow \mathbb{R}$ satisfies

$$d(x, y) \geq 0$$

$$d(x, y) = 0 \iff x = y$$

$$d(x, y) = d(y, x)$$

$$d(x, z) \leq d(x, y) + d(y, z)$$

$$\forall x, y, z \in M$$

Many Examples of Metric Spaces

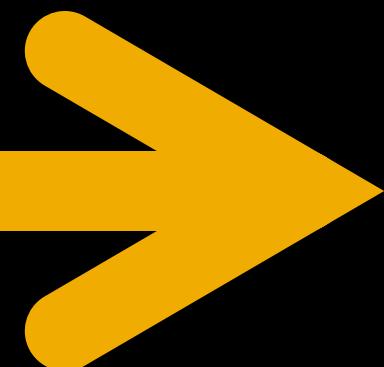
$$\mathbb{R}^n, d(x, y) := \|x - y\|_p$$

$$S \subset \mathbb{R}^3, d(x, y) := \text{geodesic}$$

$$C^\infty(\mathbb{R}), d(f, g)^2 := \int_{\mathbb{R}} (f(x) - g(x))^2 dx$$

Isometry [ahy-som-i-tree]:

A map between metric spaces
that preserves pairwise
distances.





Can you **always** embed
a metric space
isometrically in \mathbb{R}^n ?



Can you always embed
a **finite** metric space
isometrically in \mathbb{R}^n ?

Disappointing Example

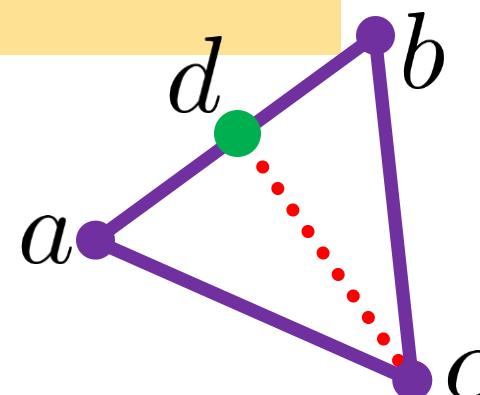
$$X := \{a, b, c, d\}$$

$$d(a, d) = d(b, d) = 1$$

$$d(a, b) = d(a, c) = d(b, c) = 2$$

$$d(c, d) = 1.5$$

Cannot be embedded in Euclidean space!



Contrasting Example

$$\begin{aligned}\ell_\infty(\mathbb{R}^n) &:= (\mathbb{R}^n, \|\cdot\|_\infty) \\ \|\mathbf{x}\|_\infty &:= \max_k |\mathbf{x}_k|\end{aligned}$$

Proposition. Every finite metric space embeds isometrically into $\ell_\infty(\mathbb{R}^n)$ for some n .

Extends to infinite-dimensional spaces!

Fréchet Embedding

Definition 7.3 (Fréchet embedding). Suppose (M, d) is a metric space that $S_1, \dots, S_r \subseteq M$. We define the Fréchet embedding of M with respect to $\{S_1, \dots, S_r\}$ to be the map $\phi : M \rightarrow \mathbb{R}^r$ given by

$$\phi(x) := (d(x, S_1), d(x, S_2), \dots, d(x, S_r)), \quad (7.2)$$

where $d(x, S) := \min_{y \in S} d(x, y)$.

Approximate Embedding

$$\text{expansion}(f) := \max_{x,y} \frac{\mu(f(x), f(y))}{\rho(x, y)}$$

$$\text{contraction}(f) := \max_{x,y} \frac{\rho(x, y)}{\mu(f(x), f(y))}$$

$$\text{distortion}(f) := \text{expansion}(f) \times \text{contraction}(f)$$

Well-Known Result

Proposition 7.2 (Bourgain's Theorem). Suppose (M, d) is a metric space consisting of n points, that is, $|M| = n$. Then, for $p \geq 1$, M embeds into $\ell_p(\mathbb{R}^m)$ with $O(\log n)$ distortion, where $m = O(\log^2 n)$. Matousek improved the distortion bound to $\log n/p$ [14].

```
m := 576 log n
for j = 1 to log n do          /* levels of density */
    for i = 1 to m do          /* repeat for high probability */
        choose set Sij by sampling each node in X
        independently with probability 2-j
    end
end
fij(x) := d(x, Sij)
f(x) := ⊕j=1log n ⊕i=1m fij(x)
```

Uses Fréchet
embedding

Euclidean Problem

Given:

$$P_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|_2^2, P \in \mathbb{R}^{n \times n}$$

Reconstruct:

$$\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^m$$

Alternative notation:

$$X \in \mathbb{R}^{m \times n}$$

Gram Matrix [gram mey-triks]:

A matrix of inner products

$$X^T X$$



Classical Multidimensional Scaling

1. Double centering: $G := -\frac{1}{2}J^\top P J$
Centering matrix $J := I_{n \times n} - \frac{1}{n}\mathbf{1}\mathbf{1}^\top$
2. Find m largest eigenvalues/eigenvectors
 $G = Q\Lambda Q^\top$
3. $\bar{X} = \sqrt{\Lambda}Q^\top$

Extension: Landmark MDS

“MDS”

Stress Majorization

$$\min_X \sum_{ij} (D_{0ij} - \|\mathbf{x}_i - \mathbf{x}_j\|_2)^2$$

Nonconvex!

SMACOF:
Scaling by **M**ajorizing a **C**omplicated **F**unction

de Leeuw, J. (1977), "Applications of convex analysis to multidimensional scaling" *Recent developments in statistics*, 133–145.

SMACOF Potential Terms

$$\min_X \sum_{ij} (D_{0ij} - \|\mathbf{x}_i - \mathbf{x}_j\|_2)^2$$

$$\sum_{ij} (D_{0ij})^2 = \text{const.}$$

$$J = I_{n \times n} - \frac{1}{n} \mathbf{1}\mathbf{1}^\top$$

$$\sum_{ij} \|\mathbf{x}_i - \mathbf{x}_j\|_2^2 = \text{tr}(X V X^\top), \text{ where } V = 2nJ$$

$$-2 \sum_{ij} D_{0ij} \|\mathbf{x}_i - \mathbf{x}_j\|_2 = -2 \text{tr}(X B(X) X^\top)$$

$$\text{where } B_{ij}(X) := \begin{cases} -\frac{2D_{0ij}}{\|\mathbf{x}_i - \mathbf{x}_j\|_2} & \text{if } \mathbf{x}_i \neq \mathbf{x}_j, i \neq j \\ 0 & \text{if } \mathbf{x}_i = \mathbf{x}_j, i \neq j \\ -\sum_{j \neq i} B_{ij} & \text{if } i = j \end{cases}$$

SMACOF Lemma

Lemma. Define

Then,

with equality exactly when $X \propto Z$.

Lemma. Define

$$\tau(X, Z) := \text{const.} + \text{tr}(X V X^\top) - 2\text{tr}(X B(Z) Z^\top)$$

Then,

$$\tau(X, X) \leq \tau(X, Z) \quad \forall Z$$

with equality exactly when $X \propto Z$.

Proof using Cauchy-Schwarz.

SMACOF: Single Step

$$X^{k+1} \leftarrow \min_X \tau(X, X^k)$$

$$\tau(X, Z) := \text{const.} + \text{tr}(X V X^\top) - 2\text{tr}(X B(Z) Z^\top)$$

$$\implies 0 = \nabla_X [\tau(X, X^k)]$$

$$= 2XV - 2X^k B(X^k)$$

$$\implies X^{k+1} = X^k B(X^k) \left(I_{n \times n} - \frac{\mathbf{1}\mathbf{1}^\top}{n} \right)$$

**Majorization-Minimization
(MM) algorithm**

Objective convergence:
 $\tau(X^{k+1}, X^{k+1}) \leq \tau(X^k, X^k)$

Visualization

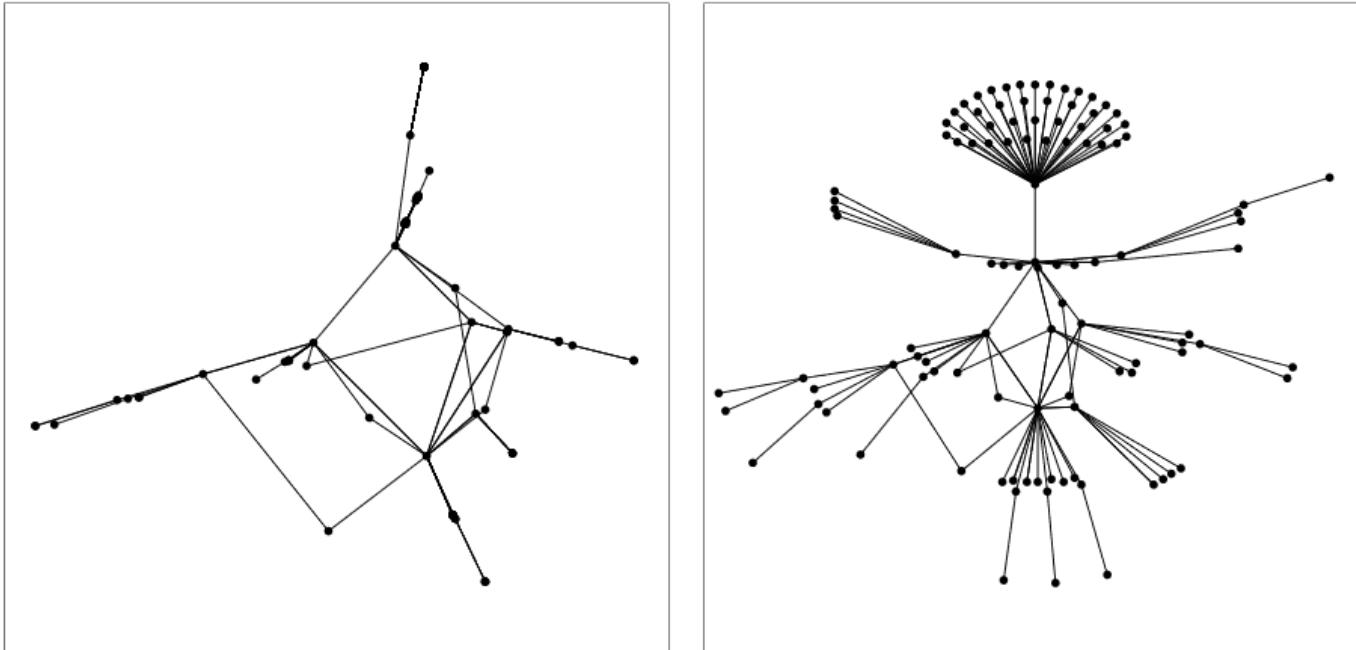


Figure 9: A Telephone Call Graph, Layed Out in 2-D. Left: classical scaling (Stress=0.34); right: distance scaling (Stress=0.23). The nodes represent telephone numbers, the edges represent the existence of a call between two telephone numbers in a given time period.

Also Useful for 3D Tasks

DOI: 10.1111/cgf.12558

EUROGRAPHICS 2015 / O. Sorkine-Hornung and M. Wimmer
(Guest Editors)

Volume 34 (2015), Number 2

Shape-from-Operator: Recovering Shapes from Intrinsic Operators

Davide Boscaini, Davide Eynard, Drosos Kourounis, and Michael M. Bronstein

Università della Svizzera Italiana (USI), Lugano, Switzerland

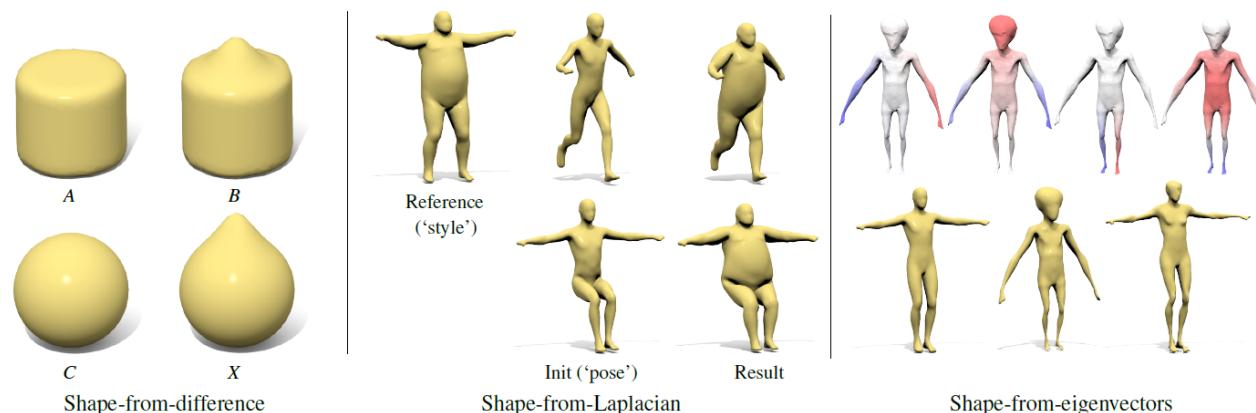
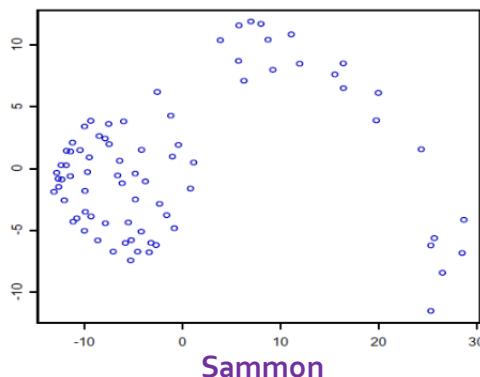
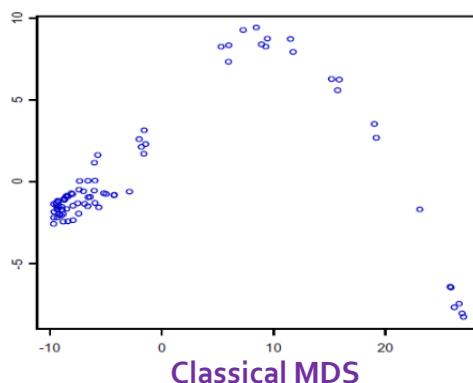


Figure 1: Examples of three different shape-from-operator problems considered in the paper. Left: shape analogy synthesis as shape-from-difference operator problem (shape X is synthesized such that the intrinsic difference operator between C,X is as close as possible to the difference between A,B). Center: style transfer as shape-from-Laplacian problem. The Laplacian of the reference shape (from the left) is used as initial value for the heat flow (red/blue) on the target shape. Right: shape-from-eigenvectors.

Related Method

$$\min_X \sum_{ij} \frac{(D_{0ij} - \|\mathbf{x}_i - \mathbf{x}_j\|_2)^2}{D_{0ij}}$$

Cares more about preserving small distances



“Sammon mapping”

Sammon (1969). “A nonlinear mapping for data structure analysis.” IEEE Transactions on Computers 18.

t-SNE

1. Construct joint distribution using high-dimensional Gaussians

$$p_{j|i} := \frac{\exp(-\|\mathbf{x}_i - \mathbf{x}_j\|_2^2/2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|\mathbf{x}_i - \mathbf{x}_k\|_2^2/2\sigma_i^2)}$$

$$p_{ij} = \frac{p_{j|i} + p_{i|j}}{2N}$$

2. Optimize KL from using low-dimensional heavy-tailed Student t-distribution

$$q_{ij} := \frac{(1 + \|\mathbf{y}_i - \mathbf{y}_j\|_2^2)^{-1}}{\sum_{k \neq l} (1 + \|\mathbf{y}_k - \mathbf{y}_l\|_2^2)^{-1}}$$

$$\min_{\{\mathbf{y}_i\}} \text{KL}(P\|Q) = \sum_{i \neq j} p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

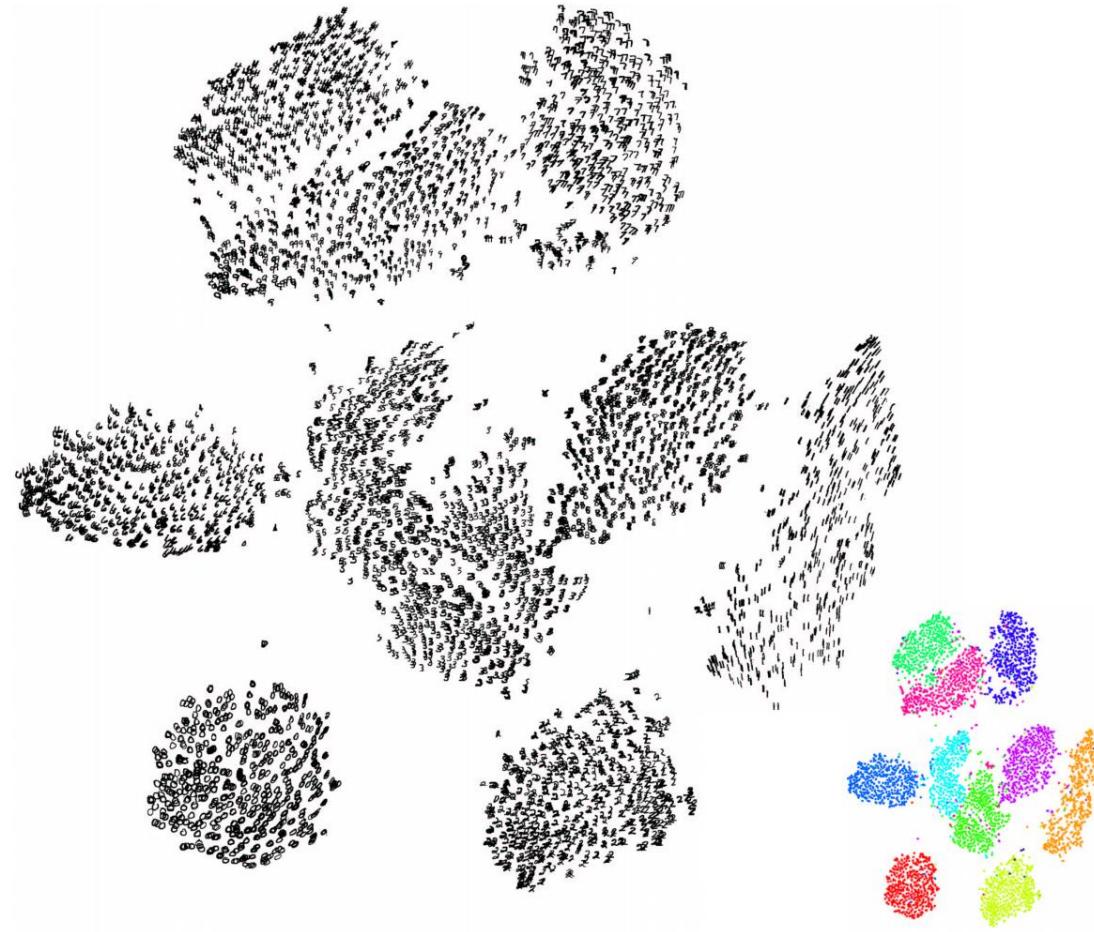
Designed to allow dissimilar objects to move far apart

"Mismatched Tails can Compensate for Mismatched Dimensionalities"

van der Maaten & Hinton 2008

t-distributed stochastic neighbor embedding

Typical t-SNE Output: MNIST



Ignores large distances

Nice Web Page

How to Use t-SNE Effectively

Although extremely useful for visualizing high-dimensional data, t-SNE plots can sometimes be mysterious or misleading. By exploring how it behaves in simple cases, we can learn to use it more effectively.

The screenshot shows a web page titled "How to Use t-SNE Effectively". The main content features a large, colorful t-SNE plot of data points arranged in a grid-like pattern. Below the plot is a control panel with several components:

- A grid of 12 smaller t-SNE plots, each showing a different configuration of data points.
- A blue circular button with a white play icon and a blue circular button with a white 'C' icon.
- A text input field labeled "Step 280".
- A slider labeled "Points Per Side 20".
- A text area with the following description:

A square grid with equal spacing between points. Try convergence at different sizes.
- A URL at the bottom: <https://distill.pub/2016/misread-tsne/>
- A small footer note: "Epsilon 5"

Recall:

Extrinsic vs. Intrinsic Embedding

“Swiss roll”



Extrinsic: 3D



Intrinsic: 2D

Intrinsic-to-Extrinsic: Theory

Theorem 7.1 (Whitney embedding theorem). *Any smooth, real k -dimensional manifold maps smoothly into \mathbb{R}^{2k} .*

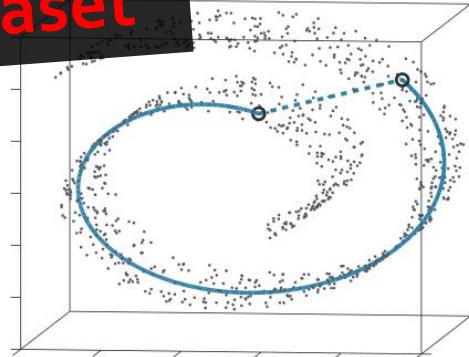
Theorem 7.2 (Nash–Kuiper embedding theorem, simplified). *Any k -dimensional Riemannian manifold admits an isometric, differentiable embedding into \mathbb{R}^{2k} .*



Intrinsic-to-Extrinsic: ISOMAP

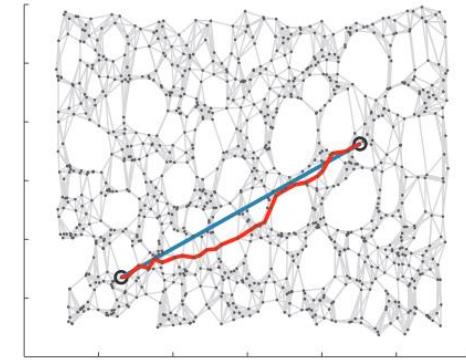
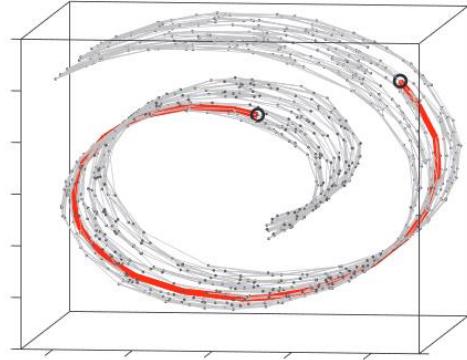
- Construct neighborhood graph
 k -nearest neighbor graph or ε -neighborhood graph
- Compute shortest-path distances
Floyd-Warshall algorithm or Dijkstra

Swiss roll dataset



- Classical MDS

Eigenvalue problem



Tenenbaum, de Silva, Langford.

"A Global Geometric Framework for Nonlinear Dimensionality Reduction." Science (2000).

Floyd-Warshall Algorithm

```
let dist be a |V| × |V| array of minimum distances initialized to ∞ (infinity)
for each vertex v
    dist[v][v] ← 0
for each edge (u,v)
    dist[u][v] ← w(u,v) // the weight of the edge (u,v)
for k from 1 to |V|
    for i from 1 to |V|
        for j from 1 to |V|
            if dist[i][j] > dist[i][k] + dist[k][j]
                dist[i][j] ← dist[i][k] + dist[k][j]
            end if
```

Landmark ISOMAP

- **Construct neighborhood graph**
 k -nearest neighbor graph or ε -neighborhood graph
- **Compute some shortest-path distances**
Dijkstra: $O(kn N \log N)$, n landmarks, N points
 - **MDS on landmarks**
Smaller $n \times n$ problem
 - **Closed-form embedding formula**
 $\delta(x)$ vector of squared distances from x to landmarks

$$\text{Embedding}(x)_i = -\frac{1}{2} \frac{v_i^\top}{\sqrt{\lambda_i}} (\delta(x) - \delta_{\text{average}})$$

Landmark MDS

Locally Linear Embedding (LLE)

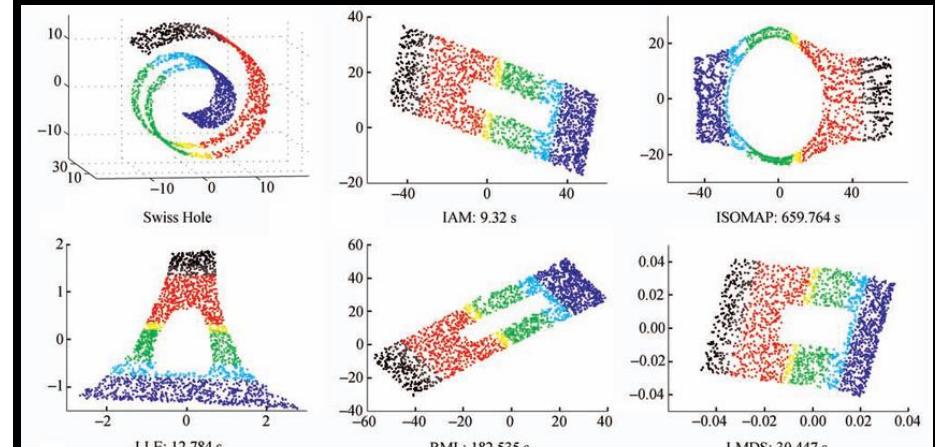
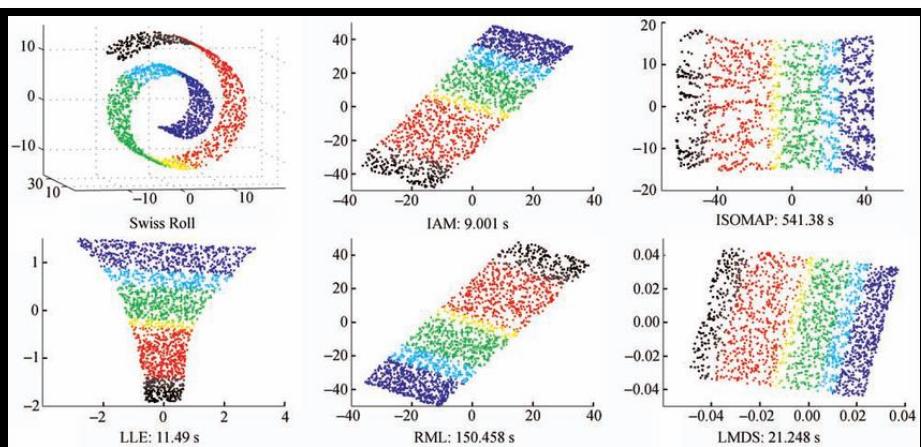
- **Construct neighborhood graph**
 k -nearest neighbor graph or ε -neighborhood graph
- **Analysis step: Compute weights W_{ij}**
$$\min_{\omega^1, \dots, \omega^k} \left\| \mathbf{x}_i - \sum_j \omega^j \mathbf{n}_j \right\|_2$$

subject to $\sum_j \omega^j = 1$
- **Embedding step: Minimum eigenvalue problem**
$$\min_Y \quad \|Y - YW^\top\|_{\text{Fro}}^2$$

subject to $YY^\top = I_{p \times p}$
 $Y\mathbf{1} = \mathbf{0}$

Comparison: ISOMAP vs. LLE

ISOMAP	LLE
Global distances	Local averaging
k -NN graph distances	k -NN graph weighting
Largest eigenvectors	Smallest eigenvectors
Dense matrix	Sparse matrix



Other option:

Diffusion Maps

- **Construct similarity (kernel) matrix**

Example: $K(x, y) := e^{-\|x-y\|^2/\varepsilon}$

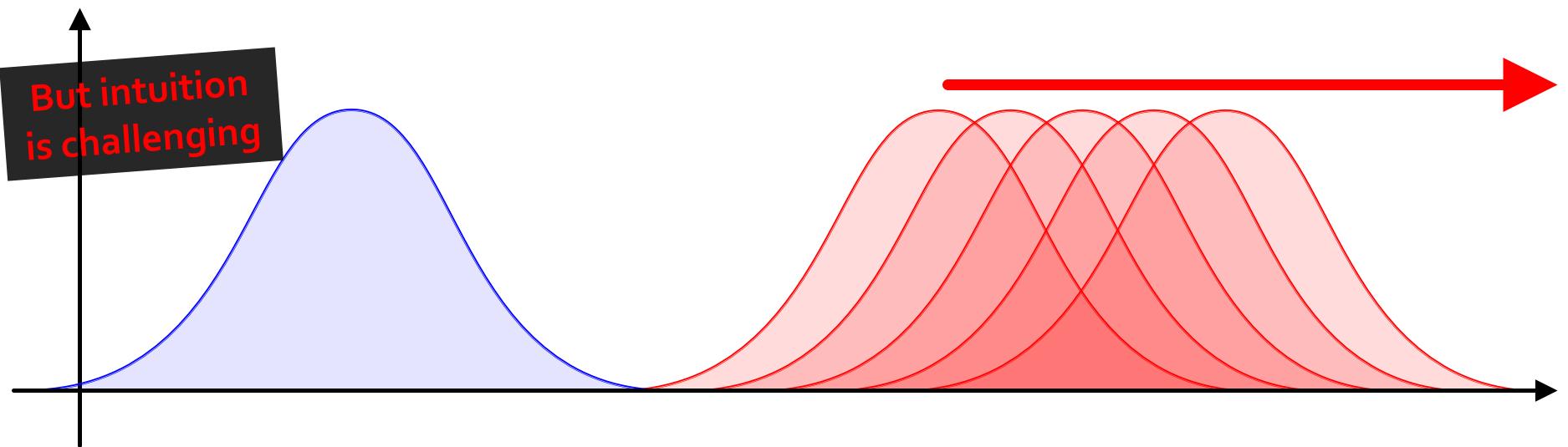
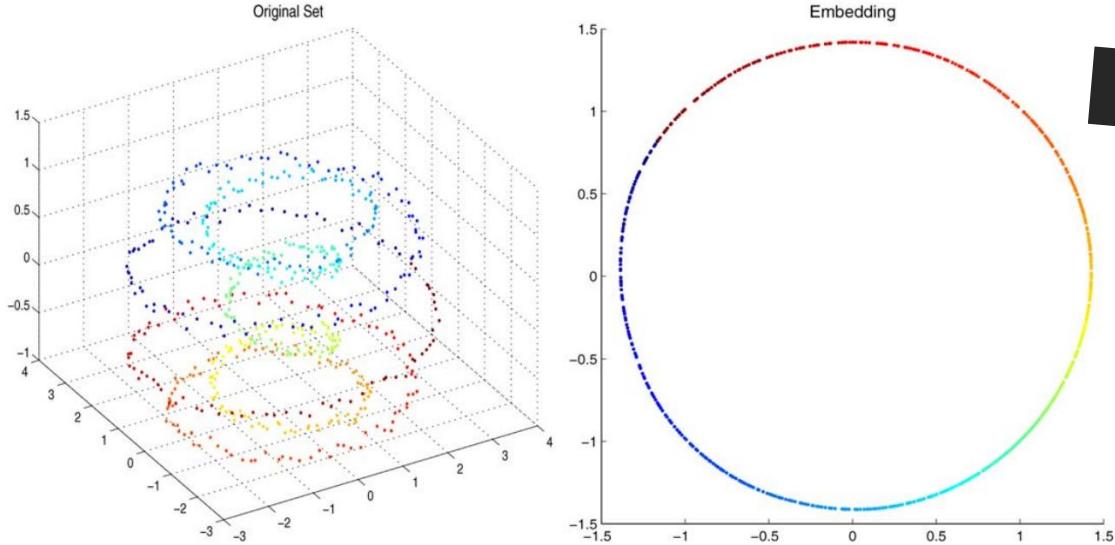
- **Normalize rows**
 $M := D^{-1}K$

- **Embed from k largest eigenvectors**

$$(\lambda_1 \psi_1, \lambda_2 \psi_2, \dots, \lambda_k \psi_k)$$

More when we
discuss Laplacians

Diffusion Distance



Take-Away

Huge zoo
of embedding techniques.

Each with different theoretical properties: Try them all!

But what if the distance matrix is incomplete or noisy?

More General: Metric Nearness

$$\min_{X \in \mathcal{M}_{N \times N}} \|X - D\|_{\text{Fro}}^2$$

`TRIANGLE_FIXING(D, ϵ)`

Input: Input dissimilarity matrix D , tolerance ϵ

Output: $M = \operatorname{argmin}_{X \in \mathcal{M}_N} \|X - D\|_2$.

for $1 \leq i < j < k \leq n$

$(z_{ijk}, z_{jki}, z_{kij}) \leftarrow 0$

for $1 \leq i < j \leq n$

$e_{ij} \leftarrow 0$

$\delta \leftarrow 1 + \epsilon$

while ($\delta > \epsilon$) {convergence test}

foreach triangle (i, j, k)

$b \leftarrow d_{ki} + d_{jk} - d_{ij}$

$\mu \leftarrow \frac{1}{3}(e_{ij} - e_{jk} - e_{ki} - b)$

$\theta \leftarrow \min\{-\mu, z_{ijk}\}$ {Stay within half-space of constraint}

$e_{ij} \leftarrow e_{ij} - \theta, e_{jk} \leftarrow e_{jk} + \theta, e_{ki} \leftarrow e_{ki} + \theta$

$z_{ijk} \leftarrow z_{ijk} - \theta$ {Update correction term}

end foreach

$\delta \leftarrow \text{sum of changes in the } e$

end while

return $M = D + E$

In other words, the vector e is projected orthogonally onto the constraint set $\{e' : e'_{ij} - e'_{jk} - e'_{ki} \leq b_{ijk}\}$. This is tantamount to solving

$$\begin{aligned} \min_{e'} \quad & \frac{1}{2}[(e'_{ij} - e_{ij})^2 + (e'_{jk} - e_{jk})^2 + (e'_{ki} - e_{ki})^2], \\ \text{subject to} \quad & e'_{ij} - e'_{jk} - e'_{ki} = b_{ijk}. \end{aligned} \quad (3.2)$$

It is easy to check that the solution is given by

$$e'_{ij} \leftarrow e_{ij} - \mu_{ijk}, \quad e'_{jk} \leftarrow e_{jk} + \mu_{ijk}, \quad \text{and} \quad e'_{ki} \leftarrow e_{ki} + \mu_{ijk}, \quad (3.3)$$

$$\text{where } \mu_{ijk} = \frac{1}{3}(e_{ij} - e_{jk} - e_{ki} - b_{ijk}) > 0.$$

Iterative
projection

Dhillon, Sra, Tropp. "Triangle Fixing Algorithms for the Metric Nearness Problem." NIPS.

Euclidean Matrix Completion

$$\begin{aligned} \min_G \quad & \|H \circ (\mathcal{D}(G) - D_{\text{input}})\|_{\text{Fro}}^2 \\ \text{s.t. } & G \succeq 0 \end{aligned}$$

Convex program

Alfakih, Khandani, and Wolkowicz. "Solving Euclidean distance matrix completion problems via semidefinite programming." *Comput. Optim. Appl.*, 12 (1999).

Maximum Variance Unfolding

(on the board, if there's time)

Network Embedding

Distributed Representations of Words and Phrases and their Compositionality

Tomas Mikolov
Google Inc.
Mountain View
mikolov@google.com

Ilya Sutskever
Google Inc.
Mountain View
ilyasu@google.com

Kai Chen
Google Inc.
Mountain View
kai@google.com

Greg Corrado
Google Inc.
Mountain View
gcorrado@google.com

Jeffrey Dean
Google Inc.
Mountain View
jeff@google.com

Abstract

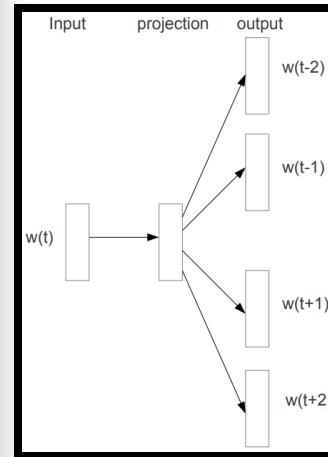
The recently introduced continuous Skip-gram model is an efficient method for learning high-quality distributed vector representations that capture a large number of precise syntactic and semantic word relationships. In this paper we present several extensions that improve both the quality of the vectors and the training speed. By subsampling of the frequent words we obtain significant speedup and also learn more regular word representations. We also describe a simple alternative to the hierarchical softmax called negative sampling.

An inherent limitation of word representations is their indifference to word order and their inability to represent idiomatic phrases. For example, the meanings of “Canada” and “Air” cannot be easily combined to obtain “Air Canada”. Motivated by this example, we present a simple method for finding phrases in text, and show that learning good vector representations for millions of phrases is possible.

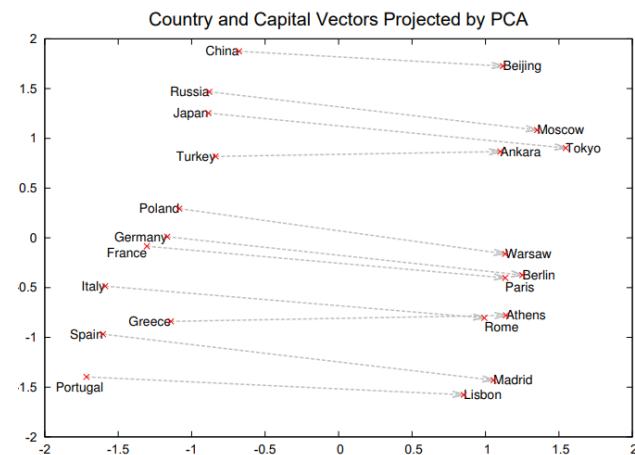
1 Introduction

Well-known example: Word2Vec

Download the embedding!



Skip-gram architecture:
Predict neighborhood of a word



Challenging Computational Problems

- Is my data **embeddable**?
- Can you compute intrinsic **dimensionality**?
- Are two metric spaces **isometric**?
- How **similar** are two metric spaces?
- What is the **average** of two metric spaces?
- Can I embed into **non-Euclidean** spaces?

NP-Hardness Result

Robust Euclidean Embedding

Lawrence Cayton
Sanjoy Dasgupta

Department of Computer Science and Engineering, University of California, San Diego
9500 Gilman Dr. La Jolla, CA 92093

LCAYTON@CS.UCSD.EDU
DASGUPTA@CS.UCSD.EDU

Abstract

We derive a robust Euclidean embedding procedure based on semidefinite programming that may be used in place of the popular classical multidimensional scaling (cMDS) algorithm. We motivate this algorithm by arguing that cMDS is not particularly robust and has several other deficiencies. General-purpose semidefinite programming solvers are too memory intensive for medium to large sized applications, so we also describe a fast subgradient-based implementation of the robust algorithm. Additionally, since cMDS is often used for dimensionality reduction, we provide an in-depth look at reducing dimensionality with embedding procedures. In particular, we show that it is NP-hard to find optimal low-dimensional embeddings under a variety of cost functions.

choice for embedding seems to be classical scaling (cMDS). Its popularity is due to its being relatively fast, parameter-free, and optimal for its cost function. In this work, we look carefully at the algorithm and argue that cMDS has some problematic features as well. We argue that the cost function is not conceptually awkward.

We propose a robust alternative to classical Euclidean embedding (REE), that retains the desirable features of cMDS, but avoids its pitfalls. We show that the global minimum of the REE cost function can be found using a semidefinite program (SDP). Though this is hard for standard SDP-solvers can only manage the program for around 100 points. So the used on more reasonably sized data sets a subgradient-based implementation of

ℓ_1 EUCLIDEAN EMBEDDING

Input: A dissimilarity matrix $D = (d_{ij})$.

Output: An embedding into the line: $x_1, x_2, \dots \in \mathbf{R}$

Goal: Minimize $\sum_{i,j} |d_{ij} - \|x_i - x_j\||$.

We show that this problem is NP-hard by reducing from a variant of not-all-equal 3SAT.

The hardness result can be extended to distortion functions of the form $\sum_{i,j} g(f(d_{ij}) - f(\|x_i - x_j\|))$. We assume that f, g are

1. symmetric;
2. monotonically increasing in the absolute values of their arguments;
3. Lipschitz on $[0, 1]$ with constant λ_U , that is, for $x, y \in [0, 1]$, $|f(x) - f(y)| \leq \lambda_U|x - y|$; and
4. similarly lower-bounded: for some $\lambda_L > 0$, for any $x, y \in [0, 1]$, $|f(x) - f(y)| \geq \lambda_L|x - y| \max\{x, y\}$.

Notice that $f(x), g(x) \in \{x, x^2\}$ satisfy these conditions with $\lambda_U = 2, \lambda_L = 1$, meaning that $\|D - D^*\|_1$ and $\|D - D^*\|_2$ are both hard to minimize over one-dimensional embeddings.

Dimensionality reduction is an important application of MDS, though it will not be NP-hard if



Embedding

Justin Solomon
MLSS 2019

