A New Technique for Dimension Reduction for Data Visualization



Joint work with Yingfan Wang, Haiyang Huang, and Yaron Shaposhnik

[PDF] Visualizing data using t-SNE. L Van der Maaten, G Hinton - Journal of machine learning research, 2008 - jmlr.org ... t-SNE is better than existing techniques at creating a single ... large data sets, we show how t-SNE can use random walks on ... We illustrate the performance of t-SNE on a wide variety of ... ☆ Save 55 Cite Cited by 25211 Related articles All 47 versions ↔

t-SNE is a dimension reduction algorithm.

Input: high-dimensional data Output: low-dimensional data that preserves...

- the graph structure?
- local neighborhoods?
- global structure?



t-SNE on MNIST



[PDF] Visualizing data using t-SNE.

L Van der Maaten, <u>G Hinton</u> - Journal of machine learning research, 2008 - jmlr.org ... **t-SNE** is better than existing techniques at creating a single ... large data sets, we show how **t-SNE** can use random walks on ... We illustrate the performance of **t-SNE** on a wide variety of ... ☆ Save 切 Cite Cited by 25211 Related articles All 47 versions ≫

How to Use t-SNE Effectively

MARTIN WATTENBERG	FERNANDA VIÉGAS	IAN JOHNSON	Oct. 13	
Google Brain	Google Brain	Google Cloud	2016	

1. Those hyperparameters really matter



2. Cluster sizes in a t-SNE plot mean nothing



SAMSI Interpretable Deep Learning Working Group



Yingfan Wang PhD student, Duke



Haiyang Huang PhD student, Duke



Yaron Shaposhnik Asst. Prof., U Rochester

Original Mammoth



Task: 3d to 2d.





Local vs Global

- Local structure: local neighborhood graph, nearest neighbors
- Global structure: relationships between clusters, respect relative distances between points in high-dimensional space.



Global Methods

- PCA (Pearson, 1901)
- MDS (Torgerson, 1952)

Local Methods

- LLE (Roweis and Saul, 2000),
- Isomap (Tenenbaum et al., 2000)
- Hessian Local Linear Embedding (Donoho and Grimes, 2003)
- Laplacian Eigenmaps (Belkin and Niyogi, 2001)
- Stochastic Neighborhood Embedding (SNE) (Hinton and Roweis, 2003)
- t-SNE (van der Maaten and Hinton, 2008)
- LargeVis (Tang et al., 2016)
- UMAP (McInnes et al., 2018)

Crowding problem

Preserve neighborhoods

Preserve distances,

not neighborhoods

Global Methods

- PCA (Pearson, 1901)
- MDS (Torgerson, 1952)

Article | Open Access | Published: 28 November 2019

The art of using t-SNE for single-cell transcriptomics

Dmitry Kobak 🖂 & Philipp Berens 🖂

Nature Communications 10, Article number: 5416 (2019) | Cite this article 36k Accesses | 67 Citations | 259 Altmetric | Metrics

Local Methods

- LLE (Roweis and Saul, 2000),
- Isomap (Tenenbaum et al., 2000)
- Hessian Local Linea arXiv.org > cs > arXiv:1708.03229
- Laplacian Eigenmap
- Stochastic Neighbori Yanshuai Cao, Luyu Wang
- t-SNE (van der Maat
- LargeVis (Tang et al., 2016)
- UMAP (McInnes et al., 2018)

How to Use t-SNE Effectively

2003)

MARTIN WATTENBERG FERNANDA V Google Brain Google Brain

FERNANDA VIÉGAS IAN JOHNSON Google Brain Google Cloud Oct. 13 2016

d Roweis, 2003)

t-Distributed stoch Automated optimal parameters for T-distributed stochastic neighbor embedding improve dimensionality red visualization and allow analysis of large datasets

October 2018 DOI: <u>10.1101/451690</u>

Project: Automated Analysis of Flow Cytometry Multidimensional Datasets

Authors:

Computer Science > Artificial Intelligence

[Submitted on 10 Aug 2017]



Algorithm	Graph components and Loss function		
	Graph components: Edges (i, j)		
t-SNE	Loss ^{t-SNE} _{<i>i,j</i>} = $p_{ij} \log \frac{p_{ij}}{q_{ij}}$, where $q_{ij} = \frac{\left(1 + \ \mathbf{y}_i - \mathbf{y}_j\ ^2\right)^{-1}}{\sum_{k \neq l} (1 + \ \mathbf{y}_k - \mathbf{y}_l\ ^2)^{-1}}$		
& Hinton, 2008)	where p_{ij} is a function of \mathbf{x}_i , \mathbf{x}_j and other \mathbf{x}_ℓ 's.		
	Graph components: Edges (i, j)		
UMAP	$\text{Loss}_{i,j}^{\text{UMAP}} = \begin{cases} \bar{w}_{i,j} \log \left(1 + a \left(\ \mathbf{y}_i - \mathbf{y}_j \ _2^2 \right)^b \right)^{-1} \\ (1 - \bar{w}_{i,j}) \log \left(1 - \left(1 + a \left(\ \mathbf{y}_i - \mathbf{y}_j \ _2^2 \right)^b \right)^{-1} \right) \end{cases}$	i, j neighbors otherwise.	
(McInnes et al., 2018)	where $\bar{w}_{i,j}$ is a function of \mathbf{x}_i , \mathbf{x}_j and nearby \mathbf{x}_ℓ 's.		
	Graph components: Triplets (i, j, k) where $\text{Distance}_{i,j} \leq \text{Distance}_{i,k}$		
TriMAP	Loss TM _{i,j,k} = $\omega_{i,j,k} \frac{s(\mathbf{y}_i, \mathbf{y}_k)}{s(\mathbf{y}_i, \mathbf{y}_j) + s(\mathbf{y}_i, \mathbf{y}_k)}$, where $s(\mathbf{y}_i, \mathbf{y}_j) = (1 + \ \mathbf{y}_i - \mathbf{y}_j\ ^2)^{-1}$		
(Amid & Warmuth, 2019)	and $\omega_{i,j,k}$ is a function of \mathbf{x}_i , \mathbf{x}_j , \mathbf{x}_k and nearby points.		

Hard to understand what's important here...

Start from the obvious:

- Attraction: high-dimensional neighbors should be attracted.
- Repulsion: points far in original space should be far in low-dim space.

But that's not enough...



Start from the obvious:

- Attraction: high-dimensional neighbors should be attracted.
- Repulsion: points far in original space should be far in low-dim space.

After a huge amount of experimentation, we found that:

- Certain specific properties of the loss function are important for local structure.
- The choice of which graph components to exert forces on is important for global structure.



Algorithm	Graph components and Loss function		
	Graph components: Edges (i, j)		
t-SNE (van der Maaten	Loss ^{t-SNE} _{<i>i,j</i>} = $p_{ij} \log \frac{p_{ij}}{q_{ij}}$, where $q_{ij} = \frac{(1 + \ \mathbf{y}_i - \mathbf{y}_j\ ^2)^{-1}}{\sum_{k \neq l} (1 + \ \mathbf{y}_k - \mathbf{y}_l\ ^2)^{-1}}$		
& Hinton, 2008)	where p_{ij} is a function of \mathbf{x}_i , \mathbf{x}_j and other \mathbf{x}_{ℓ} 's.		
	Graph components: Edges (i, j)		
UMAP	$\text{Loss}_{i,j}^{\text{UMAP}} = \begin{cases} \bar{w}_{i,j} \log \left(1 + a \left(\ \mathbf{y}_i - \mathbf{y}_j \ _2^2 \right)^b \right)^{-1} \\ (1 - \bar{w}_{i,j}) \log \left(1 - \left(1 + a \left(\ \mathbf{y}_i - \mathbf{y}_j \ _2^2 \right)^b \right)^{-1} \right) \end{cases}$	i, j neighbors otherwise,	
(McInnes et al., 2018)	where $\bar{w}_{i,j}$ is a function of \mathbf{x}_i , \mathbf{x}_j and nearby \mathbf{x}_{ℓ} 's.		
	Graph components: Triplets (i, j, k) where $\text{Distance}_{i,j} \leq \text{Distance}_{i,k}$		
TriMAP	$\text{Loss}_{i,j,k}^{\text{TM}} = \omega_{i,j,k} \frac{s(\mathbf{y}_i, \mathbf{y}_k)}{s(\mathbf{y}_i, \mathbf{y}_j) + s(\mathbf{y}_i, \mathbf{y}_k)}, \text{ where } s(\mathbf{y}_i, \mathbf{y}_j) = \left(1 + \ \mathbf{y}_i - \mathbf{y}_j\ ^2\right)^{-1}$		
(Amid & Warmuth, 2019)	and $\omega_{i,j,k}$ is a function of \mathbf{x}_i , \mathbf{x}_j , \mathbf{x}_k and nearby points.		

Hard to understand what's important here...

The "rainbow" plot

Triple *i*, *j* (neighbor), *k* (further)





Distance to neighbor j







Our principles for a good loss function



1. Monotonicity

Except at bottom, gradient should go mainly to the left.
 (if further point is sufficiently far, should focus on pulling neighbor closer.)

At bottom, gradient goes up.
 (push further points away really hard)

4. At left, gradient has small magnitude. (don't crowd, relax when close enough)

5. At bottom, gradient has large magnitude. (push farther point away)

6. Gradient fades as neighbor gets farther away. (give up on neighbors when they are too far)





supposed to give up on far neighbors...



Our principles for a good loss function

Very close to left, small gradient



VS.

Algorithm	Graph components and Loss function		
t-SNE	Graph components: Edges (i, j) $\operatorname{Loss}_{i,j}^{t-\operatorname{SNE}} = p_{ij} \log \frac{p_{ij}}{q_{ij}}$, where $q_{ij} = \frac{(1+\ \mathbf{y}_i-\mathbf{y}_j\ ^2)^{-1}}{\sum_{k \neq l} (1+\ \mathbf{y}_k-\mathbf{y}_l\ ^2)^{-1}}$ where p_{ij} is a function of \mathbf{x}_i , \mathbf{x}_j and other \mathbf{x}_ℓ 's.		
UMAP	Graph components: Edges (i, j) $ \operatorname{Loss}_{i,j}^{\mathrm{UMAP}} = \begin{cases} \bar{w}_{i,j} \log \left(1 + a \left(\ \mathbf{y}_i - \mathbf{y}_j \ _2^2 \right)^b \right)^{-1} & i, j \text{ neighbors} \\ \left(1 - \bar{w}_{i,j} \right) \log \left(1 - \left(1 + a \left(\ \mathbf{y}_i - \mathbf{y}_j \ _2^2 \right)^b \right)^{-1} \right) & \text{otherwise,} \end{cases} $ where $\bar{w}_{i,j}$ is a function of \mathbf{x}_i , \mathbf{x}_j and nearby \mathbf{x}_ℓ 's.	S	
TriMAP	Graph components: Triplets (i, j, k) where $\text{Distance}_{i,j} \leq \text{Distance}_{i,k}$ $\text{Loss}_{i,j,k}^{\text{TM}} = \omega_{i,j,k} \frac{s(\mathbf{y}_i, \mathbf{y}_k)}{s(\mathbf{y}_i, \mathbf{y}_j) + s(\mathbf{y}_i, \mathbf{y}_k)}$, where $s(\mathbf{y}_i, \mathbf{y}_j) = (1 + \ \mathbf{y}_i - \mathbf{y}_j\ ^2)^{-1}$ and $\omega_{i,j,k}$ is a function of $\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_k$ and nearby points.		

PaCMAP's loss is simpler.



- Certain specific properties of the loss function are important for local structure.
- The choice of which graph components to exert forces on is important for global structure.





PaCMAP's Loss



Mid near pair for i: sample 6 observations, choose the second closest of the 6, pair it with i.

Mid-near pairs have an effect!



- Certain specific properties of the loss function are important for local structure.
- The choice of which graph components to exert forces on is important for global structure.

How to evaluate DR?

How to Evaluate DR Algorithms?

- Local structure DR then supervised classification in 2D
- Global structure triplet loss
- Sensitivity to parameter choices
- Sensitivity to preprocessing choices
- Computational efficiency











PBMC Data, PaCMAP projection, Colored 3 Ways

(PBMC is Peripheral Blood Mononuclear Cell)



Data from Ding J, Adiconis X, Simmons SK, Kowalczyk MS et al. Systematic comparison of single-cell and singlenucleus RNA-sequencing methods. Nat Biotechnol 2020 Jun;38(6):737-746. PMID: 32341560 PaCMAP: 23.90 seconds

figure credit: Carla Moelbert

PBMC Data, UMAP projection, Colored 3 Ways



Colored by Cell Type

Colored by Experiment

Colored by Method

UMAP: 75.82 seconds figure credit: Carla Moelbert

Data from Ding J, Adiconis X, Simmons SK, Kowalczyk MS et al. Systematic comparison of single-cell and singlenucleus RNA-sequencing methods. Nat Biotechnol 2020 Jun;38(6):737-746. PMID: 32341560

How to Evaluate DR Algorithms?

- Local structure DR then supervised classification in 2D
- Global structure triplet loss
- Sensitivity to parameter choices
- Sensitivity to preprocessing choices
- Computational efficiency





Data from: Kazer, S. W. et al. Integrated single-cell analysis of multicellular immune dynamics during hyperacute HIV-1 infection. Nat. Med.26, 511–518 (2020).

How to Evaluate DR Algorithms?

- Local structure DR then supervised classification in 2D
- Global structure triplet loss
- Sensitivity to parameter choices
- Sensitivity to preprocessing choices
- Computational efficiency



Run time

Dataset (Size)	T-SNE	LARGEVIS	UMAP	TRIMAP	PACMAP
OLIVETTI FACES (0.4K)	00:00:04	00:08:13	00:00:02	00:00:01	00:00:01
COIL-20 (1.4K)	00:00:08	00:10:18	00:00:05	00:00:02	00:00:01
COIL-100(7.2K)	00:00:49	00:09:53	00:00:10	00:00:06	00:00:03
S-Curve with Hole $(9.5K)$	00:01:17	00:10:09	00:00:15	00:00:08	00:00:05
USPS $(9.5K)$	00:01:14	00:10:15	00:00:15	00:00:07	00:00:05
Mammoth (10K)	00:00:58	00:10:36	00:00:16	00:00:08	00:00:05
20Newsgroups (18K)	00:03:29	00:11:40	00:00:19	00:00:18	00:00:12
Mouse scRNA-seq $(20K)$	00:04:43	00:12:52	00:00:24	00:00:20	00:00:13
MNIST (70K)	00:14:02	00:20:19	00:01:09	00:01:14	00:00:52
F-MNIST (70K)	00:12:43	00:17:11	00:00:59	00:01:13	00:00:47
FLOW CYTOMETRY (3M)	-	-	-	02:10:27	00:58:28
KDD CUP99 (4M)	-	-	-	03:34:57	02:05:19

Name-Ethnicity Classification (helpful for assessing fairness)



EEG Monitoring



joint work with Alina Barnett, Zhicheng Guo, Jing Jing and Brandon Westover



Home Equity Line of Credit (HELOC) Dataset

This competition focuses on an anonymized dataset of Home Equity Line of Credit (HELOC) applications made by real homeowners. A HELOC is a line of credit typically offered by a bank as a percentage of home equity (the difference between the current market value of a home and its purchase price). The customers in this dataset have requested a credit line in the range of \$5,000 - \$150,000. The fundamental task is to use the information about the applicant in their credit report to predict whether they will repay their HELOC account within 2 years. This prediction is then used to decide whether the homeowner qualifies for a line of credit and, if so, how much credit should be extended.

PaCMAP result on FICO

higher PercentTradesWBalance, NetFractionRevolvingBurden and NetFractionInstallBurden



credit: Yingfan Wang

Summary







- What makes DR algorithms succeed/fail?
 - Local structure: A good loss function.
 - Should obey 6 principles.
 - Rainbow plot allows to compare across algorithms.
 - Global structure:
 - We suggest forces on non-neighbors. Mid-near pairs.
- PaCMAP
 - A simpler loss function involving 3 terms.
 - Preserves local and global structure, fast run time.



