

PDEs and Graph Based Learning

Summer School on Random Structures in Optimizations and Related Applications

Lecture 3: t-SNE

Instructor: Jeff Calder (jcalder@umn.edu)

Web: <http://www-users.math.umn.edu/~jwcalder>

Lecture Notes: <http://www-users.math.umn.edu/~jwcalder/5467S21>

Embeddings of high dimensional data

High dimensional data is hard to visualize and work with. Embeddings to low dimensional spaces help us visualize data and improve the performance of data analysis algorithms.

- PCA (linear dimension reduction)
- Spectral embedding (today)
- t-Distributed Stochastic Neighbor Embedding (t-SNE) (also today)

The key is to embed the data while still preserving important structures.

Spectral embeddings

Let v_1, v_2, v_3, \dots be the normalized eigenvectors of L , in order of increasing eigenvalues $0 = \lambda_1 \leq \lambda_2 \leq \dots$. The spectral embedding corresponding to L is the map $\Phi : I_m \rightarrow \mathbb{R}^k$ (recall $I_m = \{1, 2, \dots, m\}$ are the indices of our datapoints) given by

$$(1) \quad \Phi(i) = (v_1(i), v_2(i), \dots, v_k(i)).$$

Since the first eigenvector v_1 is the trivial constant eigenvector, it is also common to omit this to obtain the embedding

$$\Phi(i) = (v_2(i), v_3(i), \dots, v_{k+1}(i)).$$

There are other normalizations of the graph Laplacian that are commonly used, such as the symmetric normalization $L = D^{-1/2}(D - W)D^{-1/2}$, and the spectral embedding for a normalized Laplacian is defined analogously.

Spectral embedding of MNIST

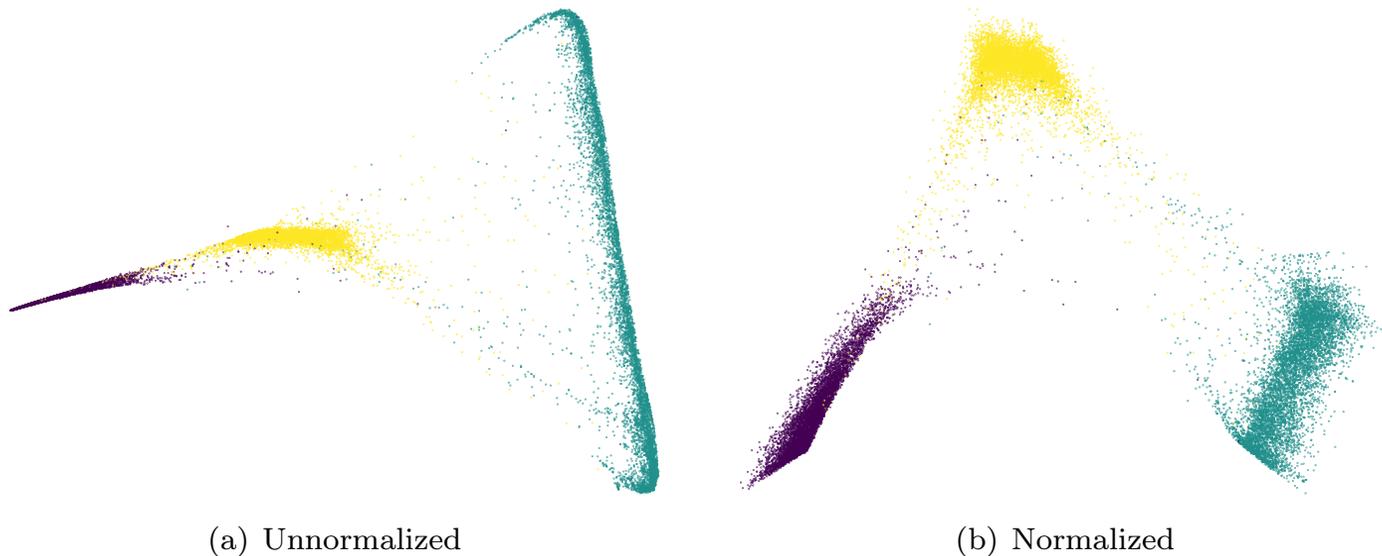


Figure 1: Example of spectral embeddings in the plane $k = 2$ of the 0, 1, and 2 digits of the MNIST dataset using the unnormalized $L = D - W$ and symmetric normalized $L = D^{-1/2}(D - W)D^{-1/2}$ graph Laplacians.

t-SNE

The t-Stochastic Neighbor Embedding (t-SNE) tries to find embedded points whose pairwise similarities match as closely as possible the given weight matrix W for the graph.

From the weight matrix W , t-SNE constructs a probability weight matrix

$$(2) \quad P = \frac{1}{2m} (D^{-1}W + W^T D^{-1}),$$

where D is the diagonal matrix of degrees $d(i) = \sum_{j=1}^m W(i, j)$. We say *probability* since all entries of P sum to one, i.e., $\mathbf{1}^T P \mathbf{1} = 1$.

t-SNE

t-SNE aims to find embedded points $y_1, y_2, \dots, y_m \in \mathbb{R}^k$, where usually $k = 2$ or $k = 3$, so that the similarity between y_i and y_j matches $P(i, j)$ as closely as possible. The similarity matrix for the y_i , denoted Q , is the $m \times m$ matrix defined by

$$(3) \quad Q(i, j) = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_{\ell \neq s} (1 + \|y_s - y_\ell\|^2)^{-1}}.$$

The discrepancy between P and Q is measured with the Kullback-Leibler divergence

$$(4) \quad E(y_1, y_2, \dots, y_k) = D(P \| Q) := \sum_{i \neq j} P(i, j) \log \left(\frac{P(i, j)}{Q(i, j)} \right).$$

t-SNE finds the embedded points y_i by minimizing E with gradient descent.

Why Kullback-Leibler?

$$E(y_1, y_2, \dots, y_k) = D(P||Q) := \sum_{i \neq j} P(i, j) \log \left(\frac{P(i, j)}{Q(i, j)} \right).$$

- When $P(i, j) \gg 0$, forces $Q(i, j) \sim P(i, j)$; i.e., preserve local structure.
- When $P(i, j) \sim 0$ we don't care what $Q(i, j)$ does; i.e., allow global structure to change.
- We cannot preserve all information in a dimension reduction.

Gradient descent

The gradient of the Kullback-Leibler divergence

$$E(y_1, y_2, \dots, y_k) = D(P\|Q) := \sum_{i \neq j} P(i, j) \log \left(\frac{P(i, j)}{Q(i, j)} \right)$$

in the variable y_i is given by

$$\nabla_{y_i} E = \underbrace{4Z \sum_{j:j \neq i} P(i, j) Q(i, j) (y_i - y_j)}_{\text{Attraction}} - \underbrace{4Z \sum_{j:j \neq i} Q(i, j)^2 (y_i - y_j)}_{\text{Repulsion}}.$$

where $Z = \sum_{i \neq j} (1 + \|y_i - y_j\|^2)^{-1}$.

Gradient descent is

$$y_i^{k+1} = y_i^k - h \nabla_{y_i} E(y_1^k, y_2^k, \dots, y_m^k),$$

where h is the time-step.

t-SNE embedding of MNIST

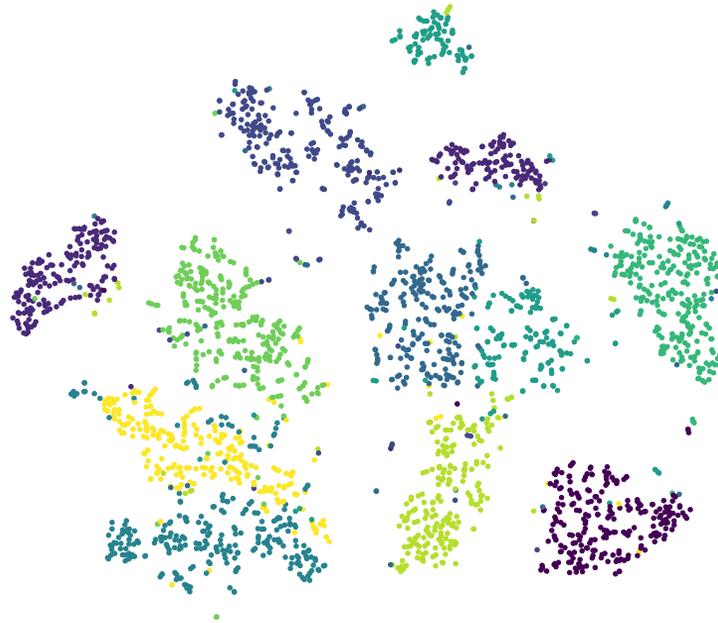


Figure 2: A t-SNE embedding of 2500 images from the MNIST dataset, with colors corresponding to the digit labels of each image.

$$E = \sum_{i \neq j} P(i, j) \log \left(\frac{P(i, j)}{Q(i, j)} \right)$$

we may as well write

$$E = - \sum_{i \neq j} P(i, j) \log(Q(i, j))$$

Recall

$$Q(i, j) = \frac{(1 + |y_i - y_j|^2)^{-1}}{\sum_{k \neq s} (1 + |y_k - y_s|^2)^{-1}}$$

$$E = \underbrace{\sum_{i \neq j} P(i, j)}_A \log(1 + |y_i - y_j|^2)$$

$$+ \left(\sum_{i \neq j} P(i, j) \right) \log \left(\sum_{l \neq s} (1 + |y_l - y_s|^2) \right)^{-1}$$

$= 1$
}
B

$$\nabla_{y_l} A = \sum_{i \neq j} P(i, j) \nabla_{y_l} \log (1 + |y_i - y_j|^2)$$

$$= \sum_{i \neq j} P(i, j) (1 + |y_i - y_j|^2)^{-1} \nabla_{y_l} |y_i - y_j|^2$$

$$= 2 \sum_{i \neq j} P(i, j) (1 + |y_i - y_j|^2)^{-1} (y_i - y_j) (\delta_{il} - \delta_{jl})$$

If $i=l$, $\nabla_{y_l} |y_i - y_j|^2 = 2(y_i - y_j)$

$$\text{If } j=l, \quad \nabla_{y_l} |y_i - y_j|^2 = 2(y_j - y_i) \\ = -2(y_i - y_j)$$

$$\delta_{ij} = \begin{cases} 1, & i=j \\ 0, & i \neq j \end{cases}$$

$\nabla_{y_l} A$

$$= 2 \sum_{i \neq j} P(i, j) (1 + |y_i - y_j|^2)^{-1} (y_i - y_j) (\delta_{il} - \delta_{jl})$$

$$= 2 \sum_{j \neq l} P(l, j) (1 + |y_l - y_j|^2)^{-1} (y_l - y_j)$$

$$- 2 \sum_{j \neq l} P(j, l) (1 + |y_j - y_l|^2)^{-1} (y_j - y_l)$$

$$= 4 \sum_{j \neq l} P(l, j) (1 + |y_l - y_j|^2)^{-1} (y_l - y_j)$$

Define $Z = \sum_{i \neq j} (1 + |y_i - y_j|^2)^{-1}$

so that $Q(i, j) = \frac{(1 + |y_i - y_j|^2)^{-1}}{Z}$

$$= 4 Z \sum_{j \neq l} P(l, j) Q(l, j) (y_l - y_j)$$

Attraction term.

$$\beta = \log \left(\sum_{i \neq j} (1 + |y_i - y_j|^2)^{-1} \right)$$

$$\nabla_{y_\ell} \beta = \nabla_{y_\ell} \log \left(\sum_{i \neq j} (1 + |y_i - y_j|^2)^{-1} \right)$$

$$= z^{-1} \nabla_{y_\ell} \sum_{i \neq j} (1 + |y_i - y_j|^2)^{-1}$$

$$= -z^{-1} \sum_{i \neq j} \underbrace{(1 + |y_i - y_j|^2)^{-2}}_{\alpha(i,j)^2 z^2} \nabla_{y_\ell} |y_i - y_j|^2$$

$$= -z \sum_{i \neq j} \alpha(i,j)^2 2 (y_i - y_j) (d_{i\ell} - d_{j\ell})$$

$$= -\tau \sum_{i \neq j} \alpha(r_{ij})^2 2 (y_i - y_j) \delta_{ie}$$

$$+ \tau \sum_{i \neq j} \alpha(r_{ij})^2 2 (y_i - y_j) \delta_{je}$$

$$= -2\tau \sum_{j \neq e} \alpha(r_{ej})^2 (y_e - y_j)$$

$$+ 2\tau \sum_{i \neq e} \alpha(r_{ie}) (y_i - y_e)$$

$$= -4\tau \sum_{j \neq e} \alpha(r_{ej})^2 (y_e - y_j)$$

Repulsion term.

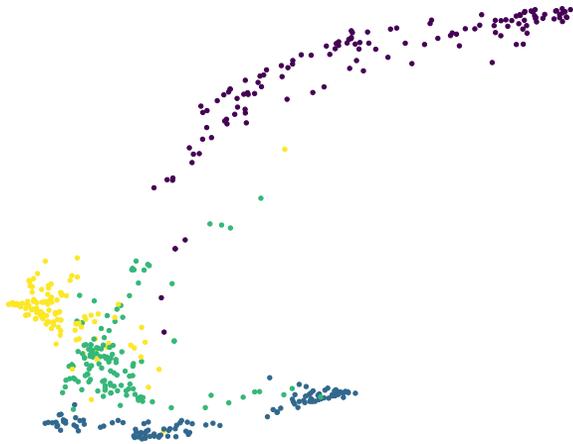
Early exaggeration

Gradient descent for t-SNE is *very* slow. To speed it up, *early exaggeration* is used, which amplifies the attraction forces for the first few hundred iterations:

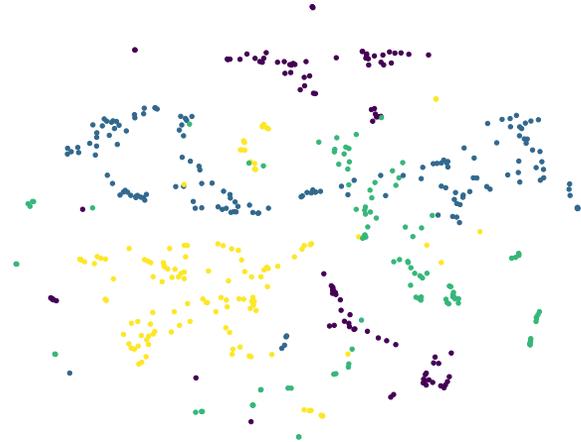
$$\nabla_{y_i} E = 4Z\alpha \sum_{j:j \neq i} P(i,j)Q(i,j)(y_i - y_j) - 4Z \sum_{j:j \neq i} Q(i,j)^2(y_i - y_j),$$

The parameter α is the amplification factor, often $\alpha \approx 10$.

Early exaggeration



(a) After early exaggeration



(b) Final embedding

Figure 3: An example of the t-SNE embedding after the early exaggeration phase and the final embedded, for a small version of MNIST with only 500 images from the digits 0, 1, 2, and 3.

Perplexity

The construction of the weight matrix W is important for the performance of t-SNE. The *perplexity* construction has the form

$$W(i, j) = \exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma_i^2}\right),$$

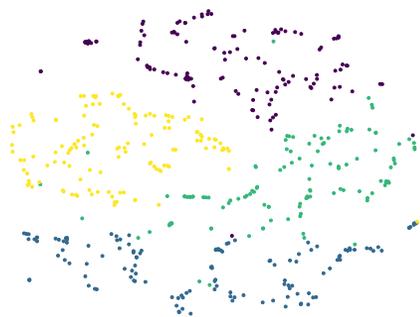
where σ_i is tuned independently for each x_i depending on a specified *perplexity* level (usually in the range 5 to 50).

The perplexity of the i^{th} row of $W(i, j)$ is $2^{H(i)}$, where

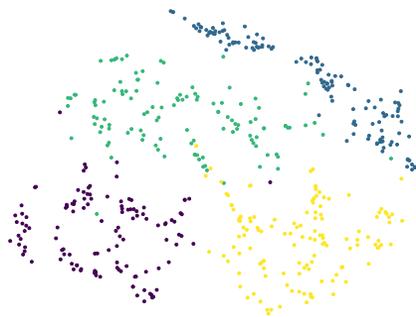
$$H(i) = -\sum_{j=1}^m p(j) \log p(j), \quad p(j) = \frac{W(i, j)}{\sum_{k=1}^m W(i, k)}.$$

The value of σ_i is determined so that the perplexity $2^{H(i)}$ equals a desired user-specified value.

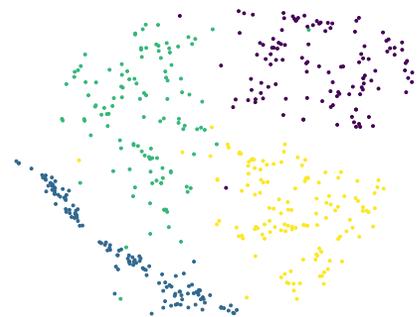
MNIST



Perplexity 5

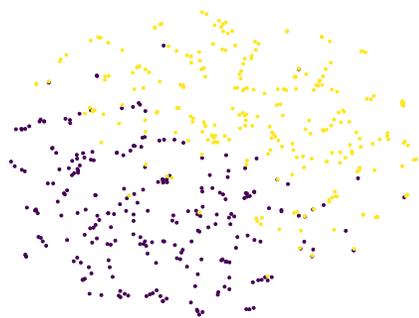


Perplexity 30



Perplexity 50.

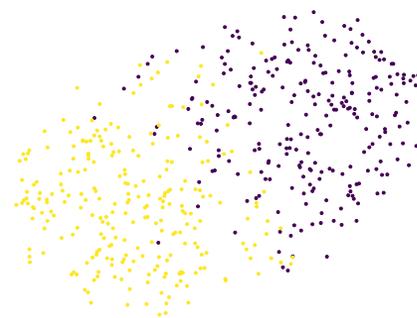
Gaussian mixture in 10 dimensions



Perplexity 5

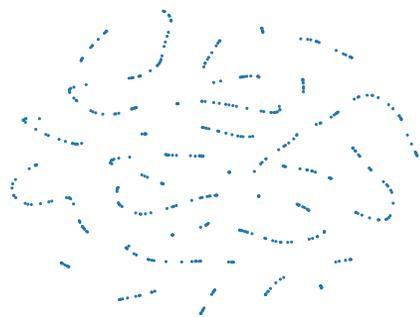


Perplexity 30

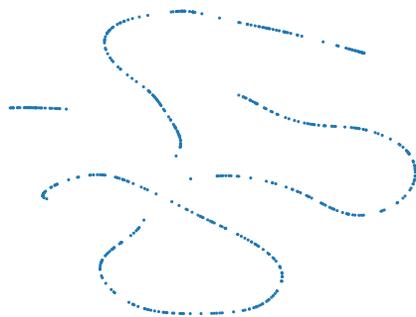


Perplexity 50.

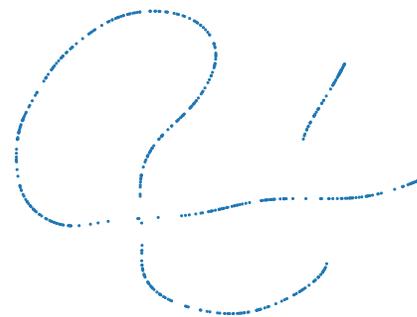
Parabolic curve in 5 dimensions



Perplexity 5



Perplexity 30



Perplexity 50.

Graph-based embeddings ([.ipynb](#))