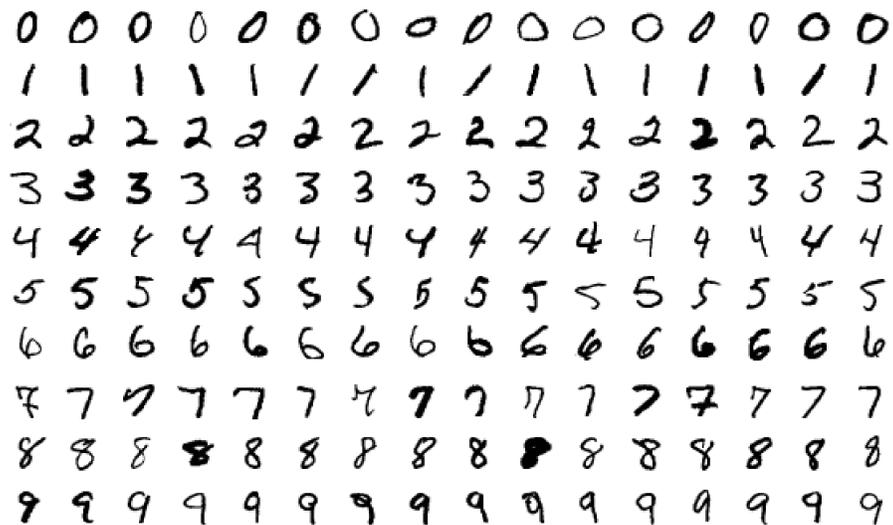


Explorative Analysis via t-SNE Visualization of High-Dimensional Data

Data Example: MNIST

The MNIST data consists of 70,000 images with gray scale, anti-aliased 28×28 pixel images, i.e., elements in \mathbb{R}^{784} , image by [1].



Goal

Provide a two-dimensional visualization of the data in a way such that:

1. Points that are close in the high-dimensional space are close in the two-dimensional embedding.
2. Points in the two-dimensional embedding are sufficiently far away from each other (no "crowding").

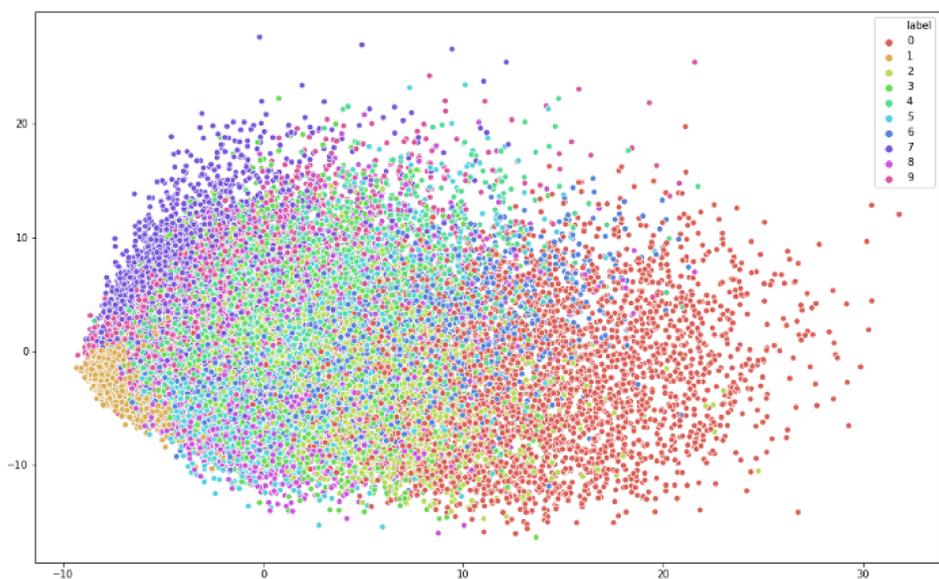
Method

Take the similarities P of the high-dimensional points and find similarities Q of the two-dimensional points such that

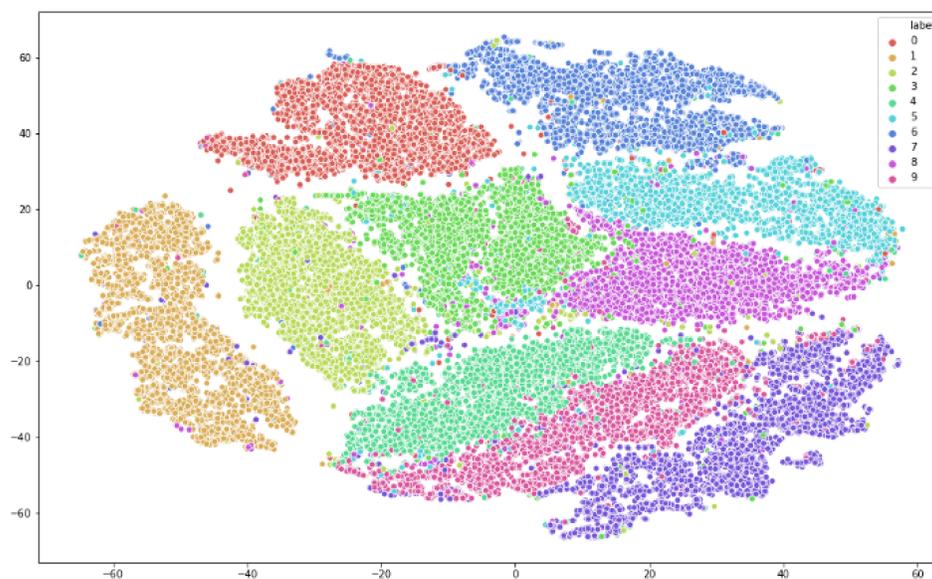
$$KL(P||Q) = \sum_i \sum_{j \neq i} p_{i|j} \log \left(\frac{p_{i|j}}{q_{ij}} \right)$$

is minimized [2].

Embedding the MNIST data via PCA [3]



Embedding the MNIST data via t-SNE [3]



Definition of the Similarities

If we choose

$$p_{j|i} = \frac{\exp(-\|x_i - x_j\|^2 / 2\sigma_i^2)}{\sum_{j' \neq j} \exp(-\|x_i - x_{j'}\|^2 / 2\sigma_i^2)},$$

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

for some specific bandwidth σ_i , then:

- ▶ large $p_{i|j}$ and small q_{ij} give big penalty,
- ▶ small $p_{i|j}$ and large q_{ij} give small penalty.

Minimization

Perform gradient descent on

$$\frac{KL(P||Q)}{\partial y_i} = 4 \sum_{j \neq i} p_{i|j} q_{ij} Z(y_i - y_j) - 4 \sum_{j \neq i} q_{ij}^2 Z(y_i - y_j),$$

where $Z = \sum_{k \neq \ell} (1 + \|y_k - y_\ell\|^2)^{-1}$ normalizes and

- ▶ the first term represents attractive forces, while
- ▶ the second term represents repulsive forces.

One exact evaluation of the gradient takes $\mathcal{O}(n^2)$ steps for n points. Thus, utilize acceleration data structures.

Project Questions:

- ▶ Development of multi-grid approaches to accelerate the optimization.
- ▶ Investigation of alternative embeddings spaces, e.g., hyperbolic space.
- ▶ Derivation of corresponding acceleration structures.
- ▶ Optimized user-interaction with the embeddings.

References

- [1] Stepan, L. (2017). A few samples from the MNIST test dataset. Wikimedia commons.
- [2] Van der Maaten, L., & Hinton, G. (2008). Visualizing data using t-SNE. Journal of machine learning research, 9(11).
- [3] Zhang, D. (2020). Dimensionality Reduction using t-Distributed Stochastic Neighbor Embedding (t-SNE) on the MNIST Dataset. Towards Data Science.