

t-Distributed Stochastic Neighbor Embedding (t-SNE)

Jerome Chou

January 31, 2025

Introduction to t-SNE

- ▶ t-SNE is a nonlinear dimensionality reduction technique used for visualization.
- ▶ It maps high-dimensional data to a lower-dimensional space while preserving local structure.
- ▶ Often used to explore and cluster high-dimensional datasets.

High-Dimensional Probability Distributions

- ▶ Given a dataset $X = \{x_1, x_2, \dots, x_n\}$, we define pairwise similarities using a Gaussian distribution:

$$p_{j|i} = \frac{\exp(-\|x_i - x_j\|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|x_i - x_k\|^2 / 2\sigma_i^2)} \quad (1)$$

- ▶ The joint probability distribution is then:

$$p_{ij} = \frac{p_{j|i} + p_{i|j}}{2n} \quad (2)$$

Low-Dimensional Probability Distributions

- ▶ In the lower-dimensional space, we model similarities with a Student's t-distribution:

$$q_{ij} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_{k \neq l} (1 + \|y_k - y_l\|^2)^{-1}} \quad (3)$$

- ▶ This distribution allows better separation of clusters.

KL Divergence Optimization

- ▶ The objective function minimizes the Kullback-Leibler (KL) divergence:

$$C = \sum_{i \neq j} p_{ij} \log \frac{p_{ij}}{q_{ij}} \quad (4)$$

- ▶ This is optimized using gradient descent to update the low-dimensional representation Y .

Machine Learning Applications of t-SNE

- ▶ Visualizing high-dimensional datasets (e.g., MNIST, gene expression data).
- ▶ Discovering clusters and patterns in data.
- ▶ Understanding feature embeddings in deep learning.