

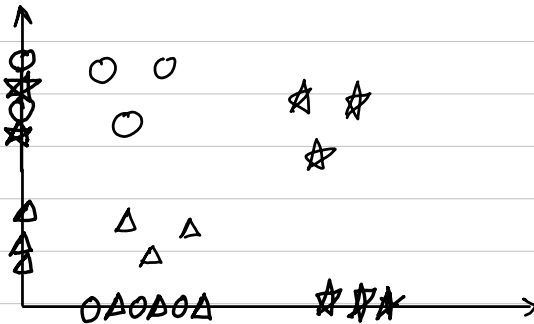
t-SNE

t-SNE preserves the original clustering of the data.

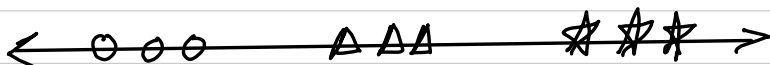
Consider data \rightarrow



Projecting on either x axis or y axis will not give the correct clustering



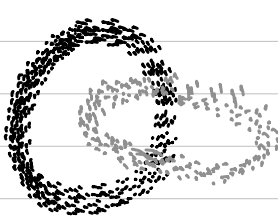
We want a projection that preserves the data like this



t - Stochastic Neighbourhood Embedding

PCA separates data using a straight line
However, not all data can be separated so.
hence it fails to interpret complex datasets.

t - SNE performs nonlinear dimensionality reduction & overcomes this limitation.



Linearly Non
separable
data.

t - SNE is a form of Manifold learning

It is mainly used for data visualization in lower dimensions (1D to 2D)

It can also be used for dimensionality reduction

Useful for revealing clusters in data

Preserves local structure of data - points that are close together in higher dimensions remain close in the lower dimensions as well

How t-SNE operates

t-SNE works on minimizing the divergence between the two distributions — one in the higher dimensions and one in the lower dimensions.

It is a probabilistic model that uses gradient descent. This may give different results when run different times.

How t-SNE operates →

1. Measure similarity between points in a Matrix
2. Map random to lower dimension & calculate similarity
3. Adjust the points to have same distribution of similarity as original high dimensional matrix
4. Visualize final point positions.

Unlike PCA, t-SNE is based on probability distributions with random walk or neighborhood graphs to find the structure within the data

t-SNE has a non-convex objective function that is minimized using gradient descent. It has random initiation, which means that different runs give different solutions.

t-SNE produces well-separated clusters, providing clear insights into the grouping of similar datapoints.

Unlike LDA, t-SNE does not assume that the covariances are equal.

Applications →

Data Exploration - Understand structure & relationship in high dimensional data

NLP → Visualize word embeddings.

Market segmentation, Social Network analysis, recommender systems.

Anomaly detection → By visualizing the data in low dimensional space t-SNE can help identify clusters of datapoints that are different from rest of data. Useful for potential fraud, anomaly detection; outliers.

How t-SNE Works



There are 3 distinct classes that we recognize

First step is to find the similarity between points

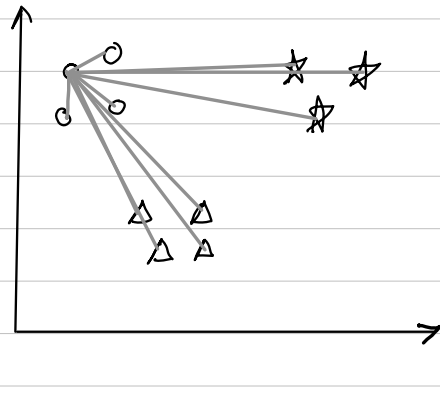
Similarity of a datapoint x_i to datapoint x_j is the datapoint $p_{j|i}$ that x_i would pick up x_j as its neighbour in a random walk under a gaussian centered at x_i

Hence for nearby datapoints $p_{j|i}$ is relatively high, whereas for widely spaced points $p_{j|i}$ is low and reduces to infinitesimal

Let the variance of the gaussian centered around x_i be σ_i^2 . σ_i will be calculated later.

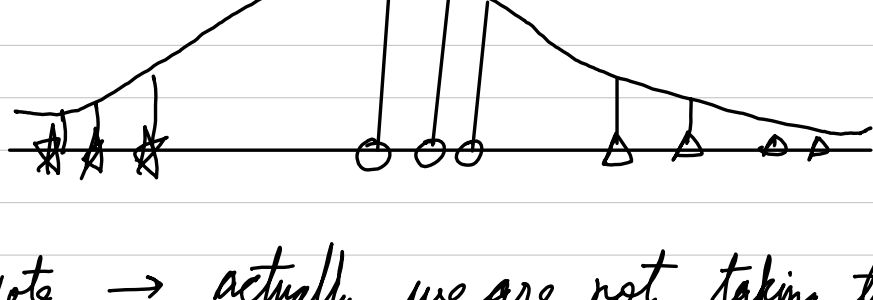
$$p_{j|i} = \frac{e^{-\frac{\|x_i - x_j\|^2}{2\sigma_i^2}}}{\sum_{k \neq i} e^{-\frac{\|x_i - x_k\|^2}{2\sigma_i^2}}}$$

$p_{j|i}$ is the similarity of point i & j measured from i

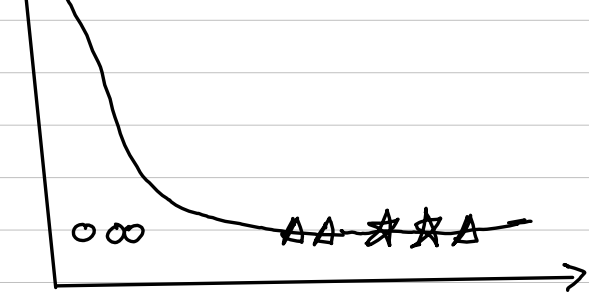


Similarly for every point j we calculate similarity

$$\sum_j p_{ij} = 1$$



Note \rightarrow actually we are not taking the gaussian values but slightly different function $e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}}$ hence graph becomes



We do this for every point & make a table of similarity

But p_{ij} & p_{ji} are not same because they have different σ values.

In order to make them same, average value is taken (symmetric SNE)

$$p'_{ij} = p'_{ji} = \frac{p_{ij} + p_{ji}}{N} \quad N \leftarrow \text{No of dimension}$$

Our similarity table \rightarrow

	x_1	x_2	x_3	...	x_k
x_1	0	0.3	0.3	...	-
x_2	0.3	0	0.4
x_3	0.2	0.2	0	.	.
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
x_k	0

The aim is that when we reduce the dimensions, the same similarity table must be followed

The next step is to create a low dimensional space with the same number of points y_i

1. Randomly place the points
2. Calculate similarity
3. Move points with gradient descent to have same similarity as the original table

New similarity using gaussian

$$q_{j|i} = \frac{e^{-\frac{\|y_i - y_j\|^2}{2\sigma^2}}}{\sum_{k \neq i} e^{-\frac{\|y_i - y_k\|^2}{2\sigma^2}}}$$

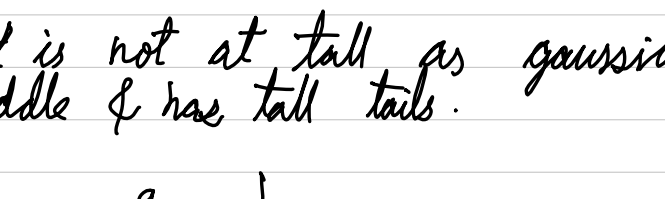
Here $\sigma = 1/\sqrt{2}$

This is the Method for SNE. But for t-SNE the method is different.

Using a gaussian leads to a "crowding problem".

Clusters all clump up in the middle & are hard to see.

Hence the t distribution is used instead with a single degree of freedom



It is not as tall as gaussian in the middle & has tall tails.

Hence q_{ij} becomes

$$q_{ij} = \frac{1}{1 + \|y_i - y_j\|^2}$$

$$\sum_{k=1}^n \frac{1}{1 + \|y_i + y_k\|^2}$$

The objective function is using the "kullback leiber divergence"

to check how similar P & Q are

$$C = D_{KL}(P||Q) = \sum_{x \in X} P(x) \log\left(\frac{P(x)}{Q(x)}\right)$$

$$= \sum_i \sum_j p_{ij} \log\left(\frac{p_{ij}}{q_{ij}}\right)$$

Gradient descent

$$\frac{\partial C}{\partial y_i} = \sum_j \frac{(p_{ij} - q_{ij})(y_i - y_j)}{(1 + \|y_i - y_j\|^2)}$$

This can be interpreted as springs that attract or repel.

t-SNE gradient strongly repels dissimilar datapoints.

Similar datapoints are attracted

This way, after random shuffling, the rearrangement occurs

Calculation of σ_i

σ_i is found using binary search to match the parameter perplexity.

$$- \sum_j p_{ij} \log_2 p_{ij}$$

$$\text{perp} = 2$$

Perplexity Parameter

Perplexity is a knob that can be used to adjust the number of effective nearest neighbours.

The most appropriate value depends on the density of the data.

Larger / denser dataset requires a large perplexity

Typical values of perplexity range between 5 & 50

If the perplexity value is too high then all points become equidistant & a strange ball with uniformly distributed points occur

If the dataset contains very large numbers then the binary search may fail & t-SNE reports a very low but results look bad. Scaling must be performed.

t -SNE	PCA
Visualization of high dimensional dataset	Dimensionality reduction
Non Linear Manifold learning	Linear transformation
Probabilistic algorithm	Algebraic algorithm
Preserves local structure	Maximizes variance
Data Reduced to 2D or 3D Visualization	Lower dimensions, may be depending on % of loss
Less scalable for large datasets	More scalable to large datasets
Robust to Noise	Sensitive to noise
Computationally Intensive	Better computation than t -SNE
Depends on random initialization	Same result everytime