

# A Gentle Introduction to Stochastic Neighbor Embedding (SNE)

Based on the Tutorial by B. Ghogogh, A. Ghodsi, et al.

Data Science University

September 27, 2025

# The Big Picture

[cite<sub>s</sub> tart] We have high – dimensional data,  $\{x_i\}_{i=1}^n$  where each  $x_i \in \mathbb{R}^d$  [cite: 46]. Think of data with hundreds of features. It's impossible to "see".

[cite<sub>s</sub> tart] Our goal is to create a meaningful low – dimensional "map" of this data [cite : 46].

We want to find a new set of points,  $\{y_i\}_{i=1}^n$  where each  $y_i \in \mathbb{R}^p$  [cite: 46]. [cite<sub>s</sub> tart]

Usually, the new dimension  $p$  is 2 or 3 so we can make a scatter plot [cite: 47].

# Stochastic Neighbor Embedding (SNE)

[cite<sub>s</sub> tart] **SNE** is a *manifold learning* *method for dimensionality reduction* [cite : 27].

## The Core Idea

[cite<sub>s</sub> tart] Instead of thinking about neighbors in a strict "yes/no" way, SNE uses probabilities [cite<sub>s</sub> tart] Every point considers every other point its neighbor, but with a certain probability [cite<sub>s</sub> tart] Closer points get a higher probability [cite : 31].

[cite<sub>s</sub> tart] The algorithm's main job is to arrange the points in the low-dimensional map so that these neighborhood probabilities are preserved as closely as possible [cite : 33].

## ① Part 1: Stochastic Neighbor Embedding (SNE)

- How do we define "neighborhood probability"?
- How do we measure if our map is good?
- How do we optimize the map?

## ② Part 2: Symmetric SNE (In the next section!)

## ③ Part 3: The Crowding Problem & t-SNE (Later!)

Let's Dive In!

## Step 1: Probabilities in the Original Space

[cite<sub>s</sub> tart] *SNE places a Gaussian distribution centered on every data point,  $x_i$*  [cite: 53].

*ite<sub>s</sub> tart This Gaussian measures the probability density for any other point,  $x_j$ , to be a neighbor of  $x_i$*  [cite: 53].

[cite<sub>s</sub> tart] *ite<sub>s</sub> tart Points that are farther away will naturally have a lower probability* [cite: 53].

[cite<sub>s</sub> tart] *The key idea is to convert the Euclidean distances between points into conditional probabilities* [cite : 60].

# The Math: Calculating $p_{j|i}$

The probability that point  $x_j$  is a neighbor of  $x_i$  is given by  $p_{j|i}$ .

## Similarity Score

[cite<sub>s</sub> tart] First, we calculate a similarity score based on the Gaussian function:  $\exp(-||x_j - x_i||^2 / 2\sigma_i^2)$  [cite: 61, 62].

## Softmax Normalization

[cite<sub>s</sub> tart] To turn these scores into real probabilities that sum to 1, we use a "softmax" [cite: 59]. [cite<sub>s</sub> tart] We divide each score by the sum of all possible scores (for a given  $i$ ) [cite: 61].

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

[cite<sub>s</sub> tart] This is the conditional probability of  $j$  given  $i$  [cite: 60]. [cite<sub>s</sub> tart] Note that  $p_{i|i}$  is set to 0 [cite: 92].

# What is that $\sigma_i^2$ (Variance)?

[cite<sub>s</sub> tart] The term  $\sigma_i^2$  is the variance of the Gaussian centered on point  $x_i$  [cite: 73].

- It controls the "width" of the Gaussian curve.
- A small  $\sigma_i$  means only very close points are considered neighbors.
- A large  $\sigma_i$  means even distant points can be neighbors.

## How is it chosen?

[cite<sub>s</sub> tart] Crucially, each point  $x_i$  gets its own  $\sigma_i$  [cite: 73].

[cite<sub>s</sub> tart] It's found using a binary search to match a user-defined value called *\*perplexity\**, which you can think of as the desired number of effective neighbors [cite : 73].

This allows the algorithm to adapt to regions of varying data density.



## Step 2: Probabilities in the Map

Now, we do the exact same thing for our low-dimensional points,  $y_i$ !

- Our goal is to arrange the  $y_i$  points on a 2D or 3D map. [cite<sub>s</sub>tart]
- We also define a conditional probability,  $q_{j|i}$ , that  $y_j$  is a neighbor of  $y_i$  in this new map[cite: 75].

### A Simpler Gaussian

We use the same softmax approach, but with a much simpler Gaussian.

[cite<sub>s</sub>tart] *We fix the variance to be  $\sigma^2 = 1/2$*  [cite: 81].

$$q_{j|i} = \frac{\exp(-\|y_i - y_j\|^2)}{\sum_{k \neq i} \exp(-\|y_i - y_k\|^2)}$$

Why fixed?

[cite<sub>s</sub>tart] *Because the scale of the map is something we can choose, so we don't need to* [cite: 81].

# Summarizing Our Two Worlds

We now have two sets of probabilities.

## High-D World (Input)

Probabilities  $P_i = \{p_{1|i}, p_{2|i}, \dots\}$

- Based on original data  $x_i$ .
- Variances  $\sigma_i^2$  are learned.

## Low-D World (Map)

Probabilities  $Q_i = \{q_{1|i}, q_{2|i}, \dots\}$

- Based on map points  $y_i$ .
- Variance is fixed.

## The Main Goal

*[cite<sub>s</sub> tart]* We want to arrange the  $y_i$  points so that the probabilities  $Q_i$  become as similar to the probabilities  $P_i$  as possible for all data points *[cite: 82]*.

## Step 3: How Do We Measure "Similarity"?

How do we measure the difference between two probability distributions?

### The Kullback-Leibler (KL) Divergence

[cite<sub>5</sub> tart] The KL Divergence is a standard way to measure how one probability distribution differs from another [cite : 83].

[cite<sub>5</sub> tart] Our total cost function

\*(C) is the sum of all the KL divergences for every point [cite : 83].  $C =$

$$\sum_{i=1}^n KL(P_i || Q_i) \text{ Which expands to } : C = \sum_{i=1}^n \sum_{j \neq i} p_{j|i} \log \frac{p_{j|i}}{q_{j|i}}$$

**Our task:** Find the arrangement of  $\{y_i\}$  that **minimizes this cost**!

# Intuition on the Cost Function

$$C = \sum_i \sum_j p_{j|i} \log \frac{p_{j|i}}{q_{j|i}}$$

Let's break this down.

*ite\_s tart The cost function heavily penalizes putting points far apart in the map (small  $p_{j|i}$ ) when they were close in the original data (large  $p_{j|i}$ ) [cite: 613].*  
*ite\_s tart A large mismatch where  $p_{j|i}$  is large creates a large cost.*  
*ite\_s tart It cares less about putting points close together in the map (large  $q_{j|i}$ ) when they were far apart in the data (small  $q_{j|i}$ ).*  
*ite\_s tart This means SNE focuses more on preserving the \*local structure\* of the data.*

## Step 4: Finding the Best Map

How do we find the positions of  $y_i$  that minimize our cost function  $C$ ?

### Gradient Descent!

This is a classic optimization problem.

[cite<sub>start</sub>]We'll use the gradient descent algorithm [cite : 68].

The process:

- 1 Start with a random arrangement of the  $y_i$  points on the map.
- 2 Calculate the "gradient" of the cost function. The gradient tells us, for each point  $y_i$ , which direction to move it to decrease the cost the most.
- 3 Take a small step in that direction.
- 4 Repeat steps 2 and 3 many times until the map settles down.

# The SNE Gradient

[cite<sub>start</sub>] The derivative (or gradient) of the cost function  $C$  with respect to the position  $y_i$  is surprisingly elegant [cite: 89].

## The Gradient Formula

$$\frac{\partial C}{\partial y_i} = 2 \sum_j ((p_{j|i} - q_{j|i}) + (p_{i|j} - q_{i|j}))(y_i - y_j)$$

This equation looks complex, but the intuition is simple. It's like a system of springs connecting all the points!

# Intuition: A System of Springs

$$\frac{\partial C}{\partial y_i} = 2 \sum_j (\dots)(y_i - y_j)$$

Think of the term  $(\dots)$  as the "stiffness" of a spring between points  $y_i$  and  $y_j$ .

- If points  $i$  and  $j$  should be closer ( $p_{j|i}$  is large but  $q_{j|i}$  is small), the "stiffness" is positive. This creates an **attractive force**, pulling  $y_i$  and  $y_j$  together.
- If points  $i$  and  $j$  are too close ( $p_{j|i}$  is small but  $q_{j|i}$  is large), the "stiffness" is negative. This creates a **repulsive force**, pushing them apart.

The algorithm finds the equilibrium state of this complex spring system!

T

- $y_i^{(t-1)}$  is the old position.
- $\frac{\partial C}{\partial y_i}$  is the gradient (the direction of "steepest ascent"). [cite<sub>s</sub> *tart*]
- $\eta$  is the **learning rate**: a small number that controls how big of a step we take[cite: 149].



b make the optimization faster and avoid getting stuck in bad local minima, we can add a **momentum** term[cite: 143, 148].

## The Idea

[cite<sub>start</sub>] The update for the current step should also depend a little on the update from previous \* step[cite : 140]. It's like giving the ball rolling down the hill some momentum, so it doesn't get stuck

The update rule with momentum is:

$$\Delta y_i^{(t)} = -\eta \frac{\partial C}{\partial y_i} + \alpha(t) \Delta y_i^{(t-1)}$$

$$y_i^{(t)} = y_i^{(t-1)} + \Delta y_i^{(t)}$$

[cite<sub>start</sub>] Here,  $\alpha(t)$  is the momentum term, which can change over time[cite: 140, 145].

# Another Trick: A Bit of Randomness

## Avoiding Local Optima

[cite<sub>start</sub>] *In the early stages of optimization, it's helpful to add a small amount of randomness after each iteration* [cite: 150].

**Why?** *This helps "shake" the points out of poor local arrangements and* [cite: 151].

# SNE: The Story So Far

- For every pair of points in the high-D space, calculate the conditional neighborhood probability  $p_{j|i}$  using Gaussians with adaptive variances[cite: 61, 73]. [cite<sub>s</sub>tart]
- Start with a random scatter plot of points  $\{y_i\}$  in the low-D space[cite: 139]. [cite<sub>s</sub>tart]
- Calculate the conditional neighborhood probability  $q_{j|i}$  for these map points using Gaussians with a fixed variance[cite: 77, 81]. [cite<sub>s</sub>tart]
- Calculate the gradient of the KL divergence cost function[cite: 90]. [cite<sub>s</sub>tart]
- Update the positions of all  $y_i$  points by taking a small step in the direction of the negative gradient, using momentum[cite: 140].
- Repeat for many iterations until the map is stable.

# Next Steps

SNE is a fantastic idea, but it has a few wrinkles.

- The probabilities  $p_{j|i}$  are not symmetric. This makes the gradient a bit messy.
- It suffers from something called the "Crowding Problem."

C

Coming up next: We'll see how **Symmetric SNE** and the revolutionary **t-SNE** solve these problems!