SNE and t-SNE

Statistical Machine Learning (ENGG*6600*02)

School of Engineering, University of Guelph, ON, Canada

Course Instructor: Benyamin Ghojogh Summer 2023 Introduction

Introduction

- Stochastic Neighbor Embedding (SNE) (2003) [1] is a manifold learning and dimensionality reduction method which can be used for feature extraction [2].
- It has a probabilistic approach. It fits the data in the embedding space locally hoping to preserve the global structure of data [3].
- The idea of SNE is to <u>consider every point as neighbors of other points with some</u> probability where the closer points are neighbors with higher probability. Therefore, rather than considering <u>k nearest neighbors in a binary manner</u> (whether being neighbors or not), it considers <u>neighbors in a stochastic way</u> (for how probable it is to be neighbors).
- It tries to preserve the probability of neighborhoods in the low-dimensional embedding space.
- There exist some other similar probabilistic dimensionality reduction methods which make use of <u>Gaussian distribution</u> for neighborhood. Some examples are <u>Neighborhood</u>
 <u>Component Analysis (NCA)</u> [4], <u>deep NCA</u> [5], and <u>Proxy-NCA</u> [6].
- SNE uses the Gaussian distribution for neighbors in both the input and embedding spaces. The Student-t distributed SNE, or so-called t-SNE [7], considers the Student-t and Gaussian distributions in the input and embedding spaces, respectively. The reason of using Student-t distribution in t-SNE is because of its heavier tails so it can include more information from the high-dimensional data.
- <u>t-SNE</u> is one of the <u>state-of-the-art methods</u> for <u>data visualization</u>; for example, it has been used for <u>DNA and single-cell data visualization</u> [8].
- The goal of SNE is to embed the high-dimensional data {x_i}_{i=1}ⁿ into the lower dimensional data {y_i}_{i=1}ⁿ where n is the number of data points. We denote the dimensionality of high- and low-dimensional spaces by d and p, respectively, i.e. x_i ∈ ℝ^d and y_i ∈ ℝ^p. We usually have p ≪ d. For data visualization, we have p ∈ {2,3}.

• In **SNE** (2003) [1], we consider a **Gaussian probability** around every point x_i where the distribution is for probability of accepting any other point as the neighbor of x_i ; the farther points are neighbors with less probability. Hence, the variable is distance, denoted by $d \in \mathbb{R}$, and the Gaussian probability is:



where the mean of distribution is assumed to be zero.

• The fixed multiplier $\frac{1}{\sqrt{2\pi\sigma^2}}$ can be **(roppe)**, however, $\exp(-d^2/2\sigma^2)$ does not add (integrate) to one and thus it is not a probability density function. In order to tackle this problem, we can do a trick and divide $\exp(-d^2/2\sigma^2)$ by the summation of all possible values of $\exp(-d^2/2\sigma^2)$ to have a softmax function. Therefore, the probability that the point $x_i \in \mathbb{R}^d$ takes $x_j \in \mathbb{R}^d$ as its neighbor is:

where:

SNE and t-SNE

 $\mathbb{R} \ni d_{ij}^2 := \frac{||\boldsymbol{x}_i - \boldsymbol{x}_j||_2^2}{|\boldsymbol{x}_i|_2}$

(3)

(1)

- The σ_i² is the variance which we consider for the Gaussian distribution used for the x_i. It can be set to a fixed number or determined by a binary search to make the entropy of distribution some specific value [1]. Note that according to the distribution of data in the input space, the best value for the variance of Gaussian distributions can be found.
- In the low-dimensional embedding space, we again consider a Gaussian probability distribution for the point $y_i \in \mathbb{R}^p$ to take $y_i \in \mathbb{R}^p$ as its neighbor:

$$\left[\mathbb{R} \ni \left[\overline{q_{ij}} \right] := \frac{\exp(-z_{ij}^2)}{\sum_{k \neq i} \exp(-z_{ik}^2)}, \right]$$
(4)

where:

$$\boxed{\mathbb{R} \ni z_{ij}^2 := ||\mathbf{y}_i - \mathbf{y}_j||_2^2} \qquad (5)$$

It is noteworthy that the variance of distribution is not used (or is set to σ_i² = 0.5 to cancel 2 in the denominator) because the variance of distribution in the embedding space is the choice of algorithm.

 We want the probability distributions in both the input and embedded spaces to be as similar as possible; therefore, the cost function to be minimized can be summation of the Kullback-Leibler (KL) divergences [9] over the n points:

$$\mathbb{R} \ni c_1 := \sum_{i=1}^n \mathsf{KL}(P_i||Q_i) = \sum_{i=1}^n \sum_{j=1, j \neq i}^n p_{ij} \log(\frac{p_{ij}}{q_{ij}}), \tag{6}$$

where p_{ij} and q_{ij} are the Eqs. (2) and (4).

- Note that divergences other than the KL divergence can be used for the SNE optimization; e.g., see [10].
- The gradient of c_1 with respect to (\mathbf{y}_i) is:

$$\underbrace{\mathbb{R}^{p}}_{p} \ni \frac{\partial c_{1}}{\partial \mathbf{y}_{i}} = 2 \sum_{j=1}^{n} (p_{ij} - q_{ij} + p_{ji} - q_{ji}) (\mathbf{y}_{i} - \mathbf{y}_{j}), \quad \boldsymbol{\leftarrow} \qquad (7)$$

where p_{ij} and q_{ij} are the Eqs. (2) and (4), and $p_{ii} = q_{ii} = 0$.

 For proof of this, refer to our tutorial "<u>Stochastic neighbor embedding with Gaussian and</u> student-t distributions: Tutorial and survey" [11] or our textbook.

• The update of the embedded point y_i is done by gradient descent. Every iteration is:

$$\begin{cases} \Delta \mathbf{y}_{i}^{(t)} := -\eta \frac{\partial c_{1}}{\partial \mathbf{y}_{i}} + \alpha(t) \Delta \mathbf{y}_{i}^{(t-1)}, \\ \mathbf{y}_{i}^{(t)} := \mathbf{y}_{i}^{(t-1)} + \Delta \mathbf{y}_{i}^{(t)}, \end{cases}$$
(8)

where momentum is used for better convergence [12].

 The α(t) is the momentum. It can be smaller for initial iterations and larger for further iterations. For example, we can have [7]:

$$\alpha(t) := \begin{cases} 0.5 & t < 250, \\ 0.8 & t \ge 250. \end{cases}$$
(9)

In the original paper of SNE [1], the momentum term is not mentioned but it is suggested in [7].

- The η is the learning rate which can be a small positive constant (e.g., η = 0.1) or can be updated adaptively according to [13].
- Moreover, in both [1] and [7], it is mentioned that in SNE we <u>should add some Gaussian</u> noise (random jitter) to the solution of the first iterations before going to the next iterations. It helps avoiding the local optimum solutions.

 In <u>symmetric SNE (2008)</u> [7], we consider a Gaussian probability around every point x_i. The probability that the point x_i ∈ ℝ^d takes x_i ∈ ℝ^d as its neighbor is:

where:

$$\mathbb{R} \ni p_{ij} := \frac{\exp(-d_{ij}^2)}{\sum_{k \neq j} \exp(-d_{kj}^2)}, \quad (10)$$

$$\mathbb{R} \ni d_{ij}^2 := \frac{||\mathbf{x}_i - \mathbf{x}_j||_2^2}{2\sigma_i^2}. \quad (11)$$

• Note that the **denominator** of Eq. (10) for all points is fixed and thus it is symmetric for *i* and *j*. Compare this with Eq. (2):

which is not symmetric.

• The Eq. (10):

$$\mathbb{R} \ni p_{ij} := \underbrace{\frac{\exp(-d_{ij}^2)}{\sum_{k \neq l} \exp(-d_{kl}^2)}}_{\text{exp}(-d_{kl}^2)},$$

has a problem with outliers. If the point x_i is an outlier, its p_{ij} will be extremely small because the denominator is fixed for every point and numerator will be small for the outlier

• However, If we use Eq. (2) for p_{ij} :

$$\mathbb{R}
i p_{ij} := rac{\exp(-d_{ij}^2)}{\sum_{k
eq i} \exp(-d_{ik}^2)},$$

the denominator for all the points is not the same and therefore, the denominator for an outlier will also be small waving out the problem of small numerator.

• For this mentioned problem, we do not use Eq. (10) and rather we use:

where:

$$\mathbb{R} \ni p_{ij} := \underbrace{\frac{p_{i|j} + p_{j|i}}{2n}}_{\sum_{k \neq j} \exp(-d_{ij}^{2})}, \quad (12)$$
is the probability that $x_i \in \mathbb{R}^d$ takes $x_i \in \mathbb{R}^d$ as its neighbor.
$$(13)$$

where.

 In the low-dimensional embedding space, we consider a Gaussian probability distribution for the point y_i ∈ ℝ^p to take y_j ∈ ℝ^p as its neighbor and we make it symmetric (fixed denominator for all points):

where:

$$\mathbb{R} \ni z_{ij}^2 := ||\boldsymbol{y}_i - \boldsymbol{y}_j||_2^2.$$
(15)

• Note that the Eq. (14) does not have the problem of outliers as in Eq. (10) because even for an outlier, the embedded points are initialized close together and not far.

 We want the probability distributions in both the input and embedded spaces to be as similar as possible; therefore, the cost function to be minimized can be summation of the Kullback-Leibler (KL) divergences [9] over the n points:

$$\longrightarrow \mathbb{R} \ni c_2 := \sum_{i=1}^n \mathsf{KL}(P_i||Q_i) = \sum_{i=1}^n \sum_{j=1, j \neq i}^n p_{ij} \log(\frac{p_{ij}}{q_{ij}}), \tag{16}$$

where p_{ij} and q_{ij} are the Eqs. (12) and (14).

The gradient of c₂ with respect to y_i is:

$$\longrightarrow \mathbb{R}^{p} \ni \frac{\partial c_{2}}{\partial \mathbf{y}_{i}} = 4 \sum_{j=1}^{n} (p_{ij} - q_{ij}) (\mathbf{y}_{i} - \mathbf{y}_{j}), \qquad (17)$$

where p_{ij} and q_{ij} are the Eqs. (12) and (14), and $p_{ii} = q_{ii} = 0$.

 For proof of this, refer to our tutorial "<u>Stochastic neighbor embedding with Gaussian and</u> student-t distributions: Tutorial and survey" [11] or our textbook.



The Crowding Problem

- In SNE [1], we are considering Gaussian distribution for both input and embedded spaces.
- That is okay for the input space because it already has a high dimensionality.
- However, when we embed the high-dimensional data into a low-dimensional space, it is very hard to fit the information of all the points in the same neighborhood area.
- For better clarification, suppose the dimensionality is like the size of a room, as depicted in this figure. In high dimensionality, we have a large hall including a huge crowd of people. Now, we want to fit all the people into a small room; of course, we cannot! This problem is referred to as the **crowding problem**.



- The main idea of <u>t-SNE</u> (2008) [7] is addressing the <u>crowding problem</u> which exists in SNE [1].
- In the example of fitting people in a room, t-SNE enlarges the room to solve the crowding problem (see the figure).
- Therefore, in the formulation of t-SNE, we use <u>Student-t distribution [14] rather than</u> <u>Gaussian distribution</u> for the <u>low-dimensional embedded space</u>.
- This is because the Student-t distribution has <u>heavier tails</u> than Gaussian distribution, which is like a larger room, and can fit the information of high dimensional data in the low dimensional embedding space.
- As we will see later, the q_{ij} in t-SNE is:

$$(q_{j}) = \underbrace{(1 + z_{jj}^{2})^{-1}}_{\sum_{k \neq l} (1 + z_{kl}^{2})^{-1}}$$

which is based on the standard Cauchy distribution:

$$f(z) = \frac{1}{\pi(1+z^2)}, \qquad = f(z) + f$$

where π is canceled from the numerator and the normalizing denominator in q_{ij} (similar to the technique of **softmax**).

• If the Student-t distribution [14] with the general degrees of freedom δ is used, we would have:

$$f(z) = \frac{\Gamma(\frac{\delta+1}{2})}{\sqrt{\delta} \times \pi} \Gamma(\frac{\delta}{2}) (1 + \frac{z^2}{\delta})^{-\frac{\delta+1}{2}},$$
(19)

where Γ is the gamma function.

• Cancelling out the scaling factors from the numerator and denominator, we would have [15]:

$$q_{ij} = \frac{(1 + z_{ij}^2/\delta)^{-(\delta+1)/2}}{\sum_{k \neq l} (1 + z_{kl}^2/\delta)^{-(\delta+1)/2}}.$$
(20)

• However, as the first degree of freedom has the heaviest tails amongst different degrees of freedom, it is the most suitable for the crowding problem; hence, we use the first degree of freedom which is the Cauchy distribution. Note that the t-SNE algorithm, which uses the Cauchy distribution, may also be called the Cauchy-SNE.

• Later, t-SNE with general degrees of freedom was proposed [15].

In t-SNE [7], we consider a Gaussian probability around every point x_i in the input space because the crowding problem does not exist in the high dimensional data. The probability that the point x_i ∈ ℝ^d takes x_i ∈ ℝ^d as its neighbor is:

$$\mathbb{R} \ni p_{j|i} := \frac{\exp(-d_{ij}^2)}{\sum_{k \neq i} \exp(-d_{ik}^2)},$$
(21)

where:

$$\mathbb{R} \ni d_{ij}^2 := \frac{||\mathbf{x}_i - \mathbf{x}_j||_2^2}{2\sigma_i^2}.$$
(22)

Note that Eq. (21) is not symmetric for i and j because of the denominator. We take the symmetric p_{ij} as the scaled average of p_{i|j} and p_{i|i}:

$$\mathbb{R} \ni p_{ij} := \frac{p_{i|j} + p_{j|i}}{2n}.$$
(23)

In the low-dimensional embedding space, we consider a Student's t-distribution with one degree of freedom (Cauchy distribution) for the point y_i ∈ ℝ^p to take y_j ∈ ℝ^p as its neighbor:

$$\mathbb{R} \ni q_{ij} := \frac{(1+z_{ij}^2)^{-1}}{\sum_{k \neq l} (1+z_{kl}^2)^{-1}},$$
(24)

where:

$$\mathbb{R} \ni z_{ij}^2 := || \mathbf{y}_i - \mathbf{y}_j ||_2^2.$$
(25)

SNE and t-SNE

 We want the probability distributions in both the input and embedded spaces to be as similar as possible; therefore, the cost function to be minimized can be summation of the Kullback-Leibler (KL) divergences [9] over the n points:

$$\mathbb{R} \ni c_3 := \sum_{i=1}^n \mathsf{KL}(P_i||Q_i) = \sum_{i=1}^n \sum_{j=1, j \neq i}^n p_{ij} \log(\frac{p_{ij}}{q_{ij}}), \tag{26}$$

where p_{ij} and q_{ij} are the Eqs. (23) and (24).

The gradient of c₃ with respect to y_i is:

$$\longrightarrow \underbrace{\frac{\partial c_3}{\partial \mathbf{y}_i}}_{ij} = 4 \sum_{j=1}^n (\mathbf{p}_{ij} - \mathbf{q}_{ij}) (1 + ||\mathbf{y}_i - \mathbf{y}_j||_2^2)^{-1} (\mathbf{y}_i - \mathbf{y}_j), \qquad (27)$$

where p_{ij} and q_{ij} are the Eqs. (23) and (24), and $p_{ii} = q_{ii} = 0$.

• Proof: Proof is according to [7]. Let:

By changing y_i, we only have change impact in z_{ij} and z_{ji} for all j's. According to chain rule, we have:

 $\mathbb{R}^{p} \ni \underbrace{\frac{\partial c_{3}}{\partial \mathbf{y}_{i}}}_{j} = \underbrace{\sum_{j} \underbrace{\frac{\partial c_{j}}{\partial r_{ij}} + \frac{\partial c_{3}}{\partial r_{jj}} \frac{\partial r_{jj}}{\partial \mathbf{y}_{i}}}_{j} + \underbrace{\frac{\partial c_{3}}{\partial r_{jj}} \frac{\partial r_{jj}}{\partial \mathbf{y}_{i}}}_{j}$ • According to Eq. (28), we have: $\underbrace{r_{ij} = ||\mathbf{y}_{i} - \mathbf{y}_{j}||_{2}^{2}}_{ij} \Longrightarrow \frac{\frac{\partial r_{ij}}{\partial \mathbf{y}_{i}}}_{ij} = 2(\mathbf{y}_{i} - \mathbf{y}_{j}),$ $\underbrace{r_{ji} = ||\mathbf{y}_{j} - \mathbf{y}_{i}||_{2}^{2}}_{ij} = ||\mathbf{y}_{i} - \mathbf{y}_{j}||_{2}^{2} \Longrightarrow \frac{\partial r_{ji}}{\partial \mathbf{y}_{i}} = 2(\mathbf{y}_{i} - \mathbf{y}_{j}).$ • Therefore:

$$\frac{\partial c_3}{\partial \mathbf{y}_i} = 2 \sum_j \left(\frac{\partial c_3}{\partial r_{ij}} \right) + \frac{\partial c_3}{\partial r_{ji}} \right) (\mathbf{y}_i - \mathbf{y}_j).$$
(29)

. .

• The cost function can be re-written as:

 $c_{3} = \sum_{k} \sum_{l \neq k} p_{kl} \log(\frac{p_{kl}}{q_{kl}}) = \sum_{k \neq l} p_{kl} \log(\frac{p_{kl}}{q_{kl}})$ $= \sum_{k \neq l} (p_{kl} \log(p_{kl}) - p_{kl} \log(q_{kl})),$

whose first term is a constant with respect to q_{kl} and thus to r_{kl} . • We have:

$$\mathbb{R} \ni \frac{\partial c_3}{\partial r_{ij}} = -\sum_{k \neq l} p_{kl} \frac{\partial (\log(q_{kl}))}{\partial r_{ij}}$$

• According to Eq. (24):

$$\mathbb{R} \ni q_{j} := \frac{(1+z_{ij}^{2})^{-1}}{\sum_{k \neq j} (1+z_{kl}^{2})^{-1}}$$

$$q_{kl} := \frac{(1 + c_{kl}^2)^{-1}}{\sum_{m \neq f} (1 + c_{mf}^2)^{-1}}, = \frac{(1 + c_{kl})^{-1}}{\sum_{m \neq f} (1 + c_{mf})^{-1}}.$$

We take the denominator of q_{kl} as:

$$\beta := \sum_{\substack{m \neq f}} (1 + z_{mf}^2)^{-1} = \sum_{\substack{m \neq f}} (1 + r_{mf})^{-1}.$$
(30)

SNE and t-SNE



• The $q_{kl}\beta$ is:



• Therefore, we have:

$$\therefore \quad \frac{\partial c_3}{\partial r_{ij}} = -\sum_{k\neq l} p_{kl} \left[\frac{1}{q_{kl}\beta} \frac{\partial ((1+r_{kl})^{-1})}{\partial r_{ij}} - \frac{1}{\beta} \frac{\partial \beta}{\partial r_{ij}} \right].$$

SNE and t-SNE



• We found:

where (a)

gives us:

٠

which is the gradient mentioned before. Q.E.D.

SNE and t-SNE

• The update of the embedded point y_i is done by gradient descent whose every iteration is as Eq. (8) where c_1 is replaced by c_3 :

$$\begin{cases} \Delta \mathbf{y}_i^{(t)} := -\eta \frac{\partial c_1}{\partial \mathbf{y}_i} + \alpha(t) \Delta \mathbf{y}_i^{(t-1)}, \\ \mathbf{y}_i^{(t)} := \mathbf{y}_i^{(t-1)} + \Delta \mathbf{y}_i^{(t)}. \end{cases}$$

- For t-SNE, there is <u>no need to add jitter</u> to the solution of initial iterations [7] because it is more robust than SNE.
- In t-SNE, it is better to **multiply all** *p_{ij}*'s by a constant (e.g., 4) in the **initial iterations**:

$$p_{ij} := p_{ij} \times 4, \tag{31}$$

which is called <u>early exaggeration</u>. This heuristic helps the optimization focus on the large p_{ij} 's (close neighbors) more in the early iterations.

- This is because large p_{ij} 's are affected more by multiplying by 4 than the small p_{ij} 's.
- After the neighbours are embedded close to one another, we are free not to do this
 multiplication any more and let far-away points be embedded using the probabilities
 without multiplication. Note that the early exaggeration is optional and not mandatory.

- We can have general degrees of freedom for Student-t distribution in t-SNE [15].
- As we saw in Eqs. (19) and (20), we can have any degrees of freedom for q_{ij} (note that α is a positive integer). We repeat Eq. (20) here for more convenience:

$$q_{ij} = \frac{(1+z_{ij}^2/\delta)^{-(\delta+1)/2}}{\sum_{k \neq l} (1+z_{kl}^2/\delta)^{-(\delta+1)/2}}.$$
(32)

- If $\delta \to \infty$, the Student-t distribution formulated in Eq. (19) tends to Gaussian distribution used in SNE [1].
- SNE and t-SNE use degrees $\delta \to \infty$ and $\delta = 1$ in Eq. (32), respectively.



Acknowledgment

- Some slides are based on our tutorial paper: "Stochastic neighbor embedding with Gaussian and student-t distributions: Tutorial and survey" [11]
- For more information on SNE and t-SNE, refer to our tutorial paper [11].
- Some slides of this slide deck are inspired by teachings of Prof. Ali Ghodsi at University of Waterloo, Department of Statistics.
- The code of SNE and t-SNE in my GitHub: https://github.com/bghojogh/SNE-tSNE
- t-SNE in sklearn: https: //scikit-learn.org/stable/modules/generated/sklearn.manifold.TSNE.html

References

- G. E. Hinton and S. T. Roweis, "Stochastic neighbor embedding," in Advances in neural information processing systems, pp. 857–864, 2003.
- [2] B. Ghojogh, M. N. Samad, S. A. Mashhadi, T. Kapoor, W. Ali, F. Karray, and M. Crowley, "Feature selection and feature extraction in pattern analysis: A literature review," arXiv preprint arXiv:1905.02845, 2019.
- [3] L. K. Saul and S. T. Roweis, "Think globally, fit locally: unsupervised learning of low dimensional manifolds," *Journal of machine learning research*, vol. 4, no. Jun, pp. 119–155, 2003.
- [4] J. Goldberger, G. E. Hinton, S. T. Roweis, and R. R. Salakhutdinov, "Neighbourhood components analysis," in *Advances in neural information processing systems*, pp. 513–520, 2005.
- [5] X. Liu, X. Yang, M. Wang, and R. Hong, "Deep neighborhood component analysis for visual similarity modeling," ACM Transactions on Intelligent Systems and Technology (TIST), vol. 11, no. 3, pp. 1–15, 2020.
- [6] Y. Movshovitz-Attias, A. Toshev, T. K. Leung, S. loffe, and S. Singh, "No fuss distance metric learning using proxies," in *Proceedings of the IEEE International Conference on Computer Vision*, pp. 360–368, 2017.
- [7] L. van der Maaten and G. Hinton, "Visualizing data using t-SNE," *Journal of machine learning research*, vol. 9, no. Nov, pp. 2579–2605, 2008.

References (cont.)

- [8] D. Kobak and P. Berens, "The art of using t-SNE for single-cell transcriptomics," *Nature communications*, vol. 10, no. 1, pp. 1–14, 2019.
- [9] S. Kullback, *Information theory and statistics*. Courier Corporation, 1997.
- [10] D. J. Im, N. Verma, and K. Branson, "Stochastic neighbor embedding under f-divergences," arXiv preprint arXiv:1811.01247, 2018.
- [11] B. Ghojogh, A. Ghodsi, F. Karray, and M. Crowley, "Stochastic neighbor embedding with gaussian and student-t distributions: Tutorial and survey," arXiv preprint arXiv:2009.10301, 2020.
- [12] N. Qian, "On the momentum term in gradient descent learning algorithms," Neural networks, vol. 12, no. 1, pp. 145–151, 1999.
- [13] R. A. Jacobs, "Increased rates of convergence through learning rate adaptation," Neural networks, vol. 1, no. 4, pp. 295–307, 1988.
- [14] W. S. Gosset (Student), "The probable error of a mean," Biometrika, pp. 1–25, 1908.
- [15] L. van der Maaten, "Learning a parametric embedding by preserving local structure," in Artificial Intelligence and Statistics, pp. 384–391, 2009.
- [16] N. Pezzotti, "Dimensionality-reduction algorithms for progressive visual analytics," 2019.