

# Unit 18: t-SNE Algorithm

#### Final Topic of the Course:

We will wrap up the new content in this course by discussing the algorithm and the theory behind the t-SNE algorithm.

# Purpose of this Lecture:

We will introduce the algorithm and the theory behind the t-SNE algorithm.

In this lecture we will cover the following topics.

- Stochastic neighbor embedding.
- Drawbacks of stochastic neighbor embedding.
- What's different in SNE algorithm vs. t-SNE algorithm?
- T-SNE algorithm general
- Full t-sne algorithm with gradient descent algorithm.

#### **Additional Resources**

van der Maaten, L.J.P.; Hinton, G.E. Visualizing High-Dimensional Data Using t-SNE. Journal of Machine Learning Research 9:2579-2605, 2008. <u>https://www.jmlr.org/papers/volume9/vandermaaten08a/vandermaaten08a.pdf</u>

# STOCHASTIC NEIGHBOR EMBEDDING (SNE)

## Input:

- <u>Dataset</u>:  $X_{m \times n} = \begin{bmatrix} x_{1*} \\ x_{2*} \\ \vdots \\ x_{m*} \end{bmatrix}$  comprised of *m* objects, where each object has a complete set of n attributes.
- <u>Perplexity Values</u>: *Perplexity*(*P<sub>i</sub>*) for *i=1,...,m*

## <u>Algorithm</u>

• Step 1: Define a projected/mapped 2-d coordinates matrix of decision

 $\underline{variables}: Y_{m \times 2} = \begin{bmatrix} y_{1*} \\ y_{2*} \\ \vdots \\ y_{m*} \end{bmatrix}.$ 

• <u>Step 2</u>: Create a similarity matrix  $P = [p_{j|i}]$  between each of the objects in  $X_{m \times n}$ .

#### (i,j) entries represents:

similarity between  $x_{i*}$  to  $x_{j*}$ 

Read ahead	
for how we	
pick $\sigma_i$ .	

#### (i,j) entries mathematically:

$$p_{j|i} = \frac{\exp\left(-\frac{\|x_{i^*} - x_{j^*}\|^2}{2\sigma_i^2}\right)}{\sum_{k \neq i} \exp\left(-\frac{\|x_{i^*} - x_{k^*}\|^2}{2\sigma_i^2}\right)}$$

#### (j,i) entries graphical interpretation:

$$p_{j|i} = \frac{\exp\left(-\frac{\left|\left|x_{i^{*}} - x_{j^{*}}\right|\right|^{2}}{2\sigma_{i}^{2}}\right)}{\sum_{k \neq i} \exp\left(-\frac{\left|\left|x_{i^{*}} - x_{k^{*}}\right|\right|^{2}}{2\sigma_{i}^{2}}\right)} = \frac{\frac{1}{\sigma_{i}\sqrt{2\pi}} \exp\left(-\frac{\left(dist(x_{i^{*}}, x_{j^{*}}) - 0\right)^{2}}{2\sigma_{i}^{2}}\right)}{\sum_{k \neq i} \frac{1}{\sigma_{i}\sqrt{2\pi}} \exp\left(-\frac{\left(dist(x_{i^{*}}, x_{k^{*}}) - 0\right)^{2}}{2\sigma_{i}^{2}}\right)} = \frac{f(dist(x_{i^{*}}, x_{j^{*}}))}{\sum_{k \neq i} f(dist(x_{i^{*}}, x_{k^{*}}))}$$

Interpretation of f(x):

- f(x) is the pdf of the \_\_\_\_\_ distribution with:
  - mean = \_\_\_\_\_
  - Standard deviation = \_\_\_\_\_

Interpretation of  $f(dist(x_{i*}, x_{j*}))$ :

 $f(dist(\mathbf{x}_{i*}, \mathbf{x}_{j*})) = g(dist(\mathbf{x}_{i*}, \mathbf{x}_{j*}))$ 

g(dist(x<sub>i\*</sub>, x<sub>j\*</sub>)) is the \_\_\_\_\_ curve, centered at \_\_\_\_\_ with standard deviation \_\_\_\_\_

#### **<u>Ex</u>**: Generate the similarity matrix for the three 1-D data points shown below.



### Properties of $P = [p_{i|i}]$

- This matrix \_\_\_\_\_\_ to be symmetric.
- The rows in this matrix form a \_\_\_\_\_\_.

# <u>Step 3</u>: Create a similarity matrix $Q = [q_{j|i}]$ between each of the <u>decision</u> <u>variable objects</u> in $Y_{m\times 2}$ .

- (j,i) entries represents:
   similarity between y<sub>i</sub>\* to y<sub>i</sub>\*
- (j,i) entries mathematically:  $q_{j|i} = \frac{\exp\left(-\left||y_{i*}-y_{j*}|\right|^{2}\right)}{\sum_{k\neq i} \exp\left(-\left||y_{i*}-y_{k*}|\right|^{2}\right)}$
- (j,i) entries graphical interpretation:

$$q_{j|i} = \frac{\exp\left(-\left|\left|\boldsymbol{y}_{i^{*}} - \boldsymbol{y}_{j^{*}}\right|\right|^{2}\right)}{\sum_{k \neq i} \exp\left(-\left|\left|\boldsymbol{y}_{i^{*}} - \boldsymbol{y}_{k^{*}}\right|\right|^{2}\right)} = \frac{\frac{1}{\left(\frac{1}{2}\right)\sqrt{2\pi}} \exp\left(-\frac{\left(dist(\boldsymbol{y}_{i^{*}}, \boldsymbol{y}_{j^{*}}) - 0\right)^{2}}{2 \cdot \left(\frac{1}{2}\right)^{2}}\right)}{\sum_{k \neq i} \frac{1}{\left(\frac{1}{2}\right)\sqrt{2\pi}} \exp\left(-\frac{\left(dist(\boldsymbol{y}_{i^{*}}, \boldsymbol{y}_{k^{*}}) - 0\right)^{2}}{2 \cdot \left(\frac{1}{2}\right)^{2}}\right)} = \frac{h(dist(\boldsymbol{y}_{i^{*}}, \boldsymbol{y}_{j^{*}}))}{\sum_{k \neq i} h(dist(\boldsymbol{y}_{i^{*}}, \boldsymbol{y}_{k^{*}}))}$$

Interpretation of h(x):

- h(x) is the pdf of the \_\_\_\_\_ distribution with:
  - mean = \_\_\_\_\_
  - Standard deviation = \_\_\_\_\_

#### Interpretation of $h(dist(y, y_{j*}))$ :

 $h(dist(\mathbf{y}_{i*}, \mathbf{y}_{j*})) = \overline{h}(dist(\mathbf{y}_{i*}, \mathbf{y}_{j*}))$ 

*h*(*dist*(*y*<sub>*i*\*</sub>, *y*<sub>*j*\*</sub>)) is the \_\_\_\_\_ curve, centered at \_\_\_\_\_ with standard deviation \_\_\_\_\_

### Properties of $Q = [q_{j|i}]$

- This matrix \_\_\_\_\_\_ to be symmetric.
- The rows in this matrix form a \_\_\_\_\_\_.

#### Local Structure Preservation Idea:

• If the mapped points in  $Y_{m\times 2}$  correctly model the *similarity* between the high dimensional data points in  $X_{m\times n}$ , then we would expect:

• <u>Step 4</u>: Find the optimal values of the coordinates in  $Y_{m\times 2}$  that minimize the sum of the Kullback-Leibler (KL) divergence over all data points.

<u>Optimization Problem:</u>  $min_{Y_{m\times 2}}C = min_{Y_{m\times 2}}\sum_{i} KL(\boldsymbol{P}_{i*}||\boldsymbol{Q}_{i*}) = \sum_{i}\sum_{j} p_{j|i}\log \frac{p_{j|i}}{q_{j|i}}$ 

Common Algorithm: Gradient descent methods



• Scenario a: Widely separated data points  $x_{i*}$  and  $x_{j*}$  and nearby map points  $y_{i*}$  and  $y_{j*}$ , Or

• Scenario b: nearby data points  $x_{i*}$  and  $x_{j*}$  and widely separated map points  $y_{i*}$  and  $y_{j*}$ ?

<u>Additional Consideration</u>: How to pick  $\sigma_i$ ?



#### <u>Goal</u>: Pick $\sigma_i$ to Adapt to Different Cluster Sparsities

Clusters of different sparsities need different values of  $\sigma_i$  to effectively preserved the local pairwise distances of all the objects in the cluster.

- If we have a cluster that is more dense, then we want the objects in this cluster to have values of σ<sub>i</sub> that are \_\_\_\_\_\_ when calculating p<sub>j|i</sub>.
- If we have a cluster that is more sparse, then we want the objects in this cluster to have values of σ<sub>i</sub> that are \_\_\_\_\_\_ when calculating p<sub>i|i</sub>.

<u>Solution</u>: Perform a binary search for  $\sigma_i$  that satisfies:

 $= Perplexity(P_i) = 2^{-\sum_j p_{j|i} \log_2 p_{j|i}}$ 

What *Perplexity*(*P<sub>i</sub>*)represents:

<u>What  $-\sum_{j} p_{j|i} \log_2 p_{j|i}$  represents:</u>

$$\underline{\qquad} \leq -\sum_{j} p_{j|i} \log_2 p_{j|i} \leq \underline{\qquad}$$

<u>Ex</u>: Let's assume for now, that we have chosen a perplexity of 1. Use the method described above to estimate a value for  $\sigma_i$  corresponding to an observation in cluster 1 below.



<u>Ex</u>: Let's assume for now, that we have chosen a perplexity that is equal to the total number of observation m in the dataset above. Use the method described above to estimate a value for  $\sigma_i$  corresponding to an observation in cluster 2 below.



Good guesses for  $Perplexity(P_i)$ :

### DRAWBACKS OF STOCHASTIC NEIGHBOR EMBEDDING (SNE)

- 1. Because of non-symmetric KL-Divergence Function:
  - a. Focuses more on preserving the local structure.
- 2. Because of non-symmetric distances matrices:
  - a. Cost function of SNE is computationally inefficient. The gradient descent algorithm needs to be run several times to select the right initial parameters such that the algorithm is less likely to get stuck in poor local minimum.
- 3. Because of Gaussian probability in mapped distance matrix Q:
  - a. "Crowding problem": Some types of datasets **X** have the property that the available area (in 2dimensions) to map \_\_\_\_\_\_ distant points isn't large enough to as the amount of available area that accurately maps \_\_\_\_\_\_ points. So \_\_\_\_\_\_ distances tend to be mapped much \_\_\_\_\_\_ than in the original dataset.
  - b. The exponential function in the cost function of the gradient descent algorithm makes for slower run time.



https://scikit-learn.org/stable/modules/generated/sklearn.datasets.make swiss roll.html

# WHAT'S DIFFERENT IN T-SNE ALGORITHM VS. SNE ALGORITHM?

1. T-SNE uses a version of P and Q (as used in SNE).

2. Uses a \_\_\_\_\_\_ (as opposed to Gaussian distribution as used in SNE) to measure distance between two points in the low dimensional space.

### T-SNE ALGORITHM – GENERAL

### Input:

- <u>Dataset</u>:  $X_{m \times n} = \begin{bmatrix} x_{1*} \\ x_{2*} \\ \vdots \\ x_{m*} \end{bmatrix}$  comprised of *m* objects, where each object has a complete set of n attributes.
- Parameters Used to Generate Similarity Matrix P:
   *Perplexity*

### <u>Algorithm</u>

- <u>Step 1</u>: Define projected/mapped 2-d coordinates matrix of decision variables:  $Y_{m\times 2} = [y_{1*}]$ 
  - *y*<sub>2∗</sub> : *y*<sub>m∗</sub>
- <u>Step 2</u>: Create a similarity matrix  $P = [p_{ij}]$  between each of the objects in  $X_{m imes n}$ .
  - (j,i) entries represents: similarity between  $x_{j*}$  to  $x_{i*}$
  - (j,i) entries mathematically:

• 
$$p_{j|i} = rac{\exp\left(-\frac{\left|\left|x_{i*}-x_{j*}\right|\right|^2}{2\sigma_i^2}\right)}{\sum_{k\neq i} \exp\left(-\frac{\left|\left|x_{i*}-x_{k*}\right|\right|^2}{2\sigma_i^2}\right)}$$

• 
$$p_{i|j} = \frac{\exp\left(-\frac{\left\||x_{i*} - x_{j*}\|\right\|^2}{2\sigma_j^2}\right)}{\sum_{k \neq j} \exp\left(-\frac{\left\||x_{j*} - x_{k*}\|\right\|^2}{2\sigma_j^2}\right)}$$

• 
$$p_{ij} = \frac{1}{2m} \left( p_{j|i} + p_{i|j} \right)$$

• Is  $P = [p_{ij}]$  symmetric?

Finding  $\sigma_i$ :

Perform a binary search for  $\sigma_i$  that satisfies:  $Perplexity = 2^{-\sum_j p_{j|i} \log_2 p_{j|i}}$ 

- Step 3: Create a similarity matrix  $m{Q} = m{[q_{ij}]}$  between each of the objects in  $Y_{m imes 2}$ .
  - (j,i) entries represents: similarity between  $y_{j*}$  to  $y_{i*}$

• (j,i) entries mathematically: 
$$q_{ij} = \frac{\left(1 + \left\|y_{i*} - y_{j*}\right\|^2\right)^{-1}}{\sum_{k \neq l} \left(1 + \left\|y_{k*} - y_{l*}\right\|^2\right)^{-1}}$$

• <u>Is</u>  $Q = [q_{ij}]$  symmetric?

Why student t-distribution with df=1?

- \_\_\_\_\_tails than Gaussian.
- Pairs of mapped points with moderate distances will have higher similarity (than if you had used Gaussian.) Helps with the crowding problem.

• (j,i) entries graphical interpretation:

$$q_{ij} = \frac{\left(1 + \left|\left|y_{i*} - y_{j*}\right|\right|^{2}\right)^{-1}}{\sum_{k \neq l} \left(1 + \left|\left|y_{k*} - y_{l*}\right|\right|^{2}\right)^{-1}} = \frac{\frac{\Gamma(\frac{1+1}{2})}{\sqrt{1\pi}\Gamma(\frac{1}{2})} \left(1 + \frac{\left|\left|y_{i*} - y_{j*}\right|\right|^{2}}{1}\right)^{-\frac{1+1}{2}}}{\sum_{k \neq l} \frac{\Gamma(\frac{1+1}{2})}{\sqrt{1\pi}\Gamma(\frac{1}{2})} \left(1 + \frac{\left|\left|y_{i*} - y_{j*}\right|\right|^{2}}{1}\right)^{-\frac{1+1}{2}}} = \frac{f(dist(y_{i*}, y_{j*}))}{\sum_{k \neq l} f(dist(y_{k*}, y_{l*}))}$$

One Interpretation of f(x):

• f(x) is the pdf of the \_\_\_\_\_\_ distribution with df = \_\_\_\_\_.

Another Interpretation of f(x):

- f(x) is the curve of the \_\_\_\_\_\_ distribution with df = \_\_\_\_\_\_, shifted to be centered at \_\_\_\_\_\_.
- **<u>Ex:</u>** Generate the similarity matrix for the three 1-d MAPPED data points shown below.



• <u>Step 4</u>: Find the optimal values of the coordinates in  $Y_{m\times 2}$  that minimize the sum of the Kullback-Leibler (KL) divergence over all data points.

<u>Optimization Problem:</u>  $min_{Y_{m\times 2}}C = min_{Y_{m\times 2}}KL(P||Q) = \sum_{i}\sum_{j}p_{ij}\log\frac{p_{ij}}{q_{ij}}$ 

Common Algorithm: Gradient descent methods

### FULL T-SNE ALGORITHM - WITH GRADIENT DESCENT ALGORITHM

### Input:

- <u>Dataset</u>:  $X_{m \times n} = \begin{vmatrix} x_{1*} \\ x_{2*} \\ \vdots \\ x_{m*} \end{vmatrix}$  comprised of *m* objects, where each object has a complete set of n attributes.
- Parameters Used to Generate Similarity Matrix P:
  - $\circ$  Perplexity
- <u>Parameters Used to Generate Similarity Matrix P</u>:
  - Number of iterations T
  - $\circ$  Learning rate  $\eta$
  - Momentum  $\alpha(t)$

### **Algorithm**

• <u>Step 1</u>: Create a similarity matrix  $P = [p_{ij}]$  between each of the objects in  $X_{m \times n}$ .

$$p_{ij} = \frac{1}{2} \left( p_{j|i} + p_{i|j} \right)$$

$$p_{j|i} = \frac{\exp\left(-\frac{\left\|x_{i*} - x_{j*}\right\|^{2}}{2\sigma_{i}^{2}}\right)}{\sum_{k \neq i} \exp\left(-\frac{\left\|x_{i*} - x_{k*}\right\|^{2}}{2\sigma_{i}^{2}}\right)}$$

- Perform a binary search for  $\sigma_i$  that satisfies:  $Perplexity = 2^{-\sum_j p_{j|i} \log_2 p_{j|i}}$
- <u>Step 2</u>: Randomly sample projected/mapped 2-d coordinates matrix:  $Y_{m \times 2} = \begin{bmatrix} y_{1*} \\ y_{2*} \\ \vdots \\ y \end{bmatrix}$  from

 $N(0, 0.0001 \cdot I_{2 \times 2}).$ 

- For *t*=1 to *T* do:
  - Create a similarity matrix  $Q = [q_{ij}]$  between each of the objects in  $Y_{m \times 2}$ .

• 
$$q_{ij} = \frac{\left(1 + \left||y_{i*} - y_{j*}|\right|^2\right)^{-1}}{\sum_{k \neq l} \left(1 + \left||y_{k*} - y_{l*}|\right|^2\right)^{-1}}$$

Compute cost function gradient as follows (for each i=1,...,m)

• 
$$\frac{\delta C}{\delta y_{i*}} = 4 \sum_{j} (p_{ij} - q_{ij}) (y_{i*} - y_{j*}) (1 + ||y_{i*} - y_{j*}||^2)^{-1}$$
  
Set  $Y_{m \times 2}^{(t)} := Y_{m \times 2}^{(t-1)} + \gamma \frac{\delta C}{\delta Y} + \alpha(t) (Y_{m \times 2}^{(t-1)} - Y_{m \times 2}^{(t-2)})$ 

End do