

T-SNE (week 2 - NL2)

17.1.23.

- ① Biết đổi lại công thức toán t-SNE, t-SNE, có tilt loss và parameter

SNE basic idea:

- "Encode" high dimensional neighborhood information as a distribution
- Find low dimensional points such that their neighborhood distribution is similar
- Intuition: Random walk between data points.
 - High probability to jump to a close point
- How do you measure distance between distribution?
 - Most common measure: KL divergence.

SNE Implementation:

- Consider the neighborhood around an input data point $x_i \in \mathbb{R}^d$
- Imagine that we have a Gaussian distribution centered around x_i
- Then the probability that x_i chooses some other data point " x_j " as its neighbor is proportional with the density under this Gaussian.
- A point closer to x_i will be more likely than one further away

we have the probability that point x_i chooses x_j as its neighbor:

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

As $j \rightarrow i$, $P_{i|i} = 0$, $P_{i|j} \neq P_{j|i}$

Final distribution over pairs is symmetrized:

$$P_{ij} = \frac{1}{2N} (P_{i|j} + P_{j|i}) \rightarrow \begin{matrix} \text{The distribution} \\ \text{of original input} \\ \text{data. [1]} \end{matrix}$$

Perplexity:

- The parameter σ_i sets the size of the neighborhood
 - \rightarrow very low σ_i - all the prob. is in the nearest neighbor
 - \rightarrow very high σ_i - uniform weights.
- Here we set σ_i differently for each data point
- Results depend heavily on σ_i - it depends the neighborhoods we are trying to preserve the local structure.



- For end distribution $p_{j|i}$ (depend on σ_i) we define the perplexity

$$\text{perp}(p_{j|i}) = 2^{-H(p_{j|i})}$$

with $H(p) = -\sum p_i \log(p_i)$ is the entropy.

- If p is uniform over K elements - perplexity = K .
 - low perplexity = small σ
 - high perplexity = large σ
- Values between 5-80 usually work well
- Important parameter - different perplexities can capture different scales in the data.

SNE objective

- Given $x^{(0)}, x^{(1)}, \dots, x^{(n)} \in \mathbb{R}^D$ we define the distribution p_{ij}
 - Goal: Find good embedding $y^{(0)}, \dots, y^{(n)} \in \mathbb{R}^d$ for some $d < D$. (normally 2 or 3) \rightarrow distribution q of reduced dimension.
- $x_1 \xrightarrow{\text{embed}} \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$
- $x_2 \xrightarrow{\text{embed}} \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$
- \vdots
- $x_n \xrightarrow{\text{embed}} \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$
- parameter as defined \rightarrow random initialization
 \downarrow
distribution Q .
- For point $y^{(0)}, y^{(1)}, \dots, y^{(n)} \in \mathbb{R}^d$ ($d < D$) we can define distribution Q similarly the form (notice no σ_i and not symmetric)

$$Q_{ij} = \frac{e^{-\|y^{(i)} - y^{(j)}\|^2}}{\sum_k e^{-\|y^{(i)} - y^{(k)}\|^2}}$$

\Rightarrow Optimize Q close P . Measure KL-divergence \rightarrow to find the embedding (parameter) $y^{(1)}, \dots, y^{(n)} \in \mathbb{R}^d$

Measure the distance between two distribution P and Q (loss func)

$$L(Q) = KL(Q \parallel P) = \sum_{ij} Q_{ij} \log \left(\frac{Q_{ij}}{P_{ij}} \right)$$

KL properties: $KL(Q \parallel P) \geq 0$ and zero only when $Q = P$.
 $KL(Q \parallel P)$ is a convex function

Now, we have P , are looking for $y^{(1)}, \dots, y^{(n)} \in \mathbb{R}^d$ such that
~~at Q will minimize $L(Q) = KL(P \parallel Q)$~~

$$L(Q) = KL(P \parallel Q) = \sum_{ij} P_{ij} \log \left(\frac{P_{ij}}{Q_{ij}} \right) = - \sum_{ij} P_{ij} \log(Q_{ij}) + \text{const}$$

Taking gradient w.r.t $y^{(i)}$

$$\frac{\partial L}{\partial y^{(i)}} = \sum_j (P_{ij} - Q_{ij})(y^{(j)} - y^{(i)})$$

Recall: we defined the similarity between point i and j as a conditional probability:

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

$P_{j|i}$ represents the probability that x_j would pick x_i as its neighbor if neighbors were picked in proportion to their probability density under a Gaussian centered at x_i . We let $P_{i|i} = 0$ [3]

Low dimensional mapping

We wish to obtain $(y_i)_i$ belong to a lower dimensional space, where each y_i represent a point in the plane corresponding to x_i .

We define the similarity between point i and j in this space

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

$$\text{let } q_{i|i} = 0$$

A good map from $(x_i)_i$ to $(y_i)_i$ is a map for which distribution $p_i : j \mapsto p_{j|i}$ and $Q_i : j \mapsto q_{j|i}$ are equal (for all i)

~~to measure the distance~~

$$\text{Let define : } q_{j|i} = \frac{e^{-\|y_i - y_j\|^2}}{\sum_{k \neq i} e^{-\|y_i - y_k\|^2}} = \frac{E_{ij}}{\sum_{k \neq i} E_{ik}} = \frac{E_{ij}}{Z_i}$$

Notice that $E_{ij} = E_{ji}$. The loss function is defined as.

$$C = \sum_{k, l \neq i} p_{elik} \log \frac{p_{elik}}{q_{elik}} = \sum_{k, l \neq i} p_{elik} \log p_{elik} - p_{elik} \log q_{elik}$$

$$= \sum_{k, l \neq i} p_{elik} \log p_{elik} - p_{elik} \log E_{kl} + p_{elik} \log Z_k.$$

We derive wrt to y_i . To make the derivation less cluttered, I will omit the δy_i term at the denominator

$$\frac{\partial C}{\partial y_i} = \left(\sum_{k, l \neq i} p_{elik} \log p_{elik} - p_{elik} \log E_{kl} + p_{elik} \log Z_k \right) \frac{\partial}{\partial y_i}$$

$$= \sum_{k, l \neq i} -p_{elik} \delta \log E_{kl} + \sum_{k, l \neq i} p_{elik} \delta \log Z_k$$

We start with the first term, noting that the derivative is non zero when $j \neq i$, $k=i$ or $l=i$.

$$\sum_{l \neq k} -p_{lk} \delta \log E_{kl} = \sum_{j \neq i} -p_{ji} \delta \log E_{ij} - p_{ii} \delta \log E_{ii}$$

Since $\delta E_{ij} = E_{ij} (-2(y_i - y_j))$ we have:

$$\begin{aligned} \sum_{j \neq i} -p_{ji} \frac{E_{ij}}{E_{ii}} (-2(y_i - y_j)) - p_{ii} \frac{E_{ii}}{E_{ii}} (2(y_j - y_i)) \\ = 2 \sum_{j \neq i} (p_{ji} + p_{ii}) (y_i - y_j) \end{aligned} \quad (4)$$

We conclude with the second term. Since $\sum_{l \neq j} p_{lj} = 1$ and Z_j does not depend on k , we can write (changing variable from l to j to make it more similar to the already computed term)

$$\sum_{j \neq i} p_{kj} \delta \log Z_j = \sum_j \delta \log Z_j$$

The derivative is non-zero when $k=i$ or $j=i$ (also, in the latter case we can move Z_i inside the summation become constant)

$$\begin{aligned} &= \sum_j \frac{1}{Z_j} \sum_{k \neq j} \delta E_{jk} \\ &= \sum_{j \neq i} \frac{E_{ij}}{Z_j} (2(y_j - y_i)) + \sum_{j \neq i} \frac{E_{ij}}{Z_i} (-2(y_i - y_j)) \\ &= 2 \sum_{j \neq i} (-q_{ji} - q_{ii}) (y_i - y_j) \end{aligned} \quad (5)$$

Combine (4) + (5) we arrive at the final result,

$$\frac{\delta C}{\delta y_i} = 2 \sum_{j \neq i} (p_{ji} - q_{ji} + p_{ii} - q_{ii}) (y_i - y_j) \quad \square$$

t-distributed Stochastic Neighbor Embedding (*t-SNE*)

Define :

$$q_{ji} = q_{ij} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_{k, l \neq k} (1 + \|y_k - y_l\|^2)^{-1}} = \frac{E_{ij}^{-1}}{\sum_{k, l \neq k} E_{ke}^{-1}} = \frac{E_{ij}^{-1}}{Z}$$

Notice that $E_{ij} = E_{ji}$. The loss function is defined as

$$C = \sum_{k,l \neq i} p_{ek} \log \frac{p_{ek}}{q_{ek}} = \sum_{k,l \neq i} p_{ek} \log p_{ek} - p_{ek} \log q_{ek}$$

$$= \sum_{k,l \neq i} p_{ek} \log p_{ek} - p_{ek} \log E_{kl}^{-1} + p_{ek} \log Z$$

We derive wrt y_i . To make the derivative less cluttered, I will omit the δy_i term at the denominator.

$$\frac{\delta C}{\delta y_i} = \sum_{k,l \neq i} -p_{ek} \delta \log E_{kl}^{-1} + \sum_{k,l \neq i} p_{ek} \log Z$$

We start with the first term, noting that the derivative is non zero when $y_j > y_i$, $k=i$ or $l=i$, that $p_{ji} = p_{ij}$ and $E_{ji} = E_{ij}$

$$\sum_{k,l \neq i} -p_{ek} \delta \log E_{kl}^{-1} = -2 \sum_{j \neq i} p_{ji} \delta \log E_{ij}^{-1}$$

Since $\delta E_{ij}^{-1} = E_{ij}^{-2} (-2(y_i - y_j))$ we have

$$-2 \sum_{j \neq i} p_{ji} \frac{E_{ij}^{-2}}{E_{ij}^{-1}} (-2(y_i - y_j)) = 4 \sum_{j \neq i} p_{ji} E_{ij}^{-1} (y_i - y_j) \quad (10)$$

We conclude the second term. Using the fact that $\sum_{k,l \neq i} p_{kl} = 1$ and that Z does not depend on k or l .

$$\begin{aligned} \sum_{k,l \neq i} p_{ek} \delta \log Z &= \frac{1}{Z} \sum_{k,l \neq i} p_{ek} \delta E_{kl}^{-1} \\ &= 2 \sum_{j \neq i} \frac{E_{ij}^{-2}}{Z} (-2(y_i - y_j)) = -4 \sum_{j \neq i} q_{ji} E_{ij}^{-1} (y_i - y_j) \end{aligned} \quad (11)$$

Combine (10) + (11), final result:

$$\begin{aligned} \frac{\delta C}{\delta y_i} &= 4 \sum_{j \neq i} (p_{ji} - q_{ji}) E_{ij}^{-1} (y_i - y_j) \\ &= 4 \sum_{j \neq i} (p_{ji} - q_{ji}) (1 + \|y_i - y_j\|^2)^{-1} (y_i - y_j) \quad \square \end{aligned}$$

EX 4 So what is t-SNE vs PCA

PCA	t-SNE
① It is a linear Dimensionality reduction technique	It is a non linear Dimensionality reduction technique.
② It tries to preserve the global structure of the data	It tries to preserve the local structure (cluster) of data
③ It doesn't work well as compared to t-SNE	It's one of the best dimensionality reduction technique.
④ It does not involve hyperparameters	It involves hyperparameters such as perplexity, learning rate and number of steps
⑤ It gets highly affected by outliers	It can handle outliers
⑥ PCA is a deterministic algorithm	It is a non deterministic or randomised algorithm.
⑦ It works by rotating the vectors for preserving variance	It works by minimising the distance between the points in a gaussian
⑧ We can fine decide on how much variance to preserve using eigen values	We can not preserve variance instead we can preserve distance using hyperparameters.