

Thứ  
Ngày  
No.

1) Biến đổi lại công thức toán t-SNE, SNE, có tính đạo hàm  
với các parameter.

Bài toán t-SNE:

The similarity of data point  $x_j$  to data point  $x_i$  is the conditional probability  $p(j|i)$ . For nearby data points,  $p(j|i)$  is high, whereas for widely separated data points,  $p(j|i)$  will be almost infinitesimal.

The conditional probability is given by:

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

where  $\sigma_i$  is the variance of the Gaussian that is centered on data point  $x_i$ . For the low-dimensional counterpart  $y_i$  &  $y_j$  of the high-dimensional data points  $x_i$  &  $x_j$ , it is possible to compute a similar conditional probability, which we denote by  $q(j|i)$ :

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

We set  $q(i|i) = 0$ . If the map points  $y_i$  and  $y_j$  correctly model the similarity between the high-dimensional data point  $x_i$  and  $x_j$ , the conditional probabilities  $p(j|i)$  and  $q(j|i)$  will be equal. SNE minimise the sum of Kullback-Leibler ~~convergence~~ divergence over all data points using a gradient descent method. The cost function  $C$  is given by:

$$C = \sum_i KL(p_i || q_i) = \sum_i \sum_j p(j|i) \log \frac{p(j|i)}{q(j|i)}$$





in which  $P_i$  represents the conditional probability distribution overall all other data points given data point  $x_i$ ,  $Q_i$  represents the conditional probability distribution overall all other map points given map point  $y_i$ . Because the Kullback-Leibler divergence is not symmetric, different types of error in the pairwise distances in the low-dimensional map are not weighted equal.

SNE performs a binary search for the value of  $\sigma_i$  that produces a  $P_i$  with a fixed perplexity that is specified by the user. The perplexity is defined as

$$\text{Per}(P_i) = \frac{1}{H(P_i)}$$

where  $H(P_i)$  is the Shannon entropy of  $P_i$  measured in:

$$H(P_i) = - \sum_j P(p_{j|i}) \log_2 P(p_{j|i})$$

The perplexity can be interpreted as a measure of the effective number of neighbors. The performance of SNE is fairly robust to changes in the perplexity & typical values are between 5 & 50. The minimization of the cost function is performed using a gradient descent method. The gradient descent has a surprisingly simple form:

$$\frac{\partial \mathcal{L}}{\partial y_i} = \alpha \sum_j (P_{ji} - Q_{ji} + P_{ij} - Q_{ij}) (y_i - y_j)$$

Mathematically, the gradient update with a momentum term is given by:

$$y^{(t)} = y^{(t-1)} + \eta \frac{\partial \mathcal{L}}{\partial y} + \alpha(t) (y^{(t-1)} - y^{(t-2)})$$

where  $y^{(t)}$  indicates the solution at iteration  $t$ ,  $\eta$  indicates the learning rate &  $\alpha(t)$  represents the momentum at iteration  $t$ .



As an alternative to ~~minimize~~ minimize the sum of the Kullback-Leibler divergence between a joint prob distribution  $P$  in the high-dimensional space & a joint distribution  $Q$  in the low-dimensional space

$$C = \sum_i KL(P||Q) = \sum_i \sum_j P_{ij} \log \frac{P_{ij}}{Q_{ij}}$$

Set  $P_{ii} \approx Q_{ii}$  to 0. We refer to the type of SNE as symmetric SNE because it has the property that  $P_{ij} = P_{ji}$  &  $Q_{ij} = Q_{ji} \forall i, j$ . In symmetric SNE, the pairwise similarities in the low-dimensional map  $q_{ij}$  are given by:

$$q_{ij} = \frac{e^{-\|y_i - y_j\|^2}}{\sum_{k \neq l} e^{-\|y_k - y_l\|^2}}$$

The obvious way to define the pairwise similarities in the high dimensional space  $P_{ij}$  is:

$$P_{ij} = \frac{e^{-\|x_i - x_j\|^2 / 2\sigma^2}}{\sum_{k \neq l} e^{-\|x_k - x_l\|^2 / 2\sigma^2}}$$

The gradient of symmetric SNE is fairly similar to that of asymmetric SNE & is given:

$$\frac{\partial C}{\partial y_i} = 4 \sum_j (P_{ij} - Q_{ij}) (y_i - y_j)$$