

t-SNE algorithm implementation

Tomas Ukkonen 2020, tomas.ukkonen@novelinsight.fi

The algorithm is based on a t-SNE paper *Visualizing Data using t-SNE*. Laurens van der Maaten and Geoffrey Hinton. *Journal of Machine Learning Research* 9 (11/2008). The implementation is in `src/neuralnetwork/TSNE.h` and `TSNE.cpp` in `dinrhiw2` repository.

We maximize KL-divergence (p_{ij} calculated from data).

$$D_{\text{KL}}(\mathbf{y}_1 \dots \mathbf{y}_N) = \sum_{i \neq j} p_{ij} \log\left(\frac{p_{ij}}{q_{ij}}\right), \quad q_{ij} = \frac{(1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2)^{-1}}{\sum_k \sum_{l \neq k} (1 + \|\mathbf{y}_k - \mathbf{y}_l\|^2)^{-1}}$$

The p_{ij} values are calculated from data using formulas.

$$p_{j|i} = \frac{e^{-\|\mathbf{x}_j - \mathbf{x}_i\|^2 / 2\sigma_i^2}}{\sum_{k \neq i} e^{-\|\mathbf{x}_k - \mathbf{x}_i\|^2 / 2\sigma_i^2}}, \quad p_{i|i} = 0, \quad \sum_j p_{j|i} = 1$$

Symmetric probability values are computed from conditional probabilities using the formula

$$p_{ij} = \frac{p_{j|i} + p_{i|j}}{2N}, \quad \sum_{i,j} p_{ij} = 1$$

The variance terms of each data point σ_i^2 is calculated using values $p_{j|i}$ to search for target perplexity $\text{perp}(P_i) = 2^{H(P_i)} = 2^{-\sum_j p_{j|i} \log_2(p_{j|i})}$. Good general perplexity value is maybe 30 which we use to solve σ_i^2 value using bisection method.

First we set minimum $\sigma_{\min}^2 = 0$ and $\sigma_{\max}^2 = \text{trace}(\mathbf{\Sigma}_{\mathbf{x}})$. We then always select $\sigma_{\text{next}}^2 = \frac{\sigma_{\min}^2 + \sigma_{\max}^2}{2}$ to half the interval and calculate perplexity at σ_{next}^2 to figure out which half contains the target perplexity value and stop if error is smaller than 0.1.

Gradient

We need to calculate gradient for each \mathbf{y}_m in D_{KL} .

$$\nabla_{\mathbf{y}_m} D_{\text{KL}} = \nabla_{\mathbf{y}_m} \sum_{i \neq j} -p_{ij} \log(q_{ij}) = -\sum_{i \neq j} \frac{p_{ij}}{q_{ij}} \nabla_{\mathbf{y}_m} q_{ij}$$

The general rule to derivate q_{ij} terms is:

$$\nabla_{\frac{f}{g}} = \nabla f g^{-1} = f' g^{-2} g - f g^{-2} g' = \frac{f' g - f g'}{g^2}$$

And when $m \neq i \neq j$ we need to derivate only the second part

$$\begin{aligned} & \nabla_{\mathbf{y}_m \neq i \neq j} \left(\frac{(1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2)^{-1}}{\sum_k \sum_{l \neq k} (1 + \|\mathbf{y}_k - \mathbf{y}_l\|^2)^{-1}} \right) \\ &= -\frac{(1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2)^{-1}}{(\sum_k \sum_{l \neq k} (1 + \|\mathbf{y}_k - \mathbf{y}_l\|^2)^{-1})^2} \nabla_{\mathbf{y}_m} \sum_k \sum_{l \neq k} (1 + \|\mathbf{y}_k - \mathbf{y}_l\|^2)^{-1} \end{aligned}$$

$$\begin{aligned} & \nabla_{\mathbf{y}_m \neq i \neq j} \sum_k \sum_{l \neq k} (1 + \|\mathbf{y}_k - \mathbf{y}_l\|^2)^{-1} \\ &= \nabla_{\mathbf{y}_m} \sum_{l \neq m} (1 + \|\mathbf{y}_m - \mathbf{y}_l\|^2)^{-1} + \nabla_{\mathbf{y}_m} \sum_{k \neq m} (1 + \|\mathbf{y}_k - \mathbf{y}_m\|^2)^{-1} \\ &= 2 \nabla_{\mathbf{y}_m} \sum_{l \neq m} (1 + \|\mathbf{y}_m - \mathbf{y}_l\|^2)^{-1} \\ &= 2 \sum_{l \neq m} \nabla_{\mathbf{y}_m} (1 + \|\mathbf{y}_m - \mathbf{y}_l\|^2)^{-1} \\ &= 4 \sum_{l \neq m} -(1 + \|\mathbf{y}_m - \mathbf{y}_l\|^2)^{-2} (\mathbf{y}_m - \mathbf{y}_l) \end{aligned}$$

And when $y = i$ or $y = j$ we need to derivate the upper part too.

$$\begin{aligned} & \nabla_{\mathbf{y}_i} \frac{(1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2)^{-1}}{\sum_k \sum_{l \neq k} (1 + \|\mathbf{y}_k - \mathbf{y}_l\|^2)^{-1}} = \frac{1}{\sum_k \sum_{l \neq k} (1 + \|\mathbf{y}_k - \mathbf{y}_l\|^2)^{-1}} \nabla_{\mathbf{y}_i} (1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2)^{-1} - \frac{f g'}{g^2} \\ & \nabla_{\mathbf{y}_i} (1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2)^{-1} = -2 (1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2)^{-2} (\mathbf{y}_i - \mathbf{y}_j) \end{aligned}$$

With these derivates we can then calculate derivate of D_{KL} for each \mathbf{y} . We just select step length for the gradient which causes increase in D_{KL} .

Optimized gradient

We can rewrite the gradient of D_{KL} by taking partial derivates of distance variables d_{ij} and d_{ji} , $d_{ij} = \|\mathbf{y}_i - \mathbf{y}_j\|$

$$\nabla_{\mathbf{y}_i} D_{KL} = \sum_j \left(\frac{\partial D_{KL}}{\partial d_{ij}} \frac{\partial d_{ij}}{\partial \mathbf{y}_i} + \frac{\partial D_{KL}}{\partial d_{ji}} \frac{\partial d_{ji}}{\partial \mathbf{y}_i} \right) = \sum_j \left(\frac{\partial D_{KL}}{\partial d_{ij}} + \frac{\partial D_{KL}}{\partial d_{ji}} \right) \frac{\partial d_{ji}}{\partial \mathbf{y}_i} = 2 \sum_j \frac{\partial D_{KL}}{\partial d_{ij}} \frac{\partial d_{ji}}{\partial \mathbf{y}_i}$$

The gradient of the distance variable d_{ji} is

$\frac{\partial d_{ji}}{\partial \mathbf{y}_i} = D \|\mathbf{y}_i - \mathbf{y}_j\| = \frac{d}{d \mathbf{y}_i} \sqrt{\|\mathbf{y}_i - \mathbf{y}_j\|^2} = \frac{1}{2\sqrt{\|\mathbf{y}_i - \mathbf{y}_j\|^2}} D \|\mathbf{y}_i - \mathbf{y}_j\|^2 = \frac{\mathbf{y}_i - \mathbf{y}_j}{\|\mathbf{y}_i - \mathbf{y}_j\|}$. Note that the research paper gives different derivate which seems to be **wrong** (?).

Gradient of the D_{KL} term is (we use auxiliary variable $Z = \sum_{k \neq l} (1 + d_{kl}^2)^{-1}$).

$$\begin{aligned} \frac{\partial D_{KL}}{\partial d_{ij}} &= -\sum_{k \neq l} p_{kl} \frac{\partial(\log(q_{kl}))}{\partial d_{ij}} = -\sum_{k \neq l} p_{kl} \frac{\partial(\log(q_{kl}Z) - \log(Z))}{\partial d_{ij}} \\ &= -\sum_{k \neq l} p_{kl} \left(\frac{1}{q_{kl}Z} \frac{\partial(1 + d_{kl}^2)^{-1}}{\partial d_{ij}} - \frac{1}{Z} \frac{\partial Z}{\partial d_{ij}} \right) = 2 \frac{p_{ij}}{q_{ij}Z} (1 + d_{ij}^2)^{-2} d_{ij} + \sum_{k \neq l} p_{kl} \frac{1}{Z} \frac{\partial Z}{\partial d_{ij}} \\ &= 2 p_{ij} (1 + d_{ij}^2)^{-1} d_{ij} - 2 \sum_{k \neq l} p_{kl} \frac{(1 + d_{ij}^2)^{-2}}{Z} d_{ij} \\ &= 2 p_{ij} (1 + d_{ij}^2)^{-1} d_{ij} - 2 \left(\frac{(1 + d_{ij}^2)^{-2}}{Z} \right) d_{ij} = 2 p_{ij} (1 + d_{ij}^2)^{-1} d_{ij} - 2 q_{ij} (1 + d_{ij}^2)^{-1} d_{ij} \\ &= 2 (p_{ij} - q_{ij}) (1 + d_{ij}^2)^{-1} d_{ij} \end{aligned}$$

So we now get a simple formula for the gradient

$$\nabla_{\mathbf{y}_i} D_{KL} = 4 \sum_j (p_{ij} - q_{ij}) (1 + d_{ij}^2)^{-1} (\mathbf{y}_i - \mathbf{y}_j).$$

To get even better results we want to use absolute value $|D|_{KL}$ (See later in this paper). This means we will compute altered gradient.

$$D \|f(\mathbf{x})\| = D \sqrt{\|f(\mathbf{x})\|^2} = \frac{1}{2\sqrt{\|f(\mathbf{x})\|^2}} D \|f(\mathbf{x})\|^2 = \frac{f(\mathbf{x})}{\|f(\mathbf{x})\|} \nabla f(\mathbf{x}) = \text{sign}(f(\mathbf{x})) \nabla f(\mathbf{x})$$

$$\frac{\partial |D|_{KL}}{\partial d_{ij}} = \sum_{k \neq l} p_{kl} \frac{1}{\partial d_{ij}} \left| \log\left(\frac{p_{kl}}{q_{kl}}\right) \right| = -\sum_{k \neq l} \text{sign}\left(\log\left(\frac{p_{kl}}{q_{kl}}\right)\right) p_{kl} \frac{\partial(\log(q_{kl}))}{\partial d_{ij}}$$

This means we only need to modify our gradient formula by multiplication of $\text{sign}(x)$ function for the first term and for the second term we need to calculate one additional term P_{sign} .

$$\begin{aligned} P_s &= \sum_{k \neq l} \text{sign}\left(\log\left(\frac{p_{kl}}{q_{kl}}\right)\right) p_{kl} \\ \nabla_{\mathbf{y}_i} |D|_{KL} &= 4 \sum_j \left(\text{sign}\left(\log\left(\frac{p_{ij}}{q_{ij}}\right)\right) p_{ij} - P_s q_{ij} \right) (1 + d_{ij}^2)^{-1} (\mathbf{y}_i - \mathbf{y}_j). \end{aligned}$$

This optimized gradient is faster because it scales as $O(N^2)$ instead of slower $O(N^3)$ of the direct method.

Optimization of computation

For large number of points the update rule is still slow with $O(N^2)$ scaling. Extra speed can be achieved by combining large away data points to a single point which is then used to calculate the divergence and gradient. This can be done by using *Barnes-Hut approximation* which changes computational complexity to nearly linear $O(N \log(N))$.

Improvement of the KL divergence based distribution comparison

The information theoretic distribution comparison metric D_{KL} can be improved by using absolute values. This also symmetrizes comparison a bit. (See my other notes about information theory/also at the end of this section.)

$$|D|_{\text{KL}}(\mathbf{y}_1 \dots \mathbf{y}_N) = \sum_{i \neq j} p_{ij} \left| \log \left(\frac{p_{ij}}{q_{ij}} \right) \right|$$

Gradient of the absolute value can be computed using a simple trick.

$$D \|f(\mathbf{x})\| = D \sqrt{\|f(\mathbf{x})\|^2} = \frac{1}{2\sqrt{\|f(\mathbf{x})\|^2}} D \|f(\mathbf{x})\|^2 = \frac{f(\mathbf{x})}{\|f(\mathbf{x})\|} \nabla f(\mathbf{x}) = \text{sign}(f(\mathbf{x})) \nabla f(\mathbf{x})$$

This means the improved gradient is:

$$\nabla_{\mathbf{y}_m} |D|_{\text{KL}} = \sum_{i \neq j} p_{ij} \nabla_{\mathbf{y}_m} \left| \log \left(\frac{p_{ij}}{q_{ij}} \right) \right| = - \sum_{i \neq j} \frac{p_{ij}}{q_{ij}} \text{sign} \left(\log \left(\frac{p_{ij}}{q_{ij}} \right) \right) \nabla_{\mathbf{y}_m} q_{ij}$$

This means we only need to add $\text{sign}(x)$ non-linearity to the gradient calculation code. The $\text{sign}(x)$ non-linearity is well defined everywhere else except at zero where we can set $\text{sign}(0) = 1$ without having much problems in practice.

Justification of the modified KL divergence

The absolute value can be justified by following calculations. Geometric mean of observed symbol string is P and the number of symbols $l = 1 \dots L$ in N symbol long string is n_l . Additionally we let the length of string to go to infinity ($N \rightarrow \infty$):

$$P = \left(\prod_k^N p(\mathbf{x}_k) \right)^{1/N} = \left(\prod_l^L p(l)^{n_l} \right)^{1/N} \approx \prod_l^L p(l)^{p(l)}$$

By taking the logarithm of P we get formula for entropy: $\log(P) = \sum_l p(l) \log(p(l)) = -H(L)$.

Comparing distributions probabilities we can write ($N \rightarrow \infty$):

$$Q_{\mathbf{x}} = \left(\frac{\prod_k^N p(\mathbf{x}_k)}{\prod_k^N p(\mathbf{y}_k)} \right)^{1/N} = \left(\prod_l^L \left(\frac{p_{\mathbf{x}}(l)}{p_{\mathbf{y}}(l)} \right)^{n_l} \right)^{1/N} \approx \prod_l^L \left(\frac{p_{\mathbf{x}}(l)}{p_{\mathbf{y}}(l)} \right)^{p_{\mathbf{x}}(l)}$$

And by taking the logarithm of Q we get Kullback-Leibler divergence:

$$\log(Q_{\mathbf{x}}) = \sum_l p_{\mathbf{x}}(l) \log \left(\frac{p_{\mathbf{x}}(l)}{p_{\mathbf{y}}(l)} \right) = D_{\text{KL}}$$

Now by always taking the maximum ratio of probabilities when computing Q we don't have the problem that multiplication (in $\prod_l^L \left(\frac{p_{\mathbf{x}}(l)}{p_{\mathbf{y}}(l)} \right)^{n_l}$ -term) of probability ratios would cancel each other reducing the usability of D_{KL} divergence when used for distribution comparison.

$$|Q_{\mathbf{x}}| = \left(\prod_l^L \max \left(\frac{p_{\mathbf{x}}(l)}{p_{\mathbf{y}}(l)}, \frac{p_{\mathbf{y}}(l)}{p_{\mathbf{x}}(l)} \right)^{n_l} \right)^{1/N} \approx \prod_l^L \max \left(\frac{p_{\mathbf{x}}(l)}{p_{\mathbf{y}}(l)}, \frac{p_{\mathbf{y}}(l)}{p_{\mathbf{x}}(l)} \right)^{p_{\mathbf{x}}(l)}$$

$$\log |Q_{\mathbf{x}}| = \sum_l p_{\mathbf{x}}(l) \log \left(\max \left(\frac{p_{\mathbf{x}}(l)}{p_{\mathbf{y}}(l)}, \frac{p_{\mathbf{y}}(l)}{p_{\mathbf{x}}(l)} \right) \right) = \sum_l p_{\mathbf{x}}(l) \left| \log \left(\frac{p_{\mathbf{x}}(l)}{p_{\mathbf{y}}(l)} \right) \right| = |D_{\mathbf{x}}|_{\text{KL}}$$

Further symmetrization can be done by taking the geometric mean:

$$|Q| = (|Q_{\mathbf{x}}| |Q_{\mathbf{y}}|)^{1/2}, \quad \log(|Q|) = \frac{1}{2} (\log |Q_{\mathbf{x}}| + \log |Q_{\mathbf{y}}|) = \frac{1}{2} (|D_{\mathbf{x}}|_{\text{KL}} + |D_{\mathbf{y}}|_{\text{KL}}).$$

Improvement of the MSE calculation code

Calculating a gradient of absolute value can be also used in minimum least squares (MSE) optimization where we can then easily use norm instead (minimum norm error - MNE) of the squared error which is then less affected by large outlier values.

$$\text{MSE}(\mathbf{w}) = E_{\mathbf{x}\mathbf{y}}\left[\frac{1}{2}\|\mathbf{y} - f(\mathbf{x})\|^2\right], \nabla_{\mathbf{w}} \text{MSE}(\mathbf{w}) = E_{\mathbf{x}\mathbf{y}}[(f(\mathbf{x}) - \mathbf{y})^T \nabla_{\mathbf{w}} \mathbf{f}(\mathbf{x})]$$

$$\text{MNE}(\mathbf{w}) = E_{\mathbf{x}\mathbf{y}}[\|\mathbf{y} - f(\mathbf{x})\|], \nabla_{\mathbf{w}} \text{MNE}(\mathbf{w}) = E_{\mathbf{x}\mathbf{y}}\left[\frac{(f(\mathbf{x}) - \mathbf{y})^T}{\|f(\mathbf{x}) - \mathbf{y}\|} \nabla_{\mathbf{w}} \mathbf{f}(\mathbf{x})\right]$$

This means we have to just to scale the backpropagation gradient of each term i by dividing with $\|\mathbf{y}_i - f(\mathbf{x}_i)\|$. This means that for the large errors the effect to gradient is now smaller and small values have equal effect to gradient.