

3.6 Calculate gradient of SNE cost function with respect to y_i .

→ The SNE cost function is defined as:

$$C = \sum_i \sum_j p_{ji} \log \frac{p_{ji}}{q_{ji}}$$

gradient of C with respect to y_i

chain rule -
$$\frac{\partial}{\partial y_i} \left(\log \frac{p_{ji}}{q_{ji}} \right) = \frac{1}{\frac{p_{ji}}{q_{ji}}} \cdot \frac{\partial}{\partial y_i} \left(\frac{p_{ji}}{q_{ji}} \right)$$

$$\frac{\partial}{\partial y_i} \left(\frac{p_{ji}}{q_{ji}} \right) = \frac{1}{q_{ji}} \cdot \frac{\partial p_{ji}}{\partial y_i} - \frac{p_{ji}}{q_{ji}^2} \cdot \frac{\partial q_{ji}}{\partial y_i}$$

using gradient descent formula,

$$\frac{\partial}{\partial y_i} \left(p_{ji} \log \frac{p_{ji}}{q_{ji}} \right) = 2 \cdot (p_{ji} - q_{ji}) \cdot (y_i - y_j)$$

summing over neighbors:

$$\frac{\partial C}{\partial y_i} = \sum_j \frac{\partial}{\partial y_i} \left(p_{ji} \log \frac{p_{ji}}{q_{ji}} \right)$$

$$\frac{\partial C}{\partial y_i} = \sum_j 2 \cdot (p_{ji} - q_{ji}) \cdot (y_i - y_j)$$

finally
rule. update y_i using gradient descent update

c) gradient descent of symmetric SNE cost function with respect to Y_i .

for symmetric SNE, there are further simplification to be made. Both p and q matrices are symmetric, so $p_{ij} = p_{ji}$ and $q_{ij} = q_{ji}$ leading to,

$$\begin{aligned}\frac{\partial C}{\partial Y_i} &= 2 \sum_j (p_{ij} - q_{ij} + p_{ji} - q_{ji}) (Y_i - Y_j) \\ &= 2 \sum_j (2p_{ij} - 2q_{ij}) (Y_i - Y_j) \\ &= 4 \sum_j (p_{ij} - q_{ij}) (Y_i - Y_j)\end{aligned}$$

d) gradient descent of t-SNE.

$$\begin{aligned}\frac{\partial C}{\partial Y_i} &= 2 \sum_j (k_{ij} + k_{ji}) (Y_i - Y_j) \\ k_{ij} &= \frac{1}{S} \left[\frac{\partial C}{\partial q_{ij}} - \sum_k \frac{\partial C}{\partial q_{kl}} q_{kl} \right] \frac{\partial w_{ij}}{\partial f_{ij}}\end{aligned}$$

Inserting Kullback divergence.

Both SNE and t-SNE use the Kullback-Leibler divergence, which as noted above, has following gradient:

$$\frac{\partial C}{\partial q_{ij}} = -\frac{p_{ij}}{q_{ij}}$$

k_{ij} therefore becomes:

$$k_{ij} = \frac{1}{S} \left[-\frac{p_{ij}}{q_{ij}} - \sum_k \frac{-p_{kl}}{q_{kl}} q_{kl} \right] \frac{\partial w_{ij}}{\partial f_{ij}} = \frac{1}{S} \left[\frac{p_{ij}}{q_{ij}} + \sum_k p_{kl} \right] \frac{\partial w_{ij}}{\partial f_{ij}}$$

$$= \frac{1}{S} \left[-\frac{P_{ij}}{q_{ij}} + 1 \right] \frac{\delta w_{ij}}{\delta f_{ij}}$$

At this point both SNE, and t-SNE output kernel (Gaussian and t-distribution respectively) have a derivative that has general form

$$\frac{\delta w_{ij}}{\delta f_{ij}} = -w_{ij}^n$$

where, $n=1$ in case of SNE, and $n=2$ in case of t-SNE substituting that, we get:

$$K_{ij} = \frac{1}{S} \left[-\frac{P_{ij}}{q_{ij}} + 1 \right] \frac{\delta w_{ij}}{\delta f_{ij}} = -\frac{w_{ij}^n}{S} \left[-\frac{P_{ij}}{q_{ij}} + 1 \right] = -\frac{w_{ij}^{n-1}}{S} \left[-\frac{P_{ij}}{q_{ij}} + 1 \right]$$

using the fact that $w_{ij}/S = q_{ij}$, now we can move $-q_{ij}$ inside the expression in parentheses to get:

$$K_{ij} = w_{ij}^{n-1} (P_{ij} - q_{ij})$$

for SNE $w_{ij}^{n-1} = 1$ because $n=1$ we get

$$\frac{\delta C}{\delta y_i} = 2 \sum_j (P_{ij} - q_{ij} + P_{ji} - q_{ji}) (y_i - y_j)$$

for t-SNE

$$\frac{\delta C}{\delta y_i} = 2 \sum_j [w_{ij} (P_{ij} - q_{ij}) + w_{ji} (P_{ji} - q_{ji})] (y_i - y_j)$$

But t-distributed kernel also creates a symmetric ~~W~~ matrix so we simplify

$$\frac{\delta C}{\delta y_i} = 4 \sum_j w_{ij} (P_{ij} - q_{ij}) (y_i - y_j)$$

9.4*

→ Consider the fitted values, that results from performing linear regression without an intercept. In this setting the i^{th} fitted value takes the form

$$\hat{y}_i = x_i \hat{\beta}$$

where,

$$\hat{\beta} = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2}$$

now we want to express \hat{y}_i as a linear combⁿ of y_i 's. we need to find coeff. a_i such that

$$\hat{y}_i = \sum_{i=1}^n a_i y_i$$

we can rewrite expression as;

$$\hat{y}_i = \frac{x_i}{\sum_{i=1}^n x_i^2} \left(\sum_{i=1}^n x_i y_i \right)$$

$$a_i = \frac{x_i}{\sum_{i=1}^n x_i^2}$$

Therefore, we can express \hat{y}_i as a linear combination of y_i 's.

$$\hat{y}_i = \sum_{i=1}^n a_i y_i$$

the coeff a_i are given by

$$a_i = \frac{x_i}{\sum_{i=1}^n x_i^2}$$

This shows that fitted values from linear regression are indeed linear combination of response values y_i .