

Homework week 1

Huỳnh Lê Linh Đan - DSEB 62

Ngày 17 tháng 1 năm 2023

1 Problem 1

For SNE:

Given point $x_1, x_2, \dots, x_N \in R^D$ we define the distribution P_{ij} which is the probability that point x_i chooses x_j as its neighbor:

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

We need to find good embedding $y_1, y_2, \dots, y_N \in R^d$ for $d < D$:

$$Q_{ij} = \frac{\exp(-\|y_i - y_j\|^2)}{\sum_k \sum_{k \neq l} \exp(-\|y_l - y_k\|^2)} = \frac{E}{Z}$$

let $\|y_i - y_j\| = d_{ij}$

We need to optimize Q to be close to P, we do so by minimizing KL-divergence: to find the embedding $y_1, \dots, y_n \in R^d$:

$$\begin{aligned} KL(P\|Q) &= \sum_{ij} P_{ij} \log \left(\frac{P_{ij}}{Q_{ij}} \right) \\ &= \sum_{ij} P_{ij} \log(P_{ij}) - P_{ij} \log(Q_{ij}) \end{aligned}$$

Where P_{ij} can be inferred from the data, we treat this as constant so (4) can be rewritten as:

$$\begin{aligned} KL(P\|Q) &= - \sum_{ij} P_{ij} \log(Q_{ij}) + \text{Constant} \\ &= - \sum_{ij} P_{ij} \log \left(\frac{E}{Z} \right) + \text{Constant} \\ &= - \sum_{ij} P_{ij} \log(E) - P_{ij} \log(Z) + \text{Constant} \end{aligned}$$

$$\begin{aligned} \frac{\delta L}{\delta y_i} &= \left(\frac{\delta L}{\delta d_{ij}} + \frac{\delta L}{\delta d_{ji}} \right) \frac{\delta d_{ji}}{\delta y_i} = 2 \frac{\delta L}{\delta d_{ij}} \frac{\delta d_{ij}}{\delta y_i} \\ \frac{\delta L}{\delta y_i} &= -2 \sum_{ij} P_{ij} \delta \log(E) - P_{ij} \delta \log(Z) \\ \delta E &= -2 (y_i - y_j) E \end{aligned}$$

First we consider the first term:

$$\begin{aligned} \sum_{ij} P_{ij} \delta \log(E) &= \sum_{ij} P_{ij} (-2 (y_i - y_j) E) \frac{1}{E} \\ &= -2 \sum_{ij} P_{ij} (y_i - y_j) \end{aligned}$$

in the second term, $\sum_{k \neq l} P_{ij} = 1$, the derivative is non-zero when $k = i$ or $l = i$:

$$\begin{aligned}
P_{ij}\delta \log(Z) &= \sum_{i \neq j} \frac{1}{Z} \delta E \\
&= 2 \sum_{i \neq j} \frac{E}{Z} (y_i - y_j) \\
&= 2 \sum_{i \neq j} Q_{ij} (y_i - y_j)
\end{aligned}$$

Plug (12), (15) into (9) we have:

$$\begin{aligned}
\frac{\delta L}{\delta y_i} &= -2 \sum_{ij} -2P_{ij} (y_i - y_j) + 2Q_{ij} (y_i - y_j) \\
&= 2 \sum_{ij} 2P_{ij} (y_i - y_j) - 2Q_{ij} (y_i - y_j) \\
&= 4 \sum_{ij} (P_{ij} - Q_{ij}) (y_i - y_j)
\end{aligned}$$

For T-SNE, Q_{ij} is defined as:

$$Q_{ij} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_k \sum_{k \neq l} (1 + \|y_k - y_l\|^2)^{-1}} = \frac{E^{-1}}{Z^{-1}}$$

We use the same P_{ij} and our loss function is:

$$KL(P\|Q) = - \sum_{ij} P_{ij} \log(E^{-1}) - P_{ij} \log(Z^{-1}) + \text{Constant}$$

$$\begin{aligned}
\frac{\delta L}{\delta y_i} &= -2 \sum_{ij} P_{ij} \delta \log(E^{-1}) - P_{ij} \delta \log(Z^{-1}) \\
\delta E^{-1} &= -(y_i - y_j) E^{-2}
\end{aligned}$$

We consider the first term:

$$\begin{aligned}
\sum_{ij} P_{ij} \delta \log(E^{-1}) &= - \sum_{ij} P_{ij} 2 (y_i - y_j) \frac{E^{-2}}{E^{-1}} \\
&= -2 \sum_{ij} P_{ij} (y_i - y_j) E^{-1}
\end{aligned}$$

We consider the second term:

$$\begin{aligned}
\sum_{ij} P_{ij} \delta \log(Z^{-1}) &= -2 \sum_{ij} (y_i - y_j) \frac{1}{Z^{-1}} E^{-2} \\
&= -2 \sum_{ij} (y_i - y_j) \frac{E^{-1}}{Z^{-1}} E^{-1} = -2 \sum_{ij} (y_i - y_j) Q_{ij} E^{-1}
\end{aligned}$$

Plug (23), (25) into (20) we have:

$$\begin{aligned}
\frac{\delta L}{\delta y_i} &= 4 \sum_{ij} P_{ij} (y_i - y_j) E^{-1} - (y_i - y_j) Q_{ij} E^{-1} \\
&= 4 \sum_{ij} (P_{ij} - Q_{ij}) (y_i - y_j) E^{-1} \\
&= 4 \sum_{ij} (P_{ij} - Q_{ij}) (y_i - y_j) (1 + \|y_i - y_j\|)^{-1}
\end{aligned}$$

2 Problem 4

Compare T-sne and PCA:

Similar: They are both unsupervised dimensional reduction and data visualization technique for very high dimensional data

index	PCA	T-sne
1	linear	non-linear
2	preserve global structure	preserve local structure
3	Does not involves hyper parameters	Involves hyper parameters
4	Affected by outliers	Can handle outliers
5	deterministic	randomised
6	rotating vectors for preserving variance	Minimising the distance between the point
7	preserve using eigenvalues	using hyper parameters

explanation:

(1) PCA works well when there is linear relation between features while T-sne does a decent job even in non-linear dataset.

(3) T-sne hyper parameters include perplexity, learning rate, iterations. (2), (4) PCA can lead to local inconsistencies, far away point can become nearest neighbor. For t-sne, low dimensional neighborhood should be the same as original neighborhood.

(5) PCA produces the same output each time, T-sne's intuition is based on random walk between data points, therefore may produce different result on the same data. Moreover, in Problem 2, difference in hyper parameters like perplexity can produce different result.

(6) Optimization problem for PCA is maximising the variance of the projected data. For T-sne, we minimize KL-divergence so that Q (distribution for projected data points) is close to P (distribution for given data)

(7) PCA decides how much variance to preserve using eigenvalues. T-sne decides the distance to preserve using Perplexity