

## t-SNE

Date \_\_\_\_\_  
Page \_\_\_\_\_

### \* High Dimensional Data

- can be predicted by 30+ features.
- High dim data is hard to visualize & work with
- Embedding to low dim spaces helps visualize the data.

### \* Earlier Techniques - Direct Visualization

#### • Parallel co-ordinates:

- allows for comparison of multiple data records, by using parallel lines to connect points based on multiple numerical variables.

#### • Chernoff Faces:

- Symbolizing data using faces

- each emoji has an eqv emotion

Ex: LA life : emoji - emotion - ranked categories

- Direct methods does not preserve ordinal nature of features.

Distance Preservation. MDS, Isomap

Topology / geometrical preservation Isomap

Information preservation. PCA

Capturing most of the variance

MNIST dataset

images  $\rightarrow$  0 1 2 --- 9    8x8 pixels

High dim Data - A selection from the 64 dimensional digits dataset.

Ex: MDS: colors have been clustered well, performed better than PCA.

ISOMAP: clearly made clusters of all digits clearly

unsupervised  $\rightarrow$  all dimensionality techniques.

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

ISOMAP - preserves global structure & focusses on geometry of data. the geodesic distance, which can be useful for understanding the underlying manifold or shape of the data

- computationally expensive ( $\because$  size of dataset or dim  $\uparrow$ )
- ALSO requires careful parameter selection for optimal results. (like # of neighbours).

Q. what to do when understanding local relationships is crucial, such as nbd relationships?

Ans: t-SNE (only looked locally).

- t-SNE provides better visualizations than other methods.
- It helps to uncover patterns, clusters & relationships in the data

#### \* Stochastic Neighbour Embedding (SNE)

$\rightarrow$  An unsupervised technique which focusses on preserving nbd's, instead of preserving distances.

Goal: To find a low-d map  
difference b/w the p (high-d) & q (low-d)  
distributions.

CLASSMATE  
Date \_\_\_\_\_  
Page \_\_\_\_\_

\* t-SNE algo

looking at neighbors

Ex: 8x8 pixel  $\Rightarrow$  64 dimensions

1 sample.

Data: data set  $X = \{x_1, x_2, x_3, \dots, x_n\}$

cost fn parameters: perplexity  $P_{\text{perp}}$

optimization parameters: # of iterations  $T$ ,

learning rate  $\eta$ , momentum  $\alpha(t)$

result: low-dim data represent<sup>n</sup>  $y^T = \{y_1, y_2, \dots, y_n\}$

begin :

compute pairwise affinities  $P_{ij|i}$  with perplexity  $P_{\text{perp}}$  (using eqn 1)

$$\text{set } P_{ij} = \frac{P_{ij|i}}{\sum_k P_{ik}}$$

sample initial solution  $y^{(0)} = \{y_1, y_2, \dots, y_n\}$  from  $\mathcal{N}(0, 10^{-4} I)$

for  $t=1$  to  $T$  do

compute low-dimensional affinities  $q_{ij}$  (eqn 4)

compute grad  $\frac{\delta C}{\delta y}$

$$\text{set } y^{(t)} = y^{(t-1)} + \eta \frac{\delta C}{\delta y} + \alpha(t)(y^{(t-1)} - y^{(t-1)})$$

end

end

→ The difference b/w the high-d & low-d maps  
are minimized using grad descent.

Compute probabilities  $P$  that  $x_i$  &  $x_j$  are neighbours,  
in high-d space.

perplexity  $\Rightarrow$  how many neighbours

(measures the effective # of neighbours, usually b/w 5 & 50).

Optimization  $\Rightarrow$  first define how many times we want to improve (usually 5000).

Parameters to speed up optimization & avoid poor local minima.

# of iterations: Number of steps need to reach the optimum. AD is iterative.

learning rate: AD steps in the dirn the error is the min. It defines how big will be the step.

$\alpha(t)$ : Momentum encourages a step that is in the same direction as previous steps.

$$q_{ij} - p_{ij} = \textcircled{c} \leftarrow \text{difference}$$

prob distribution of samples

how far am I from original prob distr' in high-d.

slope (in which direction)

1) Random guess

2) Probability distribution

3) distances (difference as small as possible)

4) slope (dir' in which diff ↑)

5) learning rate: Step size

compute pairwise affinities  $P_{ij}$

→ The probability that  $x_i$  would choose  $x_j$  as its neighbours in the high-d Space.

$x_i$  that are close (low euclidean dist)  
returns a high value.

$x_i$  that are far (high euclidean dist)  
returns a lower value.

σ is the variance, change of variance will change the values we assign to the distances.

Var will be different (clusters points are very close to each other.)

(using eqn 1)

The probability that  $x_i$  would choose  $x_j$  as its neighbours, in the high-d Space.

set  $p_{ij}$

- The low-d points are moved around to minimize the difference b/w the two distributions.
- we find a low-d represen<sup>n</sup> that captures the high-d data after successive iterations.

perplexity  
optimization

scale divergence Kullback - Leibler

$$C = KL(P||Q) = \sum \sum p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

↳ calc diff b/w 2 prob distributions

classmate

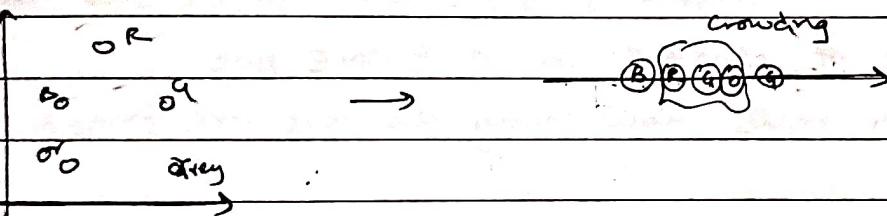
Date \_\_\_\_\_

Page \_\_\_\_\_

### \* Gaussian vs Student t-distribution

- Student's have longer tails compared to Gaussian.
- Gives high probabilities to points that are further away.
- Desirable, as we have limited low-d space & want to focus on modelling the close high-d points.

### Crowding in SNE



### probability distribution

### \* Dependence on hyperparameters

perplexity	2	5	30	50	100
step	5000	5000	5000	5000	5000

# of datapoints & steps  $\Rightarrow$  similar

(always work with values 5-50 for perplexity)

low values  $\rightarrow$  noises

high values  $\rightarrow$  more than # of points.

projecting in low dim,  $\rightarrow$  crowding  
can't give imp to distances

CLASSMATE  
Date \_\_\_\_\_  
Page \_\_\_\_\_

perplexity	30	30	30	30	30
Step	10	20	60	120	1000

separated clusters

SNE  $\rightarrow$  projecting in low dimension

clusters with different standard deviations & sizes

- ↳ has separated the into clusters
- ↳ From original graph, we cannot see relative size of clusters in a t-SNE plot.
- ↳ only tells how many clusters are formed.

SNE  $\rightarrow$ 

- good for interpreting # of clusters
- not good for understanding size, order of closeness in low & high.

\* Fine tuning perplexity

\* Take home Points:

understanding t-SNE  $\rightarrow$  Data Visualization

- Mainly used for visualization purposes, not dim reduc

Running t-SNE

- always run multiple trials
- use opt perplexity
- let the samples stabilize (iterations).

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

## Reading t-SNE

- Do not give importance to distances b/w far away points
- Do not give importance to density of clusters.
- Do not infer anything from a single o/p.