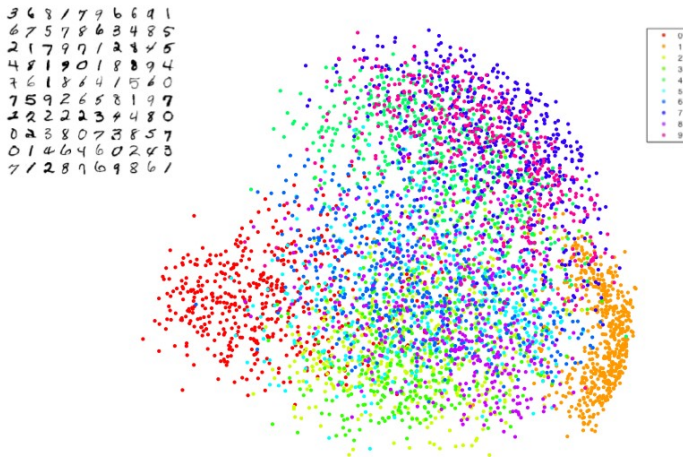


# Advanced Data Analytics

## Week 1: SNE

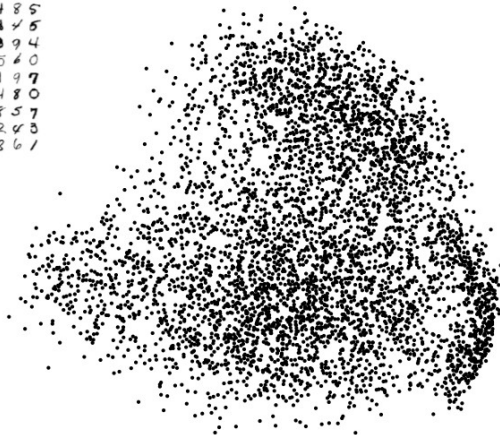
Ian T. Nabney

# Limitations of PCA: embedding MNIST



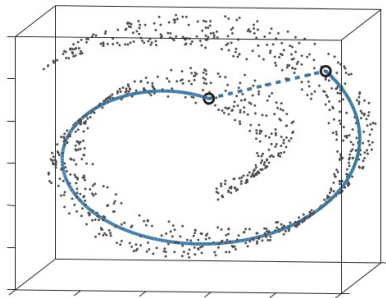
# PCA without colour coding

3 6 8 1 7 9 6 6 4 1  
6 7 5 7 8 6 3 4 8 5  
2 1 7 9 7 1 2 8 4 5  
4 8 1 9 0 1 8 8 9 4  
7 6 1 8 6 4 1 5 6 0  
7 5 9 2 6 5 8 1 9 7  
1 2 2 2 2 3 4 4 8 0  
0 2 3 8 0 7 3 8 5 7  
0 1 4 6 4 6 0 2 4 3  
7 1 2 8 1 6 9 8 6 1



# Local structure

- Stress metric measures all distances, but squared terms mean that large distances dominate.
- Similarly, techniques such as PCA consider the global representation.
- But is this how we see the data? Think back to Ware's book.
- Arguable that we actually understand neighbourhoods best.



# Stochastic Neighbourhood Embedding

- Start by converting Euclidean distances into conditional probability that represent similarities

$$p_{j|i} = \frac{\exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|\mathbf{x}_i - \mathbf{x}_k\|^2 / 2\sigma_i^2)} \quad (1)$$

where  $\sigma_i^2$  is the variance of a Gaussian centred on  $\mathbf{x}_i$ . We set  $p_{i|i} = 0$ .

- We create a similar metric in the low-dimensional space

$$q_{j|i} = \frac{\exp(-\|\mathbf{y}_i - \mathbf{y}_j\|^2)}{\sum_{k \neq i} \exp(-\|\mathbf{y}_i - \mathbf{y}_k\|^2)} \quad (2)$$

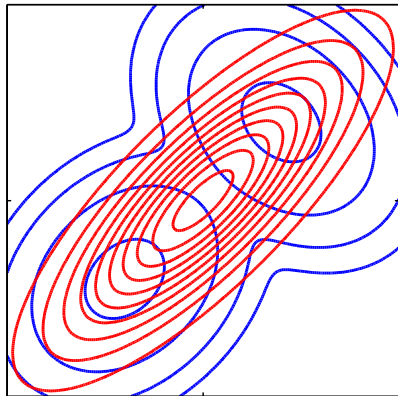
- If the mapped points  $\mathbf{y}_i$  and  $\mathbf{y}_j$  correctly model the similarity between the high-dimensional datapoints  $\mathbf{x}_i$  and  $\mathbf{x}_j$ , the conditional probabilities  $p_{j|i}$  and  $q_{j|i}$  will be equal. Measure the mismatch using the Kullback-Leibler divergence

$$C = \sum_i KL(P_i \| Q_i) = \sum_i \sum_j p_{j|i} \log \frac{p_{j|i}}{q_{j|i}}, \quad (3)$$

where  $P_i$  represents the conditional probability over all other datapoints given data point  $\mathbf{x}_i$  and similarly for  $Q_i$ . Minimise the cost function using gradient descent (and partial derivatives).

# Interpreting the cost function

- Different types of error in the pairwise distances in the low-dimensional map are not weighted equally.
- There is a large cost for using widely separated map points to represent nearby datapoints (i.e., for using a small  $q_{j|i}$  to model a large  $p_{j|i}$ ).
- There is only a small cost for using nearby map points to represent widely separated datapoints.



Blue contours show a bimodal distribution  $p(\mathbf{X})$  and the red contours correspond to a single Gaussian distribution  $q(\mathbf{Y})$  that minimises the KL-divergence  $KL(p||q)$ .

# Optimising the cost function

- The gradient has a surprisingly simple form

$$\frac{\partial C}{\partial y_i} = \sum_{j \neq i} (p_{j|i} - q_{j|i} + p_{i|j} - q_{i|j})(y_i - y_j) \quad (4)$$

and the update equation with a momentum term

$$\mathbf{Y}^{(t)} = \mathbf{Y}^{(t-1)} - \eta \frac{\partial C}{\partial \mathbf{y}_i} + \beta(t)(\mathbf{Y}^{(t-1)} - \mathbf{Y}^{(t-2)}) \quad (5)$$

- Physically, the gradient may be interpreted as the resultant force created by a set of springs between the map point  $\mathbf{y}_i$  and all other map points  $\mathbf{y}_j$ . All springs exert a force along the direction  $(\mathbf{y}_i - \mathbf{y}_j)$ .
- The force exerted by the spring between  $\mathbf{y}_i$  and  $\mathbf{y}_j$  is proportional to its length, and also proportional to its stiffness, which is the mismatch  $(p_{j|i} - q_{j|i} + p_{i|j} - q_{i|j})$  between the pairwise similarities of the data points and the map points.

- It is not likely that there is a single value of  $\sigma_i$  that is optimal for all datapoints in the data set because the density of the data is likely to vary.
- In dense regions, a smaller value of  $\sigma_i$  is usually more appropriate than in sparser regions.
- 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.
- Perplexity is defined as

$$\text{Perp}(P_i) = 2^{H(P_i)} \quad \text{where} \quad H(P_i) = - \sum_j p_{j|i} \log_2 p_{j|i}. \quad (6)$$

- The perplexity can be interpreted as a smooth measure of the effective number of neighbours. The performance of SNE is fairly robust to changes in the perplexity, and typical values are between 5 and 50.



# Effect of perplexity



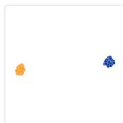
*Original*



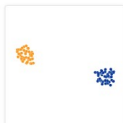
Perplexity: 2  
Step: 5,000



Perplexity: 5  
Step: 5,000



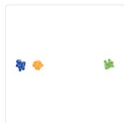
Perplexity: 30  
Step: 5,000



Perplexity: 50  
Step: 5,000



Perplexity: 100  
Step: 5,000



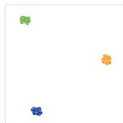
*Original*



Perplexity: 2  
Step: 5,000



Perplexity: 5  
Step: 5,000



Perplexity: 30  
Step: 5,000



Perplexity: 50  
Step: 5,000



Perplexity: 100  
Step: 5,000

- In high dimension we have more room, so points can have a lot of different neighbours while this is more difficult in lower dimension.
- This is the **crowding problem** – we don't have enough room to accommodate all neighbours. This is one of the biggest problems with SNE.
- t-SNE solution: Change the Gaussian in  $Q$  to a heavy-tailed distribution.
- If  $Q$  changes slower, we have more 'wiggle room' to place points at.

- Introduce two elements to address the crowding problem.
- Replace conditional probabilities  $p_{j|i}$  by joint probabilities  $p_{ji}$ .

$$p_{ij} = \frac{p_{j|i} + p_{i|j}}{2n} \quad \text{and} \quad q_{ij} = \frac{\exp(-\|\mathbf{y}_i - \mathbf{y}_j\|^2)}{\sum_{k \neq l} \exp(-\|\mathbf{y}_k - \mathbf{y}_l\|^2)} \quad (7)$$

- The gradient is simpler and faster to compute

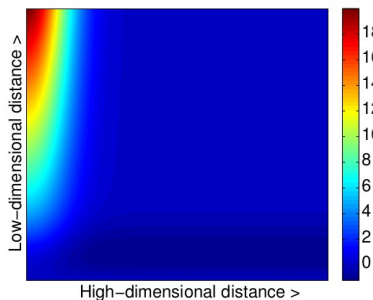
$$\frac{\partial C}{\partial \mathbf{y}_i} = 4 \sum_j (p_{ij} - q_{ij})(\mathbf{y}_i - \mathbf{y}_j). \quad (8)$$

## t-distribution Stochastic Neighbourhood Embedding

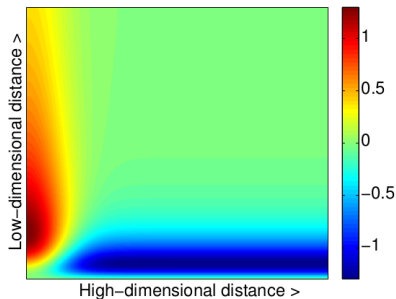
- Student-t density  $p(x) \propto (1 + \frac{x^2}{\nu})^{-(\nu+1)/2}$  which goes to zero much slower than a Gaussian. Choice of  $\nu = 1$  is equivalent to the Cauchy distribution.
- Equivalent to averaging Gaussians with a specific prior over  $\sigma^2$ . Removes need to optimise  $\sigma_i$ .
- Redefine  $q_{ij}$  but leave  $p_{ij}$  the same. This allows a moderate distance in the high-dimensional space to be faithfully modeled by a much larger distance in the map and so it eliminates the unwanted attractive forces between map points that represent moderately dissimilar datapoints.

$$q_{ij} = \frac{(1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2)^{-1}}{\sum_k \sum_{l \neq k} (1 + \|\mathbf{y}_k - \mathbf{y}_l\|^2)^{-1}} \quad (9)$$

$$\frac{\partial C}{\partial \mathbf{y}_i} = 4 \sum_{j \neq i} (p_{ij} - q_{ij}) (1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2)^{-1} (\mathbf{y}_i - \mathbf{y}_j) \quad (10)$$



Gradient of SNE



Gradient of t-SNE

# Experiments

## MNIST

- Randomly selected 6,000 images
- $28 \times 28 = 784$  pixels

## Olivetti faces

- 400 images (10 per individual)
- $92 \times 112 = 10,304$  pixels

## COIL-20

- 20 different objects and 72 equally spaced orientations, yielding a total of 1,440 images
- $32 \times 32 = 1024$  pixels

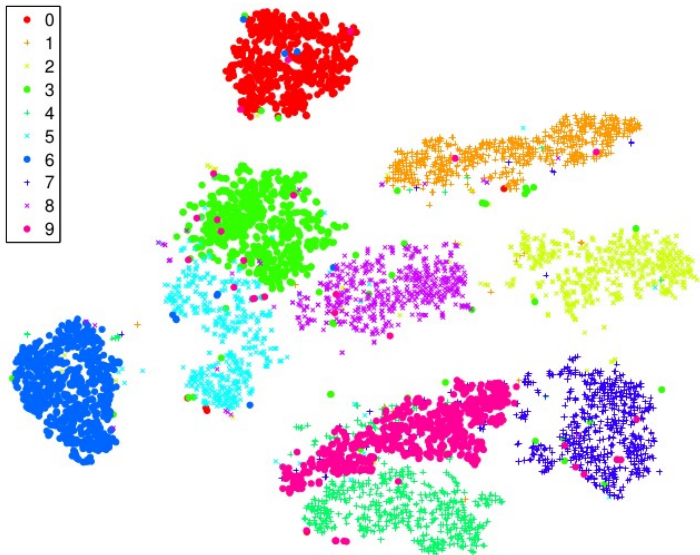
Use PCA to reduce dimensionality of data to 30.

# Experimental parameters

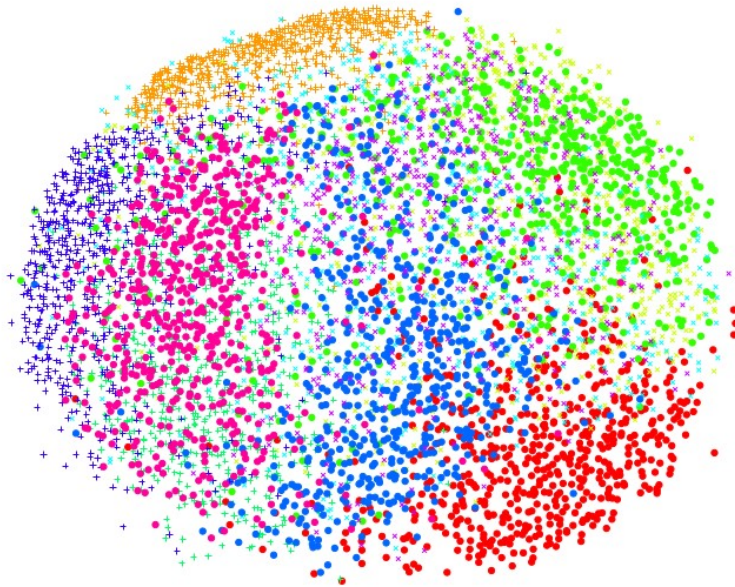
| <i>Technique</i> | <i>Cost function parameters</i> |
|------------------|---------------------------------|
| t-SNE            | $Perp = 40$                     |
| Sammon mapping   | none                            |
| Isomap           | $k = 12$                        |
| LLE              | $k = 12$                        |

Table 1: Cost function parameter settings for the experiments.

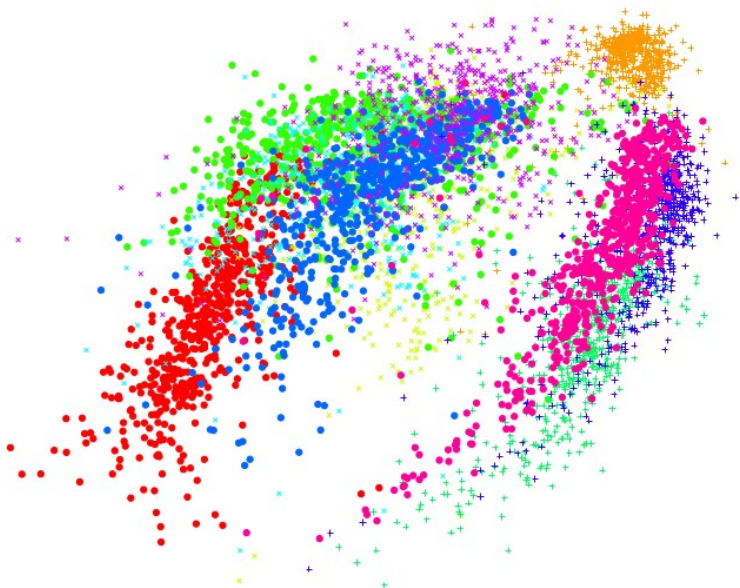
# MNIST t-SNE



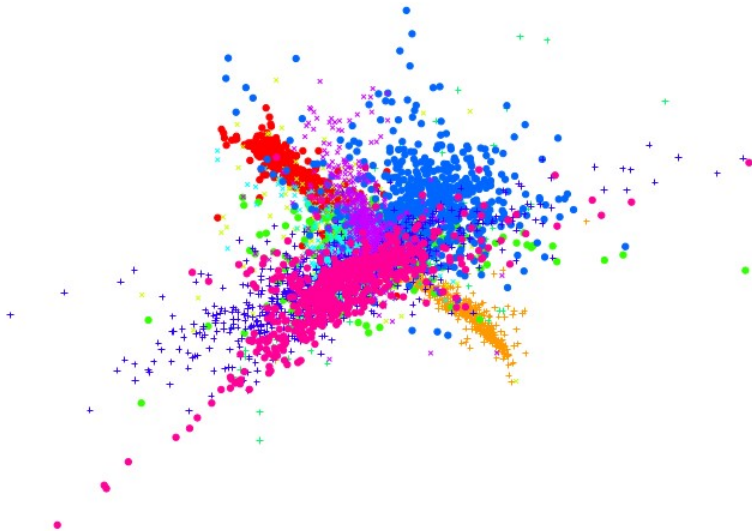




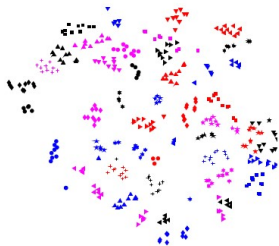
# MNIST Isomap



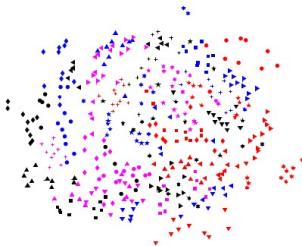
# MNIST local linear embedding (LLE)



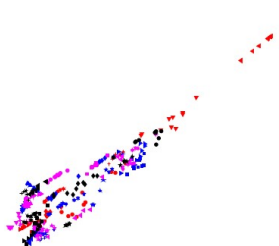
# Olivetti faces



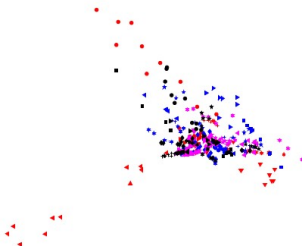
(a) Visualization by t-SNE.



(b) Visualization by Sammon mapping.



(c) Visualization by Isomap.

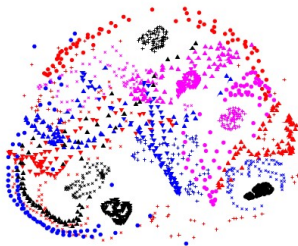


(d) Visualization by LLE.

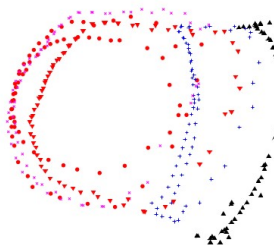
# COIL-20 results



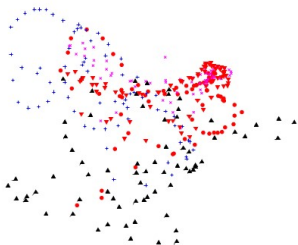
(a) Visualization by t-SNE.



(b) Visualization by Sammon mapping.



(c) Visualization by Isomap.



(d) Visualization by LLE.

- Both the computational and the memory complexity of t-SNE are  $O(n^2)$ , but the technique makes it possible to successfully visualize large real-world data sets with limited computational demands.
- We can reduce complexity from  $O(N^2)$  to  $O(N \log N)$  via Barnes Hut (tree-based) algorithm.
- Experiments on a variety of data sets show that t-SNE outperforms existing state-of-the-art techniques for visualizing a variety of real-world data sets.
- No functional mapping.