

# On t-SNE

Alex-Antoine Fortin

American Family Insurance

June 20, 2016

t-SNE stands for: t-distributed Stochastic Neighbor Embedding

## **desideratum**

We wish to have a representation in D-dim ( $R_D^*$ ) of an input in N-dim ( $R_N$ ), with  $N \gg D$

We wish that  $R_D^*$  preserves as much information as possible about  $R_N$

To do that, we need that  $R_D^*$  preserves the locality of the inputs in  $\mathbb{R}^N$

## Intuition

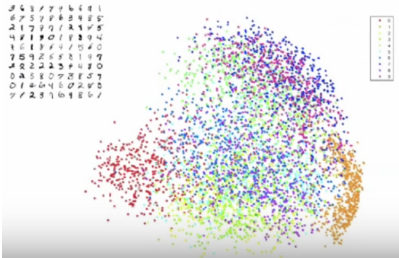
t-SNE achieves this in 3 steps:

- 1 fit a multivariate gaussian  $\mathcal{N}(\mu, \Sigma)$  to input  $R_N$  by Maximum Likelihood Estimation (MLE)
  - preserves locality of inputs by a cute trick (see *detailed calculations*)
- 2 Pose a multivariate student  $t_\nu(\mu, \Sigma)$ , with  $\nu = 1$  as the distribution of  $R_D^*$
- 3 fit  $R_D^*$  by minimizing the Kullback–Leibler divergence,  $D_{KL}(P||Q)$  between  $\mathcal{N}(\mu, \Sigma)$  and  $t_\nu(\mu, \Sigma)$

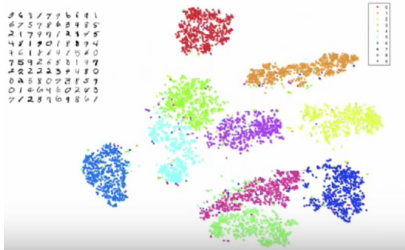
$D_{KL}(P||Q)$  is a measure of distance between 2 distributions and does not require the 2 distributions to be in the same dimension.

## t-SNE vs PCA: MNIST

Principal Components Analysis



t-Distributed Stochastic Neighbor Embedding



**Detailed calculations**

Step 1: Construct distribution of pairs of high-dim objects

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

Trick:

$$p_{ij} = \frac{p_{i|j} + p_{j|i}}{2N} \quad (2)$$

Step2: Define function in low-dim ( $\mathbf{y}_i \in \mathbb{R}^D$ , with  $D \ll N$ )

$$q_{ij} = \frac{(1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2)^{-1}}{\sum_k \sum_{k \neq l} (1 + \|\mathbf{y}_k - \mathbf{y}_l\|^2)^{-1}} \quad (3)$$

Step3: Minimizing the distance between  $p_{ij}$  and  $q_{ij}$   
loss function:

$$\sum_{i \neq j} p_{ij} (\log(p_{ij})) - p_{ij} (\log(q_{ij})) = \sum_{i \neq j} p_{ij} \cdot \log \left( \frac{p_{ij}}{q_{ij}} \right) = KL(p||q) \quad (4)$$

Using the  $KL(p||q)$  is great since it penalizes a lot when  $p_{ij}$  is big and very little when  $p_{ij}$  is small.