

$$p_{j|i} = \frac{\exp(-\|x_i - x_j\|^2 / 2\sigma^2)}{\sum_{k \neq i} \exp(-\|x_i - x_k\|^2 / 2\sigma^2)}, \quad p_{i|i} = 0, \quad p_{ij} = \frac{p_{i|j} + p_{j|i}}{2}$$

$\exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right) = \frac{1}{\sqrt{2\pi}^D} \exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right) \Rightarrow$  consider an increase in  $\sigma$  to be a change of scale  
 in the probabilities, the neighbour graph plots  
 should be identical but with a darker shade of  
 gray when  $\sigma = 5$

$\Rightarrow$  high similarity should result in darker groups  $\Rightarrow$  we should see 3 blocks in the neighbour graph plots

$\Rightarrow$  only two of them have darker blocks along the diagonal (b) & f))

$\Rightarrow$  since  $\sigma = 5$  should be darker  $\Rightarrow$ 

$$\begin{cases} \text{b)} \rightarrow \sigma = 5 \\ \text{f)} \rightarrow \sigma = 2 \end{cases}$$

② if  $k < D$ , we usually have less degrees of freedom to capture all of the relevant features in the input data.

So when applying this bottleneck in the reconstruction network, the most prominent features are learned and the rest are discarded. Since we lose some amount of information, although we are able to reconstruct most of the information, the data which was discarded can't be reconstructed, therefore we incur some reconstruction loss in our model.

• if there exists correlated data in our input we can model the full information even with the bottleneck layer.

Consider two parts of the input  $x_1$  and  $x_2$  where we have  $x_2 = \lambda x_1$  with some arbitrary  $\lambda \in \mathbb{R}$ ; the autoencoder, even with less neurons will be able to reconstruct the original information since two data points can be reconstructed with a learnt feature and different weights.

So, as long as the number of correlated data points  $\geq D - k$  the autoencoder should be able to have zero loss

formally:  $f(x) = xW_1W_2$  and  $x \in \mathbb{R}^D$ ,  $W_1 \in \mathbb{R}^{D \times k}$ ,  $W_2 \in \mathbb{R}^{k \times D}$

so as long as  $x$  can be fully represented in  $\mathbb{R}^k$  with no overlap between points

the autoencoder will have zero loss.