

t-SNE: Stochastic Neighbour Embedding

- Dim reduction that focuses on capturing local clusters

Visualisation of whole process:
<https://askadvaith-t-sne-visualisation.vercelapp/>

① Calculate pair-wise distances for every point

$$d_{ij} = \|x_i - x_j\| = \sqrt{\sum_k (x_{ik} - x_{jk})^2}$$

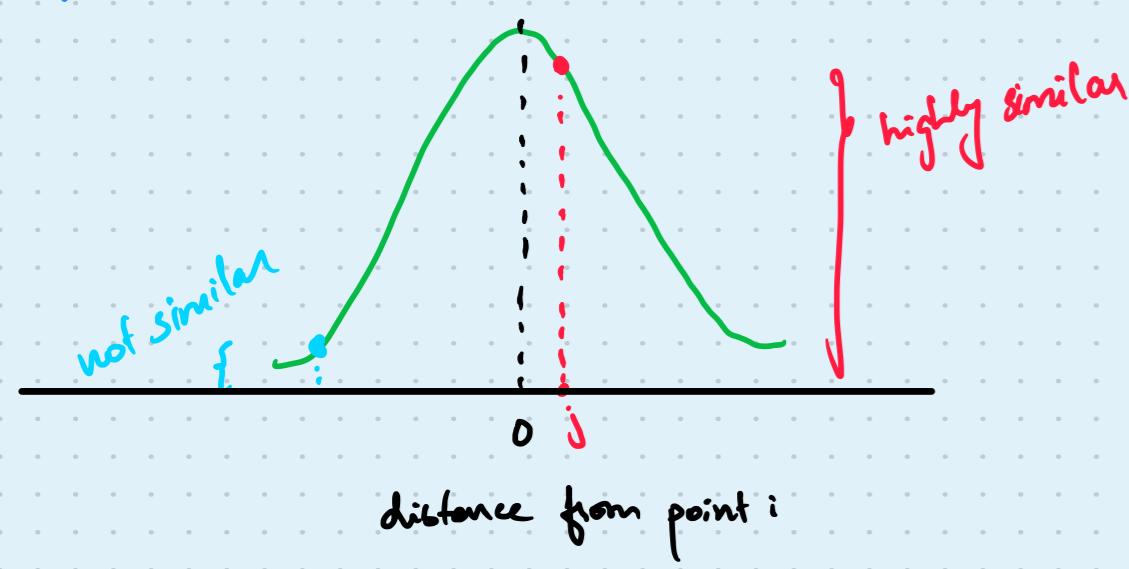
② Create normal distributions for every point.

- Determine σ for every point
- $\sigma = 5d$ (lower spread to weight similarity correctly)

Hyperparameter: Perplexity
 Defines local neighbourhood size, i.e.,
 how many neighbours are expected for
 a given point.

③ Calculate similarity to every other point using distance

for given point i :



$$P(j|i) = \frac{\exp\left(-\frac{d_{ij}^2}{2\sigma_i^2}\right)}{\sum_{k \neq i} \exp\left(-\frac{d_{ik}^2}{2\sigma_i^2}\right)}$$

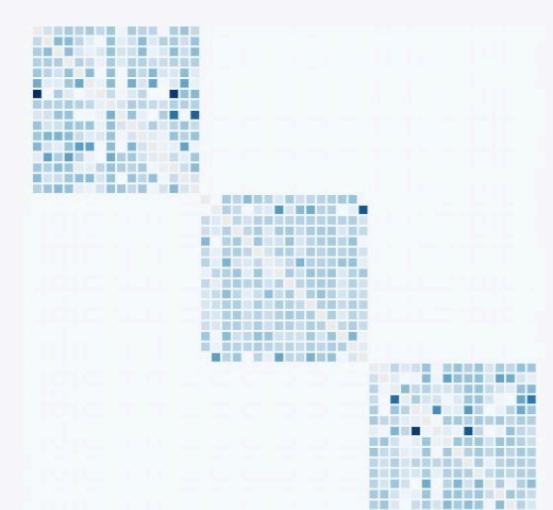
RBF for Gaussian:

$$\phi(x) = \exp\left(-\frac{x^2}{2\sigma^2}\right)$$

value determined by distance of x from origin

We thus construct a matrix P after doing this for every point. This matrix contains the similarity of every point to each other

P_{ij} - Asymmetric



④ Symmetrize P

$$\begin{matrix} & \vdots \\ & P(i|i) \\ \vdots & \\ j & \end{matrix} \quad \begin{matrix} P(j|i) \\ \vdots \\ j \end{matrix}$$

From earlier, not necessary that

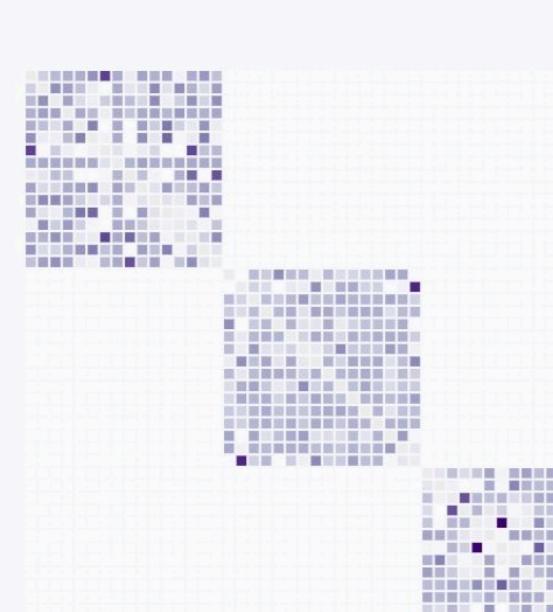
$$P(j|i) = P(i|j)$$

Depends on density etc.

Symmetrize:

$$P_{ij} = \frac{P(i|j) + P(j|i)}{2N}$$

P_{ij} - Symmetrized



⑤ Randomly init lower dimension representation

$$\begin{matrix} & \vdots \\ & \vdots \\ \vdots & \end{matrix} \quad \Rightarrow \quad \begin{matrix} & \dots \\ & \dots \end{matrix}$$

No bias.

⑥ Do the same similarity calculation in the lower dim with t-distribution

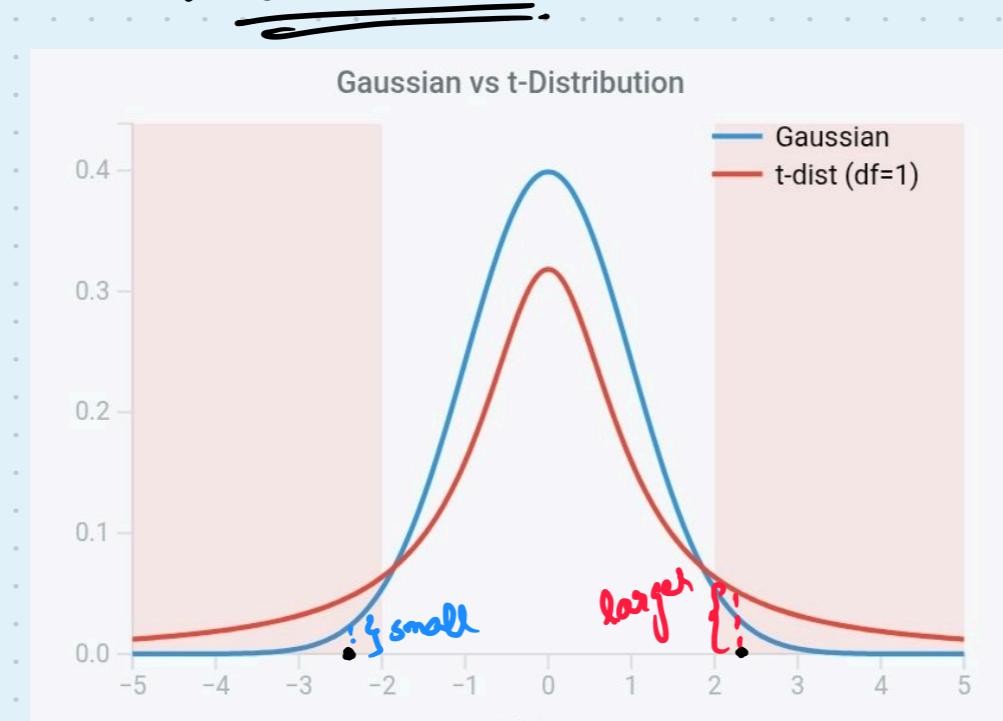
$$q_{ij} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_k (1 + \|y_i - y_k\|^2)^{-1}}$$

Why t-distribution?

Heavier tails \rightarrow less similar points given higher probability

prevents distant points from all bunching together

"Crowding problem" \downarrow not a problem in high D because more "space"/"volume" for points



⑦ Take matrix Q and make it match P

KL divergence between the two \rightarrow gradient gives direction to minimise difference between P and Q .

Attractive force: $P_{ij} > Q_{ij} \rightarrow$ points should be closer

Repulsive force: $P_{ij} < Q_{ij} \rightarrow$ points should be farther

$$\text{Loss} = \text{KL}(P || Q) = \sum_{i,j} P_{ij} \log \left[\frac{P_{ij}}{Q_{ij}} \right] \quad \left\{ \text{Minimise by gradient descent} \right.$$

Self Organising Maps (SOM)

- low dim (2D) map that preserves topological relationships
- Map high-D data onto a grid of neurons
- Similar inputs \rightarrow activate nearby neurons

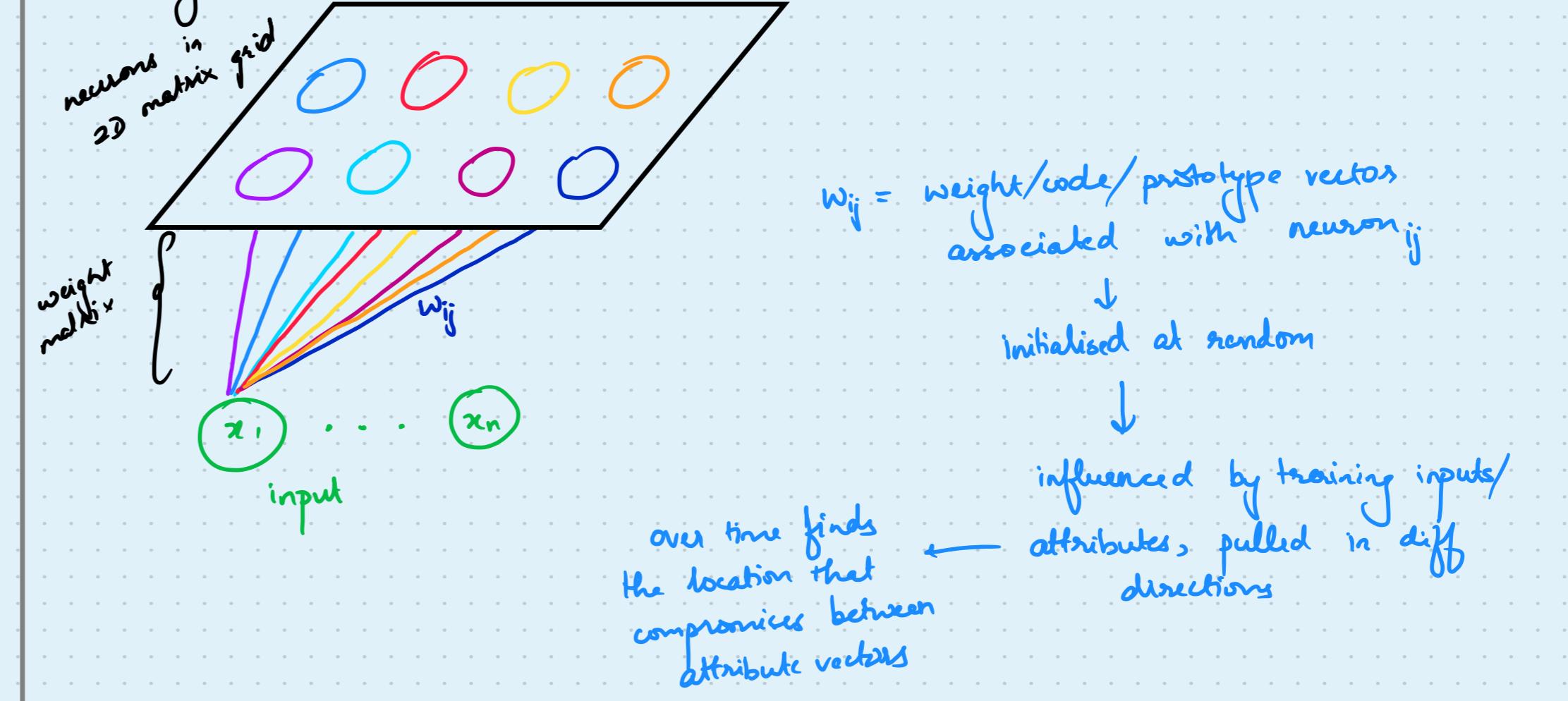
Differences

unlike PCA: non-linear

unlike t-SNE: creates regular grid structure

unlike clustering: preserves neighbourhood relationships

Working



Stochastic learning algorithm

```
for each  $w_i \in W$ 
    init  $w_i$ 
    while ! (stop condition):
        select random i/p vector from  $X$ 
        find best neuron for  $x_i$   $\rightarrow$  find closest  $w_i$  { competition
        for each  $w_i \in W$ :  $\downarrow$  cooperation + adaptation
            update  $w_i$ 
        end for
    end while
```

Batch learning algorithm

```
for each  $w_i \in W$ :
    init  $w_i$ 
    while ! (stop condition):
        for each  $w_i \in W$ :
             $X_i \leftarrow \emptyset$   $\downarrow$  set of vectors associated with  $w_i$ 
            if  $x_i \in (\text{neighbourhood of } w_i)$ :
                add  $x_i$  to  $X_i$ 
            end if
        end for
        for  $w_i \in W$ :
             $w_i \leftarrow \text{update } w_i \text{ mean over } X_i$ 
        end for
    end while
```

Outcomes

Topological map

- Winning neuron = $\arg\min_i \|x - w_i\|$
- Adapt weights:

$$w'_i = w_i + \eta(t) h_i((x - w_i))$$

η : learning rate
 t : time step
 $h_i(\cdot)$: neighborhood function
 $t \uparrow \Rightarrow \eta \downarrow$
 $h_i(\cdot) \downarrow$

- Closer neuron = stronger response to pattern \Rightarrow topological map.

Dimensionality reduction

Everything goes to 2D / low-dim grid

Phases

① Competition

$$\arg\min_i \|x - w_i\|$$

neuron closest \rightarrow Best Matching Unit (BMU)

② Cooperation

Find neighbours using h_{ij}

value as distance P
 symmetrical about winning neuron

③ Adaptation

$$w_{ij(n)} = w_{ij(n)} + \alpha h_{ij}(d(x, w_{ij(n)}))$$

for δ distance b/w i, j neurons

Additional notes

- Best lattice dim = \sqrt{N} total datapoints
- Very slow, unsupervised
- Not generative

- When to use
 - Want interpretable, visual representation
 - Need discrete clustering with spatial relationships
 - Data has natural neighborhood structure
 - Want to understand data distribution
 - Need robust method for noisy data
- When to stay away
 - Need exact distance preservation
 - Want continuous embedding space
 - Working with very high-dimensional data
 - Need fast processing for large datasets
 - Require deterministic results