

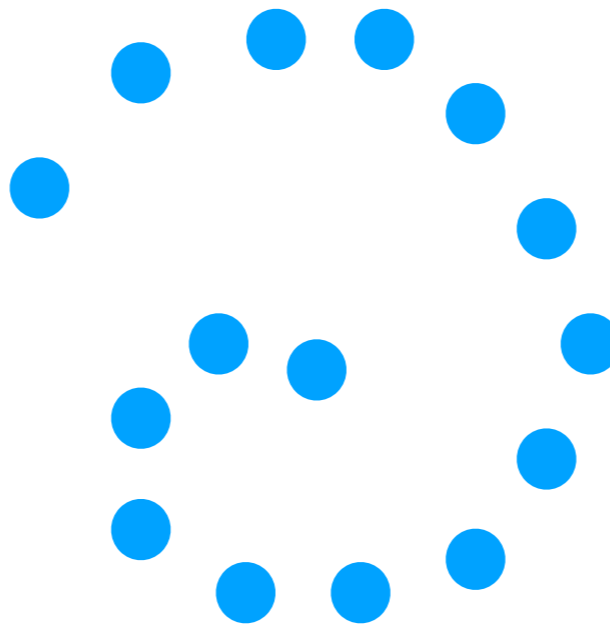
# Machine Learning for Data Science (CS4786)

## Lecture 9

Isomap + TSNE

# MANIFOLD BASED DIMENSIONALITY REDUCTION

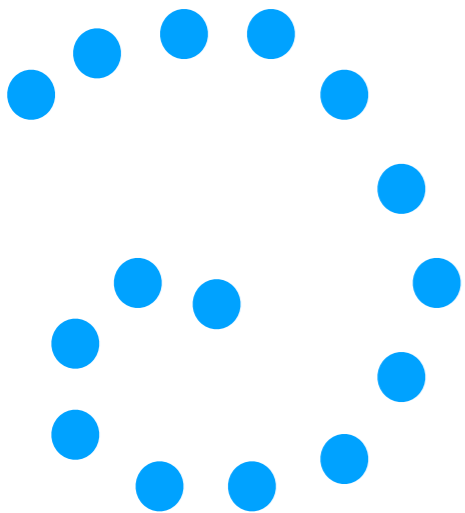
- Key Assumption: Points live on a low dimensional manifold
- Manifold: subspace that looks locally Euclidean
- Given data, can we uncover this manifold?



**Can we unfold this?**

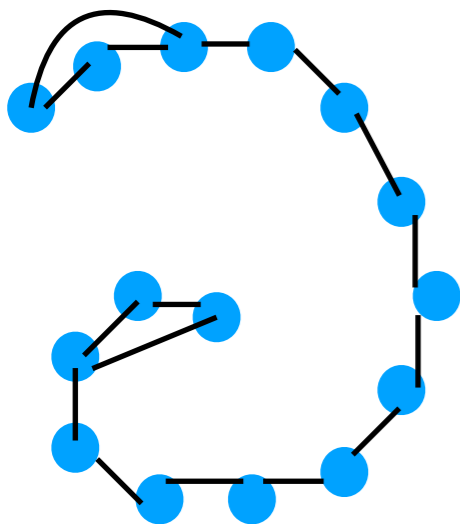
# METHOD I: ISOMAP

- 1 For every point, find its ( $k$ -) Nearest Neighbors



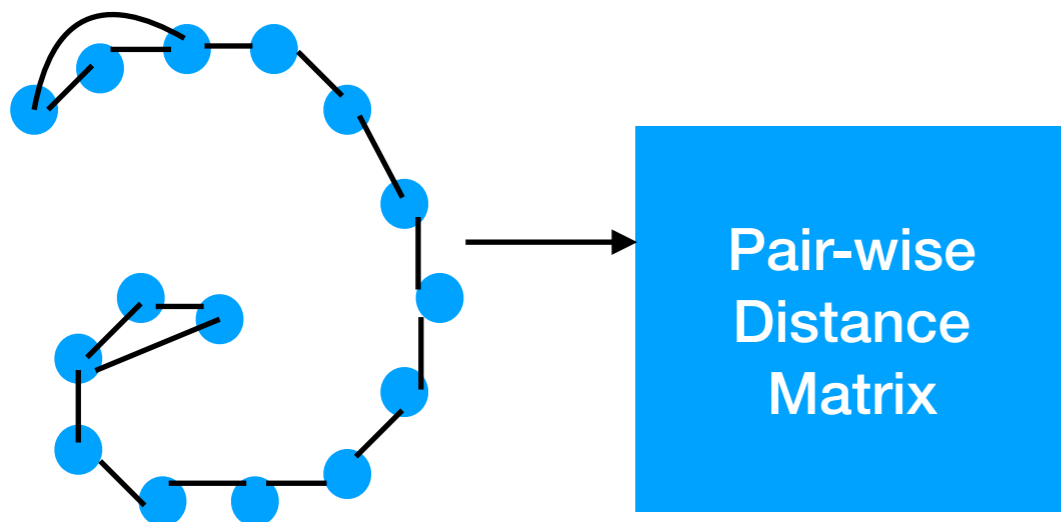
# METHOD I: ISOMAP

- 1 For every point, find its ( $k$ -) Nearest Neighbors
- 2 Form the Nearest Neighbor graph



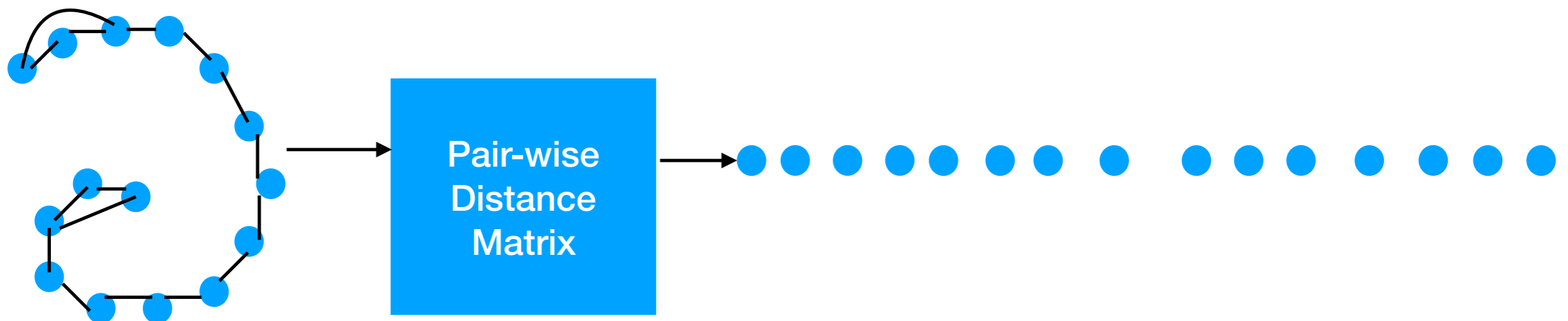
# METHOD I: ISOMAP

- 1 For every point, find its ( $k$ -) Nearest Neighbors
- 2 Form the Nearest Neighbor graph
- 3 For every pair of points  $A$  and  $B$ , distance between point  $A$  to  $B$  is shortest distance between  $A$  and  $B$  on graph



# METHOD I: ISOMAP

- 1 For every point, find its ( $k$ -) Nearest Neighbors
- 2 Form the Nearest Neighbor graph
- 3 For every pair of points  $A$  and  $B$ , distance between point  $A$  to  $B$  is shortest distance between  $A$  and  $B$  on graph
- 4 Find points in low dimensional space such that distances between points in this space is equal to distance on graph.



# ISOMAP: PITFALLS

- ① If we don't take enough nearest neighbors, then graph may not be connected
- ② If we connect points too far away, points that should not be connected can get connected
- ③ There may not be a right number of nearest neighbors we should consider!

# STOCHASTIC NEIGHBORHOOD EMBEDDING

- Use a probabilistic notion of which points are neighbors.



# STOCHASTIC NEIGHBORHOOD EMBEDDING

- Use a probabilistic notion of which points are neighbors.

Stochastic neighborhood distribution  $P$

# STOCHASTIC NEIGHBORHOOD EMBEDDING

- Use a probabilistic notion of which points are neighbors.  
    Stochastic neighborhood distribution  $P$
- Close by points are neighbors with high probability, ...



# STOCHASTIC NEIGHBORHOOD EMBEDDING

- Use a probabilistic notion of which points are neighbors.

Stochastic neighborhood distribution  $P$

- Close by points are neighbors with high probability, ...

Eg: For point  $\mathbf{x}_t$ , point  $\mathbf{x}_s$  is picked as neighbor with probability

$$p_{t \rightarrow s} = \frac{\exp\left(-\frac{\|\mathbf{x}_s - \mathbf{x}_t\|^2}{2\sigma^2}\right)}{\sum_{u \neq t} \exp\left(-\frac{\|\mathbf{x}_u - \mathbf{x}_t\|^2}{2\sigma^2}\right)}$$

Probability that points  $s$  and  $t$  are connected  $P_{s,t} = P_{t,s} = \frac{p_{t \rightarrow s} + p_{s \rightarrow t}}{2n}$

- Goal: Find  $\mathbf{y}_1, \dots, \mathbf{y}_n$  with stochastic neighborhood distribution  $Q$  such that “ $P$  and  $Q$  are similar”

# STOCHASTIC NEIGHBORHOOD EMBEDDING

- Use a probabilistic notion of which points are neighbors.

Stochastic neighborhood distribution  $P$

- Close by points are neighbors with high probability, ...  
Eg: For point  $\mathbf{x}_t$ , point  $\mathbf{x}_s$  is picked as neighbor with probability

$$p_{t \rightarrow s} = \frac{\exp\left(-\frac{\|\mathbf{x}_s - \mathbf{x}_t\|^2}{2\sigma^2}\right)}{\sum_{u \neq t} \exp\left(-\frac{\|\mathbf{x}_u - \mathbf{x}_t\|^2}{2\sigma^2}\right)}$$

Probability that points  $s$  and  $t$  are connected  $P_{s,t} = P_{t,s} = \frac{p_{t \rightarrow s} + p_{s \rightarrow t}}{2n}$

- Goal: Find  $\mathbf{y}_1, \dots, \mathbf{y}_n$  with stochastic neighborhood distribution  $Q$  such that “ $P$  and  $Q$  are similar”

i.e. minimize:

$$\text{KL}(P \parallel Q) = \sum_{s,t} P_{s,t} \log \left( \frac{P_{s,t}}{Q_{s,t}} \right) = \sum_{s,t} P_{s,t} \log (P_{s,t}) - \sum_{s,t} P_{s,t} \log (Q_{s,t})$$

# CHOICE FOR $Q$

- Just like we defined  $P$ , we can define  $Q$  for a given  $\mathbf{y}_1, \dots, \mathbf{y}_n$  by

$$q_{t \rightarrow s} = \frac{\exp\left(-\frac{\|\mathbf{y}_s - \mathbf{y}_t\|^2}{2\sigma^2}\right)}{\sum_{u \neq t} \exp\left(-\frac{\|\mathbf{y}_u - \mathbf{y}_t\|^2}{2\sigma^2}\right)}$$

and then set  $Q_{s,t} = \frac{q_{t \rightarrow s} + q_{s \rightarrow t}}{2n}$

# CHOICE FOR $Q$

- Just like we defined  $P$ , we can define  $Q$  for a given  $\mathbf{y}_1, \dots, \mathbf{y}_n$  by

$$q_{t \rightarrow s} = \frac{\exp\left(-\frac{\|\mathbf{y}_s - \mathbf{y}_t\|^2}{2\sigma^2}\right)}{\sum_{u \neq t} \exp\left(-\frac{\|\mathbf{y}_u - \mathbf{y}_t\|^2}{2\sigma^2}\right)}$$

and then set  $Q_{s,t} = \frac{q_{t \rightarrow s} + q_{s \rightarrow t}}{2n}$

- However we are faced with the crowding problem:

# CHOICE FOR $Q$

- Just like we defined  $P$ , we can define  $Q$  for a given  $\mathbf{y}_1, \dots, \mathbf{y}_n$  by

$$q_{t \rightarrow s} = \frac{\exp\left(-\frac{\|\mathbf{y}_s - \mathbf{y}_t\|^2}{2\sigma^2}\right)}{\sum_{u \neq t} \exp\left(-\frac{\|\mathbf{y}_u - \mathbf{y}_t\|^2}{2\sigma^2}\right)}$$

and then set  $Q_{s,t} = \frac{q_{t \rightarrow s} + q_{s \rightarrow t}}{2n}$

- However we are faced with the crowding problem:
  - In high dimension we have a lot of space, Eg. in  $d$  dimension we have  $d + 1$  equidistant point



# CHOICE FOR $Q$

- Just like we defined  $P$ , we can define  $Q$  for a given  $\mathbf{y}_1, \dots, \mathbf{y}_n$  by

$$q_{t \rightarrow s} = \frac{\exp\left(-\frac{\|\mathbf{y}_s - \mathbf{y}_t\|^2}{2\sigma^2}\right)}{\sum_{u \neq t} \exp\left(-\frac{\|\mathbf{y}_u - \mathbf{y}_t\|^2}{2\sigma^2}\right)}$$

and then set  $Q_{s,t} = \frac{q_{t \rightarrow s} + q_{s \rightarrow t}}{2n}$

- However we are faced with the crowding problem:
  - In high dimension we have a lot of space, Eg. in  $d$  dimension we have  $d + 1$  equidistant point
  - For  $d$  dimensional gaussians, most points are found at distance  $\sqrt{d}$  from mean!

# CHOICE FOR $Q$

- Just like we defined  $P$ , we can define  $Q$  for a given  $\mathbf{y}_1, \dots, \mathbf{y}_n$  by

$$q_{t \rightarrow s} = \frac{\exp\left(-\frac{\|\mathbf{y}_s - \mathbf{y}_t\|^2}{2\sigma^2}\right)}{\sum_{u \neq t} \exp\left(-\frac{\|\mathbf{y}_u - \mathbf{y}_t\|^2}{2\sigma^2}\right)}$$

and then set  $Q_{s,t} = \frac{q_{t \rightarrow s} + q_{s \rightarrow t}}{2n}$

- However we are faced with the crowding problem:
  - In high dimension we have a lot of space, Eg. in  $d$  dimension we have  $d + 1$  equidistant point
  - For  $d$  dimensional gaussians, most points are found at distance  $\sqrt{d}$  from mean!
  - If we use gaussians in both high and low dimensional space, all the points are squished in to a small space
  - Too many points crowd the center!

# METHOD II: T-SNE

- Instead for  $Q$  we use, student  $t$  distribution which is heavy tailed:

$$q_{t \rightarrow s} = \frac{(1 + \|\mathbf{y}_s - \mathbf{y}_t\|^2)^{-1}}{\sum_{u \neq t} (1 + \|\mathbf{y}_u - \mathbf{y}_t\|^2)^{-1}}$$

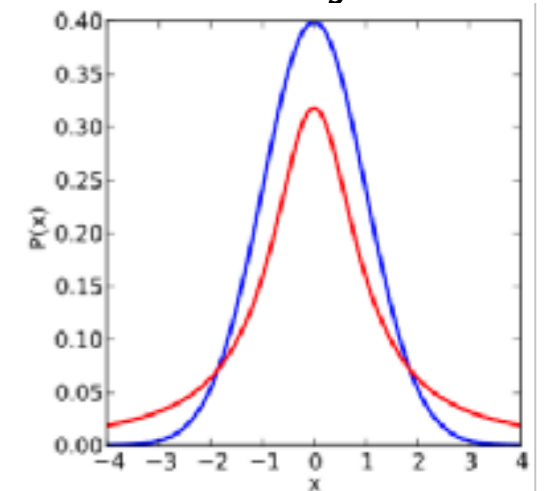
and then set  $Q_{s,t} = \frac{q_{t \rightarrow s} + q_{s \rightarrow t}}{2n}$

# METHOD II: T-SNE

- Instead for  $Q$  we use, student  $t$  distribution which is heavy tailed:

$$q_{t \rightarrow s} = \frac{(1 + \|\mathbf{y}_s - \mathbf{y}_t\|^2)^{-1}}{\sum_{u \neq t} (1 + \|\mathbf{y}_u - \mathbf{y}_t\|^2)^{-1}}$$

and then set  $Q_{s,t} = \frac{q_{t \rightarrow s} + q_{s \rightarrow t}}{2n}$

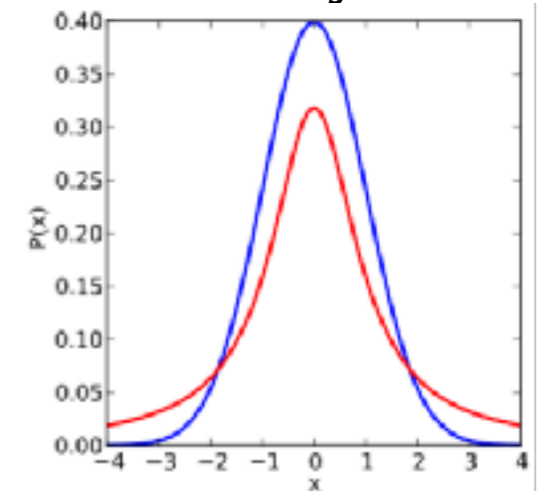


# METHOD II: T-SNE

- Instead for  $Q$  we use, student  $t$  distribution which is heavy tailed:

$$q_{t \rightarrow s} = \frac{(1 + \|\mathbf{y}_s - \mathbf{y}_t\|^2)^{-1}}{\sum_{u \neq t} (1 + \|\mathbf{y}_u - \mathbf{y}_t\|^2)^{-1}}$$

and then set  $Q_{s,t} = \frac{q_{t \rightarrow s} + q_{s \rightarrow t}}{2n}$



- It can be verified that

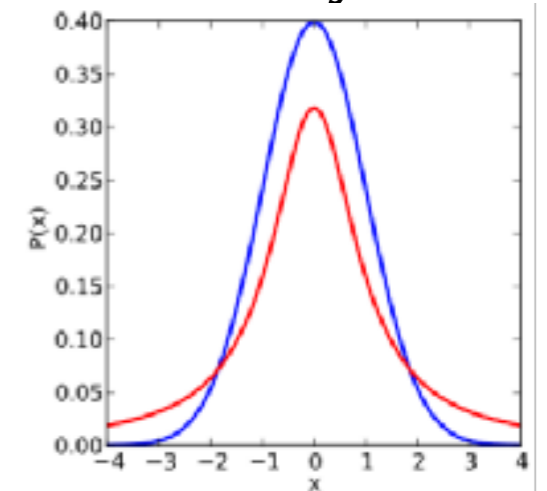
$$\nabla_{\mathbf{y}_t} \text{KL}(P \| Q) = 4 \sum_{s=1}^n (P_{s,t} - Q_{s,t}) (\mathbf{y}_t - \mathbf{y}_s) (1 + \|\mathbf{y}_s - \mathbf{y}_t\|^2)^{-1}$$

# METHOD II: T-SNE

- Instead for  $Q$  we use, student  $t$  distribution which is heavy tailed:

$$q_{t \rightarrow s} = \frac{(1 + \|\mathbf{y}_s - \mathbf{y}_t\|^2)^{-1}}{\sum_{u \neq t} (1 + \|\mathbf{y}_u - \mathbf{y}_t\|^2)^{-1}}$$

and then set  $Q_{s,t} = \frac{q_{t \rightarrow s} + q_{s \rightarrow t}}{2n}$



- It can be verified that

$$\nabla_{\mathbf{y}_t} \text{KL}(P \| Q) = 4 \sum_{s=1}^n (P_{s,t} - Q_{s,t}) (\mathbf{y}_t - \mathbf{y}_s) (1 + \|\mathbf{y}_s - \mathbf{y}_t\|^2)^{-1}$$

- Algorithm: Find  $\mathbf{y}_1, \dots, \mathbf{y}_n$  by performing gradient descent

Demo