# Similarity Graphs

Input: $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \ldots, \mathbf{x}_N$

- raw data
- flat data
- vectorial data

# Similarity Graphs

Similarity graph: $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ — **(un)weighted**

*Task 1:* For each pair $i$, $j$: define a **similarity function** $s_{ij}$

*Task 2:* Decide which edges to include

## Similarity Graphs

Similarity graph: $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ — **(un)weighted**

*Task 1:* For each pair $i$, $j$: define a **similarity function** $s_{ij}$

*Task 2:* Decide which edges to include

fully connected graphs - consider everything

$\varepsilon$-neighborhood graphs – connect the points with the distances smaller than $\varepsilon$

$k$-NN neighborhood graphs – take $k$ nearest neighbors

## Similarity Graphs

Similarity graph: $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ — **(un)weighted**

Task 1: For each pair $i$, $j$: define a **similarity function** $s_{ij}$

Task 2: Decide which edges to include

fully connected graphs - consider everything

$\varepsilon$-neighborhood graphs – connect the points with the distances smaller than $\varepsilon$

$k$-NN neighborhood graphs – take $k$ nearest neighbors

*This is art (not much theory exists).*

http://www.informatik.uni-hamburg.de/ML/contents/people/luxburg/
publications/Luxburg07_tutorial.pdf

# Similarity Graphs: Fully connected graphs

Edges connect everything.

# Similarity Graphs: Fully connected graphs

Edges connect everything.

- choose a "meaningful" similarity function $s$
- default choice:

$$s_{ij} = \exp\left(\frac{-\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right)$$

- why the exponential decay with the distance?

# Similarity Graphs: Fully connected graphs

Edges connect everything.

- choose a "meaningful" similarity function $s$
- default choice:

$$s_{ij} = \exp\left(\frac{-\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right)$$

- why the exponential decay with the distance?
- $\sigma$ controls the width of the neighborhoods
  - **a** practical rule of thumb: 10% of the average empirical std
  - possibility: learn $\sigma_i$ for each feature independently
- metric learning (a whole field of ML)

# Similarity Graphs: $\varepsilon$-neighborhood graphs

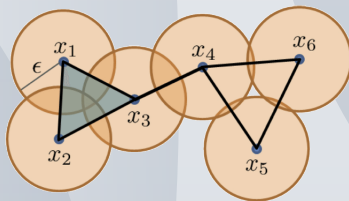Edges connect the points with the distances smaller than $\varepsilon$.



Figure: An $\varepsilon$-graph. Source: Illustration of a Rips complex

# Similarity Graphs: $\varepsilon$-neighborhood graphs

Edges connect the points with the distances smaller than $\varepsilon$.

- distances are roughly on the same scale $(\varepsilon)$
- weights may not bring additional info $\rightarrow$ unweighted
- equivalent to: similarity function is at least $\varepsilon$
- theory [Penrose, 1999]: $\varepsilon = ((\log N)/N)^d$ to guarantee connectivity $N$ nodes, $d$ dimension
- practice: choose $\varepsilon$ as the length of the longest edge in the MST - minimum spanning tree

What could be the problem with this MST approach?

# Similarity Graphs: $\varepsilon$-neighborhood graphs

Edges connect the points with the distances smaller than $\varepsilon$.

- distances are roughly on the same scale ($\varepsilon$)
- weights may not bring additional info $\rightarrow$ unweighted
- equivalent to: similarity function is at least $\varepsilon$
- theory [Penrose, 1999]: $\varepsilon = ((\log N)/N)^d$ to guarantee connectivity $N$ nodes, $d$ dimension
- practice: choose $\varepsilon$ as the length of the longest edge in the MST - minimum spanning tree

What could be the problem with this MST approach?

Anomalies can make $\varepsilon$ too large.

# Similarity Graphs: $k$-nearest neighbors graphs

Edges connect each node to its $k$-nearest neighbors.

# Similarity Graphs: $k$-nearest neighbors graphs

Edges connect each node to its $k$-nearest neighbors.

- asymmetric (or directed graph)
  - option OR: ignore the direction
  - option AND: include if we have both direction (mutual $k$-NN)

- how to choose $k$?

- $k \approx \log N$ - suggested by asymptotics (practice: up to $\sqrt{N}$)

- for mutual $k$-NN we need to take larger $k$

- mutual $k$-NN does not connect regions with different density

- why don't we take $k = N - 1$?
  - space and time
  - manifold considerations (preserving local properties)

# Similarity Graphs: Important considerations

- *calculate all $s_{ij}$ and threshold* has its limits ($N \approx 10000$)

# Similarity Graphs: Important considerations

- *calculate all $s_{ij}$ and threshold* has its limits ($N \approx 10000$)
- graph construction step can be a huge bottleneck

# Similarity Graphs: Important considerations

- *calculate all $s_{ij}$ and threshold* has its limits ($N \approx 10000$)
- graph construction step can be a huge bottleneck
- want to go higher? (we often have to)

# Similarity Graphs: Important considerations

- *calculate all $s_{ij}$ and threshold* has its limits ($N \approx 10000$)
- graph construction step can be a huge bottleneck
- want to go higher? (we often have to)
  - down-sample

# Similarity Graphs: Important considerations

- *calculate all $s_{ij}$ and threshold* has its limits ($N \approx 10000$)
- graph construction step can be a huge bottleneck
- want to go higher? (we often have to)
  - down-sample
  - approximate NN

# Similarity Graphs: Important considerations

- *calculate all $s_{ij}$ and threshold* has its limits ($N \approx 10000$)
- graph construction step can be a huge bottleneck
- want to go higher? (we often have to)
    - down-sample
    - approximate NN
        - ▶ **LSH** - Locally Sensitive Hashing

# Similarity Graphs: Important considerations

- *calculate all $s_{ij}$ and threshold* has its limits ($N \approx 10000$)
- graph construction step can be a huge bottleneck
- want to go higher? (we often have to)
  - down-sample
  - approximate NN
    - **LSH** - Locally Sensitive Hashing
    - **CoverTrees**

# Similarity Graphs: Important considerations

- *calculate all $s_{ij}$ and threshold* has its limits ($N \approx 10000$)
- graph construction step can be a huge bottleneck
- want to go higher? (we often have to)
    - down-sample
    - approximate NN
        - **LSH** - Locally Sensitive Hashing
        - **CoverTrees**
    - sometime we may not need the graph (just the final results)

# Similarity Graphs: Important considerations

- *calculate all $s_{ij}$ and threshold* has its limits ($N \approx 10000$)
- graph construction step can be a huge bottleneck
- want to go higher? (we often have to)
  - down-sample
  - approximate NN
    - ▶ **LSH** - Locally Sensitive Hashing
    - ▶ **CoverTrees**
  - sometime we may not need the graph (just the final results)
  - yet another story: when we start with a large graph and want to make it sparse (later in the course)

# Similarity Graphs: Important considerations

- *calculate all $s_{ij}$ and threshold* has its limits ($N \approx 10000$)
- graph construction step can be a huge bottleneck
- want to go higher? (we often have to)
  - down-sample
  - approximate NN
    - ▶ **LSH** - Locally Sensitive Hashing
    - ▶ **CoverTrees**
  - sometime we may not need the graph (just the final results)
  - yet another story: when we start with a large graph and want to make it sparse (later in the course)
- these rules have little theoretical underpinning

# Similarity Graphs: Important considerations

- *calculate all $s_{ij}$ and threshold* has its limits ($N \approx 10000$)
- graph construction step can be a huge bottleneck
- want to go higher? (we often have to)
  - down-sample
  - approximate NN
    - **LSH** - Locally Sensitive Hashing
    - **CoverTrees**
  - sometime we may not need the graph (just the final results)
  - yet another story: when we start with a large graph and want to make it sparse (later in the course)
- these rules have little theoretical underpinning
- similarity is very data-dependent

# Michal Valko

michal.valko@inria.fr

Inria & ENS Paris-Saclay, MVA


https://misovalko.github.io/mva-ml-graphs.html