

CS5284 : Graph Machine Learning

Lecture 2 : Introduction to Graph Science

Semester 1 2024/25

Xavier Bresson

<https://x.com/xbresson>

Department of Computer Science
National University of Singapore (NUS)



Course lectures

- Introduction to Graph Machine Learning
- Part 1: GML without feature learning (before 2014)
 - • Introduction to Graph Science
 - Graph Analysis Techniques without Feature Learning
 - Graph clustering
 - Graph SVM
 - Recommendation on graphs
 - Graph-based visualization
- Part 2 : GML with shallow feature learning (2014-2016)
 - Shallow graph feature learning
- Part 3 : GML with deep feature learning, a.k.a. GNNs (after 2016)
 - Graph Convolutional Networks (spectral and spatial)
 - Weisfeiler-Lehman GNNs
 - Graph Transformer & Graph ViT
 - Graph generation & molecular science
 - GNNs for combinatorial optimization
 - Integrating GNNs and LLMs
 - Benchmarking GNNs
 - GNNs for recommendation & knowledge graphs

Outline

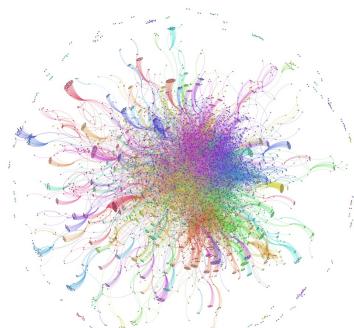
- Graph theory
- Graph categories
- Basic definitions
- Curse of dimensionality and structure
- Manifolds and graphs
- Spectral graph theory
- Graph construction
- Conclusion

Outline

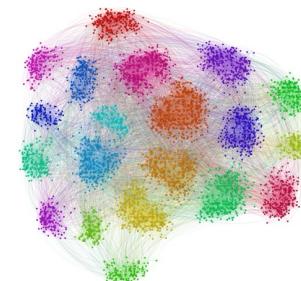
- Graph theory
- Graph categories
- Basic definitions
- Curse of dimensionality and structure
- Manifolds and graphs
- Spectral graph theory
- Graph construction
- Conclusion

Graphs

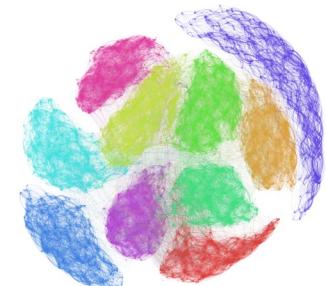
- Graphs encode complex data structures.
 - They are everywhere! Internet, social networks, customer-product relationships, etc



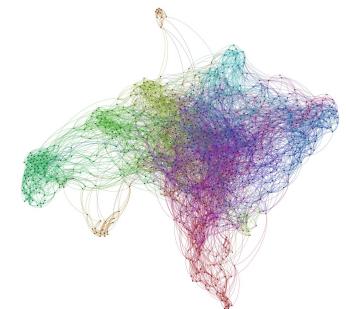
Internet Network of California



Social Network

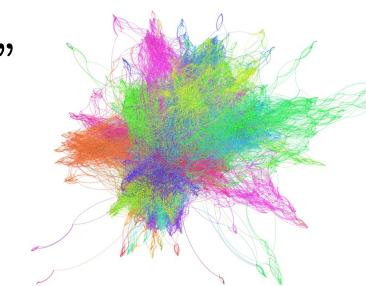


MNIST Image Network



GTZAN Music Network

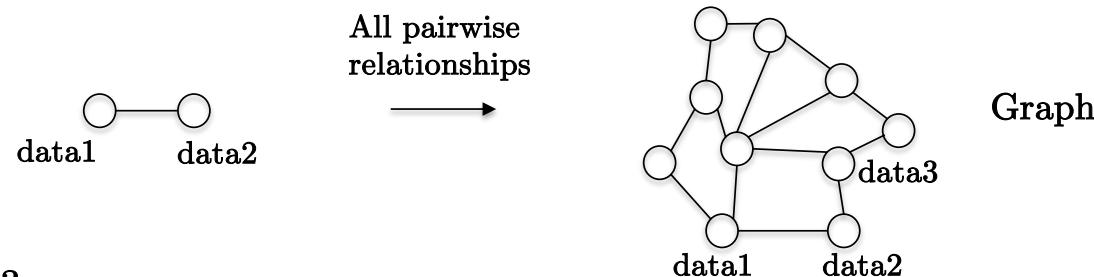
- “Graphs are the most important discrete models in the world.”
 - Gil Strang (MIT)



Network of Text Documents
20newsgroups

Graphs

- Definition : (Simple) mathematical model representing pairwise relationships between data.



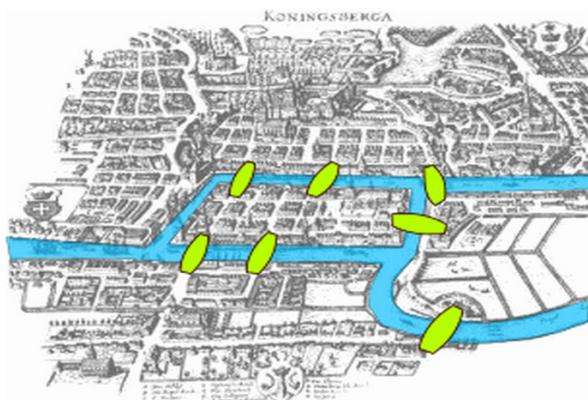
- Why are graphs useful?
 - Graphs offer a global view and analysis of data structures.
 - They possess meaningful patterns, i.e. insights about data properties.
 - Some tasks are exclusively designed for graphs, e.g. Google PageRank recommendation.
 - They can boost performance with additional priors, a.k.a. inductive bias.
 - They can benefit from GPUs as graphs are (sparse) matrices.



Handwritten note: "Helpful : many zeros → simple calculations"

Graph theory

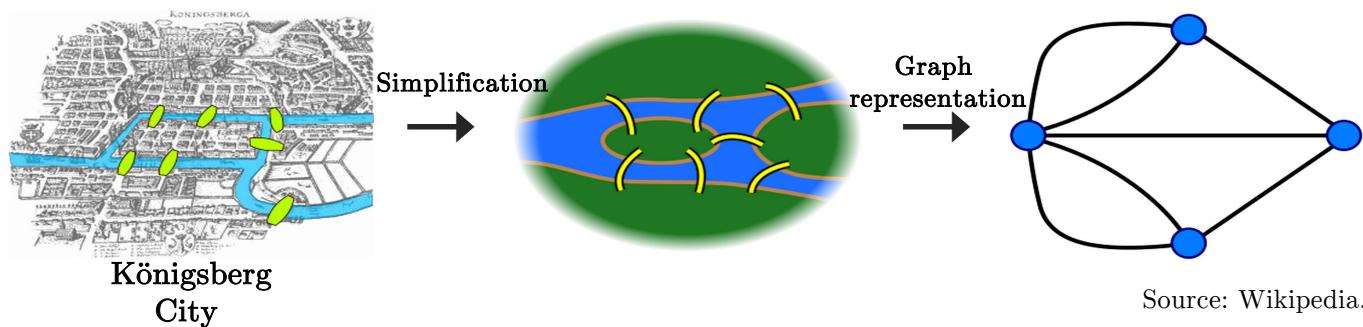
- When did it start?
 - History of graph theory : Graphs have been formally studied since 1736, starting with Mathematician Leonhard Euler and the famous problem of “Seven Bridges of Königsberg” :



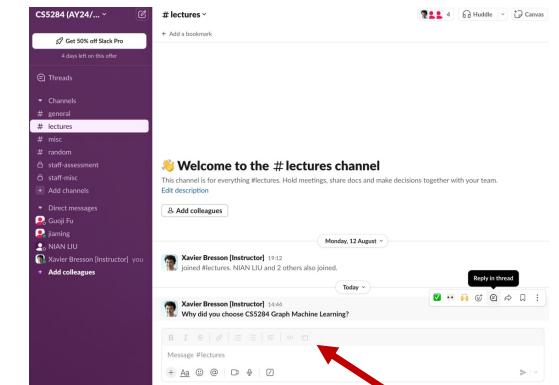
Königsberg
City

In-lecture question

- Can we find a path through the city (starting from any place) that crosses each bridge once and only once? Justify your response.



- In Slack #lectures
 - Identify the question and Reply in thread with a short response



Graph theory

- When did it start?
 - History of graph theory : Graphs have been formally studied since 1736, starting with Mathematician Leonhard Euler and the famous problem of “Seven Bridges of Königsberg” :
- Graph theory have since developed tools to analyze and process networks for all sort of applications : clustering, classification, visualization, recommendation, etc.

Outline

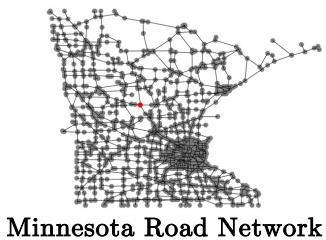
- Graph theory
- **Graph categories**
- Basic definitions
- Curse of dimensionality and structure
- Manifolds and graphs
- Spectral graph theory
- Graph construction
- Conclusion

Focus for GNNs
 (1) Artificial networks
 can now capture all
 communication & interconnect
 ion via a natural graph

Graph categories

- Natural Graphs
 - Social networks : Meta, LinkedIn, Twitter
 - Biological networks : Brain connectivity & functionality, gene regulatory networks
 - Communication networks : Internet, networking devices
 - Transportation networks : Trains, cars, airplanes, pedestrians
 - Power networks : Electricity, water
- Natural graphs mean graphs that are not artificially hand-crafted/constructed.

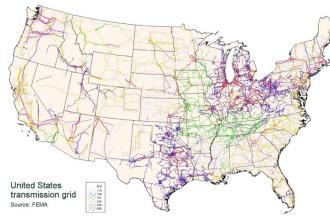
ISn't Meta hand-crafted b/c



Minnesota Road Network



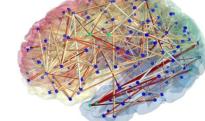
Telecommunication Network



US Electrical Network



Facebook



Brain
Connectivity

=

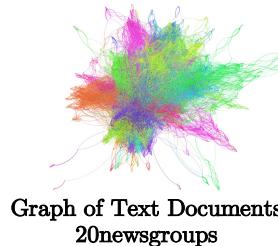


Graphs

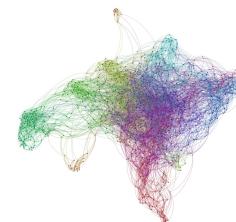
Graph categories

- Graphs constructed from data :
 - MNIST image network
 - GTZAN music network
 - 20NEWS text document network
 - 3D mesh points
- Optimal graph construction : Unfortunately, no theoretical approach is available – it is empirical and depends mostly on domain expertise knowledge and good common practice (later discussed).
- What is the computational time needed to construct a graph from data features ?
 - Exact construction : $O(n^2 \cdot d)$, n = num data, d = num features.

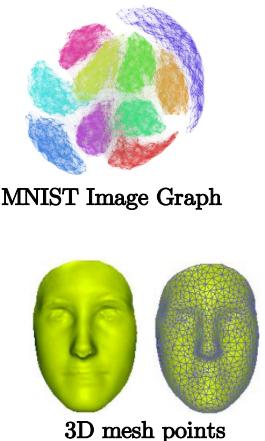
for $d = 1K$, $n = 1K$ \Rightarrow time < 1sec
 $n = 100K$ \Rightarrow time = 1 min
 $n = 1M$ \Rightarrow time > 1 hour
 - Approximate technique : kd-trees $O(n \log n \cdot d)$ with e.g. FLANN^[1], a library for fast approximate nearest neighbor search in high-dimensional spaces.



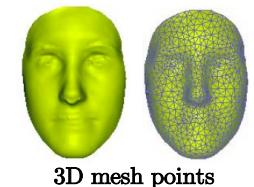
Graph of Text Documents
20newsgroups



GTZAN Music Graph



MNIST Image Graph

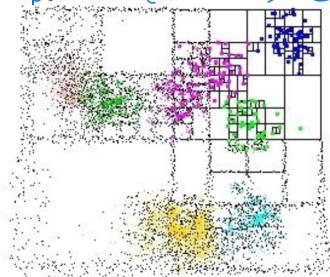


3D mesh points

edge features b/c. If node features, should be $O(n^2 + dn)$ Θ

n^2 : Going through adjacency matrix A of size n^2 .

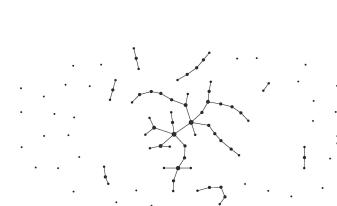
d : Adding of d features per edge



Graph categories

- Mathematical/simulated graphs

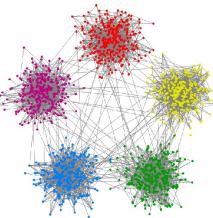
- (graphs generated by some random process)*
- Erdos-Renyi graphs^[1]
 - Stochastic block models (SBM)^[2]
 - Lancichinetti-Fortunato-Radicchi (LFR) graphs^[3]



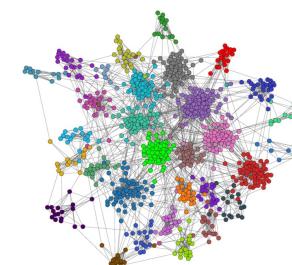
Erdos-Renyi Network
Source: Wikipedia.



Paul Erdős
1913 – 1996



SBM
Source: Abbe, JMLR'17



LFR
Source: Kojaku, 2018

- Why using artificial networks?

- Mathematical Modeling

- Advantage : Precise control of your data model (estimate best performance given some data assumptions). No need to perform extensive experiments. *can find out some graph properties given it follows certain assumptions*
 - Limitation : Most data assumptions are often too restrictive, and it is not guaranteed that real-world data follow the given model assumptions.

[1] Erdos, Renyi, On random graph, 1959

[2] Anderson, Wasserman, Faust, Building stochastic blockmodels, 1992

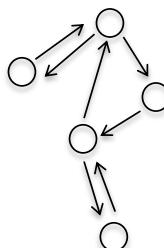
[3] Lancichinetti, Fortunato, Filippo, Benchmark graphs for testing community detection algorithms, 2008

Outline

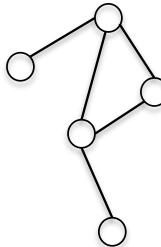
- Graph theory
- Graph categories
- **Basic definitions**
- Curse of dimensionality and structure
- Manifolds and graphs
- Spectral graph theory
- Graph construction
- Conclusion

Basic definitions

- Graphs are defined as $G = (V, E, A)$ where
 - V is the set of vertices (or nodes) w/ $|V| = n$.
 - E is the set of edges.
 - A is the adjacency similarity matrix.
- Directed or undirected graphs :



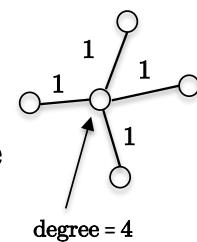
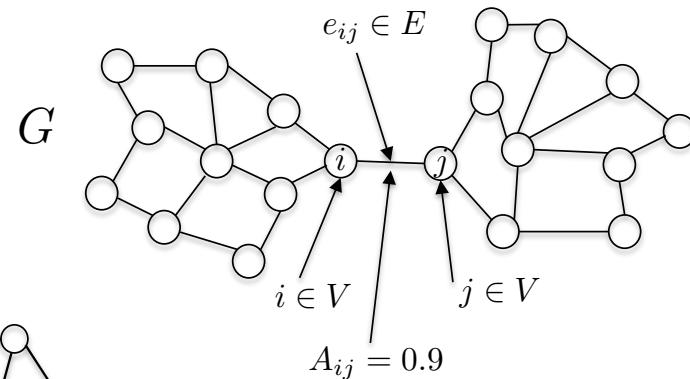
Directed graph



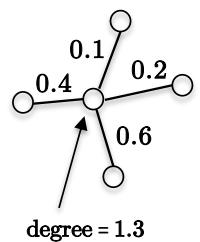
Undirected graph

- Node degree :
 - For binary graphs, $A_{ij} \in \{0,1\} \Rightarrow$ degree = num of edges connected to a node
 - For weighted graphs, A_{ij} in $[0,1]$ \Rightarrow degree is defined as $d_i = \sum_{j \in V} A_{ij}$

$$\text{Degree in both binary and weighted graphs} = d_i = \sum_{j \in V} A_{ij}$$



Binary graph

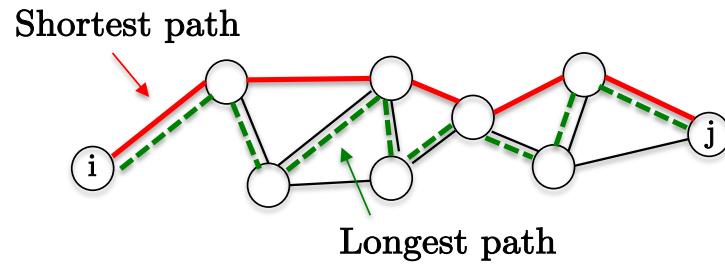


Weighted graph

Algorithms

- Standard graph algorithms (a.k.a. software engineering, no learning) :
 - Breadth-first search, depth-first search, minimum spanning tree, topological sorting, strongly connected components, graph colouring, maximum flow/minimum cut, graph matching, etc.
- Shortest path algorithm : Find a path on a graph with the smallest possible length.
 - Fast Algorithm : Dijkstra's algorithm^[1], *A* algorithm*
 - Popular application is the road navigator product, e.g. from New York to Los Angeles.

minimum # of edges s.t.
all nodes are connected by this
set of edges

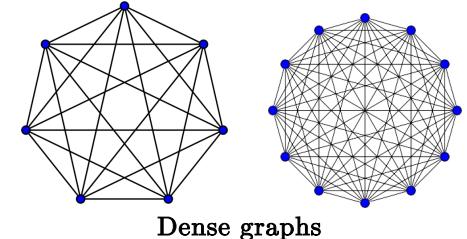


[1] Dijkstra, A note on two problems in connexion with graphs, 1959

Dense vs. sparse graphs

- Dense/complete/full graphs : Each vertex is connected to all other vertices.

$$|E| = \frac{n(n - 1)}{2} = O(n^2)$$

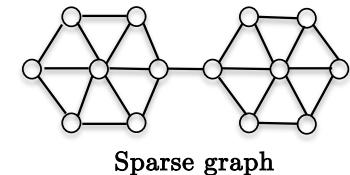



- Sparse graphs : Each vertex is connected to a few other $k \ll n$ vertices.

$$|E| = O(kn) = O(n)$$

$k = E(\deg(v)) = \text{constant};$





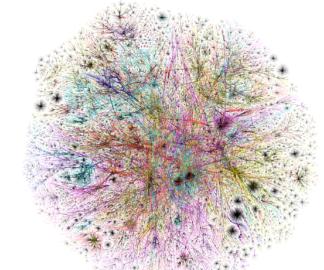
- Which graph is better -- dense or sparse?

- Sparse networks are highly desirable for memory and computational efficiency.

For instance, Internet has $n = 4.73$ billion pages (as of August 2016)

- $|E| = n^2 = 10^{18}$ if Internet was full.
- $|E| = k.n = 10^{11}$ as it is sparse with $k \approx 100$ (mean degree).

*Each page references about
100 other pages*



Internet

- Actually, most real-world networks (e.g. social, communication, etc) are sparse !

- Because sparsity induces structure (a fully connected graph has no structure).

*good
the relationships between nodes
gives no information / value*

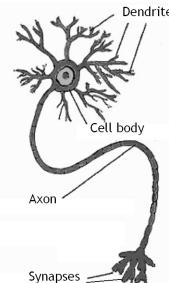
In-lecture question

- In the human brain, some 86 billion neurons form 100 trillion axon connections to each other. Can we model the human brain as a sparse graph? What is the percentage of non-zero elements? 

• $n = ? \quad 86 \times 10^9 = 8.6 \times 10^{10}$

• $|E| = ? \quad 100 \times 10^{12} = 10^{14}$

$$\text{mean degree} = \frac{10^{14}}{8.6 \times 10^{10}} = 86000 \leftarrow \text{sparse in comparison to } 8.6 \times 10^9 \text{ nodes}$$



One neuron

Inputs are dendrites
Outputs are synapses
Connections are axons

% of Non-zero elements in the adjacency matrix = $\frac{86000}{8.6 \times 10^9} \times 100\% = 0.001\% = 10^{-3}\%$

$$= \left[\frac{8.6 \times 10^9}{8.6 \times 10^9} \times 100 \right] \% = 10^{-3}\%$$

- In Slack #lectures

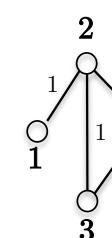
- Identify the question and Reply in thread with a short response

Adjacency matrix

- Definition : Matrix A in $G = (V, E, A)$ represents structural/topological information about the network. We have two classes of A :

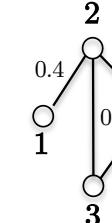
- Binary matrix : $A_{ij} \in \{0,1\}$

$$A_{ij} = \begin{cases} 1 & \text{if } (i, j) \in E \\ 0 & \text{otherwise} \end{cases}$$

$G =$  $\Rightarrow A = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 1 & 0 & 1 & 0 & 0 \\ 2 & 1 & 0 & 1 & 1 \\ 3 & 0 & 1 & 0 & 1 \\ 4 & 0 & 1 & 1 & 0 \end{bmatrix}$

- Weighted matrix : $A_{ij} \in [0,1]$ (commonly normalized to 1)

$$A_{ij} = \begin{cases} \in [0, 1] & \text{if } (i, j) \in E \\ 0 & \text{otherwise} \end{cases}$$

$G =$  $\Rightarrow A = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 1 & 0 & 0.4 & 0 & 0 \\ 2 & 0.4 & 0 & 0.7 & 0.3 \\ 3 & 0 & 0.7 & 0 & 1 \\ 4 & 0 & 0.3 & 1 & 0 \end{bmatrix}$

How is this normalization
done?

not normalized

9

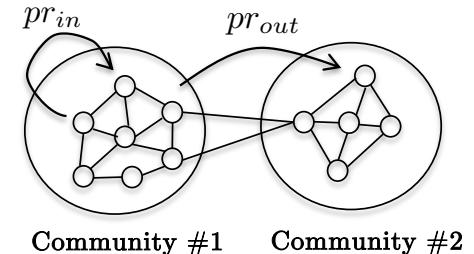
Lab 1 : LFR social networks

(27 Aug 2015)

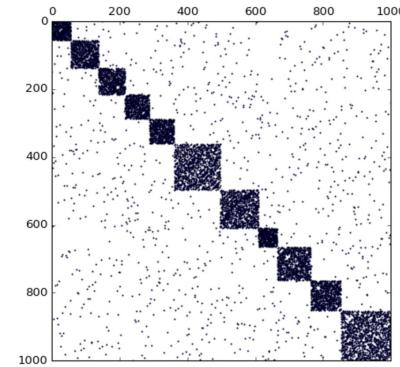
- Run code01.ipynb and synthesize LFR social networks.
 - Play with the **mixing parameter μ** :
 - μ small : Communities are well separated.
 - μ large : Communities are mixed.

$$\mu = \frac{pr_{out}}{pr_{in}}$$

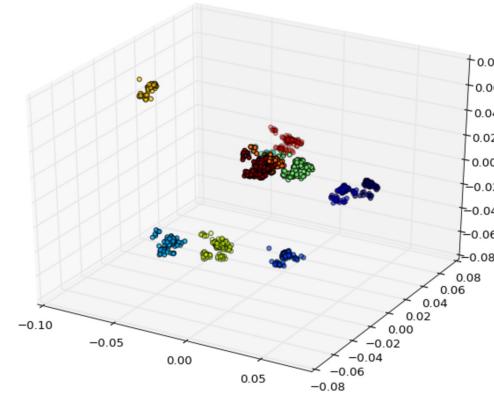
pr: probability.



```
In [15]: # Plot same W but according to communities
# Any structure?
plt.figure(2)
plt.spy(W,precision=0.01, markersize=1)
plt.show()
```



```
In [19]: # Visualize the social network in 3D
fig = pylab.figure(4)
ax = Axes3D(fig)
ax.scatter(X, Y, Z, c=C)
pyplot.show()
```



Outline

- Graph theory
- Graph categories
- Basic definitions
- **Curse of dimensionality and structure**
- Manifolds and graphs
- Spectral graph theory
- Graph construction
- Conclusion

$$d_{\max}^{\ell_2}(x_i, V \setminus x_i) = \text{max L2 dist of } x_i \text{ to another node}$$

$$d_{\min}^{\ell_2}(x_i, V \setminus x_i) = \text{min L2 dist of } x_i \text{ to another node}$$

Curse of dimensionality

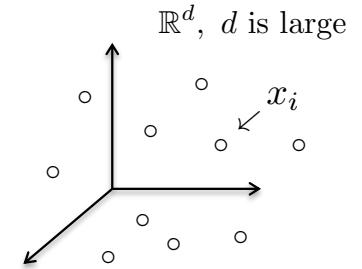
- What is the curse of dimensionality ?

- In high dimensions, Euclidean distance between data becomes meaningless.
- Theorem^[1] : Suppose data are uniformly distributed in \mathbb{R}^d ,

pick any data $x_i \in V$, we have :

$$\lim_{d \rightarrow \infty} \mathbb{E}_{x_i} \left(\frac{d_{\max}^{\ell_2}(x_i, V \setminus x_i) - d_{\min}^{\ell_2}(x_i, V \setminus x_i)}{d_{\min}^{\ell_2}(x_i, V \setminus x_i)} \right) = 0 \quad \text{Proof } \textcircled{?}$$

dimensions d
 $\rightarrow \infty$

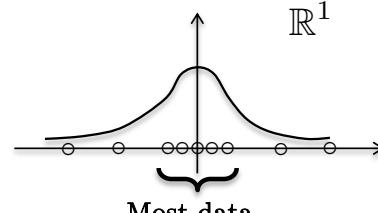


Interpretation : All data are far away to each other with the same distance value !

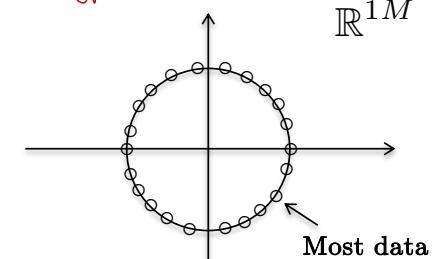
- Besides, loss of intuition in high dimensions, e.g. with the Normal distribution :

- In low-dim, most data are concentrated at the center.
- In high-dim, most data are concentrated on the surface.

surface of what $\textcircled{?}$ *the Gaussian* $\textcircled{?} \sqrt{d}$



1-D Gaussian

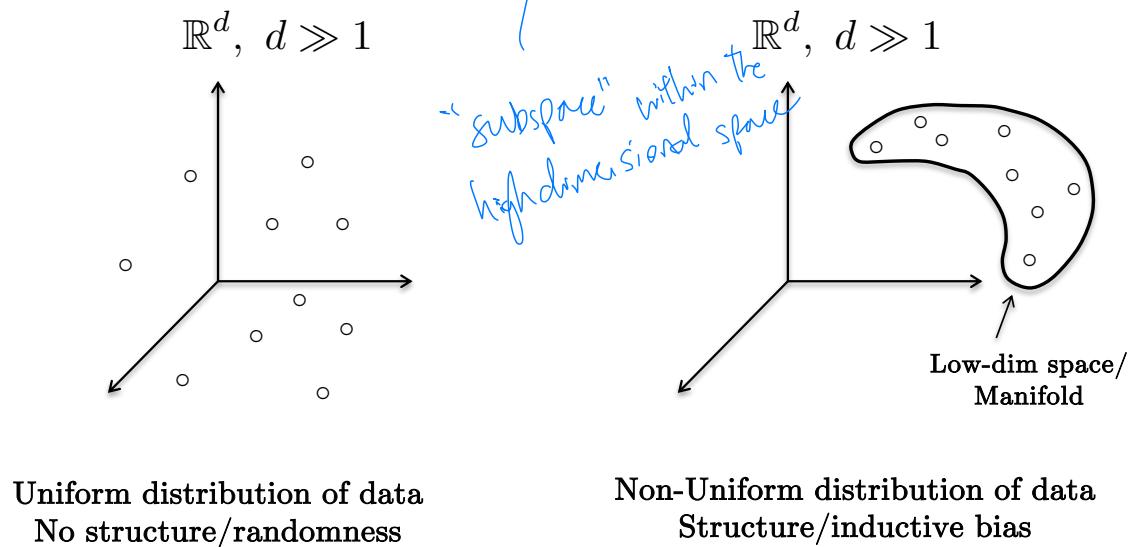


1,000,000-D Gaussian

[1] Beyer, Goldstein, Ramakrishnan, Shaft, When is “nearest neighbor” meaningful? 1999

Blessing of structure

- What is the blessing of structure ?
 - Previously, the assumption “data are uniformly” distributed is actually not true for real-world data. Data have always properties, i.e. structures or invariances, such that they belong to a low-dimensional space called **manifold** where (geodesic) distances are meaningful.



① underlying data distribution ?

“subspace” within the high-dimensional space

Low-dim space/
Manifold

Convolutional layers help to achieve shift invariance (you already know this)
↓
Act in different parts of different images

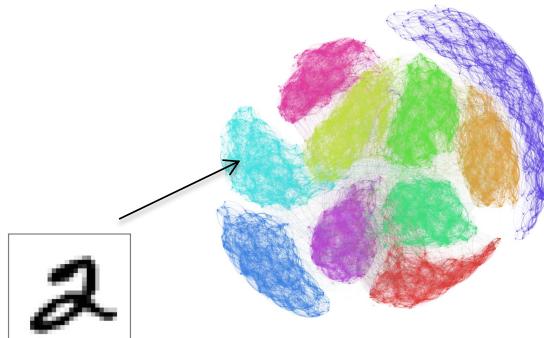
Outline

- Graph theory
- Graph categories
- Basic definitions
- Curse of dimensionality and structure
- **Manifolds and graphs**
- Spectral graph theory
- Graph construction
- Conclusion

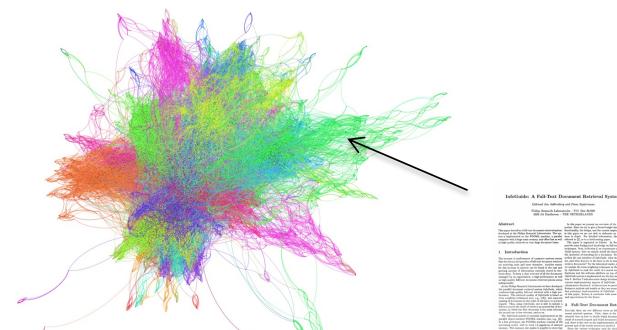
Manifold learning

- It can be challenging to identify structures hidden in data because of
 - The curse of dimensionality (i.e. high-dimensional data).
 - Some data have easy structures, but most have complex ones.
- A class of algorithms that extracts low-dimensional patterns is manifold learning (later discussed).

finding the right combination
of features that map to classes
(combination of features that map to classes)



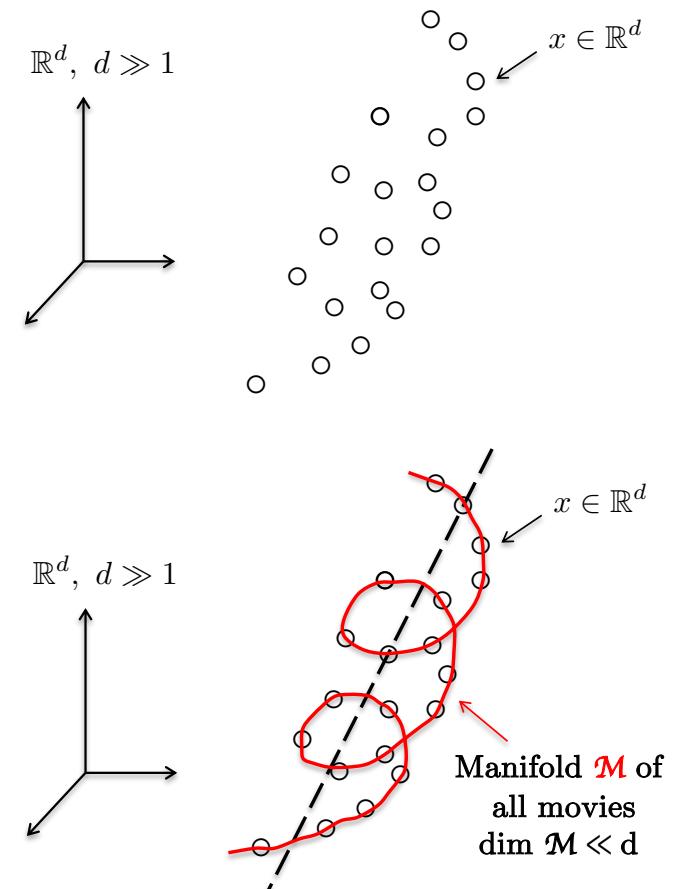
MNIST Image Graph



Graph of Text Documents
20newsgroups

From manifolds to graphs

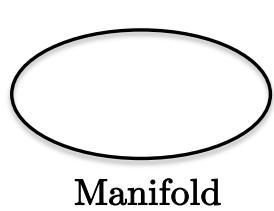
- Manifold assumption : High-dimensional data are sampled from a low-dimensional manifold.
 - Example : Let x be a movie, each movie is defined by d features/attributes like genre, actors, release year, origin country, etc such that $x \in \mathbb{R}^d$. We can make the assumption that all movies form a manifold \mathcal{M} in \mathbb{R}^d .
- Assumption validity : The manifold is a good working hypothesis for
 - Several types of data, including images, text documents, music, etc.
 - Most machine learning tasks e.g. classification, visualization, recommendation.
 - However, it can also be limited as it can be a crude approximation of the (true) data distribution.



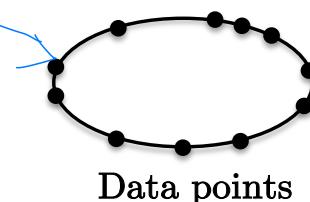
From manifolds to graphs

- Graphs can be regarded as a manifold sampling process.

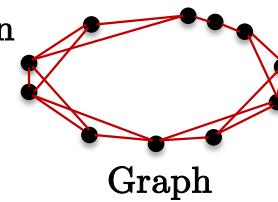
- The manifold information is represented by a neighborhood graph (observe that the manifold is never directly observed or constructed).



Sampling



Graph construction



$$G = (V, E, A)$$

- Neighborhood graphs :

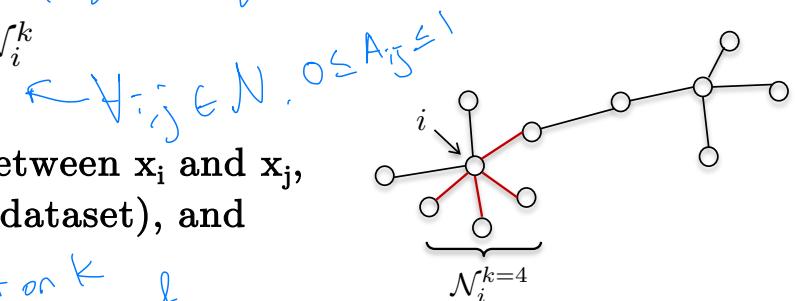
- Most popular are k-NN graphs defined as :

$$A_{ij} = \begin{cases} e^{-\frac{\text{dist}(x_i, x_j)^2}{\sigma^2}} & \text{if } j \in \mathcal{N}_i^k \\ 0 & \text{otherwise} \end{cases}$$

where $\text{dist}(x_i, x_j)$ is a distance (to be decided) between x_i and x_j , σ is the scale parameter (value depends on the dataset), and \mathcal{N}_i^k is the neighborhood of data x_i .

hyperparameter

Note: constraint on K
ensures the construction of
a sparse graph



Outline

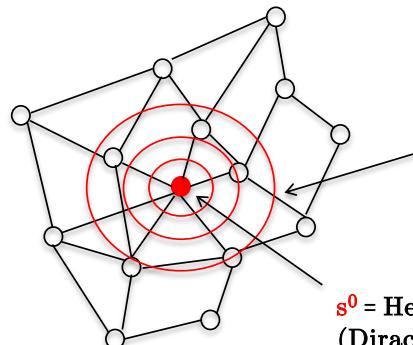
- Graph theory
- Graph categories
- Basic definitions
- Curse of dimensionality and structure
- Manifolds and graphs
- **Spectral graph theory**
- Graph construction
- Conclusion

Spectral graph theory

KIV

Global Patterns

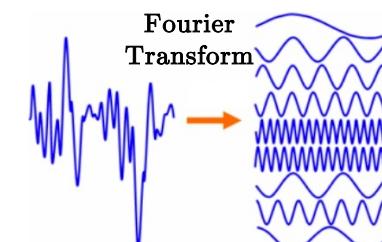
- Given a graph $G = (V, E, A)$, spectral graph theory (SGT) can
 - Find meaningful patterns that reveal multi-scale graph structures. *What is multi-scale structure?*
 - Process data defined on the graph domain (a.k.a. graph signal processing).
 - Boost performance of learning tasks s.a. clustering, classification, recommendation, etc.
- The most fundamental tool in SGT is the graph Laplacian operator L .
 - Why is the Laplacian useful?
 - It is the central operator of diffusion processes (on graphs).
 - Basis functions of this operator are the well-known Fourier modes (later discussed).



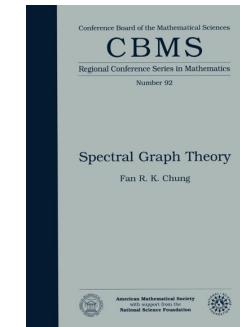
Heat propagation using
Laplacian operator :
 $s^{t+1} = L s^t$

*vector containing
the heat/temperature
of all nodes at
time t .*

s^0 = Heat source
(Dirac function)



← "eigen vector is the sine & cosine functions" ?



Fan Chung
SGT book
1997



Pierre-Simon Laplace
(1749–1827)

continuous Laplacian has a unique def'n
graph Laplacian does not

② makes the normalized Laplacian
"invariant" to the manifold
positive semi-definite

Graph Laplacian

- Un-normalized/combinatorial graph Laplacian :

$$L_{un} = D - A \quad \text{with } D \text{ is the degree matrix : } D = \text{diag}(d_1, \dots, d_n),$$

$n \times n$
 $n = |V|$

$$L_{un} : l_{ij} = \begin{cases} -A_{ij} & \text{if } i \neq j \\ \deg(V_i) & \text{otherwise} \end{cases}$$

Eigenvalues of L_{un} are real & non-negative

- Normalized Laplacian (most popular) :

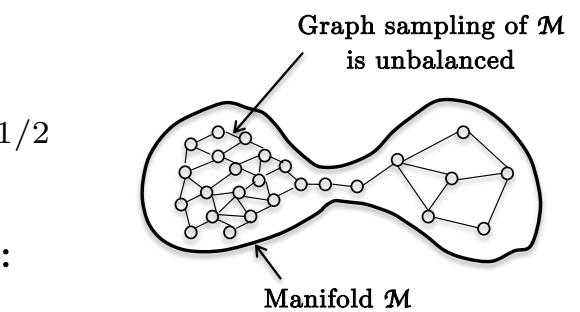
- Robust w.r.t. unbalanced sampling

$$L = D^{-1/2} L_{un} D^{-1/2} = I_n - D^{-1/2} A D^{-1/2}$$

- Random Walk Laplacian (for Google PageRank, later discussed) :

$$L = D^{-1} L_{un} = I_n - D^{-1} A$$

- All Laplacians are diffusion operators. *think*



Graph spectrum

- Spectral graph theory can extract the modes of variation of the graph system.

- How? With the eigenvalue decomposition (EVD) of the Laplacian operator L :

$$L = U \Lambda U^T \in \mathbb{R}^{n \times n}$$

with $U = [u_1, \dots, u_n] \in \mathbb{R}^{n \times n}$,
 orthonormal vectors u_1, \dots, u_n

$$U^T U = I_n, \text{ i.e. } \langle u_k, u_{k'} \rangle = \begin{cases} 1 & k = k' \\ 0 & \text{otherwise} \end{cases},$$

$$\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n) \in \mathbb{R}^{n \times n},$$

$$0 = \lambda_{\min} = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n \leq 2$$

- Interpretation:

- u_k : Laplacian eigenvectors a.k.a. Fourier functions, i.e. vibration modes of the graph.
- λ_k : Laplacian eigenvalues a.k.a. frequencies of the Fourier functions, i.e. how fast u_k vibrate.
- EVD answers the famous question 'Can One Hear the Shape of a Drum'?^[1]

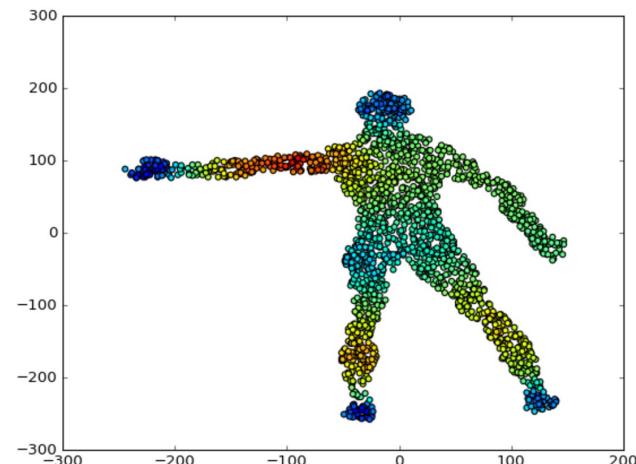
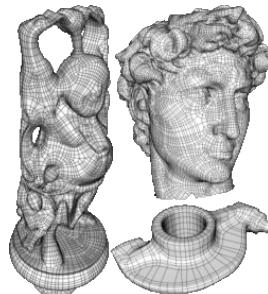


[1] Kac, Can One Hear the Shape of a Drum, 1966

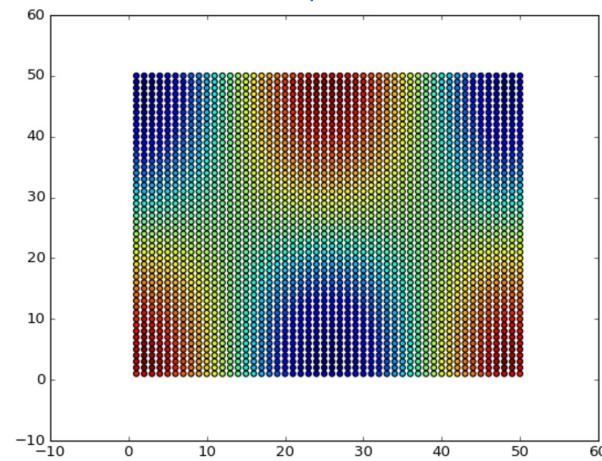
Lab 2 : Spectrum of point cloud and grid

- Run code02.ipynb and visualize the Fourier functions of a human body graph and a regular grid.
- What is the main property of the smallest and largest eigenvectors?
 - Smallest eigenvectors \Rightarrow Smoothest modes of vibration, i.e. low-frequency information.
 - Largest eigenvectors \Rightarrow Highest frequencies of the graph, i.e. details or noise.

```
In [5]:  
# Compute graph Laplacian  
L = graph_laplacian(W)  
  
# Compute modes of variations of graph system = Fourier functions  
lamb, U = scipy.sparse.linalg.eigh(L, k=9, which='SM')
```



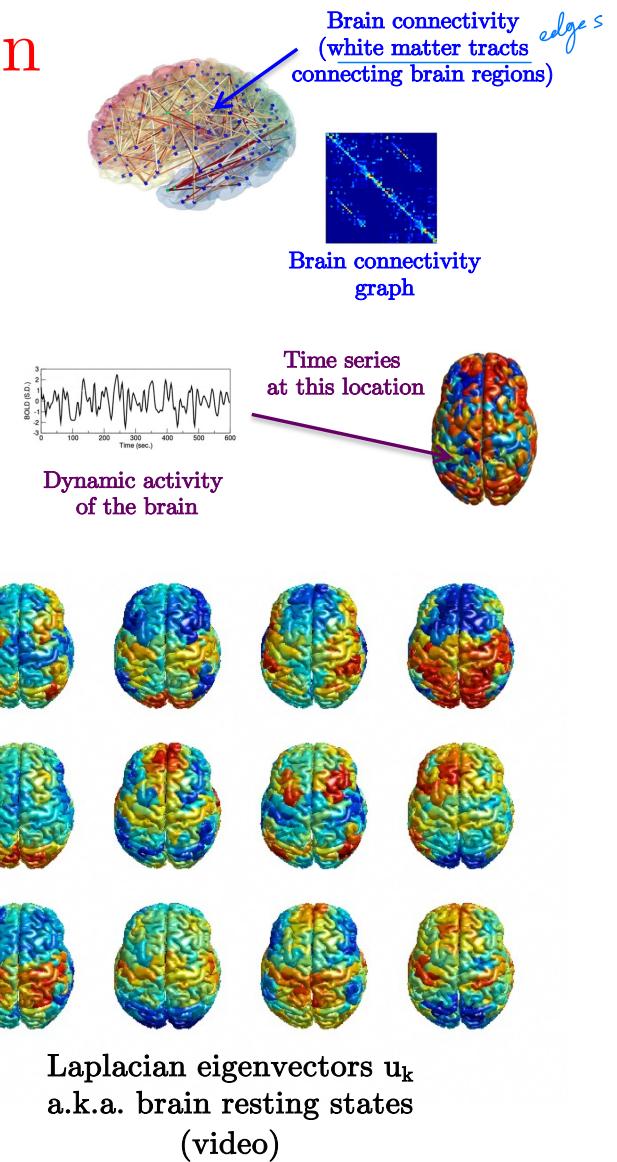
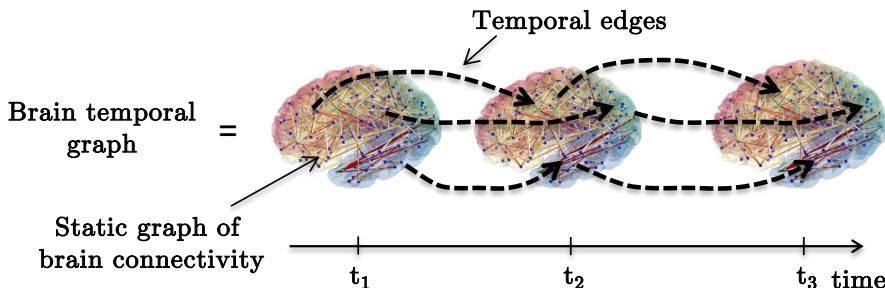
Fourier mode from a Graph =
Set of points or meshes in graphics
(e.g. 2D/3D shape recognition)



Fourier mode from a Graph =
Regular grid
(e.g. JPEG image compression,
most used technique)

Spectrum of human brain

- Goal : Find brain activation patterns from structural MRI, a.k.a. brain resting states (brain structure implies brain function).
- Methodology : Given G (brain connectivity graph) \Rightarrow design a temporal graph \Rightarrow compute Laplacian L \Rightarrow compute eigenvectors u_k \Rightarrow visualize brain temporal activation patterns.
- Result : u_k represent the dynamic patterns related to basic functional brain tasks s.a. vision, body motor, language, etc.
- How to construct the brain temporal graph?



Spectral Graph Theory

$$L = D - A$$

$$D = \text{diag}(deg(v_1), deg(v_2), \dots, deg(v_N))$$

A = Binary $(N \times N)$ matrix

$$x^T L x = \sum_{i=1}^N \sum_{j=1}^N L_{ij} x_i x_j$$

$$L_{ij} x_i x_j = \begin{cases} \text{diag}(v_i) x_i^2 & \text{if } i=j \\ -E(v_i, v_j) x_i x_j & \text{otherwise} \end{cases}$$

$$x^T L x = \sum_{i=1}^N \text{diag}(v_i) x_i^2 - \sum_{i=1}^N \sum_{j \neq i} E(v_i, v_j) x_i x_j$$

$$(x_1 \ x_2 \ x_3) \begin{pmatrix} L_{11} & L_{12} & L_{13} \\ L_{21} & L_{22} & L_{23} \\ L_{31} & L_{32} & L_{33} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = (x_1 \ x_2 \ x_3) \begin{pmatrix} L_{11} x_1 + L_{12} x_2 + L_{13} x_3 \\ L_{21} x_1 + L_{22} x_2 + L_{23} x_3 \\ L_{31} x_1 + L_{32} x_2 + L_{33} x_3 \end{pmatrix}$$

$$\begin{aligned} &= L_{11} x_1^2 + L_{12} x_1 x_2 + L_{13} x_1 x_3 \\ &\quad + L_{21} x_1 x_2 + L_{22} x_2^2 + L_{23} x_2 x_3 \\ &\quad + L_{31} x_1 x_3 + L_{32} x_2 x_3 + L_{33} x_3^2 \\ &= \sum_{i=1}^N \sum_{j=1}^N L_{ij} x_i x_j \end{aligned}$$

Outline

- Graph theory
- Graph categories
- Basic definitions
- Curse of dimensionality and structure
- Manifolds and graphs
- Spectral graph theory
- **Graph construction**
- Conclusion

Qn: Is it possible to reconstruct a graph
solely using data about each node's
degree?

Ans: No.

i	$\deg(v_i)$
1	2
2	2
3	2
4	1
5	1

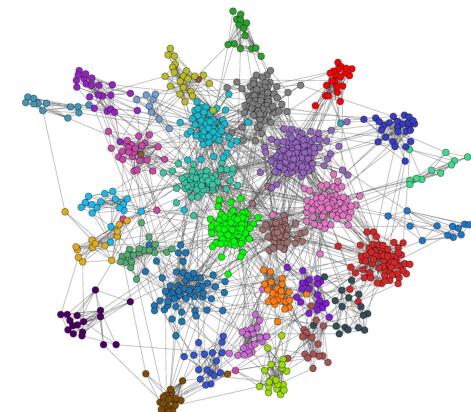
This can be two
different graphs



This is why we need either set E or matrix A

How to construct graphs from data?

- Three basic questions
 - Which type of graphs?
 - Which data distance?
 - Which data feature?
- Optimal graph construction
 - No universal recipe is available (no theory).
 - It depends on data and analysis task (empirical).
 - Domain expertise and good practice are essential.



Type of constructed graphs

complexity to construct : $O(kN)$

- Neighborhood graphs
 - k-NN graphs, i.e. sparse graphs by design.
 - Parameters
 - k : number of nearest neighbors, a common value is between {5,50}.
 - σ : positive scale parameter
 - Two options to compute scale σ
 - Global scale : $\sigma = \text{mean distance of all } k^{\text{th}} \text{ neighbors.}$
 - Local scale^[1] : $\sigma_i = \text{distance of the } k^{\text{th}} \text{ neighbor for node } i.$ (most common)

$$A_{ij} = \begin{cases} e^{-\frac{\text{dist}(x_i, x_j)^2}{\sigma^2}} & \text{if } j \in \mathcal{N}_i^k \\ 0 & \text{otherwise} \end{cases}$$

Adjacency matrix
with global scale

$$A_{ij} = \begin{cases} e^{-\frac{\text{dist}(x_i, x_j)^2}{\sigma_i \sigma_j}} & \text{if } j \in \mathcal{N}_i^k \\ 0 & \text{otherwise} \end{cases}$$

Adjacency matrix
with local scale

[1] Zelnik-Manor, Perona, Self-tuning spectral clustering, 2004

In-lecture question

- Consider an ϵ -NN graph, i.e. a graph with the following weights :

$$A_{ij}^\epsilon = \begin{cases} e^{-\frac{\text{dist}(x_i, x_j)^2}{2\epsilon^2}} & \text{if } \text{dist}(x_i, x_j) \leq \epsilon \\ 0 & \text{otherwise} \end{cases}$$

any have this here ?

- The k-NN graph's adjacency matrix is defined as :

$$A_{ij}^{\text{kNN}} = \begin{cases} e^{-\frac{\text{dist}(x_i, x_j)^2}{\sigma^2}} & \text{if } j \in \mathcal{N}_i^k \\ 0 & \text{otherwise} \end{cases}$$

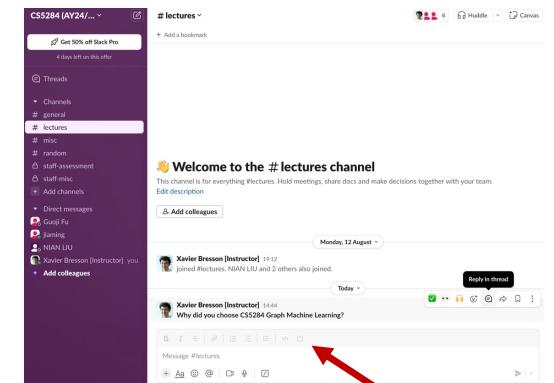
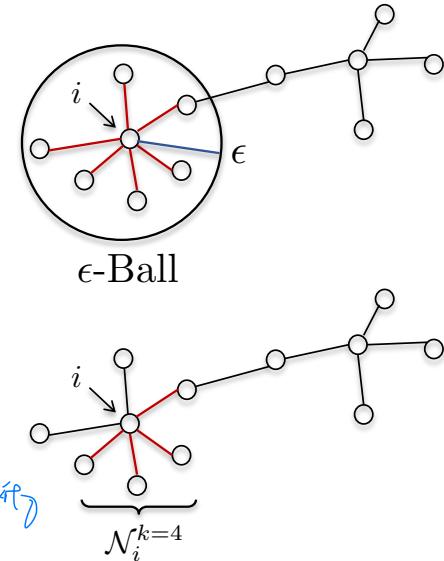
- Which graph is the most used in practice? Justify. *k -NN. Lower time complexity*

- In Slack #lectures

- Identify the question and Reply in thread with a short response

Euclidean
 In high dimensions, distance between all pairs of datapoints becomes very similar \rightarrow many datapoints within each ϵ -Ball \rightarrow graph becomes very dense \rightarrow less structural info \rightarrow Bad

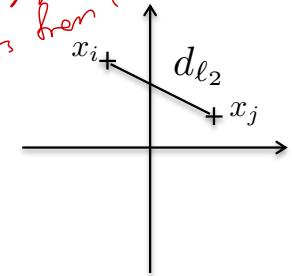
Xavier Bresson



Distance definition

- Euclidean/L2 distance :
 - Good distance for low-dim data, e.g. $d < 10$.
 - Good distance for high-dim data **with linearly separable data** (e.g. MNIST).

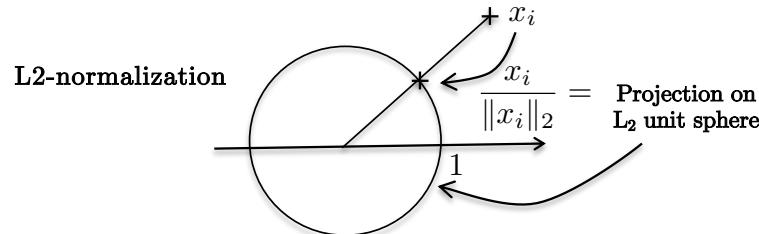
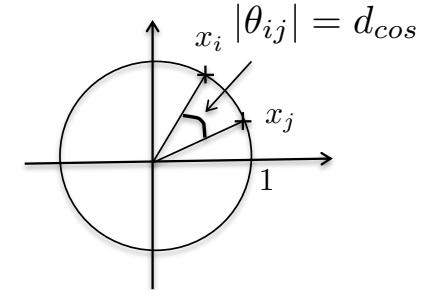
*In example in later slide,
if we use L2 distance to
construct a graph, we might get many
edges between datapoints from the 2 different
nodes*



$$d_{\ell_2}(x_i, x_j) = \|x_i - x_j\|_2 = \sqrt{\sum_{m=1}^d |x_{i,m} - x_{j,m}|^2}$$

- Cosine/dot product distance :
 - Good distance for high-dim sparse data (e.g. text documents)

$$d_{cos}(x_i, x_j) = \left| \cos^{-1} \left(\frac{\langle x_i, x_j \rangle}{\|x_i\|_2 \|x_j\|_2} \right) \right| = |\theta_{ij}|$$

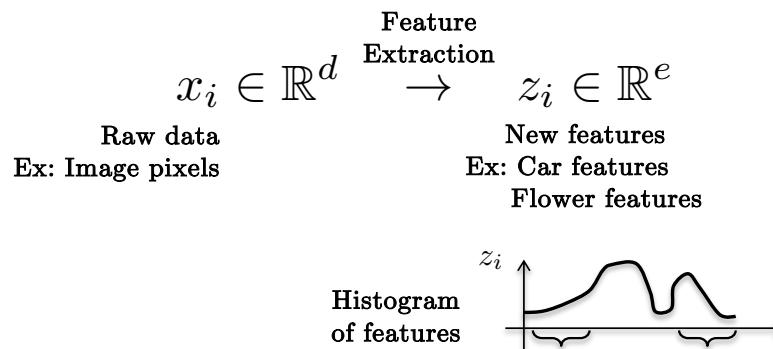


$$X = \begin{matrix} n \times d \\ \left[\begin{array}{c} \frac{x_1}{\|x_i\|_2} \\ \frac{x_2}{\|x_i\|_2} \\ \vdots \\ \frac{x_n}{\|x_i\|_2} \end{array} \right] \end{matrix} \quad \|x_i\|_2 = 1$$

- Other analytical distances : Kullback-Leibler distance (information theory), Wasserstein distance (optimal transport), etc.

Data features

- Types of data features
 - Raw features (e.g. movie features such as genre, actors, year, etc)
 - Hand-crafted features (e.g. SIFT in computer vision, etc)
 - Learned features (PCA, NMF, sparse coding, deep learning, etc)
 - should transform the raw data first*
- It is generally unsuccessful to directly use the raw features for graph construction.
 - Issues are noise, unbalanced scaling, lack of expressiveness, curse of dimensionality, etc
- For successful graph construction, raw data should be transformed into meaningful data representation by designing new features, s.a. handcrafted or learned features.



$$A_{ij} = e^{-\frac{\text{dist}(x_i, x_j)^2}{\sigma^2}}$$

\Downarrow

$A_{ij} = e^{-\frac{\text{dist}(z_i, z_j)^2}{\sigma^2}}$
transformed feature

should transform the raw data first

Standard pre-processing

- Center data (along feature dimension) : zero-mean property

$$x_i \leftarrow x_i - \text{mean}(\{x_i\}) \in \mathbb{R}^d$$

- Normalize data variance (along feature dimension) : z-scoring property

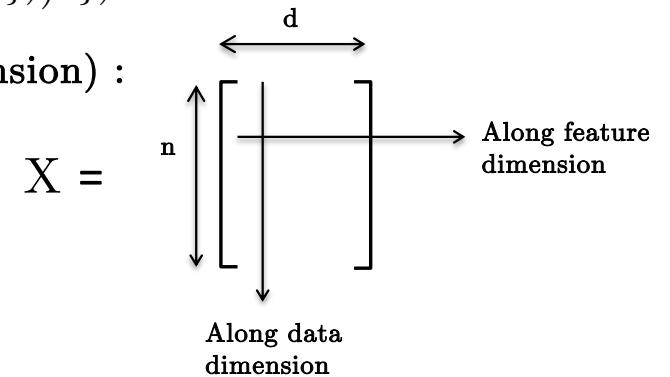
$$x_i \leftarrow x_i / \text{std}(\{x_i\}) \in \mathbb{R}^d$$

$$\text{with } \text{std}(\{x_i\})^2 = \text{mean}(\{(x_j - \text{mean}(\{x_i\}))^2\})$$

- Project data on L2-sphere (along feature dimension or data dimension) :



$$x_i \leftarrow x_i / \|x_i\|_2 \in \mathbb{R}^d$$

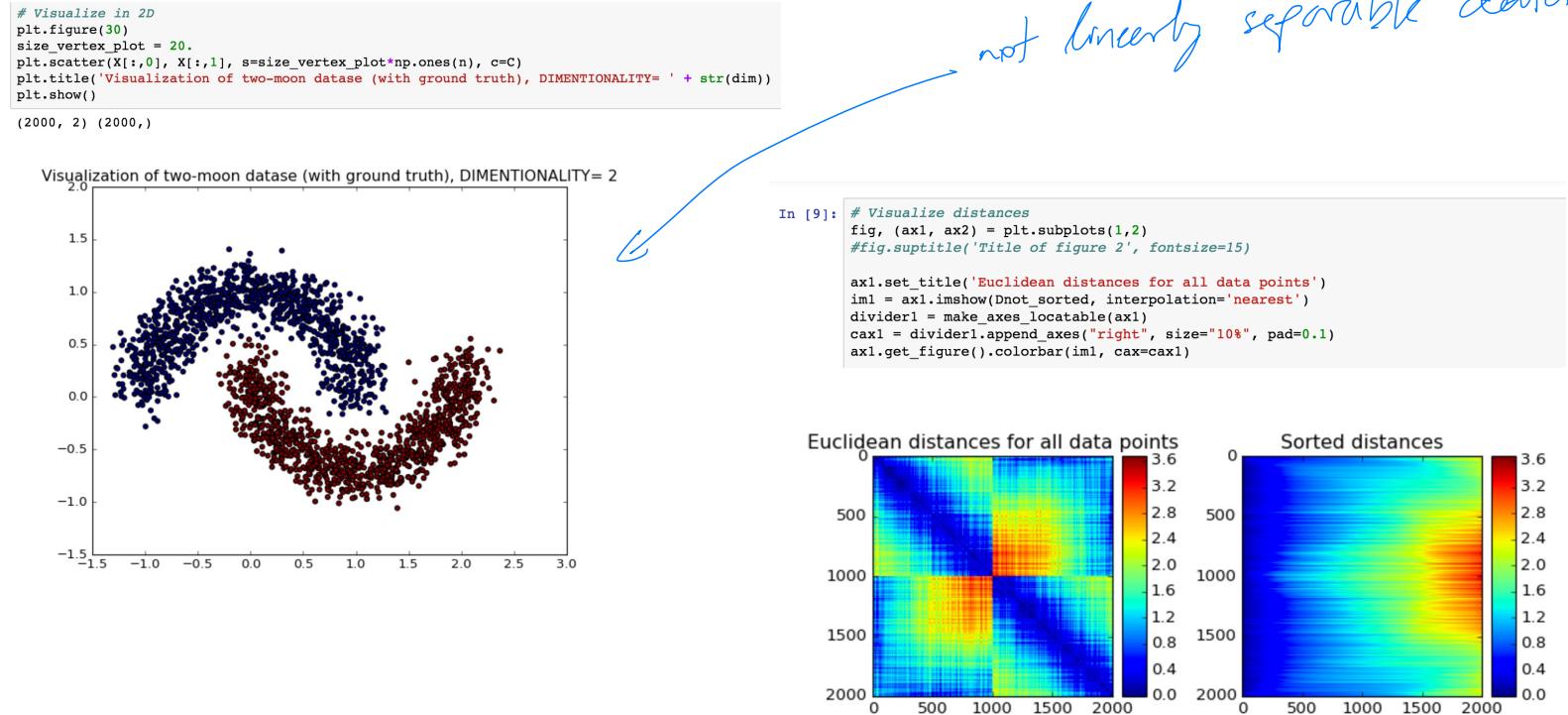


- Normalize max and min of feature value :

$$x_i \leftarrow \frac{x_i - \text{min}(\{x_i\})}{\text{max}(\{x_i\}) - \text{min}(\{x_i\})} \in [0, 1]^d$$

Lab 3 : Graph construction for two-moon

- Run code03.ipynb and study data pre-processing, construction of k-NN graphs, visualization of distances, adjacency matrix, and graph quality with clustering accuracy.

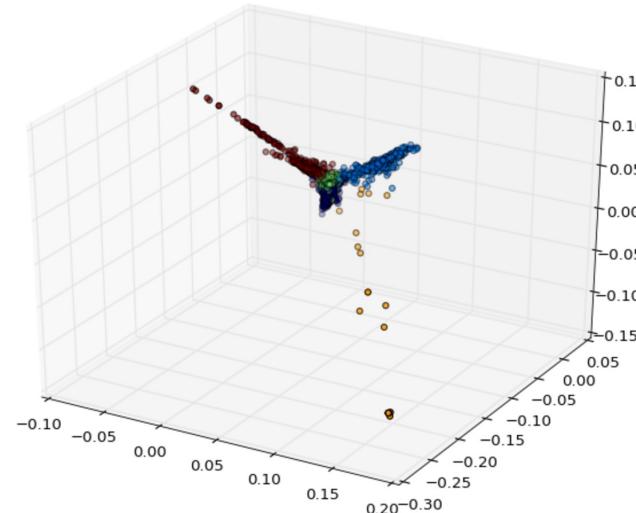
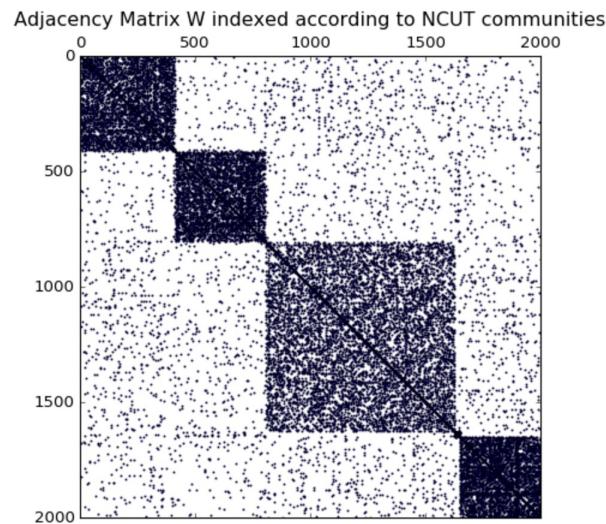


Lab 4 : Graph construction for text documents

- Run code04.ipynb and construct graph text documents.

```
In [9]: # Compute the k-NN graph with Cosine distance  
W_cosine = construct_knn_graph(X, 10, 'cosine')
```

k-NN graph with cosine distance



Outline

- Graph theory
- Graph categories
- Basic definitions
- Curse of dimensionality and structure
- Manifolds and graphs
- Spectral graph theory
- Graph construction
- Conclusion

Summary

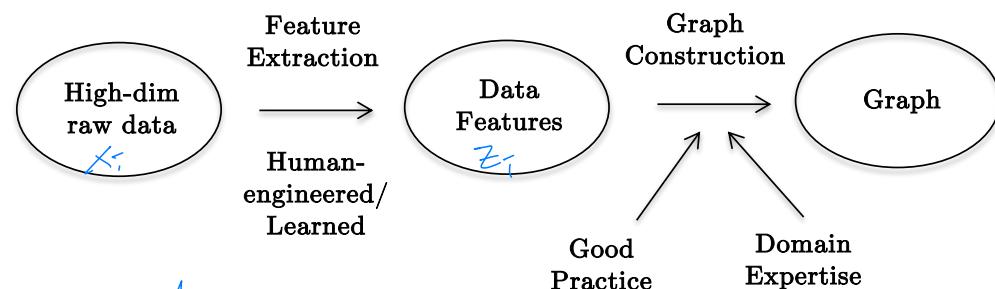
- Graphs can represent complex and heterogenous relationships between data.
 - They help to go beyond the standard assumption that data are i.i.d. (independent and identically distributed) as they explicitly leverage the connections between pairs of data.
 - Any dataset and task that use explicit relations between data is a graph-based task.
 - When we begin looking for graphs, we discover they are ubiquitous! ☺
- Graphs offer an augmented representation of data.
 - With data feature X and data relationship depicted by $G = (V, E, A)$.

Summary

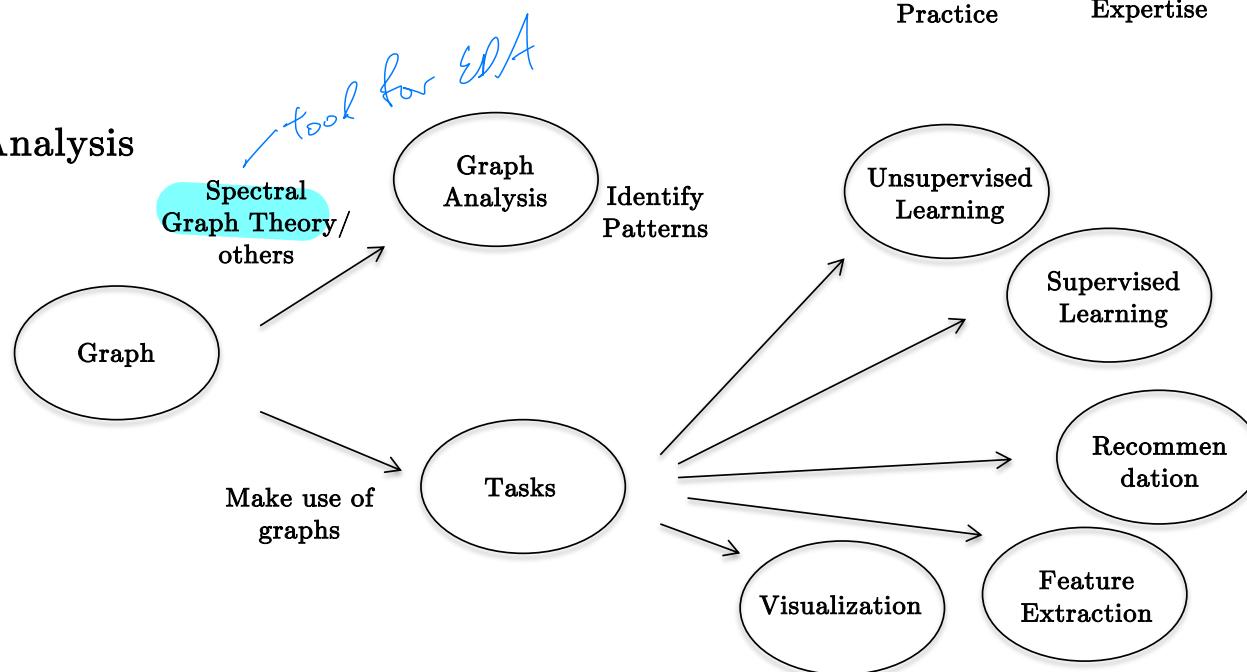
- First fundamental tool of graph science :
 - Adjacency matrix A
 - It reveals global structures in data relationship (graph spectrum).
 - It can visualize graphs in 2D or 3D Euclidean spaces (later discussed).
 - It can boost performance of machine learning techniques (later discussed).
 - Second fundamental tool of graph science :
 - ① • Graph Laplacian matrix L
 - It is a diffusion operator -- It propagates information on graphs (parabolic PDEs).
 - It is used to reveal modes of variation of the graph system.
 - It is used for image compression (jpeg), neuroscience (brain activity), positional encoding (later discussed), etc.
- and also controls the speed of propagation*

Pipeline

- Step #1: Data \Rightarrow Graph
Skip this step if graph is given,
e.g. molecule, social network, etc



- Step #2: Graph \Rightarrow Analysis





Questions?