

Comp683: Computational Biology

Lecture 3

January 18, 2024

Today

- Graph Laplacian Refresher and Nice Properties
- Building Graphs from Data (k-D trees and random projection trees)
- Graph Representations of Data and Dimension Reduction
- Modularity-Based Graph Partitioning

Truncated SVD

Here is how you can construct a rank k approximation for your matrix, $\mathbf{M} \in \mathbb{R}^{n \times m}$. Assuming that $\mathbf{M} = \mathbf{U}\mathbf{D}\mathbf{V}^T$ is the SVD of \mathbf{M} , then you can compute the rank k truncated SVD of \mathbf{M} as,

$$\mathbf{M}_{(k)} = \mathbf{U}_{(k)}\mathbf{D}_{(k)}\mathbf{V}_{(k)}^T \quad (1)$$

- entries of \mathbf{D} are sorted in descending order on the diagonal and $\mathbf{D}_{(k)}$ is the first $k \times k$ submatrix of \mathbf{D}
- $\mathbf{U}_{(k)} \in \mathbb{R}^{n \times k}$ is the first k columns of $\mathbf{U} \in \mathbb{R}^{n \times n}$
- $\mathbf{V}_{(k)} \in \mathbb{R}^{m \times k}$ is the first k columns of $\mathbf{V} \in \mathbb{R}^{m \times m}$

What Can we Do with That?

You can construct rank k approximations of \mathbf{M} for a bunch of k . You can also sequentially compute $\|\mathbf{M} - \mathbf{M}_k\|$. At some point, these two matrices will become quite close. If you don't need that high of a k for that to be the case, then the matrix doesn't have as much unique 'signal'.

An Application of SVD in Something You Use All of the Time

The Moore Penrose Pseudo Inverse can approximate the inverse of a rank k matrix, \mathbf{M} . Suppose $\mathbf{M} = \mathbf{U}\mathbf{D}\mathbf{V}^T$ is the rank- k truncated SVD. Then the pseudo inverse of \mathbf{M} is defined as,

$$\mathbf{M}^{-1} = \mathbf{V}\mathbf{D}^{-1}\mathbf{U}^T \tag{2}$$

A very exciting application of linear algebra → graphs.

Graphs

A graph is a collection of nodes joined by edges. Or generally, a data structure for capturing relational information.

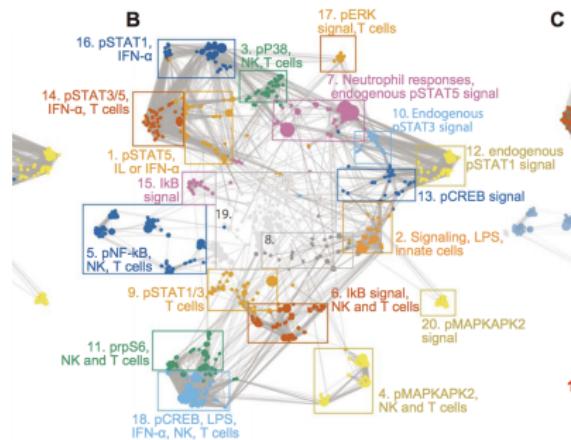


Figure: Figure from Aghapeeour *et al.* Science Immunology. 2018. A graph of function and frequency features immune cell-types (nodes), with edges representing strong correlation across patients.

Mathematical Representations of Graphs

Consider an undirected, unweighted network with N nodes. Then its adjacency matrix (\mathbf{A}) is defined as,

$$A_{ij} = \begin{cases} 1 & \text{there is an edge between nodes } i \text{ and } j \\ 0 & \text{otherwise} \end{cases}$$

In practical situations, you will want to use a sparse representation of your adjacency matrix or define an edgelist also known a $|E| \times 2$ matrix of ‘source’ (column 1) and ‘target’ (column 2) node indices.

Encoding the Strength of Connection Between Nodes

Sometimes we have a notion about how similar pairs of nodes are. For example, correlation, or another similarity measure. This extra information can be helpful in our analysis and interpretation of graph structure. The adjacency matrix of an undirected, weighted graph with N nodes is similarly defined as,

$$A_{ij} = \begin{cases} w_{ij} & \text{there is an edge between nodes } i \text{ and } j \\ 0 & \text{otherwise} \end{cases}$$

Here w_{ij} represents the strength of the connection between nodes i and j .

- The adjacency matrix can be used to quickly gain intuition about properties of the graph, which we will see as a coming attraction.

Example

Consider the adjacency matrix for an undirected graph with $N = 3$ nodes,

$$\mathbf{A} = \begin{bmatrix} 0 & 2 & 1 \\ 2 & 0 & 0.5 \\ 1 & 0.5 & 0 \end{bmatrix}$$

Here, nodes 1 and 2 have twice as strong of a connection as nodes 1 and 3. This graph is also **complete graph** in that there exists an edge connecting every pair of nodes.

The Simplest Node Ranking Strategy: Degree

Let \mathbf{A} be the adjacency matrix for a graph with N nodes. The **degree** of a node of node i (k_i) counts the number of nodes that node i is connected to or,

$$k_i = \sum_{j=1}^N A_{ij} \quad (3)$$

Every edge has two ends. So, for m edges in a graph, there are $2 \times m$ ends of edges. Therefore, the total sum of edges in the graph can be written as,

$$2m = \sum_{i=1}^N k_i \quad (4)$$

Paths in a Graph

- **Path:** A path in a graph is a sequence of nodes such that every consecutive pair of nodes in the sequence is connected by an edge. This is also known as a route across the graph.
- **Path Length:** The path length between nodes i and j is the number of edges that need to be traversed to reach node j from node i .
- **Geodesic Path:** The geodesic path or shortest path is the shortest network distance between a pair of nodes.
 - There are a lot of beautiful shortest path algorithms for doing this efficiently that we will not cover here.

Example

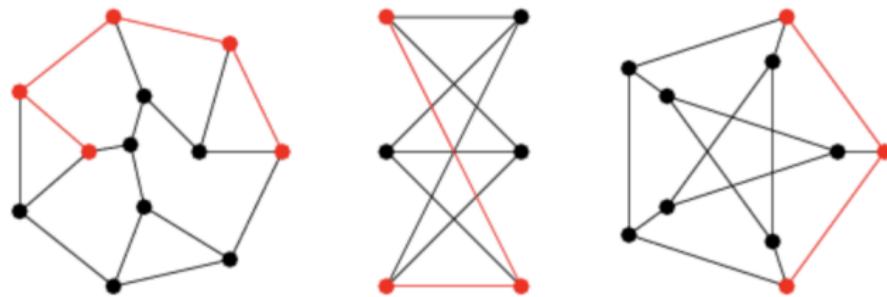


Figure: There could be multiple paths between a pair of nodes, but we are interested in the path that is the shortest.

With linear algebra and graph fundamentals established, we will move to discussion about quantifying structural properties of graphs.

Setup for Graph Diffusion

Diffusion offers the chance to study how some quantity (e.g. signal) relates to connectivity in the graph. Here we assume that some quantity can only flow to adjacent nodes.

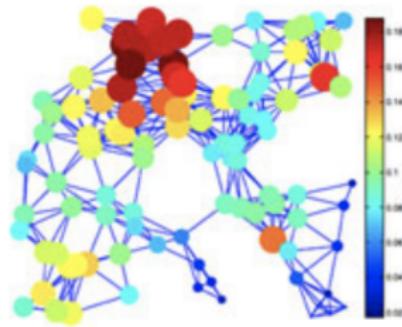


Figure: from <https://web.media.mit.edu/~xdong/paper/tsipn17.pdf>

Setup for Graph Diffusion

Understanding how long this quantity takes to reach a different part of the graph, for example, helps to understand how localized the signal is. So, if the signal is very localized and specific to a particular graph region, it will be re-enforced in that graph region and take longer to *flow* to a different part of the graph.

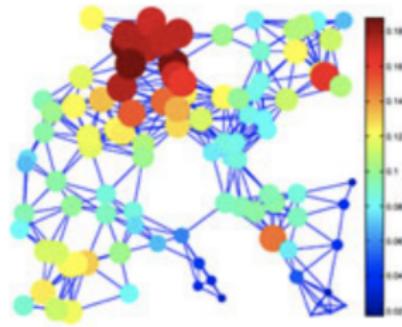


Figure: from <https://web.media.mit.edu/~xdong/paper/tsipn17.pdf>

Diffusion on Graphs

Suppose we have some quantity that will spread or *diffuse* across a graph, given the graph's connectivity. Suppose there is a quantity of ψ_i at node i . Suppose that this quantity flows along the graph from node j to an adjacent node i at a rate of $C(\psi_j - \psi_i)$. The rate at which ψ_i is changing is given as,

$$\frac{d\psi_i}{dt} = C \sum_j A_{ij}(\psi_j - \psi_i) \quad (5)$$

Here, A_{ij} enforces the idea that the only terms in a sum are pairs of nodes connected by an edge.

$$\frac{d\psi_i}{dt} = C \sum_j A_{ij}(\psi_j - \psi_i) \quad (6)$$

Splitting up terms, we can rewrite this as,

$$\frac{d\psi_i}{dt} = C \sum_j A_{ij}\psi_j - C\psi_i \sum_j A_{ij} = C \sum_j A_{ij}\psi_j - C\psi_i k_i \quad (7)$$

- Remember, $k_i = \sum_j A_{ij}$ is simply the degree of node i

Diffusion continued

$$\frac{d\psi_i}{dt} = C \sum_j A_{ij} \psi_j - C \psi_i k_i \quad (8)$$

further simplifies to,

$$\frac{d\psi_i}{dt} = C \sum_j (A_{ij} - \delta_{ij} k_j) \psi_j. \quad (9)$$

Here, k_i is the degree of node i and $\delta_{ij} = 1$ if $i = j$ and is 0 otherwise.

- That was just some massaging so that everything can be written as a \sum_j to make our next step easier.

Massaging Explained

In case you were wondering how we reworked $C\psi_i k_i$ to be written in terms of \sum_j as, $C \sum_j \delta_{ij} \psi_j k_j$

- We only need to count it for node i , so δ_{ij} is our indicator for indexing over j that ensures that we only add a non-zero value when $i = j$

Diffusion and Connection with the Graph Laplacian

$$\frac{d\psi_i}{dt} = C \sum_j (A_{ij} - \delta_{ij} k_j) \psi_j. \quad (10)$$

- This simplifies because $C\psi_i k_i = C \sum_j \delta_{ij} \psi_j k_j$

This can further be expressed in matrix form as,

$$\frac{d\psi}{dt} = C(\mathbf{A} - \mathbf{D})\psi. \quad (11)$$

Here, \mathbf{D} is an $N \times N$ diagonal matrix where k_i is the i th diagonal element and represents the degree of node i .

- The **Graph Laplacian** is the matrix $\mathbf{L} = \mathbf{D} - \mathbf{A}$, and satisfies $\frac{d\psi}{dt} + C\mathbf{L}\psi = 0$ (aka the diffusion equation).

Diffusion on a Bunny Using its Graph Representation

A small amount of signal spreads locally according to specified edges and takes many time steps to spread across the bunny.

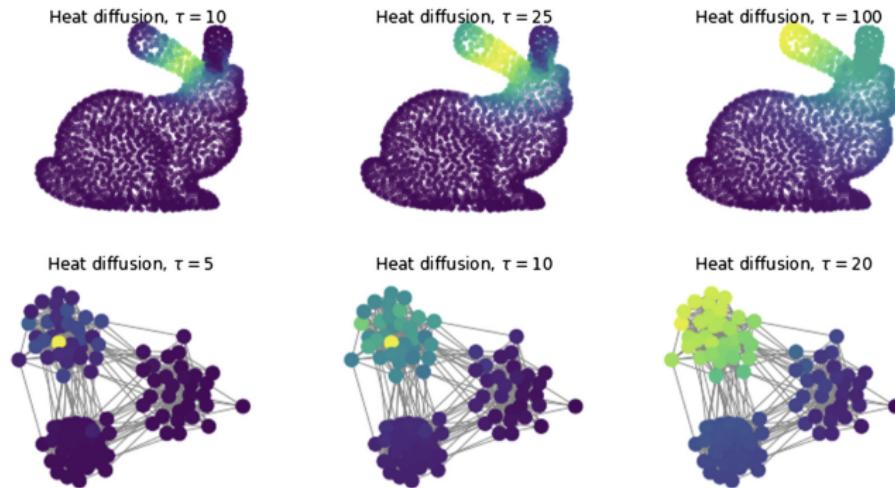


Figure: from <https://www.sciencedirect.com/science/article/pii/S1631070519301094>

Welcome Graph Laplacian

- Besides describing diffusion, the Graph Laplacian has many other nice properties.
- Notably, properties about **connectivity patterns** in the graph and can be useful for **graph partitioning**.
- We are going to be using the graph Laplacian for **many tasks!**

Laplacian Matrix in More Detail

The elements of the Laplacian matrix are defined as follows for an undirected, unweighted matrix.

- Remember: $\mathbf{L} = \mathbf{D} - \mathbf{A}$

$$L_{ij} = \begin{cases} k_i & \text{if } i = j \\ -1 & \text{if } i \neq j \text{ and there is an edge } (i, j) \\ 0 & \text{otherwise} \end{cases}$$

The Graph Laplacian Has Many Nice Linear Algebraic Properties

- In a network with c components, the number of zero eigenvalues is equal to the number of components, (c).

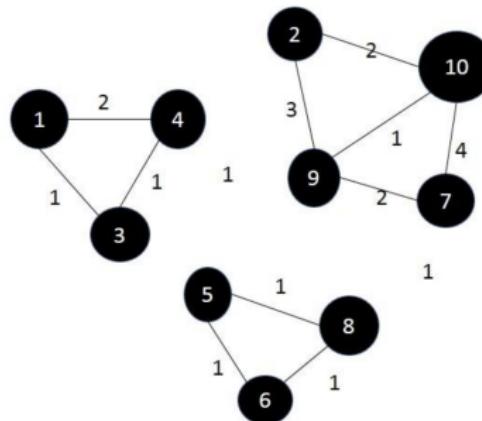


Figure: Each group of nodes makes a component.

Building graphs from data

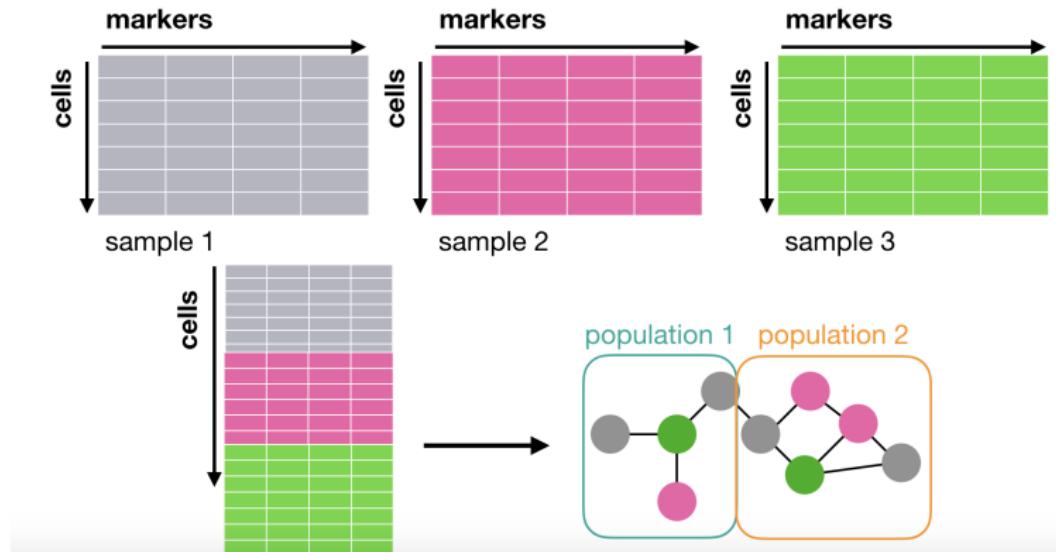


Figure: The task is to compute distances between all pairs of cells, and connect each cell with their nearest neighbors

k-Nearest Neighbor Graphs

- Suppose our dataset, \mathcal{X} is defined as $\mathcal{X} = \{\mathbf{x}_i\}_{i=1\dots N}, \mathbf{x}_i \in \mathbb{R}^d$.
- The general idea is to connect each node with its k nearest neighbors
 - Compute all pairwise similarities between all pairs of nodes
 - For a node, i , find the k nodes that are closest and draw edges.
- Computing all pairwise distances quickly becomes expensive ($O(N^2 d)$).

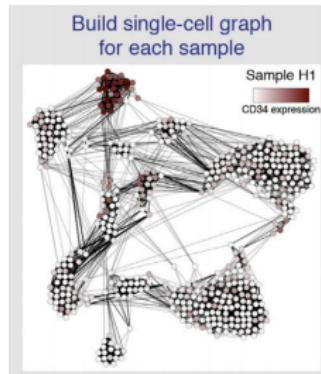


Figure: from Levine *et al.* Cell. 2015. Each cell is connected to its k -nearest

A Common way to Calculate Edge Weights in a Graph

For better or for worse, people commonly assume a Gaussian distribution on pairwise similarities between nodes. So, a weight between nodes i and j is often computed as,

$$W_{ij} = \exp\left(\frac{-||x_i - x_j||^2}{2\sigma^2}\right) \quad (12)$$

Randomized Trees for Nearest Neighbor Search

People have come up with clever ways to compute k -NN graphs much faster. A popular approach is k -d trees and more recently random projection trees.

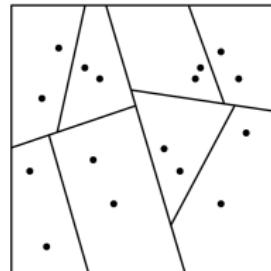
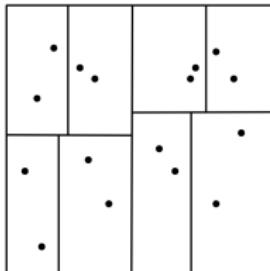


Figure: from Dasgupta *et al.* JMLR. 2013. $k - d$ tree (left), random projection tree (right).

k -d tree

This is a partition of \mathbb{R}^d hyper-rectangular cells based on the datapoints.

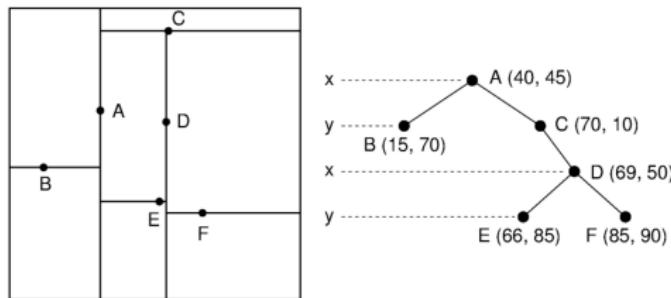


Figure: Alternate between splitting according to the 'x' and 'y' dimensions. First split according to 'A'. Based on partition by A, you can split horizontally according to C and B.

Random Projection Trees

Instead of doing axis parallel splits, split directions and split points are randomly selected.

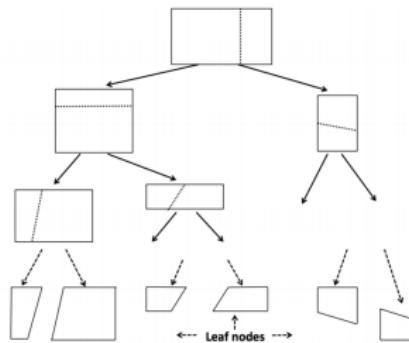


Figure: Random projection tree. (see
<https://cseweb.ucsd.edu/~dasgupta/papers/exactnn-colt.pdf> for a discussion on randomized partition trees.)

One Application of k -NN Graph is Visualization (LargeVis)

Using random projection trees, LargeVis creates an embedding for a set of points with a probabilistic model.

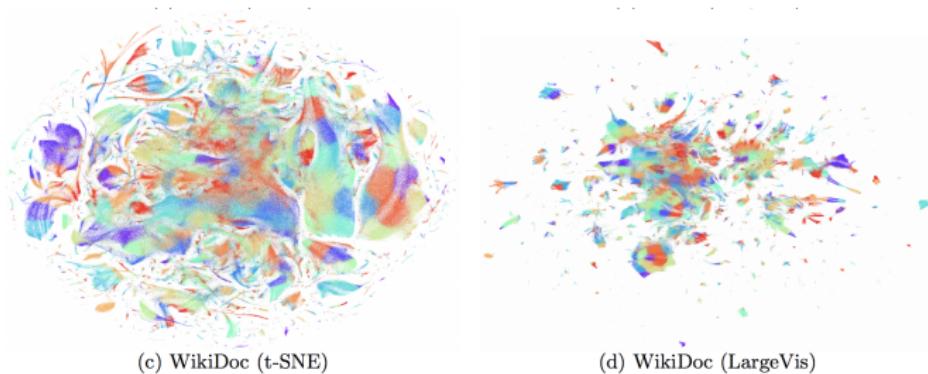


Figure: from Tang *et al.* The Web Conference (WWW). 2016. tSNE and ‘string behavior’ from being designed to preserve super local similarities.

A Probabilistic Model for Graph Embedding

For some nodes, i and j , let \mathbf{y}_i and \mathbf{y}_j be their embedding coordinates such that $\mathbf{y}_i, \mathbf{y}_j \in \mathbb{R}^d$ (for visualization purposes, $d = 2$). We can model the probability that an edge exists between nodes i and j as a function ($f(\cdot)$) of their learned embedding coordinates, or,

$$P(e_{ij} = 1) = f(||\mathbf{y}_i - \mathbf{y}_j||) \quad (13)$$

- An important criterion to be satisfied is that $f(\cdot)$ must ensure a high probability when i and j are close, and low otherwise.
- For example: $f(x) = \frac{1}{1+ax^2}$

Writing Down the Likelihood of the Observed Graph

The authors define a similarity measure between pairs of nodes where there is an edge, based on the original data, \mathbf{x}_i and \mathbf{x}_j . Here, γ being a constant fixed weight for negative edges.

The likelihood of the graph can be written as,

$$\begin{aligned} O &= \prod_{(i,j) \in E} P(e_{ij} = 1) \prod_{(i,j) \in \tilde{E}} (1 - P(e_{ij} = 1))^\gamma \\ &\propto \sum_{(i,j) \in E} \log P(e_{ij} = 1) + \sum_{(i,j) \in \tilde{E}} \gamma \log(1 - P(e_{ij} = 1)) \end{aligned}$$

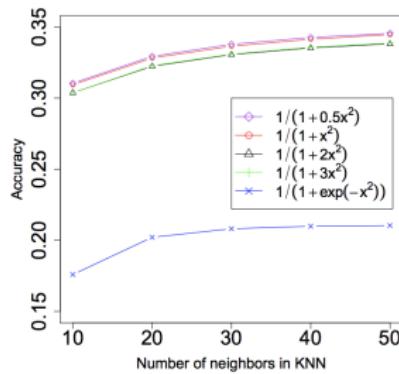
\tilde{E} represents pairs of nodes with no edge.

Brief Overview of Optimization of Parameters

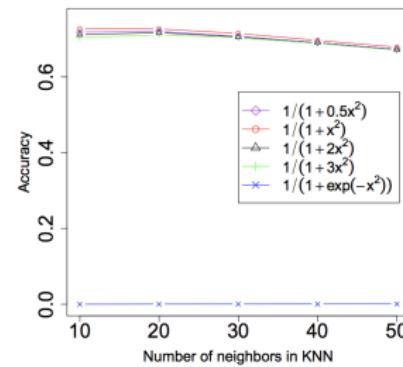
- Since these graphs are very sparse, there are a lot of ‘negative edges’ or edges $\in \tilde{E}$. So the authors used a smart negative sampling approach to only consider a subset of negative node pairs
 - For each node, i , consider sample some ‘negative’ pair nodes.
 - A node, j from the negative pair set to pair with i as a negative example is sampled with probability $\propto d_j^{0.75}$. Here d_j is the degree of node j .
 - Optimize parameters with asynchronous stochastic gradient descent (<https://papers.nips.cc/paper/2011/file/218a0aef1d1a4be65601cc6ddc1520e-Paper.pdf>)

Choice of $f(\cdot)$

The authors tested how choice of $f(\cdot)$ affected the quality of their embedding. Looks like there is one bad choice, but overall, does not make a huge difference.



(a) WikiDoc



(b) LiveJournal

Figure: from Tang *et al.* The Web Conference (WWW). 2016.

More Practical Reasons to Use LargeVis: Scales better than tSNE

Nothing would ever be as fast as PCA. But, if you need to understand more local types of similarities between data points, perhaps it's worth the wait....

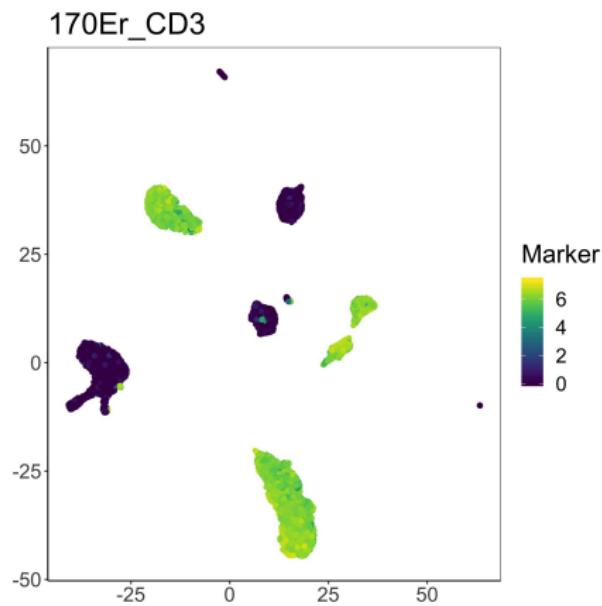
Table 2: Comparison of running time (hours) in graph visualization between the t-SNE and LargeVis.

Algorithm	20NG	MNIST	WikiWord	WikiDoc	LiveJournal	CSAuthor	DBLPaper
t-SNE	0.12	0.41	9.82	45.01	70.35	28.33	18.73
LargeVis	0.14	0.23	2.01	5.60	9.26	4.24	3.19
Speedup Rate	0	0.7	3.9	7	6.6	5.7	4.9

Figure: from Tang *et al.* The Web Conference (WWW). 2016. Comparison in run-time between tSNE and LargeVis.

Example Use Case- Single Cell Data!

Your input, \mathbf{X} is the cells \times marker matrix. Not only do you get the graph out, but we can easily find T-cells and different T-cell subsets!



Community Detection Illustrated

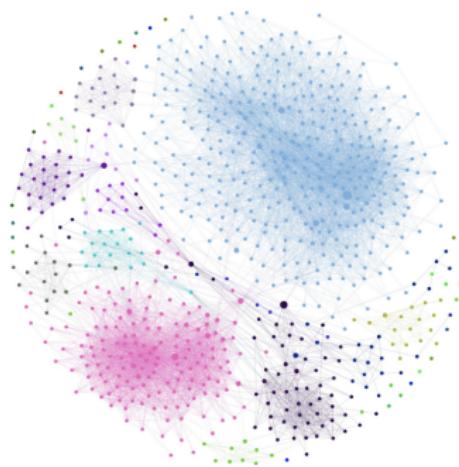


Figure: The high-level goal is to partition nodes into coherent node-subsets (communities), such that nodes within a community have on average more within group edges than between-group edges.

Community Detection (Graph Partitioning) is just Clustering

- Instead of your standard clustering problem in a matrix of N objects with P features, our input here is an adjacency matrix, \mathbf{A} , where we want to cluster nodes based on similarities in their connectivity patterns.
- Most community detection optimization problems seek a hard partition (each node assigned to a single community)
- There are some variants that learn a soft partition, or a propensity or probability that each node is assigned to each community.
- Optimization for this problem can come in many flavors

Optimization Approaches

- Quality Function with a Null Model + Heuristic for Optimization ←
 - A null model describes a graph with no structure, for example, nodes connected randomly.
- Probabilistic and Likelihood Optimization ←
- Spectral Clustering Methods (Partition based on graph Laplacian)
- Most recently: graph embedding + clustering on embedding
 - We will touch on this briefly Thursday with node2vec

Why Might a Person Waste Time Partitioning a Graph?

Build a graph between cells based on marker expression and partition into cell-populations. Members of a cell-population should be phenotypically similar and therefore express the same markers.

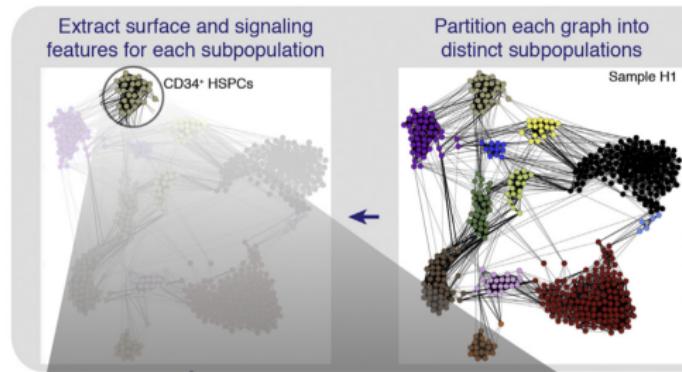


Figure: from Levine *et al.* Cell. 2015. This is the PhenoGraph algorithm.

Recap

- Matrix Rank and SVD
- Graphs: Properties, Laplacian
- LargeVis: An example of using a graph-based representation of data that was computed efficiently for visualization.

Next time:

Graph partitioning and ranking nodes (in a way that is more sophisticated than just node degree....)