# Block 2: Intro to graphs

### ELEC 573: Network Science and Analytics

Santiago Segarra

Electrical and Computer Engineering
Rice University
segarra@rice.edu

Fall 2021

| Wk. | Date | Topic | HW | Project |
|-----|------|-------|-----|---------|
| 1 | 23-Aug | Introduction to course | HW0 out | |
| 2 | 30-Aug | Graph theory / Centrality measures | HW0 solutions posted | |
| 3 | 6-Sep | LABOR DAY (no class) | HW1 out | |
| 4 | 13-Sep | Centrality measures / Community detection | | |
| 5 | 20-Sep | Community detection | | |
| 6 | 27-Sep | Signal Processing and Deep learning for graphs | HW1 due | |
| 7 | 4-Oct | Signal Processing and Deep learning for graphs | HW2 out | |
| 8 | 11-Oct | FALL BREAK (no class) | | |
| 9 | 18-Oct | Network models | HW2 due | |
| 10 | 25-Oct | Network models | HW3 out | Project proposal due |
| 11 | 1-Nov | Epidemics | | |
| 12 | 8-Nov | Inference of network topologies, features, and processes | HW3 due | |
| 13 | 15-Nov | Inference of network topologies, features, and processes | | |
| 14 | 22-Nov | Inference of network topologies, features, and processes | | Project progress report |
| 15 | 29-Nov | Inference of network topologies, features, and processes | | |
| | 13-Dec | Project presentation (video recording) and final report due | | |

(I) 3 homework sets worth 40% (plus an ungraded Homework 0)

▶ Mix of analytical problems and programming assignments

▶ Collaboration accepted, welcomed, and encouraged

▶ However, the submitted work must be your own

# Homework, project and grading

(I) 3 homework sets worth 40% (plus an ungraded Homework 0)

▶ Mix of analytical problems and programming assignments

▶ Collaboration accepted, welcomed, and encouraged

▶ However, the submitted work must be your own

(II) Research project on a topic of your choice, worth 60%

▶ Important and demanding part of this class. Three deliverables:

    1) Proposal by the end of week 9, worth 10%

    2) Progress report by the end of week 12, worth 15%

    3) Final report and recorded presentation, worth 35%

# Homework, project and grading

(I) 3 homework sets worth 40% (plus an ungraded Homework 0)

▶ Mix of analytical problems and programming assignments

▶ Collaboration accepted, welcomed, and encouraged

▶ However, the submitted work must be your own

(II) Research project on a topic of your choice, worth 60%

▶ Important and demanding part of this class. Three deliverables:

    1) Proposal by the end of week 9, worth 10%

    2) Progress report by the end of week 12, worth 15%

    3) Final report and recorded presentation, worth 35%

▶ This is a research-oriented graduate level class

    ⇒ Focus should be on thinking, reading, asking, implementing

# A few things about handing-in homework

- All submissions must be via Canvas
- Can be scanned copies of handwritten work if you are tidy

- All homework is released on Monday evenings
- All homework is due on Tuesday by midnight

- For coding exercises, ready-to-run code must be included
    - ⇒ Submit all your work in a compressed folder
    - ⇒ Jupyter notebooks highly appreciated
    - ⇒ 'lastname_homework_i'

- If anything comes up, please come talk to me with time!
    - ⇒ An ounce of prevention is worth a pound of cure

# Basic notions and definitions

▶ Graph $G(V, E) \Rightarrow$ A set $V$ of vertices or nodes
$\Rightarrow$ Connected by a set $E$ of edges or links
$\Rightarrow$ Elements of $E$ are unordered pairs $(u, v)$, $u, v \in V$

▶ In figure $\Rightarrow$ Vertices are $V = \{1, 2, 3, 4, 5, 6\}$
$\Rightarrow$ Edges $E = \{(1, 2), (1, 5), (2, 3), (3, 4), ...$
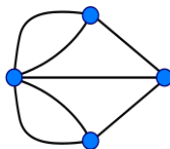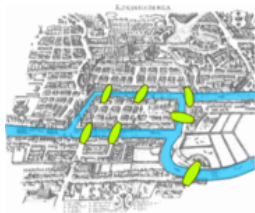$(3, 5), (3, 6), (4, 5), (4, 6)\}$

- Networks are complex systems of inter-connected components

- Graphs are mathematical representations of these systems
  - $\Rightarrow$ Formal language we use to talk about networks

# From networks to graphs

- **Networks** are complex systems of inter-connected components

- **Graphs** are mathematical representations of these systems
   - $\Rightarrow$ Formal language we use to talk about networks



- **Components:** nodes, vertices      $V$
- **Inter-connections:** links, edges      $E$
- **Systems:** networks, graphs      $G(V, E)$

| Network | Vertex | Edge |
|---|---|---|
| Internet | Computer/router | Cable or wireless link |
| Metabolic network | Metabolite | Metabolic reaction |
| WWW | Web page | Hyperlink |
| Food web | Species | Predation |
| Gene-regulatory network | Gene | Regulation of expression |
| Friendship network | Person | Friendship or acquaintance |
| Power grid | Substation | Transmission line |
| Affiliation network | Person and club | Membership |
| Protein interaction | Protein | Physical interaction |
| Citation network | Article/patent | Citation |
| Neural network | Neuron | Synapse |
| ⋮ | ⋮ | ⋮ |

▶ In general, graphs may have self-loops and multi-edges

⇒ A graph with either is called a multi-graph

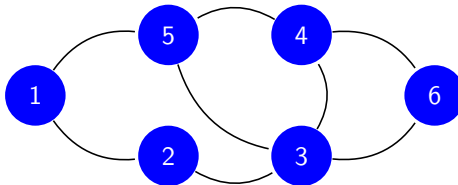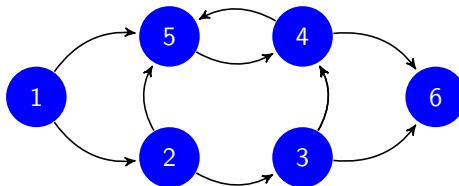▶ In general, graphs may have self-loops and multi-edges

⇒ A graph with either is called a multi-graph



▶ Mostly work with simple graphs, with no self-loops or multi-edges

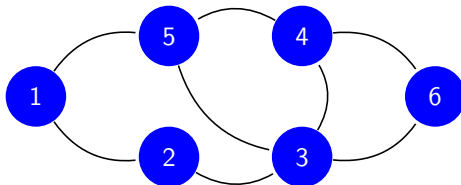# Directed graphs



- In directed graphs, elements of $E$ are ordered pairs $(u, v)$, $u, v \in V$
    - $\Rightarrow$ Means $(u, v)$ distinct from $(v, u)$

- Directed graphs often called digraphs
    - $\Rightarrow$ By convention $(u, v)$ points to $v$
    - $\Rightarrow$ If both $\{(u, v), (v, u)\} \subseteq E$, the edges are said to be mutual

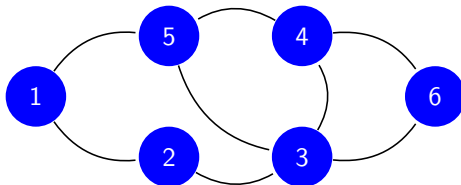- Ex: who-calls-whom phone networks, Twitter follower networks

# Subgraphs

- Consider a given graph $G(V, E)$
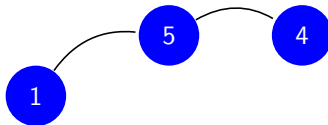


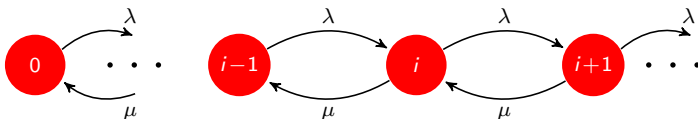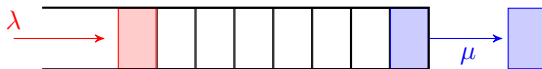- **Def:** Graph $G'(V', E')$ is an induced subgraph of $G$ if $V' \subseteq V$ and $E' \subseteq E$ is the collection of edges in G among that subset of vertices

# Subgraphs

- Consider a given graph $G(V, E)$



- **Def:** Graph $G'(V', E')$ is an induced subgraph of $G$ if $V' \subseteq V$ and $E' \subseteq E$ is the collection of edges in G among that subset of vertices

- Ex: Graph induced by $V' = \{1, 4, 5\}$

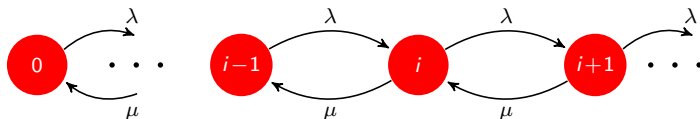# Weighted graphs

- ▶ Oftentimes one labels vertices, edges or both with numerical values
  - ⇒ Such graphs are called weighted graphs
- ▶ Useful in modeling e.g., Markov chain transition diagrams
- ▶ Ex: Single server queuing system (M/M/1 queue)

- Oftentimes one labels vertices, edges or both with numerical values
    - ⇒ Such graphs are called weighted graphs
- Useful in modeling e.g., Markov chain transition diagrams
- Ex: Single server queuing system (M/M/1 queue)



- Labels could correspond to measurements of network processes
- Ex: Node is infected or not with influenza, IP traffic carried by a link

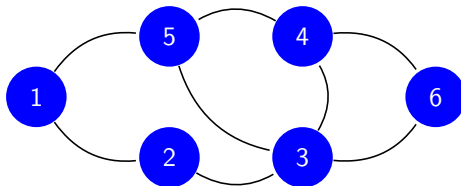| Network | Graph representation |
|---|---|
| WWW | Directed multi-graph (with loops), unweighted |
| Facebook friendships | Undirected, unweighted |
| Citation network | Directed, unweighted, acyclic |
| Collaboration network | Undirected, unweighted |
| Mobile phone calls | Directed, weighted |
| Protein interaction | Undirected multi-graph (with loops), unweighted |
| ⋮ | ⋮ |

▶ Note that multi-edges are often encoded as edge weights (counts)

# Adjacency
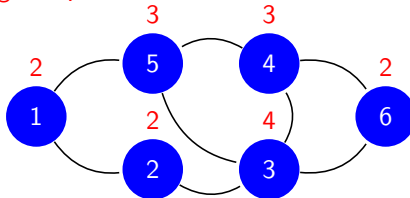
▶ Useful to develop a language to discuss the connectivity of a graph

▶ A simple and local notion is that of adjacency

⇒ Vertices $u, v \in V$ are said adjacent if joined by an edge in $E$

⇒ Edges $e_1, e_2 \in E$ are adjacent if they share an endpoint in $V$



▶ In figure ⇒ Vertices 1 and 5 are adjacent; 2 and 4 are not
⇒ Edge $(1, 2)$ is adjacent to $(1, 5)$, but not to $(4, 6)$

# Degree

- An edge $(u, v)$ is incident with the vertices $u$ and $v$

- **Def:** The degree $d_v$ of vertex $v$ is its number of incident edges



- In figure $\Rightarrow$ Vertex degrees shown in red, e.g., $d_1 = 2$ and $d_5 = 3$
- High-degree vertices likely influential, central, prominent.
- The neighborhood $\mathcal{N}_i$ of a node $i$ is the set of all its adjacent nodes
    $\Rightarrow \mathcal{N}_5 = \{1, 3, 4\}$ $\Rightarrow$ In general, $|\mathcal{N}_i| = d_i$

# Properties and observations about degrees

- ▶ Degree values range from 0 to $|V| - 1$
- ▶ The sum of the degree sequence is twice the size of the graph

$$\sum_{v=1}^{|V|} d_v = 2|E|$$

# Properties and observations about degrees

- Degree values range from 0 to $|V| - 1$
- The sum of the degree sequence is twice the size of the graph

$$\sum_{v=1}^{|V|} d_v = 2|E|$$
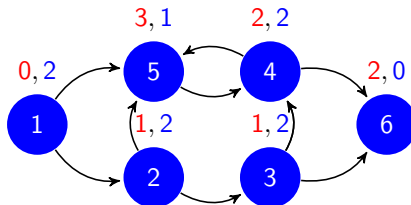
⇒ The number of vertices with odd degree is even

# Properties and observations about degrees

- Degree values range from 0 to $|V| - 1$
- The sum of the degree sequence is twice the size of the graph

$$\sum_{v=1}^{|V|} d_v = 2|E|$$

  $\Rightarrow$ The number of vertices with odd degree is even

- In digraphs, we have vertex in-degree $d_v^{in}$ and out-degree $d_v^{out}$



- In figure $\Rightarrow$ Vertex in-degrees shown in red, out-degrees in blue
  $\Rightarrow$ For example, $d_1^{in} = 0, d_1^{out} = 2$ and $d_5^{in} = 3, d_5^{out} = 1$

# Degree distribution

► Let $N(d)$ denote the number of vertices with degree $d$

  $\Rightarrow$ Fraction of vertices with degree $d$ is $P[d] := \frac{N(d)}{|V|}$
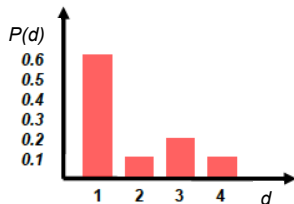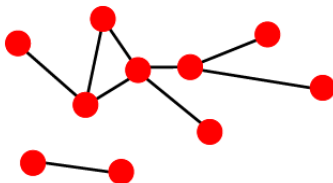
# Degree distribution

▶ Let $N(d)$ denote the number of vertices with degree $d$

  $\Rightarrow$ Fraction of vertices with degree $d$ is $P[d] := \frac{N(d)}{|V|}$

▶ **Def:** The collection $\{P[d]\}_{d \geq 0}$ is the degree distribution of $G$

  ▶ Histogram formed from the degree sequence (bins of size one)



▶ $P[d] =$ probability that randomly chosen node has degree $d$

  $\Rightarrow$ Summarizes the local connectivity in the network graph

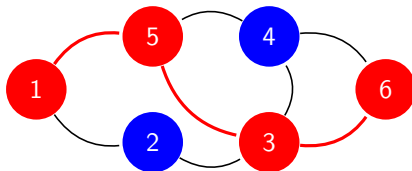- A path of length $l$ from $v_0$ to $v_l$ is a sequence

$$\{v_0, v_1, \ldots, v_{l-1}, v_l\}, \text{ where } v_i \text{ and } v_{i+1} \text{ are adjacent}$$

▶ A path of length $l$ from $v_0$ to $v_l$ is a sequence

$$\{v_0, v_1, \ldots, v_{l-1}, v_l\}, \text{ where } v_i \text{ and } v_{i+1} \text{ are adjacent}$$

▶ A simple path is a path without repeated nodes



▶ A closed path ($v_0 = v_l$) is denominated a circuit

⇒ If no other nodes are repeated, then it is a cycle

▶ A path of length $l$ from $v_0$ to $v_l$ is a sequence

$$\{v_0, v_1, \ldots, v_{l-1}, v_l\}, \text{ where } v_i \text{ and } v_{i+1} \text{ are adjacent}$$

▶ A simple path is a path without repeated nodes
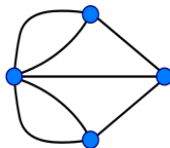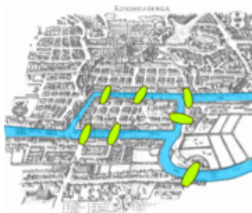


▶ A closed path ($v_0 = v_l$) is denominated a circuit

⇒ If no other nodes are repeated, then it is a cycle

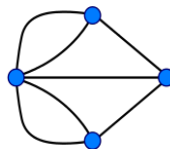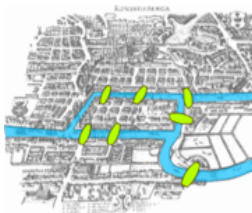▶ All these notions generalize naturally to directed graphs

▶ Can you cross each bridge exactly once in a path?

▶ Can you cross each bridge exactly once in a path?



▶ Graph matters, not physical properties of the island
▶ Suppose a walk existed, what would this imply about the degrees?
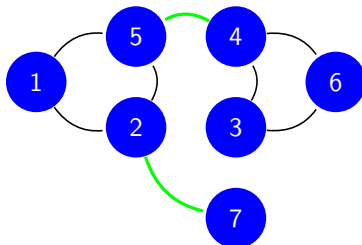▶ Can such a walk exist in Königsberg?

# Distances in a graph

▶ Assume that the edge weights represent length of walk

⇒ Length of a path ⇒ Sum weights of traversed edges

▶ The distance between two nodes $i$ and $j$ is the length of the shortest path linking $i$ and $j$

⇒ In the absence of such a path, the distance is $\infty$

⇒ The diameter of the graph is the value of the largest distance

▶ There exist efficient algorithms to compute distances in graphs
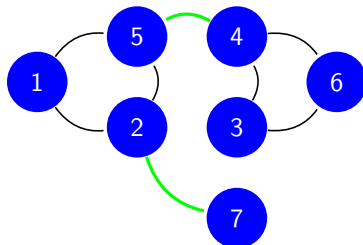
⇒ Dijkstra, Floyd-Warshall, Johnson, ...

- ▶ Vertex $v$ is reachable from $u$ if there exists a $u - v$ path
- ▶ **Def:** Graph is connected if every vertex is reachable from every other

▶ Vertex $v$ is reachable from $u$ if there exists a $u - v$ path

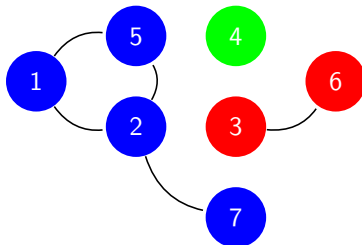▶ **Def:** Graph is connected if every vertex is reachable from every other



▶ If bridge edges are removed, the graph becomes disconnected

▶ **Def:** A component is a maximally connected subgraph

  ⇒ Maximal means adding a vertex will ruin connectivity



▶ In figure ⇒ Components are $\{1, 2, 5, 7\}$, $\{3, 6\}$ and $\{4\}$

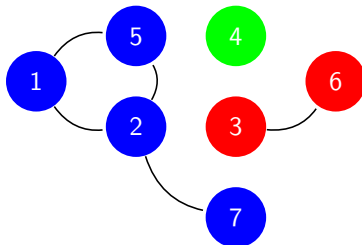     ⇒ Subgraph $\{3, 4, 6\}$ not connected, $\{1, 2, 5\}$ not maximal

▶ **Def:** A component is a maximally connected subgraph
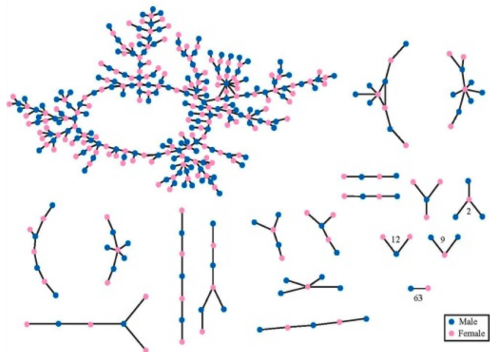
⇒ Maximal means adding a vertex will ruin connectivity



▶ In figure ⇒ Components are $\{1, 2, 5, 7\}$, $\{3, 6\}$ and $\{4\}$

⇒ Subgraph $\{3, 4, 6\}$ not connected, $\{1, 2, 5\}$ not maximal

▶ Disconnected graphs have 2 or more components

⇒ Largest component often called giant component

# Giant connected components

▶ Large real-world networks typically exhibit one giant component

▶ Ex: romantic relationships in a US high school [Bearman et al'04]

- Large real-world networks typically exhibit one giant component
- Ex: romantic relationships in a US high school [Bearman et al'04]
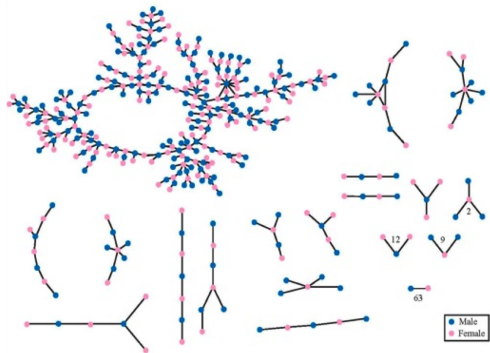


- Q: Why do we expect to find a single giant component?

- Large real-world networks typically exhibit one giant component
- Ex: romantic relationships in a US high school [Bearman et al'04]



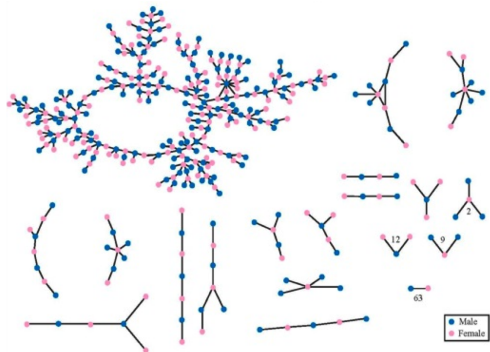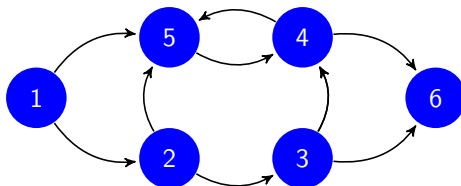- Q: Why do we expect to find a single giant component?
- A: It only takes one edge to merge two giant components

# Connectivity of directed graphs

- Connectivity is more subtle with directed graphs. Two notions

- Digraph is strongly connected if for every pair $u, v \in V$, $u$ is reachable from $v$ (via a directed path) and vice versa

- Digraph is weakly connected if connected after disregarding edge directions, i.e., the underlying undirected graph is connected
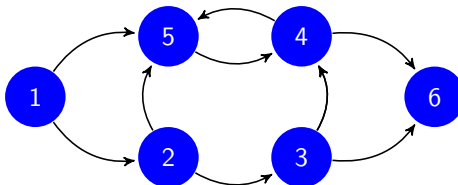
- Connectivity is more subtle with directed graphs. Two notions

- Digraph is strongly connected if for every pair $u, v \in V$, $u$ is reachable from $v$ (via a directed path) and vice versa

- Digraph is weakly connected if connected after disregarding edge directions, i.e., the underlying undirected graph is connected



- Above graph is weakly connected but not strongly connected
  $\Rightarrow$ Strong connectivity implies weak connectivity

▶ A complete graph $K_n$ of order $n$ has all possible edges



$K_2 \qquad K_3 \qquad K_4 \qquad K_5$

▶ Q: How many edges does $K_n$ have?

- ▶ A complete graph $K_n$ of order $n$ has all possible edges



$K_2$      $K_3$      $K_4$      $K_5$

- ▶ Q: How many edges does $K_n$ have?
- ▶ A: Number of edges in $K_n$ = Number of vertex pairs = $\binom{n}{2} = \frac{n(n-1)}{2}$

▶ A complete graph $K_n$ of order $n$ has all possible edges



$K_2$    $K_3$    $K_4$    $K_5$

▶ Q: How many edges does $K_n$ have?

▶ A: Number of edges in $K_n$ = Number of vertex pairs = $\binom{n}{2} = \frac{n(n-1)}{2}$

▶ Of interest in network analysis are cliques, i.e., complete subgraphs

⇒ Extreme notions of cohesive subgroups, communities

# Regular graphs

- A *d*-regular graph has vertices with equal degree *d*



- Naturally, the complete graph $K_n$ is $(n-1)$-regular
  - $\Rightarrow$ Cycles are 2-regular (sub) graphs

- Regular graphs arise frequently in e.g.,
  - Physics and chemistry in the study of crystal structures
  - Geo-spatial settings as pixel adjacency models in image processing

# Trees and directed acyclic graphs

▶ A tree is a connected acyclic graph

  ⇒ A collection of trees is denominated a forest

▶ Ex: river network, information cascades in Twitter, citation network



▶ A directed tree is a digraph whose underlying undirected graph is a tree

  ⇒ Rooted if paths from one vertex to all others

▶ Vertex terminology: parent, children, ancestor, descendant, leaf

▶ A tree is a connected acyclic graph

⇒ A collection of trees is denominated a forest

▶ Ex: river network, information cascades in Twitter, citation network



▶ A directed tree is a digraph whose underlying undirected graph is a tree

⇒ Rooted if paths from one vertex to all others

▶ Vertex terminology: parent, children, ancestor, descendant, leaf

▶ Underlying graph of a directed acyclic graph (DAG) need not be a tree

# Bipartite graphs

▶ A graph $G(V, E)$ is called bipartite when

⇒ $V$ can be partitioned in two disjoint sets, say $V_1$ and $V_2$; and

⇒ Each edge in $E$ has one endpoint in $V_1$, the other in $V_2$



▶ Useful to represent e.g., membership or affiliation networks

⇒ Nodes in $V_1$ could be people, nodes in $V_2$ clubs

⇒ Associated graph $G(V_1, E_1)$ joins members of same club

▶ The mathematics collaboration network centered at Paul Erdős



▶ Most mathematicians have an Erdős number of at most 4 or 5
  ⇒ Drawing created by R. Graham in 1979

# What is your Erdős number?

- ▶ There are resources online that can help you calculate it
- ▶ My Erdős number is 3, where a shortest path is given by



| Paul Erdös | Vance Faber | Gunnar Carlsson | me |

- ▶ What about the Erdős-Bacon number?
  - ⇒ Richard Feynman 6 $(3 + 3)$
  - ⇒ Carl Sagan 6 $(4 + 2)$
  - ⇒ Natalie Portman 7 $(5 + 2)$
  - ⇒ Santiago Segarra $\infty$ $(3 + \infty)$ ... so far

Basic notions and definitions

Algebraic graph theory

Graph data structures and algorithms

Strength of weak ties

▶ Algebraic graph theory deals with matrix representations of graphs

⇒ Leverage algebra to 'visualize' graphs as if being plotted

# Adjacency matrix

- ▶ Algebraic graph theory deals with matrix representations of graphs
    - ⇒ Leverage algebra to 'visualize' graphs as if being plotted

- ▶ Q: How can we capture the connectivity of $G(V, E)$ in a matrix?

- ▶ A: Binary, symmetric adjacency matrix $\mathbf{A} \in \{0, 1\}^{|V| \times |V|}$, with entries

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

- ⇒ Note that vertices are indexed with integers $1, \dots, |V|$

# Adjacency matrix

▶ Algebraic graph theory deals with matrix representations of graphs

    ⇒ Leverage algebra to 'visualize' graphs as if being plotted

▶ Q: How can we capture the connectivity of $G(V, E)$ in a matrix?

▶ A: Binary, symmetric adjacency matrix $\mathbf{A} \in \{0, 1\}^{|V| \times |V|}$, with entries

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

    ⇒ Note that vertices are indexed with integers $1, \ldots, |V|$

▶ In words, $\mathbf{A}$ is one for those entries whose row-column indices denote vertices in $V$ joined by an edge in $E$, and is zero otherwise

# Adjacency matrix examples

▶ Examples for undirected graphs and digraphs

⇒ There is an implicit labeling function



$$\mathbf{A}_u = \begin{pmatrix} 0 & 1 & 0 & 1 \\ 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{pmatrix}, \quad \mathbf{A}_d = \begin{pmatrix} 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 \end{pmatrix}$$

▶ If the graph is weighted, store the $(i, j)$ weight instead of 1

# Adjacency matrix properties

- Adjacency matrix useful to store graph structure.

    ⇒ Also, operations on **A** yield useful information about $G$

# Adjacency matrix properties

▶ Adjacency matrix useful to store graph structure.

　　⇒ Also, operations on **A** yield useful information about $G$

▶ Degrees: Row-wise sums give vertex degrees, i.e., $\sum_{j=1}^{|V|} A_{ij} = d_i$

▶ For digraphs **A** is not symmetric and row-, colum-wise sums differ

$$\sum_{j=1}^{|V|} A_{ij} = d_i^{out}, \qquad \sum_{i=1}^{|V|} A_{ij} = d_j^{in}$$

- Adjacency matrix useful to store graph structure.
    - ⇒ Also, operations on **A** yield useful information about $G$

- Degrees: Row-wise sums give vertex degrees, i.e., $\sum_{j=1}^{|V|} A_{ij} = d_i$
- For digraphs **A** is not symmetric and row-, colum-wise sums differ

$$\sum_{j=1}^{|V|} A_{ij} = d_i^{out}, \qquad \sum_{i=1}^{|V|} A_{ij} = d_j^{in}$$

- Paths: Let $\mathbf{A}^r$ denote the $r$-th power of **A**, with entries $A_{ij}^r$
    - ⇒ Then $A_{ij}^r$ yields the number of $i - j$ paths of length $r$ in $G$

# Adjacency matrix properties

▶ Adjacency matrix useful to store graph structure.

⇒ Also, operations on **A** yield useful information about $G$

▶ Degrees: Row-wise sums give vertex degrees, i.e., $\sum_{j=1}^{|V|} A_{ij} = d_i$

▶ For digraphs **A** is not symmetric and row-, colum-wise sums differ

$$\sum_{j=1}^{|V|} A_{ij} = d_i^{out}, \qquad \sum_{i=1}^{|V|} A_{ij} = d_j^{in}$$

▶ Paths: Let $\mathbf{A}^r$ denote the $r$-th power of **A**, with entries $A_{ij}^r$

⇒ Then $A_{ij}^r$ yields the number of $i - j$ paths of length $r$ in $G$

▶ Corollary: $\text{tr}(\mathbf{A}^2)/2 = |E|$ and $\text{tr}(\mathbf{A}^3)/6 = \#\triangle$ in $G$

# Adjacency matrix properties

- ▶ Adjacency matrix useful to store graph structure.
  - ⇒ Also, operations on **A** yield useful information about $G$

- ▶ Degrees: Row-wise sums give vertex degrees, i.e., $\sum_{j=1}^{|V|} A_{ij} = d_i$

- ▶ For digraphs **A** is not symmetric and row-, colum-wise sums differ

$$\sum_{j=1}^{|V|} A_{ij} = d_i^{out}, \qquad \sum_{i=1}^{|V|} A_{ij} = d_j^{in}$$

- ▶ Paths: Let $\mathbf{A}^r$ denote the $r$-th power of **A**, with entries $A_{ij}^r$
  - ⇒ Then $A_{ij}^r$ yields the number of $i - j$ paths of length $r$ in $G$

- ▶ Corollary: $\text{tr}(\mathbf{A}^2)/2 = |E|$ and $\text{tr}(\mathbf{A}^3)/6 = \#\triangle$ in $G$

- ▶ Spectrum: $G$ is $d$-regular if and only if **1** is an eigenvector of **A**, i.e.,

$$\mathbf{A1} = d\mathbf{1}$$

# Incidence matrix

- ▶ A graph can be also represented by its $|V| \times |E|$ incidence matrix **B**
  - ⇒ **B** is in general not a square matrix, unless $|V| = |E|$
- ▶ If the graph is undirected, we first assign arbitrary directions to edges
- ▶ We encode the direction of the edges, namely

$$B_{ij} = \left\{ \begin{array}{cl} 1, & \text{if edge } j \text{ is } (k, i) \\ -1, & \text{if edge } j \text{ is } (i, k) \\ 0, & \text{otherwise} \end{array} \right. .$$

▶ Example of undirected graph with added directions



$$\mathbf{B} = \begin{pmatrix} -1 & 0 & -1 & 0 & 1 \\ 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & -1 \\ 0 & -1 & 1 & -1 & 0 \end{pmatrix}$$

▶ If the graph is weighted, modify nonzero entries accordingly

# Graph Laplacian

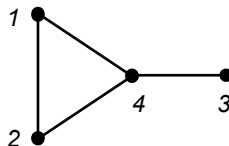▶ Vertex degrees often stored in the diagonal matrix **D**, where $D_{ii} = d_i$

$$\mathbf{D} = \begin{pmatrix} 2 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 3 \end{pmatrix}$$

# Graph Laplacian

▶ Vertex degrees often stored in the diagonal matrix $\mathbf{D}$, where $D_{ii} = d_i$

$$\mathbf{D} = \begin{pmatrix} 2 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 3 \end{pmatrix}$$



▶ The $|V| \times |V|$ symmetric matrix $\mathbf{L} := \mathbf{D} - \mathbf{A}$ is called graph Laplacian

$$L_{ij} = \begin{cases} d_i, & \text{if } i = j \\ -1, & \text{if } (i,j) \in E \\ 0, & \text{otherwise} \end{cases}, \quad \mathbf{L} = \begin{pmatrix} 2 & -1 & 0 & -1 \\ -1 & 2 & 0 & -1 \\ 0 & 0 & 1 & -1 \\ -1 & -1 & -1 & 3 \end{pmatrix}$$

# Graph Laplacian

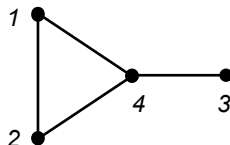▶ Vertex degrees often stored in the diagonal matrix $\mathbf{D}$, where $D_{ii} = d_i$

$$\mathbf{D} = \begin{pmatrix} 2 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 3 \end{pmatrix}$$



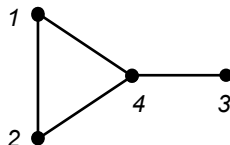▶ The $|V| \times |V|$ symmetric matrix $\mathbf{L} := \mathbf{D} - \mathbf{A}$ is called graph Laplacian

$$L_{ij} = \begin{cases} d_i, & \text{if } i = j \\ -1, & \text{if } (i,j) \in E \\ 0, & \text{otherwise} \end{cases} \quad, \quad \mathbf{L} = \begin{pmatrix} 2 & -1 & 0 & -1 \\ -1 & 2 & 0 & -1 \\ 0 & 0 & 1 & -1 \\ -1 & -1 & -1 & 3 \end{pmatrix}$$

▶ Variants of the Laplacian exist, with slightly different interpretations
  ⇒ Normalized Laplacian $\mathbf{L}_n = \mathbf{D}^{-1/2} \mathbf{L} \mathbf{D}^{-1/2}$
  ⇒ Random-walk Laplacian $\mathbf{L}_{rw} = \mathbf{D}^{-1} \mathbf{L}$

▶ Smoothness: For any vector $\mathbf{x} \in \mathbb{R}^{|V|}$ of "vertex values", one has

$$\mathbf{x}^\top \mathbf{L} \mathbf{x} = \sum_{(i,j) \in E} (x_i - x_j)^2$$

which can be minimized to enforce smoothness of functions on $G$

# Laplacian matrix properties

▶ Smoothness: For any vector $\mathbf{x} \in \mathbb{R}^{|V|}$ of "vertex values", one has

$$\mathbf{x}^\top \mathbf{L} \mathbf{x} = \sum_{(i,j) \in E} (x_i - x_j)^2$$

which can be minimized to enforce smoothness of functions on $G$

▶ Incidence relation: $\mathbf{L} = \mathbf{B}\mathbf{B}^\top$ where $\mathbf{B}$ has arbitrary orientation

▶ Positive semi-definiteness: Follows since $\mathbf{x}^\top \mathbf{L} \mathbf{x} \geq 0$ for all $\mathbf{x} \in \mathbb{R}^{|V|}$

# Laplacian matrix properties

▶ Smoothness: For any vector $\mathbf{x} \in \mathbb{R}^{|V|}$ of "vertex values", one has

$$\mathbf{x}^\top \mathbf{L} \mathbf{x} = \sum_{(i,j) \in E} (x_i - x_j)^2$$

which can be minimized to enforce smoothness of functions on $G$

▶ Incidence relation: $\mathbf{L} = \mathbf{B}\mathbf{B}^\top$ where $\mathbf{B}$ has arbitrary orientation

▶ Positive semi-definiteness: Follows since $\mathbf{x}^\top \mathbf{L} \mathbf{x} \geq 0$ for all $\mathbf{x} \in \mathbb{R}^{|V|}$

▶ Rank deficiency: Since $\mathbf{L}\mathbf{1} = \mathbf{0}$, $\mathbf{L}$ is rank deficient

# Laplacian matrix properties

▶ **Smoothness:** For any vector $\mathbf{x} \in \mathbb{R}^{|V|}$ of "vertex values", one has

$$\mathbf{x}^\top \mathbf{L} \mathbf{x} = \sum_{(i,j) \in E} (x_i - x_j)^2$$

which can be minimized to enforce smoothness of functions on $G$

▶ **Incidence relation:** $\mathbf{L} = \mathbf{B} \mathbf{B}^\top$ where $\mathbf{B}$ has arbitrary orientation

▶ **Positive semi-definiteness:** Follows since $\mathbf{x}^\top \mathbf{L} \mathbf{x} \geq 0$ for all $\mathbf{x} \in \mathbb{R}^{|V|}$

▶ **Rank deficiency:** Since $\mathbf{L}\mathbf{1} = \mathbf{0}$, $\mathbf{L}$ is rank deficient

▶ **Spectrum and connectivity:** The smallest eigenvalue $\lambda_1$ of $\mathbf{L}$ is 0
  ▶ If the second-smallest eigenvalue $\lambda_2 \neq 0$, then $G$ is connected
  ▶ If $\mathbf{L}$ has $n$ zero eigenvalues, $G$ has $n$ connected components
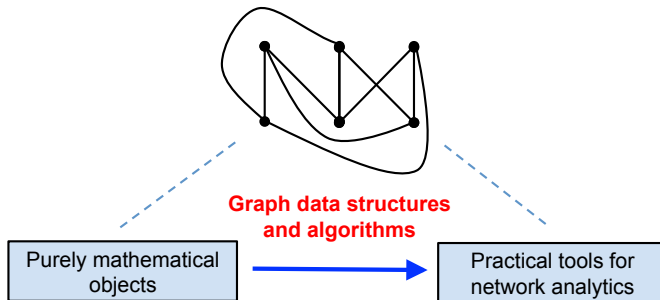
Basic notions and definitions

Algebraic graph theory

Graph data structures and algorithms

Strength of weak ties

▶ Q: How can we store and analyze a graph $G$ using a computer?



**Graph data structures and algorithms**

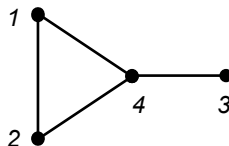| Purely mathematical objects | ⟶ | Practical tools for network analytics |

▶ Data structures: efficient storage and manipulation of a graph

▶ Algorithms: scalable computational methods for graph analytics

   ⇒ Contributions in this area primarily due to computer science

# Adjacency matrix as a data structure

► Q: How can we represent and store a graph $G$ in a computer?

► A: The $|V| \times |V|$ adjacency matrix **A** is a natural choice

$$A_{ij} = \left\{ \begin{array}{ll} 1, & \text{if } (i,j) \in E \\ 0, & \text{otherwise} \end{array} \right. .$$

$$\mathbf{A} = \left( \begin{array}{cccc} 0 & 1 & 0 & 1 \\ 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{array} \right)$$



► Matrices (arrays) are basic data objects in software environments

  ⇒ Naive memory requirement is $O(|V|^2)$

  ⇒ May be undesirable for large, sparse graphs

# Networks are sparse graphs

▶ Most real-world networks are sparse, meaning

$$|E| \ll \frac{|V|(|V|-1)}{2} \text{ or equivalently } \bar{d} := \frac{1}{|V|} \sum_{v=1}^{|V|} d_v \ll |V| - 1$$

▶ Figures from the study by Leskovec et al '09 are eloquent

| Network dataset | $|V|$ | Avg. degree $\bar{d}$ |
|---|---|---|
| WWW (Stanford-Berkeley) | 319,717 | 9.65 |
| Social network (LinkedIn) | 6,946,668 | 8.87 |
| Communication (MSN IM) | 242,720,596 | 11.1 |
| Collaboration (DBLP) | 317,080 | 6.62 |
| Roads (California) | 1,957,027 | 2.82 |
| Proteins (S. Cerevisiae) | 1,870 | 2.39 |

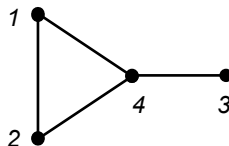▶ Graph density $\rho := \frac{|E|}{|V|^2} = \frac{\bar{d}}{2|V|}$ is another useful metric

▶ An adjacency-list representation of graph $G$ is an array of size $|V|$

⇒ The $i$-th array element is a list of the vertices adjacent to $i$

$$L_a[1] = \{2, 4\}$$
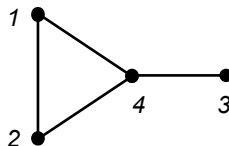$$L_a[2] = \{1, 4\}$$
$$L_a[3] = \{4\}$$
$$L_a[4] = \{1, 2, 3\}$$

▶ An adjacency-list representation of graph $G$ is an array of size $|V|$

⇒ The $i$-th array element is a list of the vertices adjacent to $i$

$L_a[1] =$ {2, 4}
$L_a[2] =$ {1, 4}
$L_a[3] =$ {4}
$L_a[4] =$ {1, 2, 3}



▶ Similarly, an edge list stores the vertex pairs incident to each edge

$$L_e[1] = \{1, 2\}$$
$$L_e[2] = \{1, 4\}$$
$$L_e[3] = \{2, 4\}$$
$$L_e[4] = \{3, 4\}$$

▶ In either case, the memory requirement is $O(|E|)$

▶ Numerous interesting questions may be asked about a given graph

▶ For few simple ones, lookup in data structures suffices

  Q1: Are vertices $u$ and $v$ linked by an edge?

  Q2: What is the degree of vertex $u$?

# Graph algorithms and complexity

▶ Numerous interesting questions may be asked about a given graph

▶ For few simple ones, lookup in data structures suffices

Q1: Are vertices $u$ and $v$ linked by an edge?

Q2: What is the degree of vertex $u$?

▶ Some others require more work. Still can tackle them efficiently

Q1: What is the shortest path between vertices $u$ and $v$?

Q2: How many connected components does the graph have?

Q3: Is a given digraph acyclic?

# Graph algorithms and complexity

- Numerous interesting questions may be asked about a given graph

- For few simple ones, lookup in data structures suffices
    - Q1: Are vertices $u$ and $v$ linked by an edge?
    - Q2: What is the degree of vertex $u$?

- Some others require more work. Still can tackle them efficiently
    - Q1: What is the shortest path between vertices $u$ and $v$?
    - Q2: How many connected components does the graph have?
    - Q3: Is a given digraph acyclic?

- Unfortunately, in some cases there is likely no efficient algorithm
    - Q1: What is the maximal clique in a given graph?

- We'll keep algorithm complexity in mind. Not focus of the course

- ▶ Goal: verify connectivity of a graph based on its adjacency list
- ▶ Idea: start from vertex $s$, explore the graph, mark vertices you visit

Output : List $M$ of marked vertices in the component
Input  : Graph $G$ (e.g., adjacency list)
Input  : Starting vertex $s$
$L := \{s\}$; $M := \{s\}$; % Initialize exploration and marking lists
% Repeat while there are still nodes to explore
**while** $L \neq \emptyset$ **do**
    choose $u \in L$; % Pick arbitrary vertex to explore
    **if** $\exists\,(u, v) \in E$ *such that* $v \notin M$ **then**
        choose $(u, v)$ with $v$ of smallest index;
        $L := L \cup \{v\}$; $M := M \cup \{v\}$; % Mark and augment
    **else**
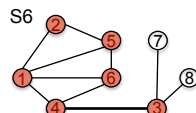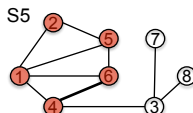        $L := L \setminus \{u\}$; % Prune
    **end**
**end**

# Graph exploration example

▶ Below we indicate the chosen and marked nodes. Initialize $s = 2$

| L | Mark |
|---|------|
| {2} | 2 |
| {2,1} | 1 |
| {2,1,5} | 5 |
| {2,1,5,6} | 6 |
| {1,5,6} | |
| {1,5,6,4} | 4 |
| {5,6,4} | |
| {5,4} | |
| {5,4,3} | 3 |
| {5,3} | |
| {5,3,7} | 7 |
| {5,3} | |
| {3} | |
| {3,8} | 8 |
| {3} | |
| {} | |



▶ Exploration takes $2|V|$ steps. Each node is added and removed once

# Breadth-first search

- Choices made arbitrarily in the exploration algorithm. Variants?
- Breadth-first search (BFS): choose for $u$ the first element of $L$

Output : List $M$ of marked vertices in the component
Input   : Graph $G$ (e.g., adjacency list)
Input   : Starting vertex $s$

$L := \{s\}$; $M := \{s\}$; % Initialize exploration and marking lists
% Repeat while there are still nodes to explore
**while** $L \neq \emptyset$ **do**
    $u := \text{first}(L)$; % Breadth first
    **if** $\exists\, (u, v) \in E$ *such that* $v \notin M$ **then**
        choose $(u, v)$ with $v$ of smallest index;
        $L := L \cup \{v\}$; $M := M \cup \{v\}$; % Mark and augment
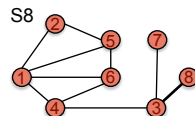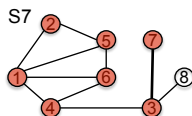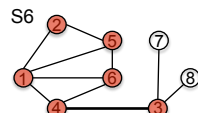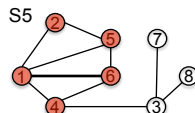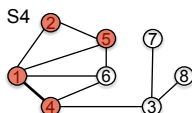    **else**
        $L := L \setminus \{u\}$; % Prune
    **end**
**end**

▶ Below we indicate the chosen and marked nodes. Initialize $s = 2$

| L | Mark |
|---|---|
| {2} | 2 |
| {2,1} | 1 |
| {2,1,5} | 5 |
| {1,5} | |
| {1,5,4} | 4 |
| {1,5,4,6} | 6 |
| {5,4,6} | |
| {4,6} | |
| {4,6,3} | 3 |
| {6,3} | |
| {3} | |
| {3,7} | 7 |
| {3,7,8} | 8 |
| {7,8} | |
| {8} | |
| {} | |



▶ The algorithm builds a wider tree (breadth first)

# Depth-first search

▶ Depth-first search (DFS): choose for $u$ the last element of $L$

Output : List $M$ of marked vertices in the component
Input   : Graph $G$ (e.g., adjacency list)
Input   : Starting vertex $s$

$L := \{s\}$; $M := \{s\}$; % Initialize exploration and marking lists
% Repeat while there are still nodes to explore
**while** $L \neq \emptyset$ **do**
$\quad$ $u := \text{last}(L)$; % Depth first
$\quad$ **if** $\exists\,(u, v) \in E$ *such that* $v \notin M$ **then**
$\quad\quad$ choose $(u, v)$ with $v$ of smallest index;
$\quad\quad$ $L := L \cup \{v\}$; $M := M \cup \{v\}$; % Mark and augment
$\quad$ **else**
$\quad\quad$ $L := L \setminus \{u\}$; % Prune
$\quad$ **end**
**end**

# DFS example

▶ Below we indicate the chosen and marked nodes. Initialize $s = 2$

| L | Mark |
|---|------|
| {2} | 2 |
| {2,1} | 1 |
| {2,1,4} | 4 |
| {2,1,4,3} | 3 |
| {2,1,4,3,7} | 7 |
| {2,1,4,3} | |
| {2,1,4,3,8} | 8 |
| {2,1,4,3} | |
| {2,1,4} | |
| {2,1,4,6} | 6 |
| {2,1,4,6,5} | 5 |
| {2,1,4,6} | |
| {2,1,4} | |
| {2,1} | |
| {2} | |
| {} | |



▶ The algorithm builds longer paths (depth first)

▶ **Def:** The distance between vertices $u$ and $v$ is the length of the shortest $u - v$ path. Oftentimes referred to as geodesic distance

⇒ In the absence of a $u - v$ path, the distance is $\infty$

⇒ The diameter of a graph is the value of the largest distance

- **Def:** The distance between vertices $u$ and $v$ is the length of the shortest $u - v$ path. Oftentimes referred to as geodesic distance
    - $\Rightarrow$ In the absence of a $u - v$ path, the distance is $\infty$
    - $\Rightarrow$ The diameter of a graph is the value of the largest distance
- Q: What are efficient algorithms to compute distances in a graph?
- A: BFS (for unit weights) and Dijkstra's algorithm

- ▶ Use BFS and keep track of path lengths during the exploration
- ▶ Increment distance by 1 every time a vertex is marked

Output : Vector $d$ of distances from reference vertex
Input   : Graph $G$ (e.g., adjacency list)
Input   : Reference vertex $s$

$L := \{s\}$; $M := \{s\}$; $d(s) = 0$; % Initialization
% Repeat while there are still nodes to explore
**while** $L \neq \emptyset$ **do**
$\quad$ $u := \text{first}(L)$; % Breadth first
$\quad$ **if** $\exists\, (u, v) \in E$ *such that* $v \notin M$ **then**
$\quad\quad$ choose $(u, v)$ with $v$ of smallest index;
$\quad\quad$ $L := L \cup \{v\}$; $M := M \cup \{v\}$;% Mark and augment
$\quad\quad$ $d(v) := d(u) + 1$ % Increment distance
$\quad$ **else**
$\quad\quad$ $L := L \setminus \{u\}$; % Prune
$\quad$ **end**
**end**

- BFS tree output for your friendship network

# Strength of weak ties

Basic notions and definitions
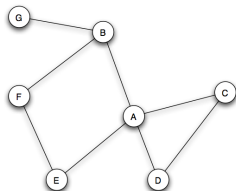
Algebraic graph theory

Graph data structures and algorithms
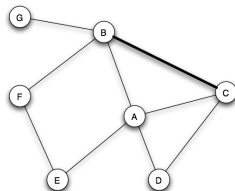
Strength of weak ties

▶ Networks are rarely static structures ⇒ Think about their evolution

  ⇒ How are edges formed? ⇒ Universal feature ⇒ Triadic closure

▶ *If two people in a social network have a friend in common, then there is an increased likelihood that they will become friends at some point in the future*
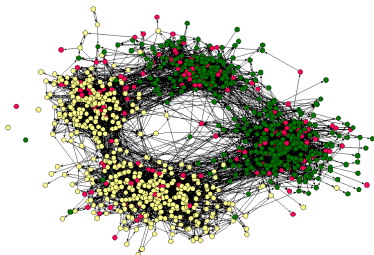


(a) *Before B-C edge forms.*　　　　(b) *After B-C edge forms.*

▶ Triadic closure is very natural ⇒ Some reasons ...

  ⇒ Opportunity: B and C have a higher chance of meeting

  ⇒ Trusting: B and C are predisposed to trusting each other
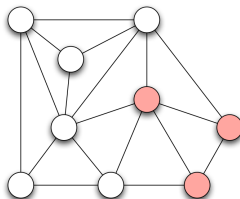
  ⇒ Incentive: A might have incentive to make B and C friends

▶ We tend to be similar to our friends ⇒ Well known for long time

⇒ Age, race, interests, beliefs, opinions, affluence, ...

▶ Contextual (as opposed to intrinsic) effect on network formation

⇒ Contextual: Friends because we attend the same school

⇒ Intrinsic: Friends because a common friend introduces us



▶ In previous slide, B and C high chance of becoming friends
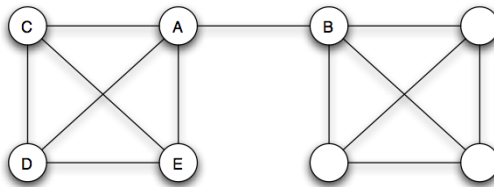
⇒ Even if they are not aware of common knowledge of A

- ▶ Is homophily present or is it an artifact of how the network is drawn?
    - ⇒ We need to formulate a precise mathematical measure
- ▶ Consider a small network of girls ($q = 3/9$) and boys ($p = 6/9$)



- ▶ If edges are agnostic to gender, portion of cross-gender edges is $2pq$
    - ⇒ Homophily Test: If the fraction of cross-gender edges is significantly less than $2pq$, then there is evidence for homophily
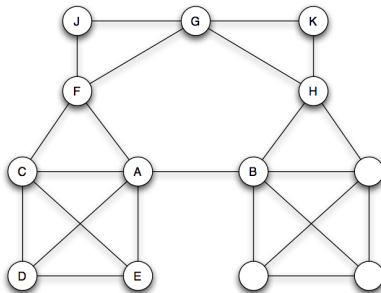    - ⇒ Cross-gender edges $5/18 < 8/18 = 2pq$ ⇒ Mild homophily

- Mark Granovetter (1960s) interviewed people that changed jobs
- Most heard about new jobs from acquaintances rather than close friends
  - ⇒ Explanation takes into account local properties and global structure



- *A*'s friends *E*, *C*, and *D* form a tightly-knit group
- *B* reaches to a different part of the network ⇒ New information
- Deleting (*A*, *B*) disconnects the network ⇒ (*A*, *B*) is a bridge
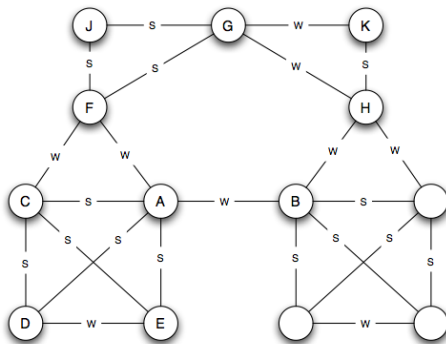  - ⇒ But bridges are rare in real-world networks

# A social network closer to reality

▶ In real life, there are other multi-step paths joining $A$ and $B$

⇒ If $(A, B)$ is deleted, distance becomes more than 2

⇒ Local bridge

⇒ An edge is a local bridge when it is not part of a triangle



▶ Closely knit group of friends are eager to help
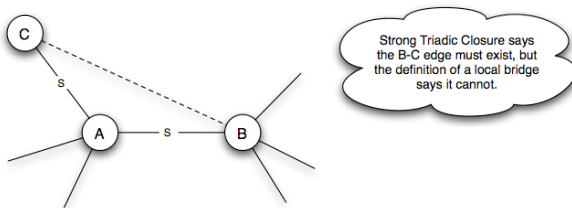
⇒ But have almost the same information as you

# Strong triadic closure

▶ How does overrepresentation of bridges relate to acquaintances?

▶ Consider two different levels of strength in the links of a social network

⇒ Strong ties correspond to friends, weak ties to acquaintances



▶ $A$ violates the Strong Triadic Closure if it has strong ties to two other nodes $B$ and $C$, and there is no edge at all (strong or weak) between $B$ and $C$

# Local bridges and weak ties

- Tie strength ⇒ Local/interpersonal feature
- Bridge property ⇒ Global/structural feature
- How do these two features relate in light of the strong triadic closure?
- *If A satisfies the strong triadic closure and is involved in at least two strong ties, then any local bridge it is involved in must be a weak tie*



Strong Triadic Closure says the B-C edge must exist, but the definition of a local bridge says it cannot.

- Acquaintances are natural sources of new information
    ⇒ Strict modeling assumptions, first-order conclusions, testable