

GRAPHS



Agenda


Morning

- General understanding of graphs
- Traversing a graph
- Importance of an node

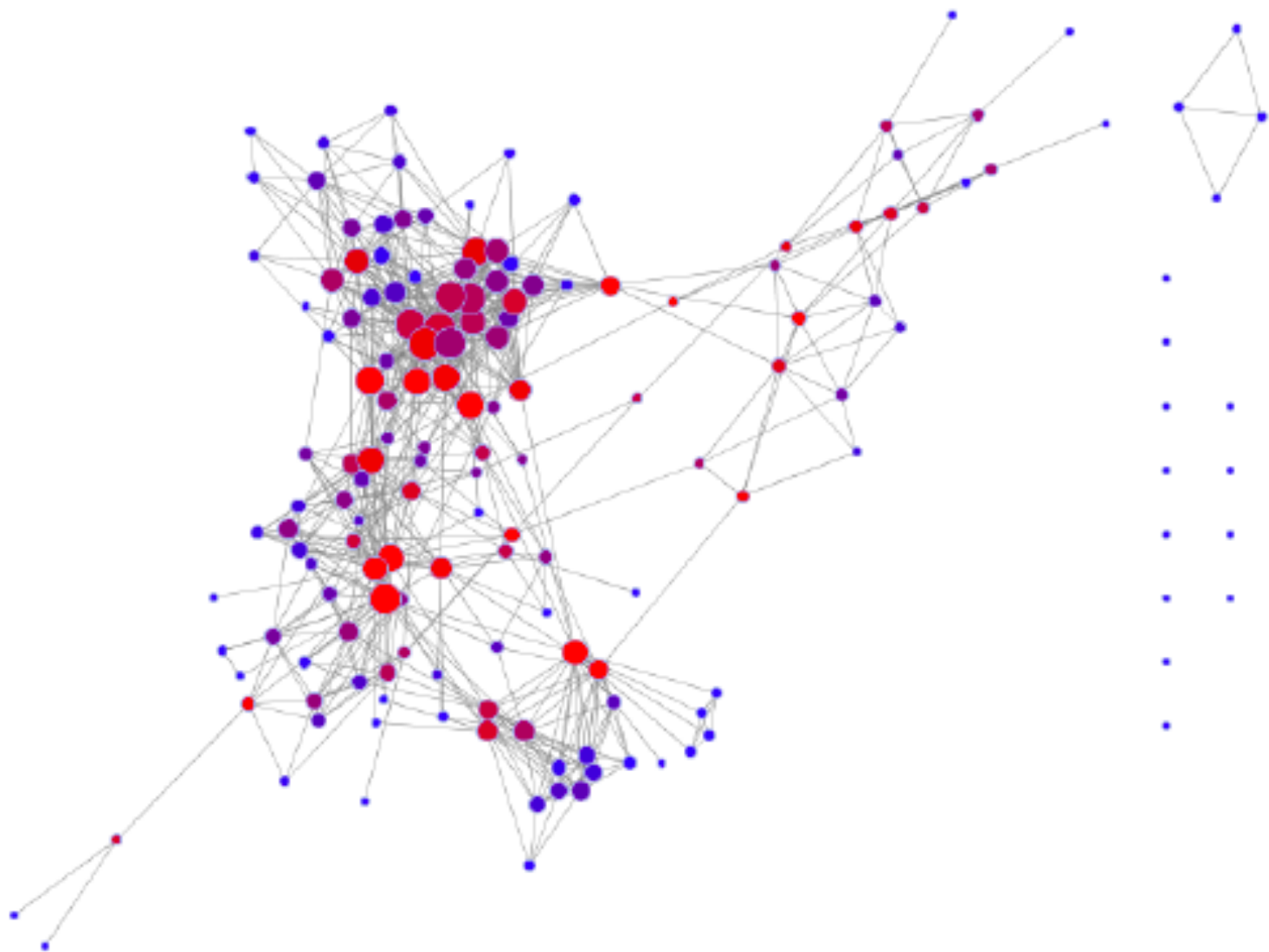
Afternoon

- Defining communities in graphs
- Finding communities

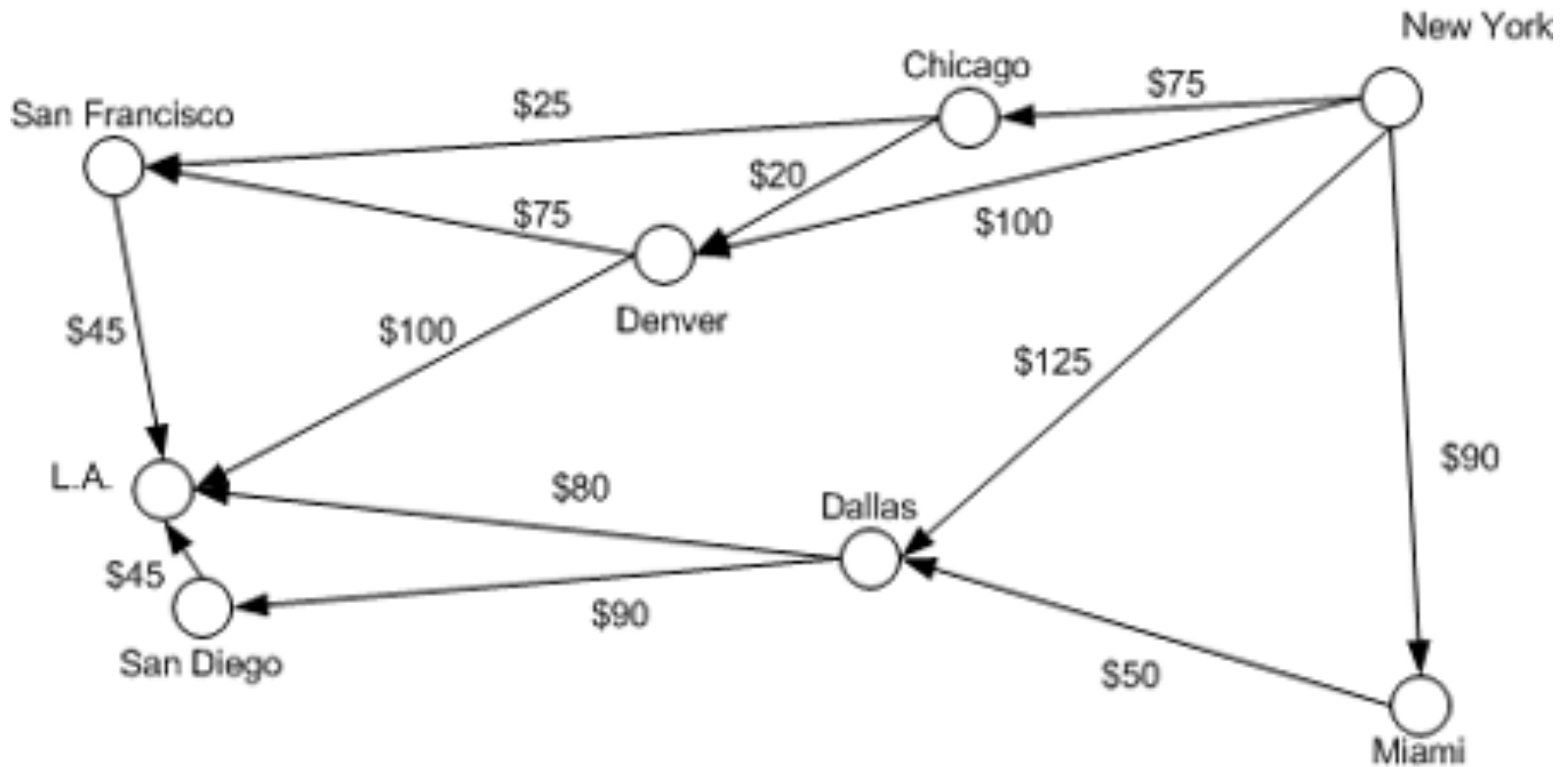
SESSION 1



General understanding of graphs
Traversing a graph
Importance of an node



Real world example



Definition of a graph

A graph is an ordered pair $G = (V, E)$ such that:

- V is a set of **vertices**
- E is a set of **relations**

Each **edge** is 2-element subset $\{V_i, V_j\}$ of V

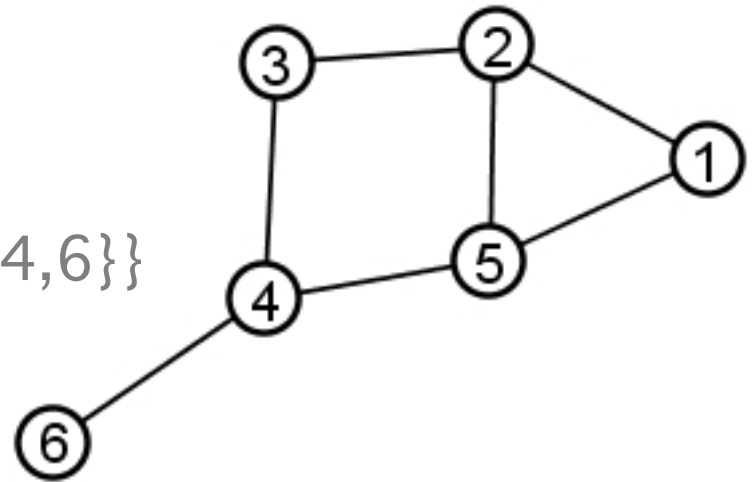
$V =$

$\{1, 2, 3, 4, 5, 6\}$

$E =$

$\{\{1, 2\}, \{1, 5\}, \{2, 3\}, \{2, 5\}, \{3, 4\}, \{4, 5\}, \{4, 6\}\}$

$G = (V, E)$



A *node* is also known as a vertex

If there's an edge between two nodes, we say that those nodes are *connected*.

Graph Terminology

Neighbors: The *neighbors* of a node are the nodes that it is connected to.

Degree: The *degree* is the number of neighbors a node has.

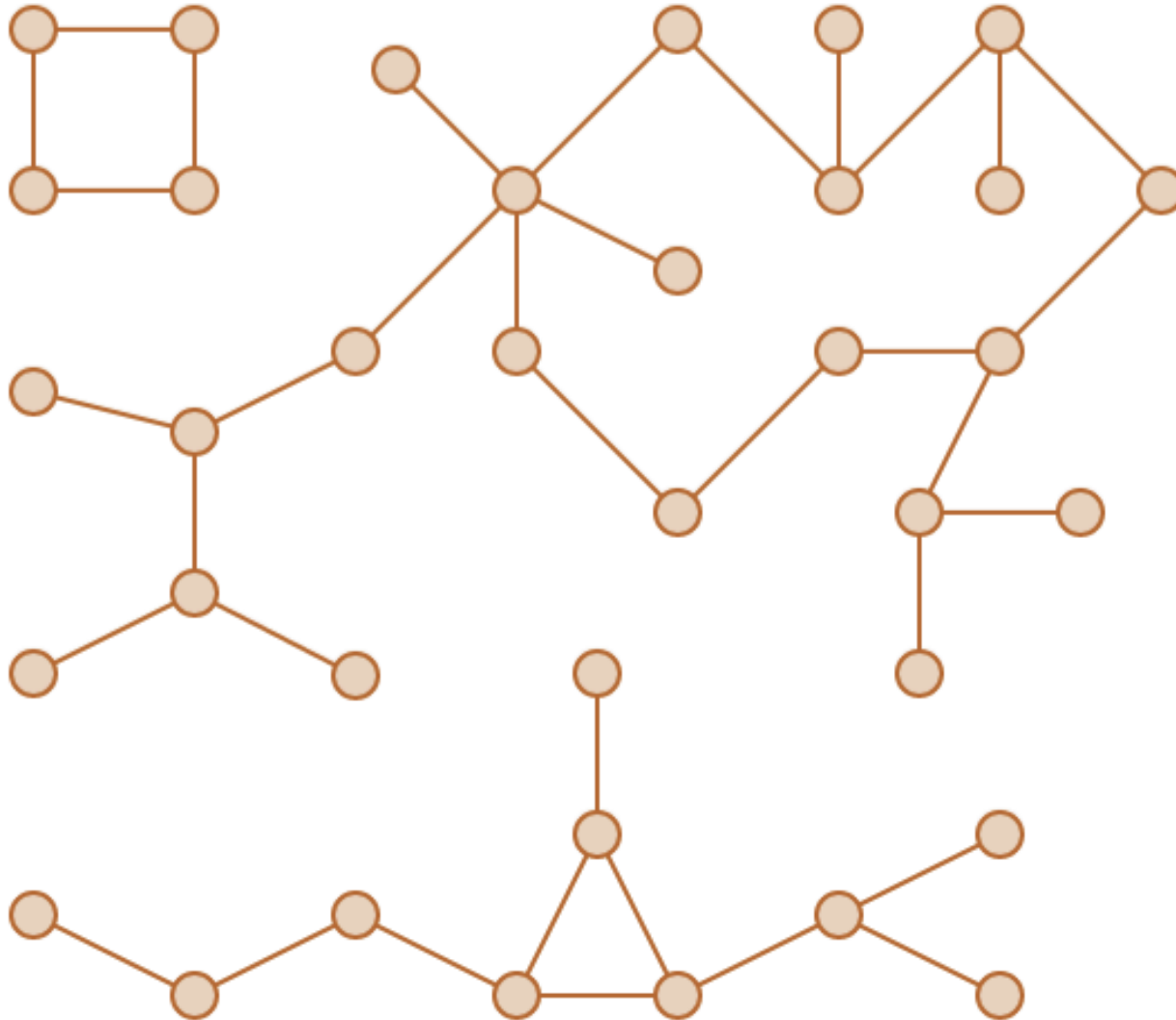
Path: A *path* is a series of nodes and the edges that connect them

Connected: A graph is *connected* if every pair of vertices is connected by some path.

Subgraph: A *subgraph* is a subset of the nodes of a graph and all the edges between them.

Connected Component: A *connected component* is a subgraph that is *connected* and which is connected to no additional vertices in the supergraph.

Graph Terminology – Pop Quiz



More Terminology – types of graphs

Directed Graph: A directed graph is a graph where edges only go in one direction.

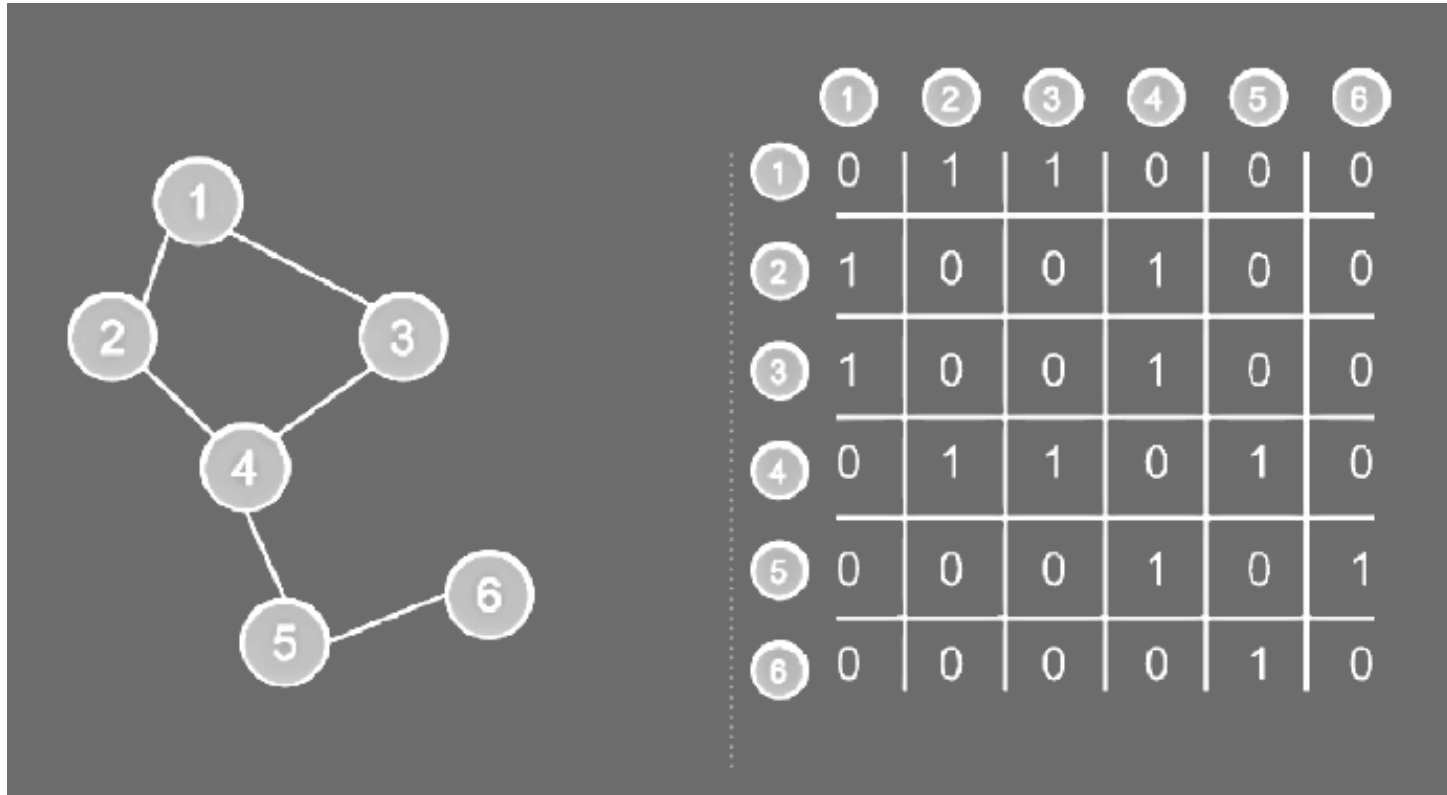
Undirected Graph: In an undirected graph, edges go both ways.

Weighted Graph: In a weighted graph, the edges have weights.

Unweighted Graph: In an unweighted graph, all edges are the same.

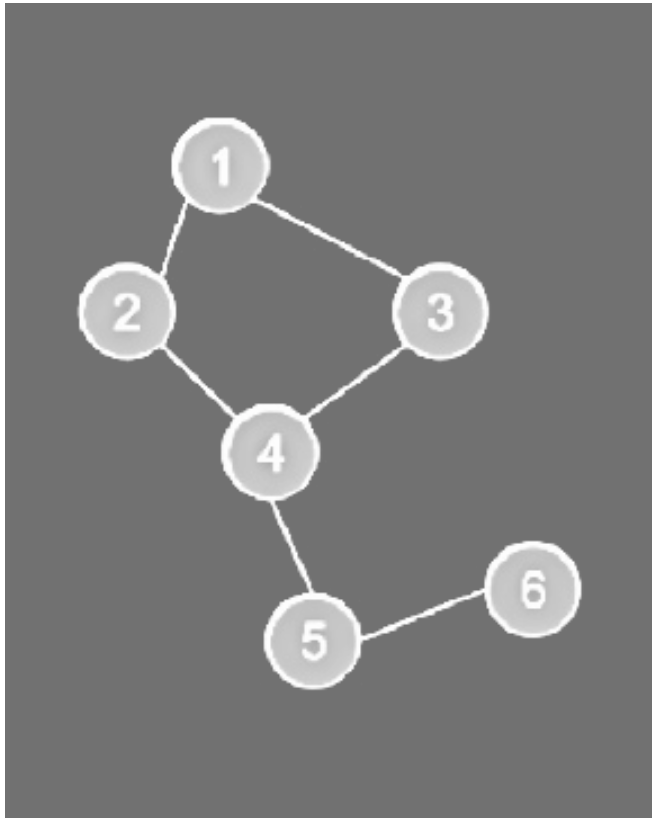
Graph Representations

Adjacency Matrix



Graph Representations

Adjacency List



$\{1\}: \{2, 3\}$

$\{2\}: \{1, 4\}$

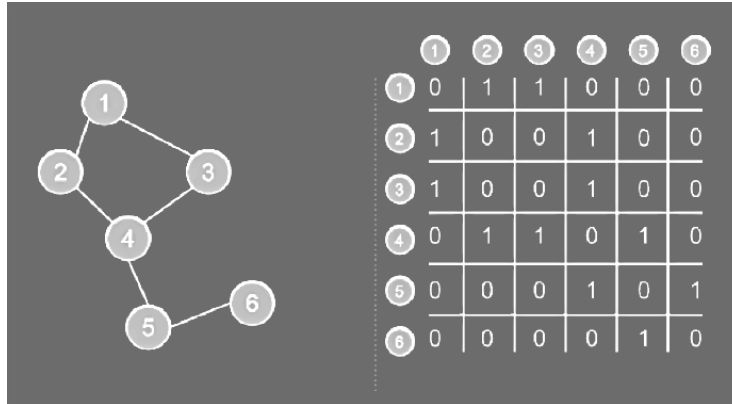
$\{3\}: \{1, 4\}$

$\{4\}: \{2, 3, 5\}$

$\{5\}: \{4, 6\}$

$\{6\}: \{5\}$

Graph Representation



{1}: {2, 3}
{2}: {1, 4}
{3}: {1, 4}
{4}: {2,3,5}
{5}: {4,6}
{6}: {5}

Storage:

- Adjacency matrix: $O(|V|^2)$
- Adjacency list: $(O|V|+O|E|)$

Lookup, Finding edge between two vertices:

- Adjacency matrix: $O(1)$, constant time
- Adjacency list: $(O|V|)$ or $(O|1|)$ depending on implementation

Graph Representation, Pop Quiz

Which representation will give you an answer faster if you wanted to check if two nodes were connected?

Which representation will give you an answer faster if you wanted to find the neighbors of a node?

What if your graph had a ton of edges and maybe even loops?

Graph Search Algorithms

- Systematically visit/traverse all nodes
- Useful to perform all sorts of graph related operations

“Breadth First Search” (BFS) is a popular graph search algorithm

Breadth First Search (BFS)

Animation:

[https://upload.wikimedia.org/wikipedia/commons/4/46/
Animated_BFS.gif](https://upload.wikimedia.org/wikipedia/commons/4/46/Animated_BFS.gif)

Starting with a given node, find all of that node's neighbors. Then, find all of those neighbors' neighbors, and so on....

Visit **ALL** neighbors of a node **BEFORE** visiting neighbors of those neighbors...

BFS, pseudo code

function BFS(graph, starting_node):

 create an empty queue Q

 initialize empty set V (set of visited nodes)

 add A to Q

 while Q is not empty:

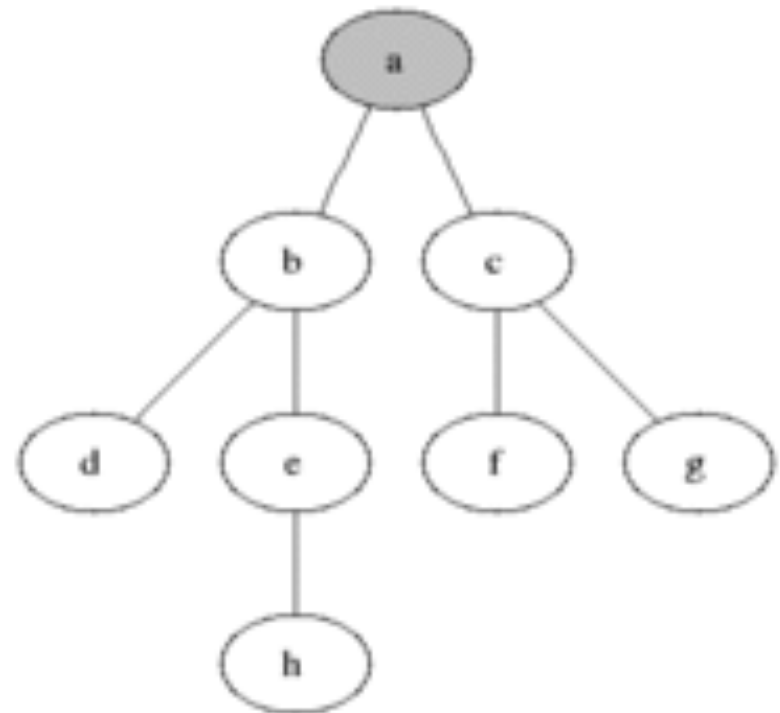
 take node off Q, call it N

 if N isn't already visited:

 add N to the visited set

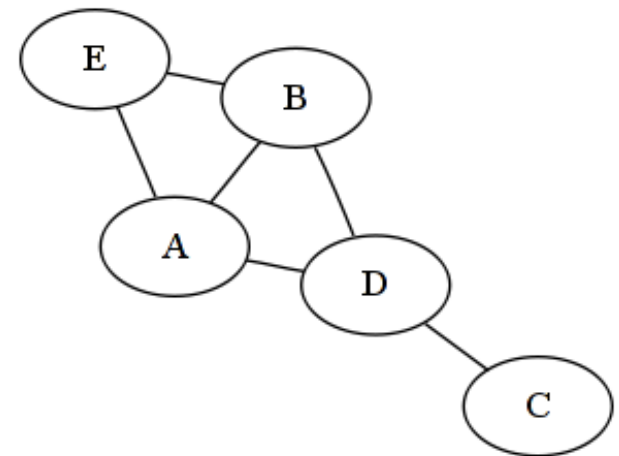
 add every neighbor of N to Q

Sprint : Modify for shortest path



	Current Node (N)	Queue (Q)	Visited (V)	explanation
		A,0		initialization
1	A,0			take first node off queue
	A,0	E,1 B,1 D,1	A	add A's neighbors to queue
2	E,1	B,1 D,1	A	take E off queue
	E,1	B,1 D,1 A,2 B,2	A,E	add E's neighbors to queue
3	B,1	D,1 A,2 B,2	A,E	take B off queue
	B,1	D,1 A,2 B,2 E,2 D,2	A,E,B	add B's neighbors to queue
4	D,1	A,2 B,2 E,2 D,2	A,E,B	take D off queue
	D,1	A,2 B,2 E,2 D,2 C,2 B,2	A,E,B,D	add D's neighbors to queue
5	A,2	B,2 E,2 D,2 C,2 B,2	A,E,B,D	take A off queue
	A,2	B,2 E,2 D,2 C,2 B,2	A,E,B,D	skip A since already visited
6	B,2	E,2 D,2 C,2 B,2	A,E,B,D	take B off queue
	B,2	E,2 D,2 C,2 B,2	A,E,B,D	skip B since already visited
7	E,2	D,2 C,2 B,2	A,E,B,D	take E off queue
	E,2	D,2 C,2 B,2	A,E,B,D	skip E since already visited
8	D,2	C,2 B,2	A,E,B,D	take D off queue
	D,2	C,2 B,2	A,E,B,D	skip D since already visited
9	C,2	B,2	A,E,B,D	take C off queue
	C,2	B,2	A,E,B,D	We're done!! :tada:

Modify for
shortest path from A to C



BFS, pop quiz

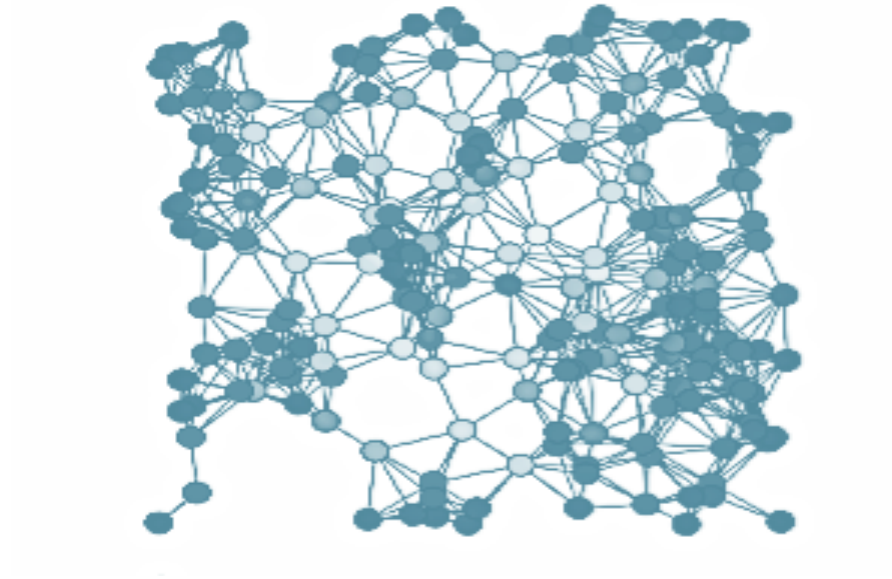
BFS implies that there is also a depth first search(DFS). Any guesses on what that would look like?

BFS is much more widely used than DFS. Why would that be?

What would be a good use case for depth first search?

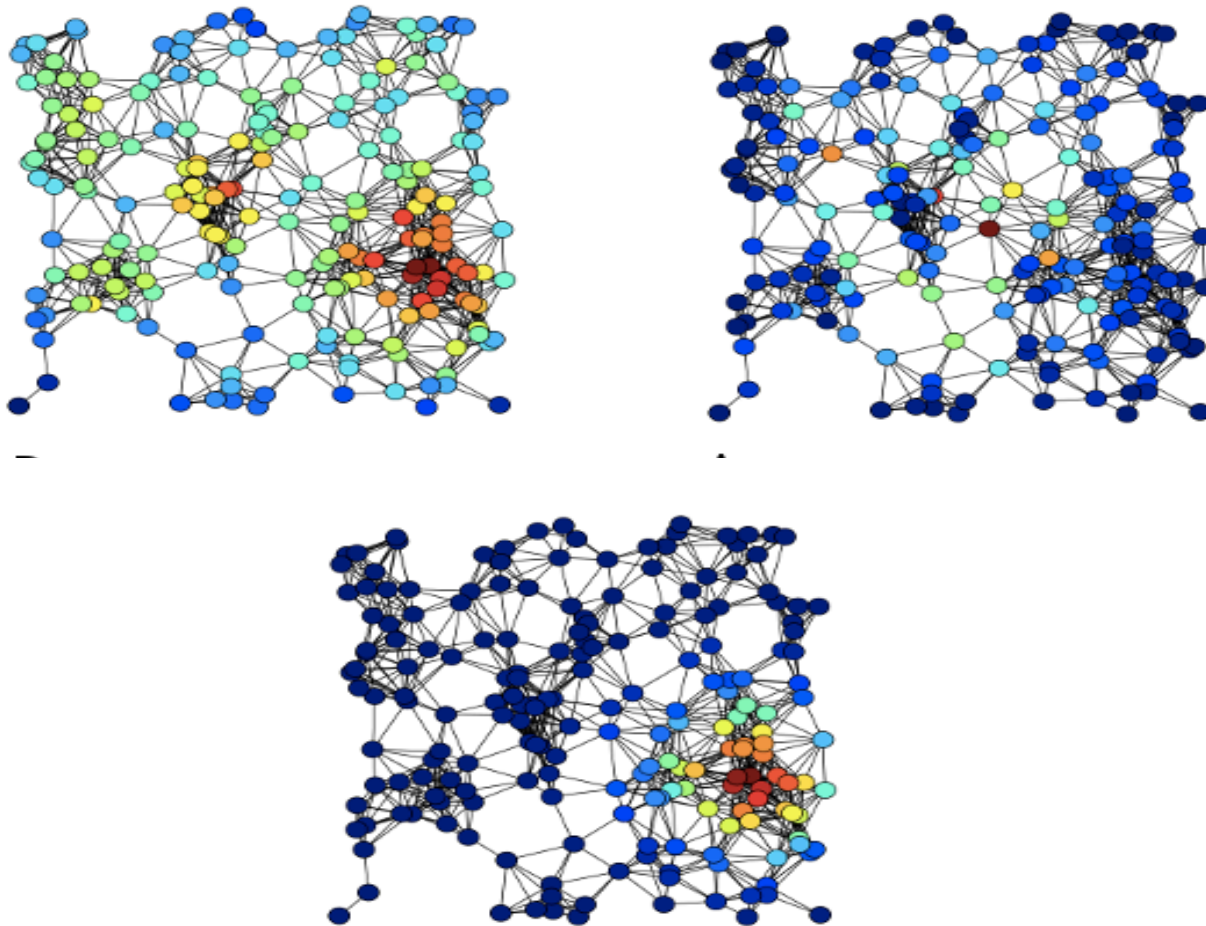
FYI, BFS pseudo-code is a popular interview question.

“Important” nodes?



What makes a node important or central?

Centrality of a node



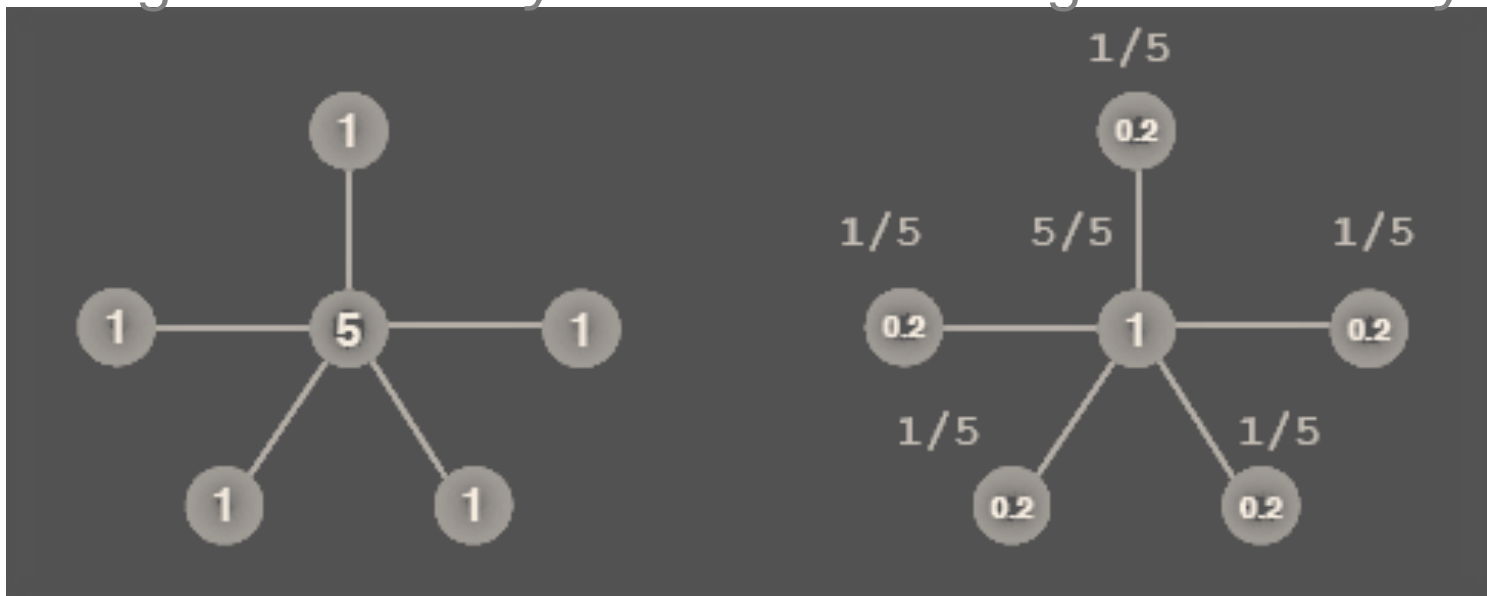
Different definitions of importance different measures of centrality

https://en.wikipedia.org/wiki/File:6_centrality_measures.png

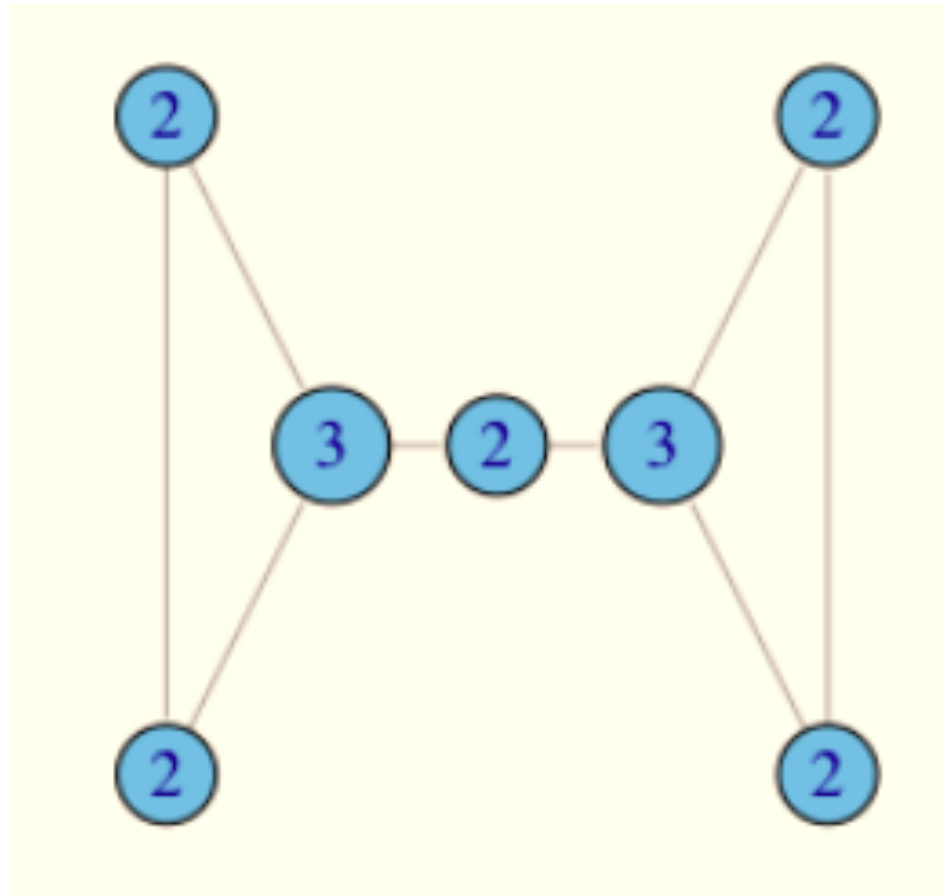
Degree Centrality

$$\text{degree centrality}(n) = \frac{d(n)}{|V| - 1} = \frac{\text{degree of } n}{\text{number of nodes in } G - 1}$$

Degree Centrality & Normalized Degree Centrality



In what way does degree fail to capture centrality here?



Betweenness Centrality

Betweenness centrality can be represented as:

$$C_B(v) = \sum_{s \neq v \neq t \in V} \frac{\sigma_{st}(v)}{\sigma_{st}}$$

σ_{st} is total number of shortest paths from node s to node t

$\sigma_{st}(v)$ is the number of those paths that pass through v .

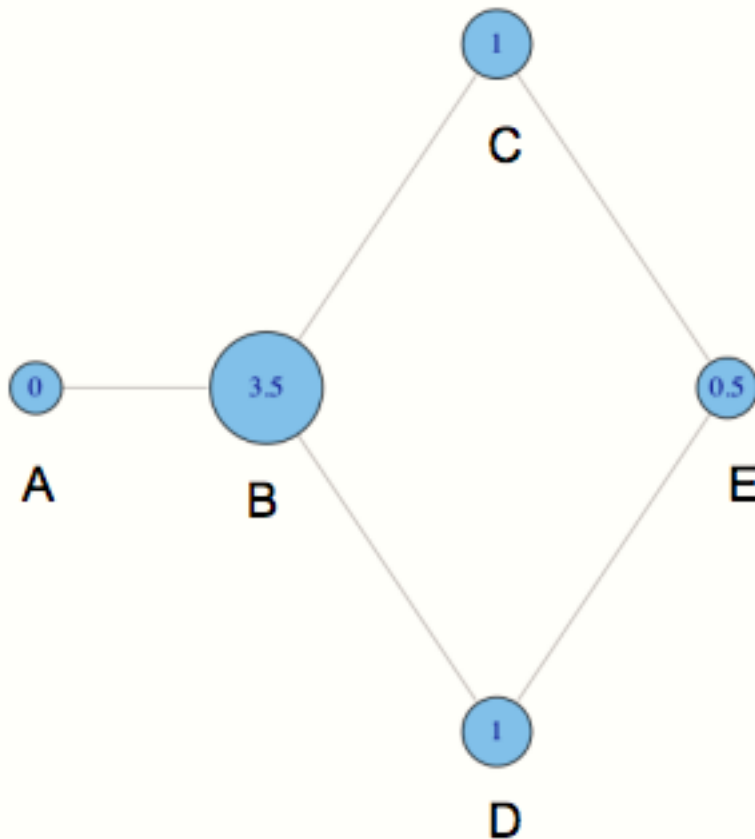
Normalized Betweenness Centrality (undirected graph)

$$C'_B(i) = C_B(i) / [(n-1)(n-2)/2]$$

number of pairs of vertices
excluding the vertex itself

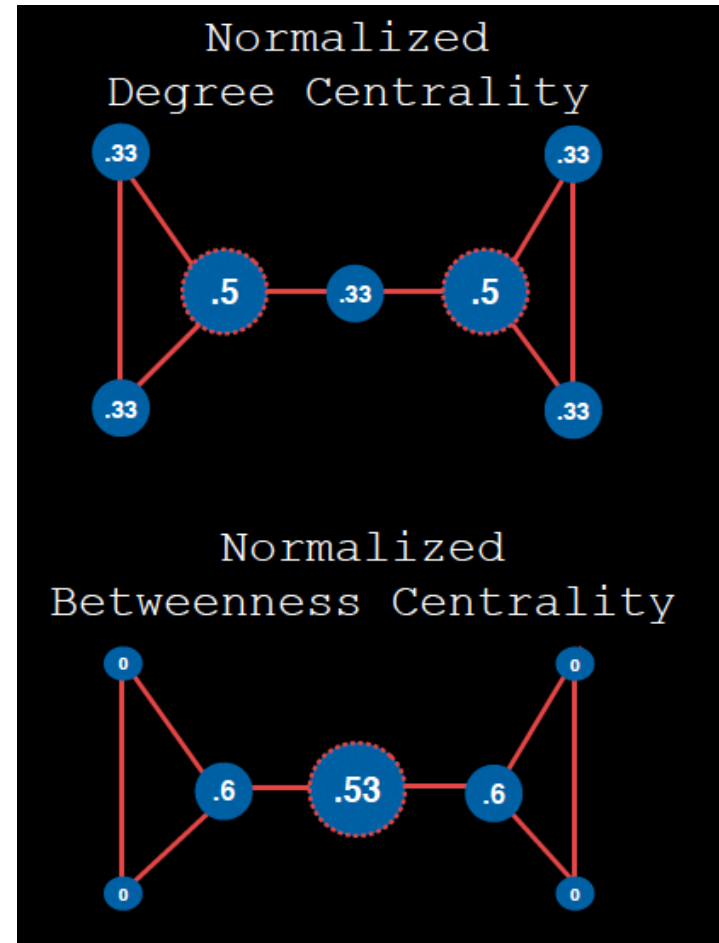
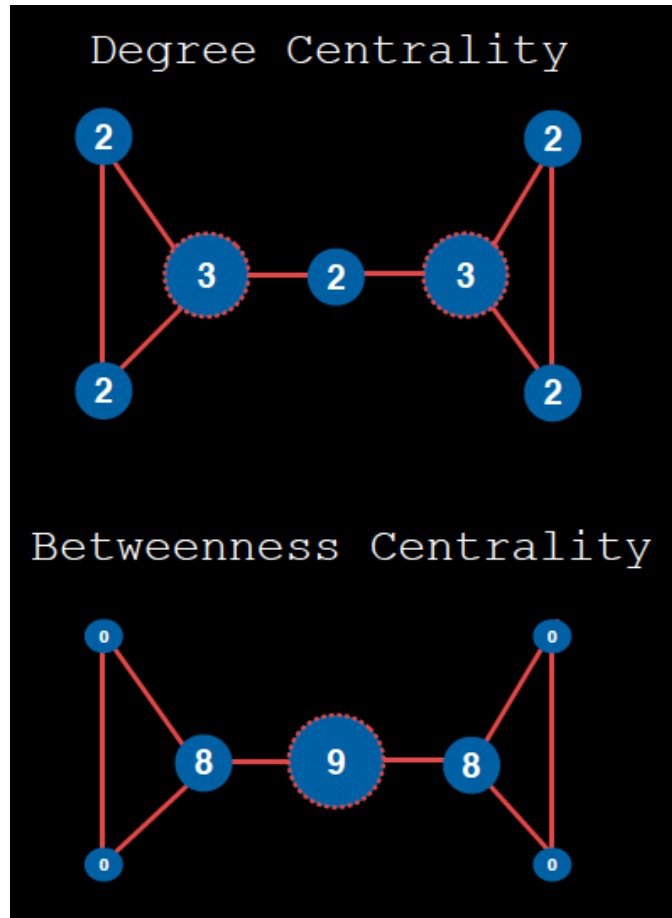
Betweenness Centrality, toy example

- non-normalized version:



- why do C and D each have betweenness 1?
- They are both on shortest paths for pairs (A,E), and (B,E), and so must share credit:
 - $\frac{1}{2} + \frac{1}{2} = 1$
- Can you figure out why B has betweenness 3.5 while E has betweenness 0.5?

In review:



Eigenvector Centrality

Important nodes

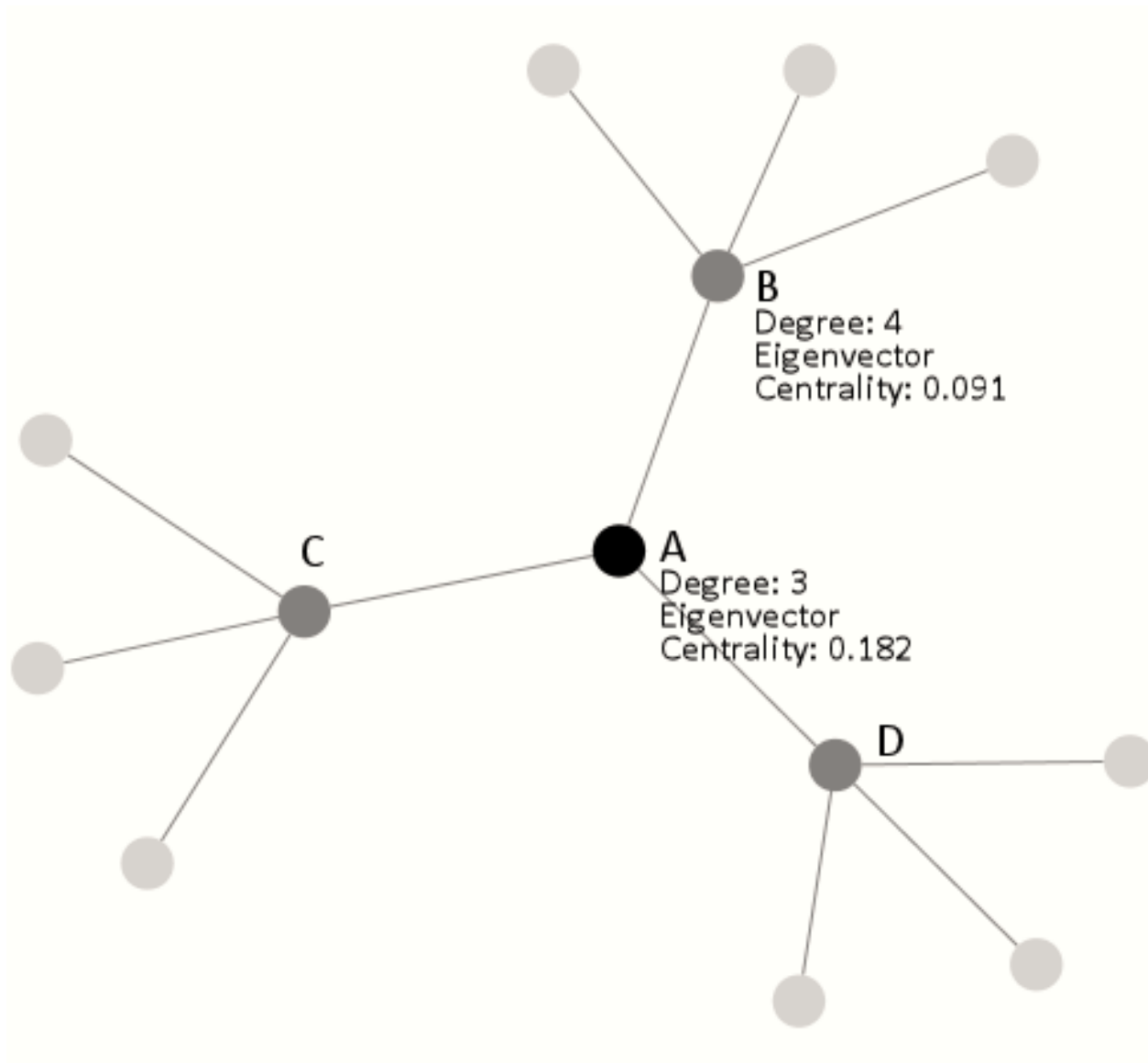
- Have important neighbors
- Connected to node with high degrees
- Themselves do not necessarily have high degree

Tied to the **idea of influence**. The centrality of a given node has is equal to sum of the centrality of my neighbors.

Page rank, a variation on this, was used to rank websites in a search result.

NOTE: We have seen this in an earlier individual sprint with stochastic matrices! An implementation of page rank was extra credit.

Eigenvector Centrality



Eigenvector Centrality

$$x_i = \frac{1}{\lambda} \sum_{j \in G} A_{ij} \cdot x_j$$

neighbors 1/0
A is adjacency matrix

x_i λ A_{ij} x_j

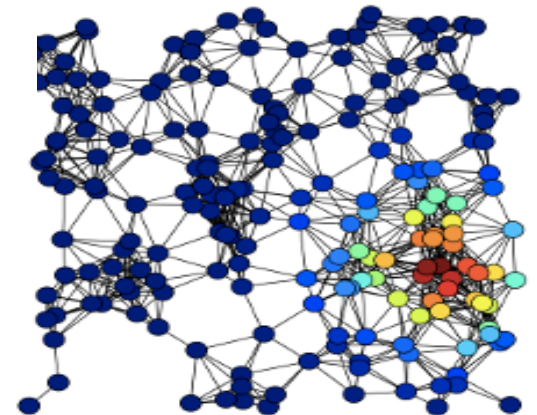
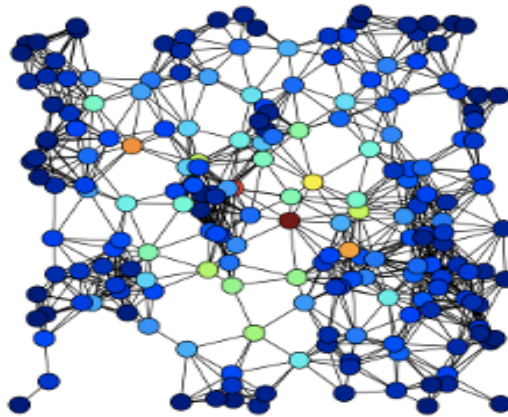
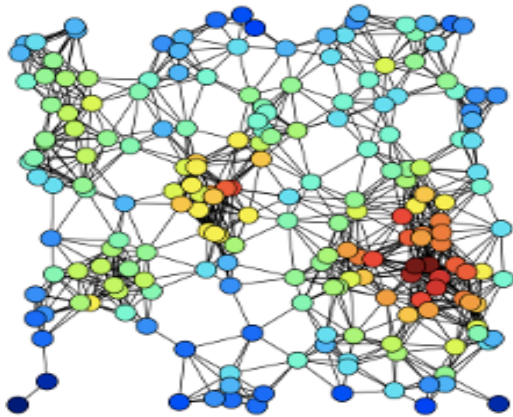
C_E of node i eigenvalue C_E of node j

$$Ax = \lambda x$$

Eigenvector with
 C_E of all nodes

“Centrality” pop quiz

Match the graph to the centrality depicted



Session Summary and SPRINT

- Formal definition of graphs and applications
- Traversing a graph with Breadth First Search
- Importance of an node

SPRINT

- Explore the IMBD dataset of movies and actors
 - Write code to find the shortest path between two actors using BFS. Kevin Bacon Number, anyone?