

Graphs

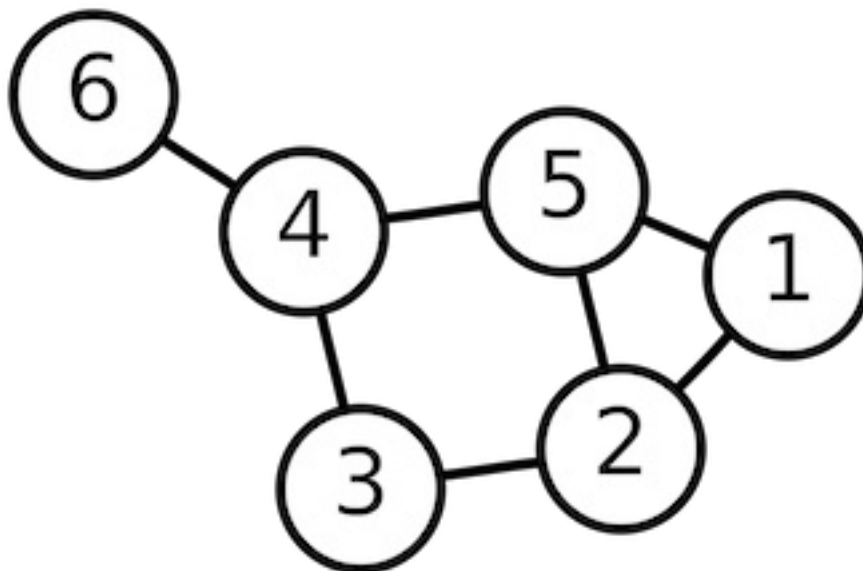
Graphs are a different way of looking at your data. The best way to understand a graph is by looking at some examples. We look at data points as *nodes* and see the connections between them (*edges*).

1. [Example Graphs](#)
2. [Terminology](#)
3. [Representing Graphs](#)
4. [Searching](#)
5. [Centrality](#)
6. [Community Detection](#)

Example Graphs

1. A Social Network (e.g. Facebook)

- A node is a person.
- Two nodes are connected with an edge if they are friends.



Here there are 6 users. These are their friend lists:

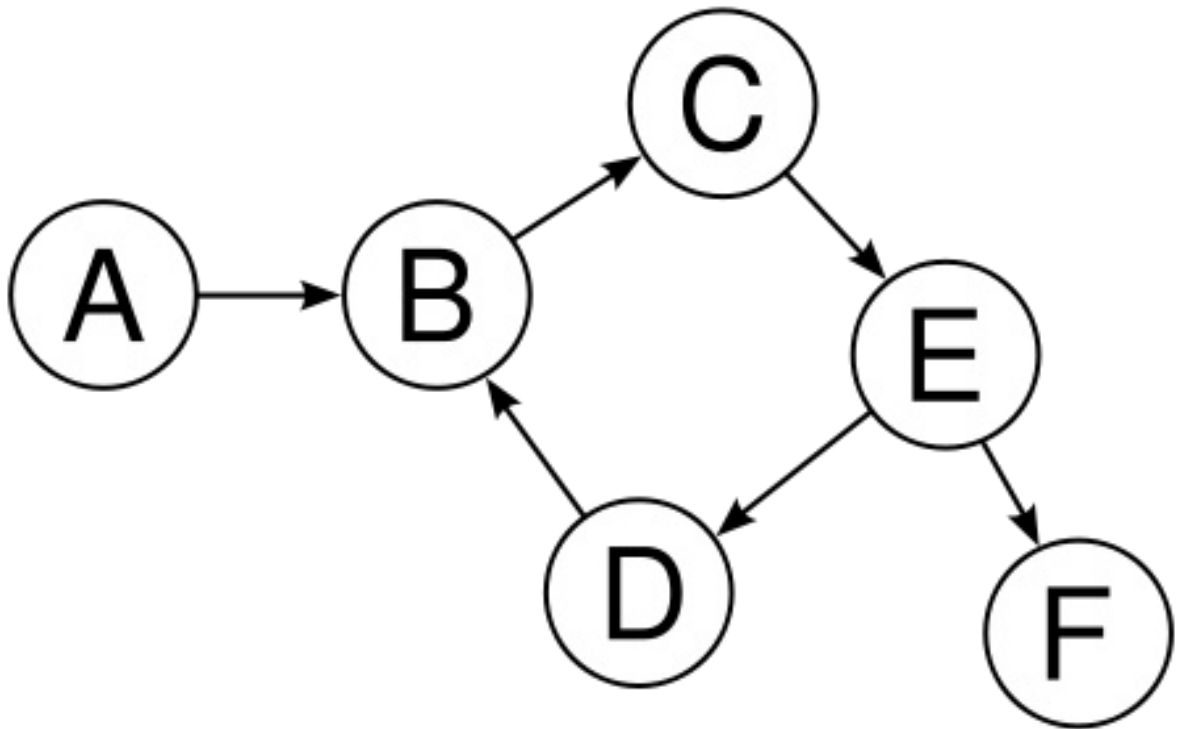
```
1: 2, 5
2: 1, 3, 5
3: 2, 4
4: 5, 6
```

5: 1, 2, 4

6: 4

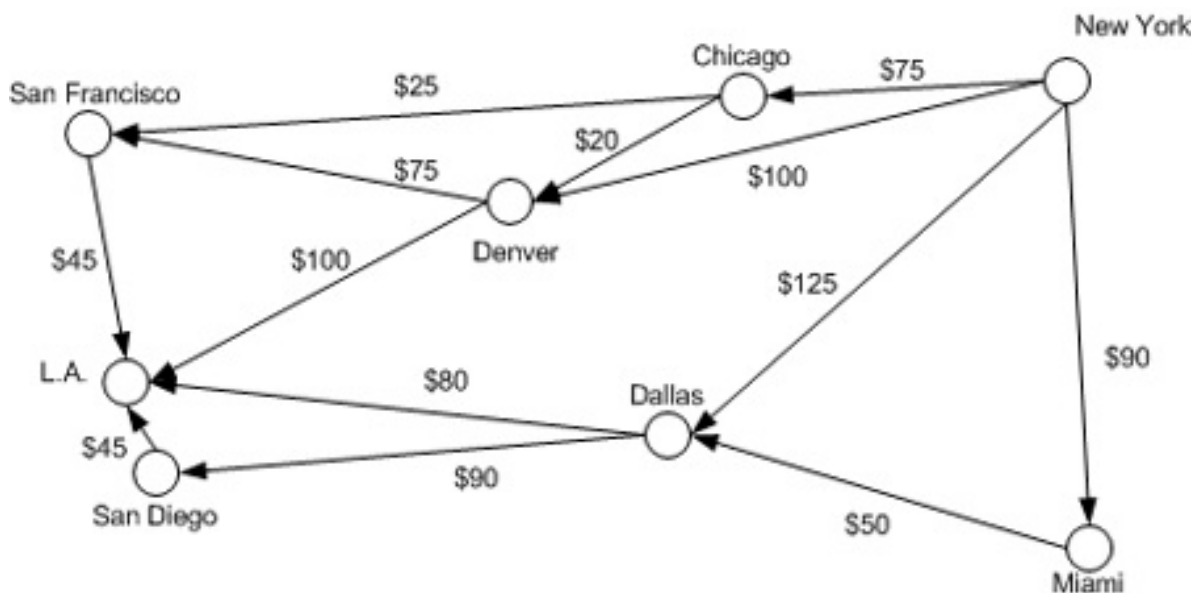
2. A One-Directional Social Network (e.g. Twitter)

- A node is a person.
- There is an edge from person A to person B if person A follows person B



3. Plane Routes

- A node is an airport.
- Two airports are connected if there is a flight between them. The edge is *weighted* according to the cost of the flight.



Terminology

Graphs come with a lot of terminology....

Node (Vertex)	A <i>node</i> (aka <i>vertex</i>) is a datapoint in the graph.
Edge	An <i>edge</i> is a connection between two nodes. If there's an edge between two nodes, we say that those nodes are <i>connected</i> . This can be that two users are friends, that there's a bus line between two cities, a link between two webpages, etc.

Types of Graphs

Directed Graph	A <i>directed graph</i> is a graph where edges only go in one direction. For example, in Twitter, user A can follow user B without user B following user A.
Undirected Graph	In an <i>undirected graph</i> , edges go both ways. For example, with Facebook, if user A is friends with user B, then user B is friends with user A.
Weighted Graph	In a <i>weighted graph</i> , the edges have weights. For example, in a graph of airports and flight routes, the edges can have weights of the cost or

	number of miles.
Unweighted Graph	In an <i>unweighted graph</i> , all edges are the same. For example, with Facebook, there aren't different levels of FB (though you could imagine them existing...)

More Terminology

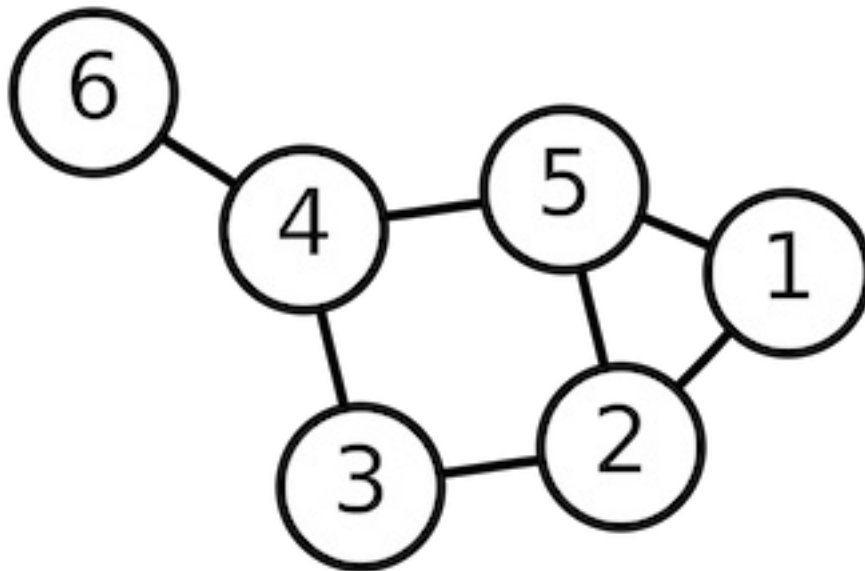
Neighbors	The <i>neighbors</i> of a node are the nodes that it is connected to. e.g., B is a neighbor of A if there is an edge from A to B.
Degree	The <i>degree</i> is the number of neighbors a node has. If the graph is directed, we have separate values for the <i>outdegree</i> (number of edges leaving the node) and <i>indegree</i> (number of edges coming into the node)
Complete	A graph is <i>complete</i> if there is an edge from every node to every other node. This is generally only seen in weighted graphs.
Path	A <i>path</i> is a series of nodes and the edges that connect them. For example, in the airport route example, a path from New York to San Diego is <code>New York, Miami, Dallas, San Diego</code> . Note that there are often multiple paths and this example isn't the <i>shortest</i> path from New York to San Diego.
Connected	A graph is <i>connected</i> if there is a path from every node to every other node.
Subgraph	A <i>subgraph</i> is a subset of the nodes of a graph and all the edges between them.
Connected Component	A <i>connected component</i> is a subgraph that is <i>connected</i> .

Representing Graphs

There are two main ways of storing graphs. Keep in mind the simple operations that are often useful:

- Getting a list of a nodes neighbors
- Determining if two nodes are connected
- Calculating the degree of a node

Take another look at the first example so we can see the representations:



Adjacency Lists

This is the most common way of storing a graph.

For each node, we store a list of its neighbors. In python, we could use a set for this (but we would still probably call it an adjacency list).

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

You can see that all of the above operations are easy! If we have an undirected graph, we do need to make sure that if we remove edge (1, 2), we remove 2 from 1's list as well as removing 1 from 2's list.

Adjacency Matrix

We create an $n \times n$ matrix (where n is the number of nodes). The value in cell i, j is 1 if the nodes are connected and 0 if they aren't connected. If it's a weighted graph, that value would

be the weight.

Here's the adjacency matrix for the above example.

```
      1 2 3 4 5 6
1 [[0 1 0 0 1 0]
2  [1 0 1 0 1 1]
3  [0 1 0 1 0 1]
4  [0 0 1 0 1 1]
5  [1 1 0 1 0 1]
6  [0 0 0 1 0 0]]
```

Searching

Graph Search Algorithms are ways of traversing all of the nodes of a graph.

These are useful for several reasons:

- Finding the shortest path between two nodes.
 - Find out an actor's Kevin Bacon number: You have a Kevin Bacon number of 1 if you've been in a film with him. Your number is 2 if you've been in a film with someone who's been in a film with him, etc...
 - Find out how many clicks it takes to get from cnn.com to galvanize.com.
- Finding an extended network.
 - Find all the users who are my friends of friends of friends).

Breadth First Search

The most common and most useful search algorithm is Breadth First Search (BFS). Here's the general algorithm:

- We start at a node
- We investigate all the neighbors of that node
- We investigate all the neighbors of neighbors
- We investigate all the neighbors of neighbors of neighbors
- etc, etc, ...

If we're trying to find the shortest path from A to B, we stop after we find B.

If we're trying to find an extended network, we stop after how every many iterations we're looking for.

Queues

In order to implement the BFS algorithm, we first need to talk about queues. A *queue* is a data structure where you have these operations:

- Add to end
- Remove from beginning
- See if it's empty

We see queues used in processor queues, since a computer can only run one program at a time, so it needs a queue of the processes that are waiting to run. It will run them in the order that they got in the line.

Note: A python list would not be good for implementing a queue. We can add to the end of the list efficiently, but to remove from the beginning, we need to slide everything over to replace the hole. They are generally implemented with what are called linked lists. Python has an implementation of queues for us, so we can just use that.

BFS Pseudocode

```
function BFS(A, B):
    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
            do whatever we want to do with N (will depend on the application)
            add every neighbor of N to Q
```

Shortest Path Pseudocode

As we mentioned above, one use case of BFS is for finding the shortest path. Here is the pseudocode. Note that along with storing the nodes we've visited, we need to store the length of the path from A to that node.

If we cared about the path and not just the distance, we'd have to store the whole path.

```
function shortest_path(A, B):
    create an empty queue Q
    initialize empty set V (set of visited nodes)
    add (A, 0) to Q
```

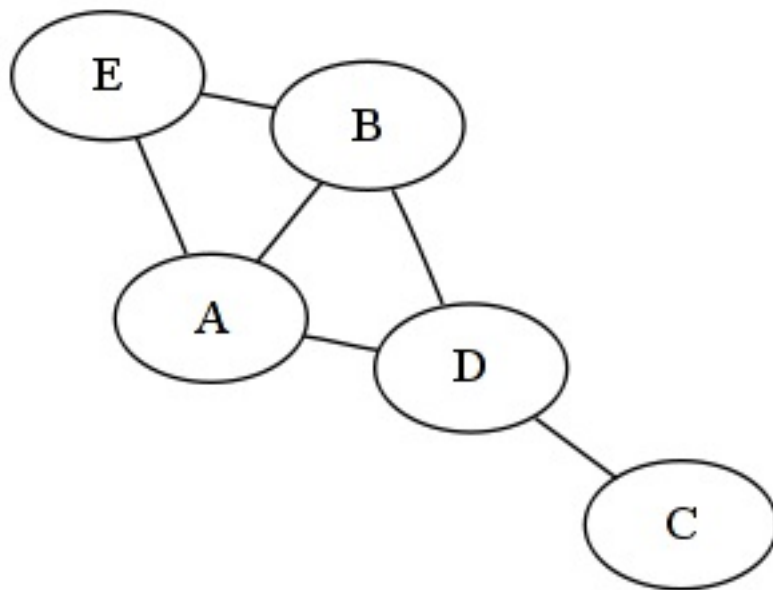
```

while Q is not empty:
    take node off Q, call it N (comes with distance d)
    if N isn't already visited:
        add N to the visited set
        if N is B: we're done!
    else: add every neighbor of N to Q with distance d+1

```

Example

Let's take a look at the following example:



Let's find the shortest path from node A to node C.

At each iteration we can look at the values for `Q`, `v` and `N`

	Current Node (<code>N</code>)	Queue (<code>Q</code>)	Visited (<code>v</code>)	explanation
	<code>A, 0</code>			<i>initialization</i>
1	<code>A, 0</code>			<i>take first node off queue</i>
	<code>A, 0</code>	<code>E, 1 B, 1 D, 1</code>	<code>A</code>	<i>add A's neighbors to queue</i>
2	<code>E, 1</code>	<code>B, 1 D, 1</code>	<code>A</code>	<i>take E off queue</i>
	<code>E, 1</code>	<code>B, 1 D, 1 A, 2 B, 2</code>	<code>A, E</code>	<i>add E's neighbors to queue</i>
3	<code>B, 1</code>	<code>D, 1 A, 2 B, 2</code>	<code>A, E</code>	<i>take B off queue</i>
	<code>B, 1</code>	<code>D, 1 A, 2 B, 2 E, 2 D, 2</code>	<code>A, E, B</code>	<i>add B's neighbors to queue</i>
4	<code>D, 1</code>	<code>A, 2 B, 2 E, 2 D, 2</code>	<code>A, E, B</code>	<i>take D off queue</i>
	<code>D, 1</code>	<code>A, 2 B, 2 E, 2 D, 2 C, 2 B, 2</code>	<code>A, E, B, D</code>	<i>add D's neighbors to queue</i>
5	<code>A, 2</code>	<code>B, 2 E, 2 D, 2 C, 2 B, 2</code>	<code>A, E, B, D</code>	<i>take A off queue</i>
	<code>A, 2</code>	<code>B, 2 E, 2 D, 2 C, 2 B, 2</code>	<code>A, E, B, D</code>	<i>skip A since already visited</i>


```
| 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!! |
```

Centrality

Centrality measures are measures of the importance of nodes. There are several ways that we compute these, but we generally standardize them all to be normalized from 0 to 1.

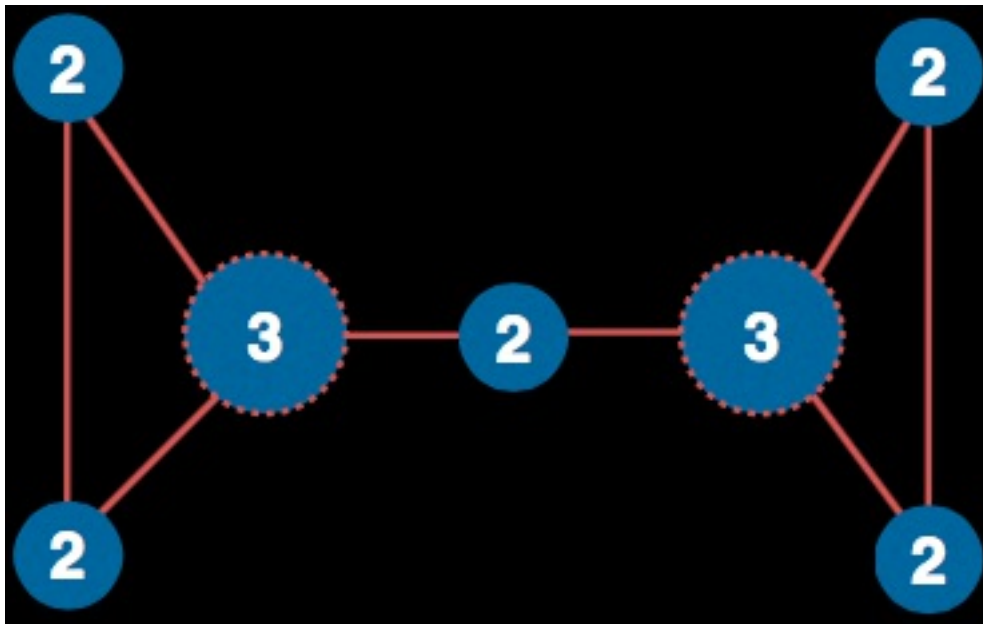
Degree Centrality

The *degree centrality* is the degree of the node divided by a normalizing factor. We divide by $n-1$ since that's the maximum possible degree.

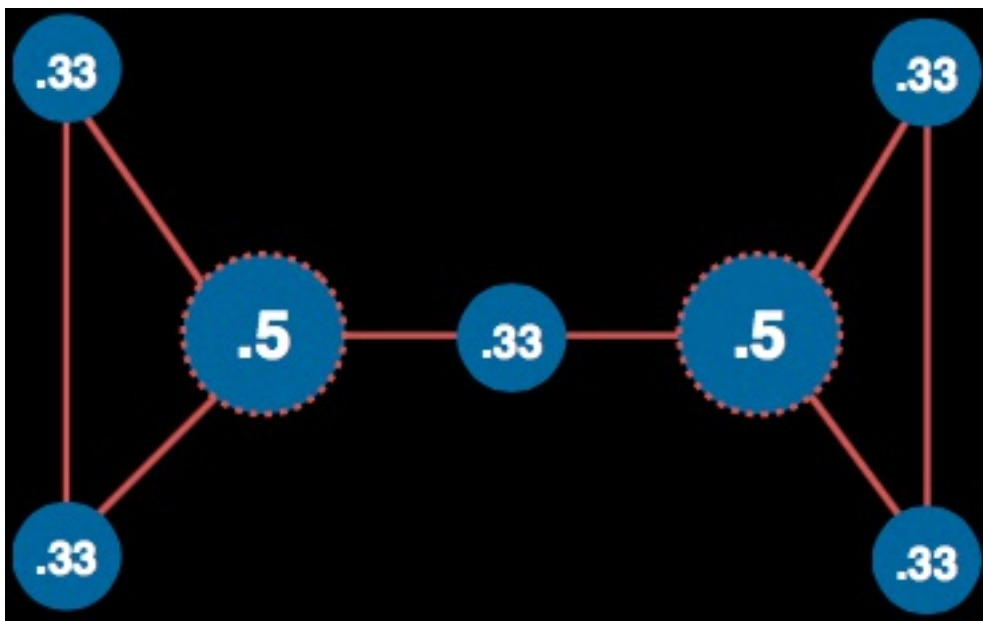
$$\text{degree centrality}(n) = \frac{d(n)}{|V| - 1} = \frac{\text{degree of } n}{\text{number of nodes in } G - 1}$$

In the same example from above, node 4 has a degree of 3 and a degree centrality of $3/(6+1) = 0.429$.

Note: Degree centrality doesn't always capture the important nodes. Take a look at the following example with the degree centralities calculated. You could argue that the middle node is the most important.



Normalized:



Betweenness Centrality

To better represent the importance of a node, we can use the *betweenness centrality*. This is a measure of how many paths the node is a part of.

Formally, here's the definition:

$$\begin{aligned}\text{betweenness}(v) &= \sum_{s \neq v \neq t} \text{percent of shortest paths from } s \text{ to } t \text{ which pass through } v \\ &= \sum_{s \neq v \neq t} \frac{\sigma_{st}(v)}{\sigma_{st}}\end{aligned}$$

where

$\sigma_{st}(v)$ = # of shortest paths from s to t which pass through v

σ_{st} = # of shortest paths from s to t

$$\text{normalized betweenness}(v) = \frac{\text{betweenness}(v)}{(n-1)(n-2)}$$

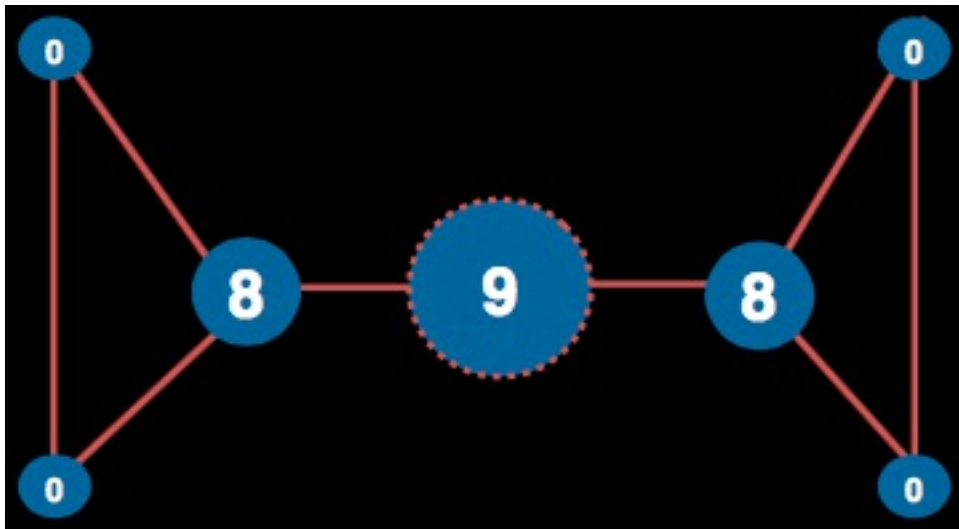
To normalize, we are dividing by $(n-1)(n-2)$ for directed graphs and $(n-1)(n-2)/2$ for undirected graphs, since that's the maximum possible value for the betweenness. This would happen if every single shortest path went through the node.

Let's calculate the betweenness centrality for node 4 above. There are $5 \cdot 4 = 20$ terms, so the 0 terms are omitted here for simplicity.

$$\begin{aligned}\text{betweenness}(4) &= \sum_{s \neq 4 \neq t} \frac{\sigma_{st}(4)}{\sigma_{st}} \\ &= \frac{\sigma_{16}(4)}{\sigma_{16}} + \frac{\sigma_{26}(4)}{\sigma_{26}} + \frac{\sigma_{36}(4)}{\sigma_{36}} + \frac{\sigma_{56}(4)}{\sigma_{56}} + \frac{\sigma_{35}(4)}{\sigma_{35}} \\ &= \frac{1}{1} + \frac{2}{2} + \frac{1}{1} + \frac{1}{1} + \frac{1}{2} \\ &= 4.5\end{aligned}$$

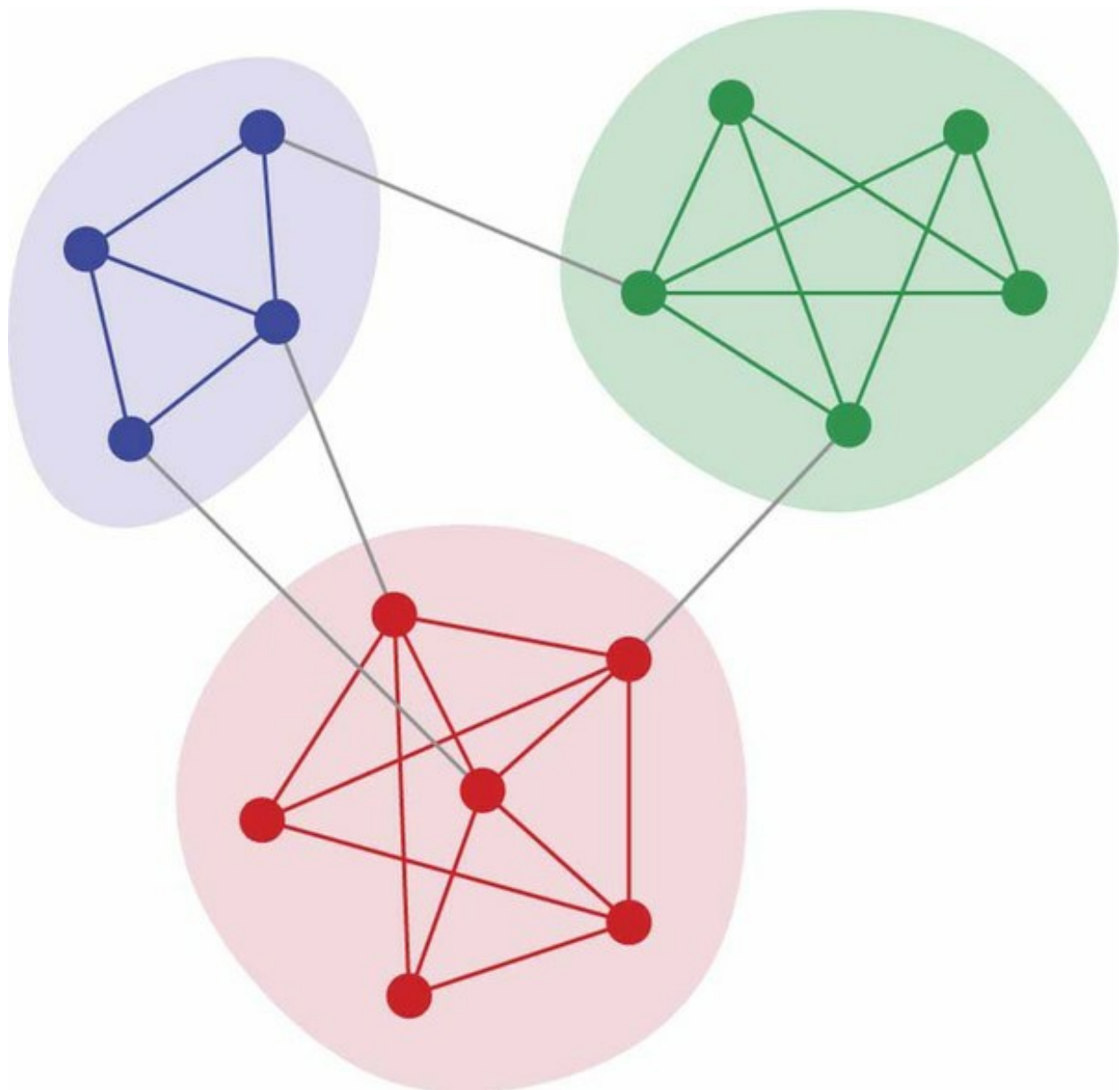
$$\text{normalized betweenness}(4) = \frac{4.5}{(6-1)(6-2)} = 0.225$$

Here's the betweenness centrality for the example from above:



Community Detection

A common task is to determine the *communities* within a graph.



If you think about the Facebook graph, you probably have a network of friends who are mostly friends with each other. This would be a community. While you have friends outside of your community, you have different connections outside of your community.

The idea to communities is that there will be more edges inside the community than outside!

Modularity

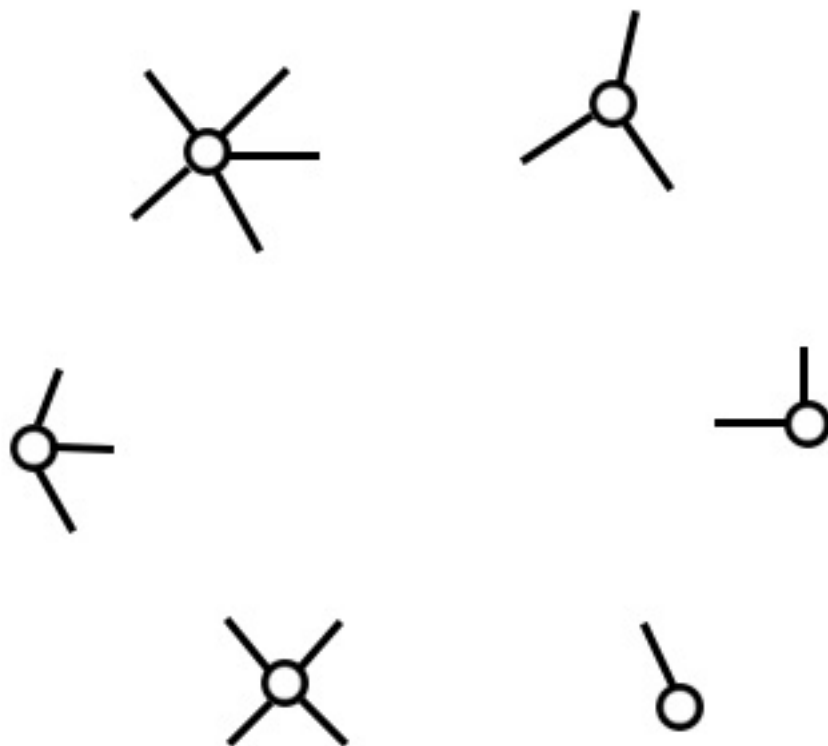
Modularity is a measure of how much of the edges are within the communities versus between communities. The higher the modularity, the better the communities.

The idea is that if we created a random graph where every node maintained its degree, how many edges would there be within the communities (and not between)? How does this

compare to how many edges there actually are within communities.

We can start by calculating the expected number of edges within the communities.

The idea is to imagine cutting every edge in half. The two halves each remain connected to one node and you end up with something that looks like this:



Now randomly connect stubs together to form a graph! We can calculate the probability that node i gets connected to node j .

$$P(\text{single edge stub gets connected to } j) = \frac{d(j)}{2m}$$

$$P(\text{edge from } i \text{ to } j) = d(i) \cdot \frac{d(j)}{2m} = \frac{d(i)d(j)}{2m}$$

where

$d(i)$ = degree of node i

m = number of edges in the graph

The expected number of edges which occur within communities is as follows:

$$E(\text{edges within communities}) = \sum_{i,j \text{ in same community}} \frac{d(i)d(j)}{2m}$$

The modularity is the true number of edges within communities minus the expected number of edges within communities. We also normalize by dividing by $2m$.

$$\text{modularity}(G, \mathcal{C}) = \frac{1}{2m} \left(\sum_{C \in \mathcal{C}} (\# \text{ edges in } C) + \sum_{i,j \in \text{same } C} \frac{d(i)d(j)}{2m} \right)$$

where \mathcal{C} is the collection of communities

We often will see it written like this:

$$\text{modularity}(G, \mathcal{C}) = \frac{1}{2m} \sum_{C \in \mathcal{C}} \sum_{i,j \in C} A_{ij} - \frac{d(i)d(j)}{2m}$$

where

\mathcal{C} = the collection of communities

m = number of edges in the graph

$$A_{ij} = \begin{cases} 1 & \text{if } (i, j) \text{ is an edge} \\ 0 & \text{if } (i, j) \text{ is not an edge} \end{cases}$$

$d(i)$ = degree of node i

Girvan-Newman Algorithm

The *Girvan-Newman Algorithm* is the most common community detection algorithm.

We iteratively remove the edge with the highest *edge betweenness*. The *edge betweenness* is a measure of how many paths an *edge* is part of.

Here's the formal definition for betweenness of an edge e (note that it's essentially the same as the node betweenness defined above).

$$\begin{aligned}\text{betweenness}(e) &= \sum_{s \neq v \neq t} \text{percent of shortest paths from } s \text{ to } t \text{ which pass through } e \\ &= \sum_{s \neq v \neq t} \frac{\sigma_{st}(e)}{\sigma_{st}}\end{aligned}$$

where

$\sigma_{st}(e)$ = # of shortest paths from s to t which pass through e

σ_{st} = # of shortest paths from s to t

$$\text{normalized betweenness}(e) = \frac{\text{betweenness}(e)}{(n-1)(n-2)}$$

Here's the pseudocode for the algorithm:

```
function GirvanNewman:
  repeat:
    repeat until a new connected component is created:
      calculate the edge betweenness centralities for all the edges
      remove the edge with the highest betweenness
```

This will iteratively create new communities. To determine the appropriate number of communities, we calculate the modularity for each set of communities and pick the one with the maximum modularity.