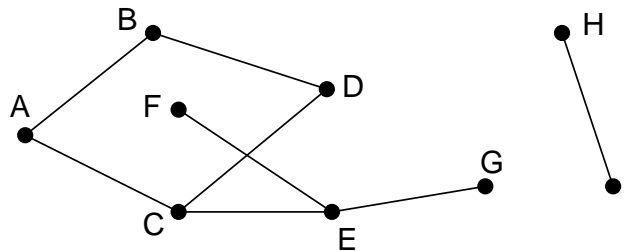


# Chapter 6

## Graph Theory

### 6.1 Introduction

Informally, a *graph* is a bunch of dots connected by lines. Here is an example of a graph:



Sadly, this definition is not precise enough for mathematical discussion. Formally, a graph is a pair of sets  $(V, E)$ , where:

- $V$  is a nonempty set whose elements are called *vertices*.
- $E$  is a collection of two-element subsets of  $V$  called *edges*.

The vertices correspond to the dots in the picture, and the edges correspond to the lines. Thus, the dots-and-lines diagram above is a pictorial representation of the graph  $(V, E)$  where:

$$V = \{A, B, C, D, E, F, G, H, I\}$$
$$E = \{\{A, B\}, \{A, C\}, \{B, D\}, \{C, D\}, \{C, E\}, \{E, F\}, \{E, G\}, \{H, I\}\}.$$

### 6.1.1 Definitions

A nuisance in first learning graph theory is that there are so many definitions. They all correspond to intuitive ideas, but can take a while to absorb. Some ideas have multiple names. For example, graphs are sometimes called *networks*, vertices are sometimes called *nodes*, and edges are sometimes called *arcs*. Even worse, no one can agree on the exact meanings of terms. For example, in our definition, every graph must have at least one vertex. However, other authors permit graphs with no vertices. (The graph with no vertices is the single, stupid counterexample to many would-be theorems— so we're banning it!) This is typical; everyone agrees more-or-less what each term means, but disagrees about weird special cases. So do not be alarmed if definitions here differ subtly from definitions you see elsewhere. Usually, these differences do not matter.

Hereafter, we use  $A-B$  to denote an edge between vertices  $A$  and  $B$  rather than the set notation  $\{A, B\}$ . Note that  $A-B$  and  $B-A$  are the same edge, just as  $\{A, B\}$  and  $\{B, A\}$  are the same set.

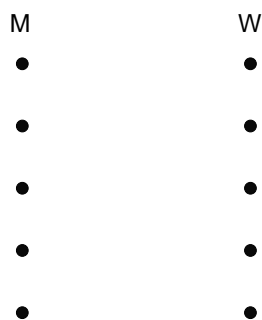
Two vertices in a graph are said to be *adjacent* if they are joined by an edge, and an edge is said to be *incident* to the vertices it joins. The number of edges incident to a vertex is called the *degree* of the vertex. For example, in the graph above,  $A$  is adjacent to  $B$  and  $B$  is adjacent to  $D$ , and the edge  $A-C$  is incident to vertices  $A$  and  $C$ . Vertex  $H$  has degree 1,  $D$  has degree 2, and  $E$  has degree 3.

Deleting some vertices or edges from a graph leaves a *subgraph*. Formally, a subgraph of  $G = (V, E)$  is a graph  $G' = (V', E')$  where  $V'$  is a nonempty subset of  $V$  and  $E'$  is a subset of  $E$ . Since a subgraph is itself a graph, the endpoints of every edge in  $E'$  must be vertices in  $V'$ .

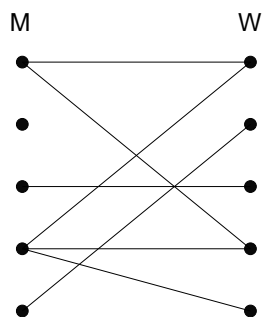
### 6.1.2 Sex in America

A 1994 University of Chicago study entitled *The Social Organization of Sexuality* found that on average men have 74% more opposite-gender partners than women.

Let's recast this observation in graph theoretic terms. Let  $G = (V, E)$  be a graph where the set of vertices  $V$  consists of everyone in America. Now each vertex either represents either a man or a woman, so we can partition  $V$  into two subsets:  $M$ , which contains all the male vertices, and  $W$ , which contains all the female vertices. Let's draw all the  $M$  vertices on the left and the  $W$  vertices on the right:



Now, without getting into a lot of specifics, *sometimes an edge appears* between an  $M$  vertex and a  $W$  vertex:



Since we're only considering opposite-gender relationships, every edge connects an  $M$  vertex on the left to a  $W$  vertex on the right. So the sum of the degrees of the  $M$  vertices must equal the sum of the degrees of the  $W$  vertices:

$$\sum_{x \in M} \deg(x) = \sum_{y \in W} \deg(y)$$

Now suppose we divide both sides of this equation by the product of the sizes of the two sets,  $|M| \cdot |W|$ :

$$\left( \frac{\sum_{x \in M} \deg(x)}{|M|} \right) \cdot \frac{1}{|W|} = \left( \frac{\sum_{y \in W} \deg(y)}{|W|} \right) \cdot \frac{1}{|M|}$$

The terms above in parentheses are the *average degree of an  $M$  vertex* and the *average degree of a  $W$  vertex*. So we know:

$$\frac{\text{Avg. deg in } M}{|W|} = \frac{\text{Avg. deg in } W}{|M|}$$

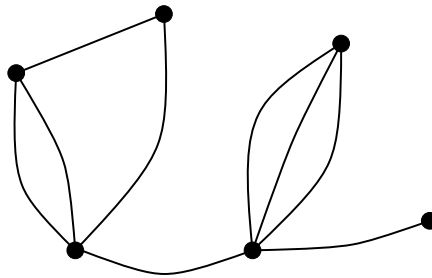
$$\text{Avg. deg in } M = \frac{|W|}{|M|} \cdot \text{Avg. deg in } W$$

Now the Census Bureau reports that there are slightly more women than men in America; in particular  $|W| / |M|$  is about 1.035. So— assuming the Census Bureau is correct—

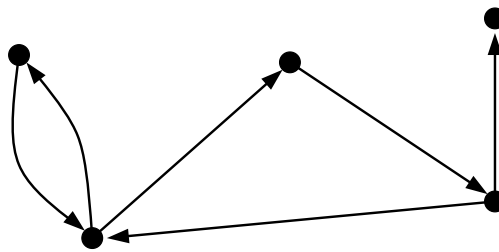
we've just proved that the University of Chicago study got bad data! On average, men have 3.5% more opposite-gender partners. Furthermore, this is totally unaffected by differences in sexual practices between men and women; rather, it is completely determined by the relative number of men and women!

### 6.1.3 Graph Variations

There are many variations on the basic notion of a graph. Three particularly common variations are described below. In a *multigraph*, there may be more than one edge between a pair of vertices. Here is an example:

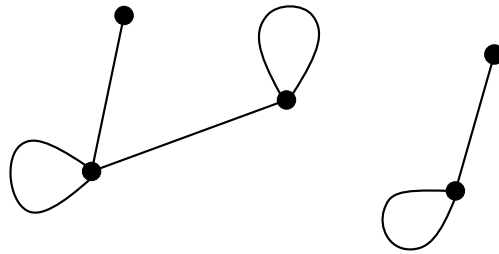


The edges in a *directed graph* are arrows pointing to one endpoint or the other. Here is an example:



Directed graphs are often called *digraphs*. We denote an edge from vertex  $A$  to vertex  $B$  in a digraph by  $A \rightarrow B$ . Formally, the edges in a directed graph are ordered pairs of vertices rather than sets of two vertices. The number of edges directed into a vertex is called the *indegree* of the vertex, and the number of edges directed out is called the *outdegree*.

One can also allow *self-loops*, edges with both endpoints at one vertex. Here is an example of a graph with self-loops:



Combinations of these variations are also possible; for example, one could work with directed multigraphs with self-loops.

*Except where stated otherwise, the word “graph” in this course refers to a graph without multiple edges, directed edges, or self-loops.*

### 6.1.4 Applications of Graphs

Graphs are the most useful mathematical objects in computer science. You can model an enormous number of real-world systems and phenomena using graphs. Once you’ve created such a model, you can tap the vast store of theorems about graphs to gain insight into the system you’re modeling. Here are some practical situations where graphs arise:

**Data Structures** Each vertex represents a data object. There is a directed edge from one object to another if the first contains a pointer or reference to the second.

**Attraction** Each vertex represents a person, and each edge represents a romantic attraction. The graph could be directed to model the unfortunate asymmetries.

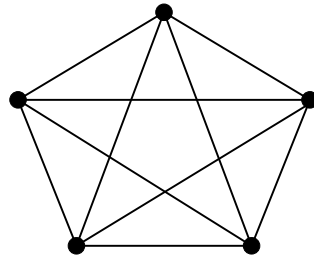
**Airline Connections** Each vertex represents an airport. If there is a direct flight between two airports, then there is an edge between the corresponding vertices. These graphs often appear in airline magazines.

**The Web** Each vertex represents a web page. Directed edges between vertices represent hyperlinks.

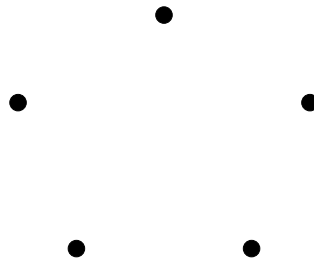
People often put numbers on the edges of a graph, put colors on the vertices, or add other ornaments that capture additional aspects of the phenomenon being modeled. For example, a graph of airline connections might have numbers on the edges to indicate the duration of the corresponding flight. The vertices in the attraction graph might be colored to indicate the person’s gender.

### 6.1.5 Some Common Graphs

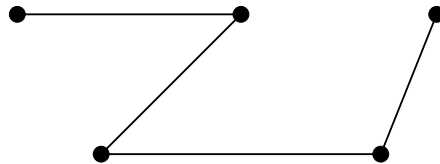
Some graphs come up so frequently that they have names. The *complete graph* on  $n$  vertices, also called  $K_n$ , has an edge between every pair of vertices. Here is  $K_5$ :



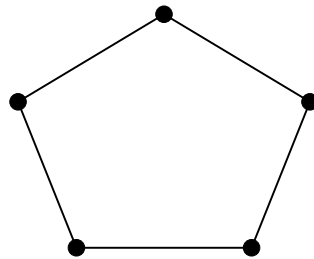
The *empty graph* has no edges at all. Here is the empty graph on 5 vertices:



Here is a *path* with 5 vertices:



And here is a *cycle* with 5 vertices, which is typically denoted  $C_5$ :



Paths and cycles are going to be particularly important, so let's define them precisely. A *path* is a graph  $P = (V, E)$  of the form

$$V = \{v_1, v_2, \dots, v_n\} \quad E = \{v_1-v_2, v_2-v_3, \dots, v_{n-1}-v_n\}$$

where  $n \geq 1$  and vertices  $v_1, \dots, v_n$  are all distinct. Vertices  $v_1$  and  $v_n$  are the *endpoints* of the path. Note that a path may consist of a single vertex, in which case both endpoints are

the same. We'll often say that there is a path from  $u$  to  $v$  in a graph  $G$ ; this is a shorthand for saying that a path with endpoints  $u$  and  $v$  is a subgraph of  $G$ .

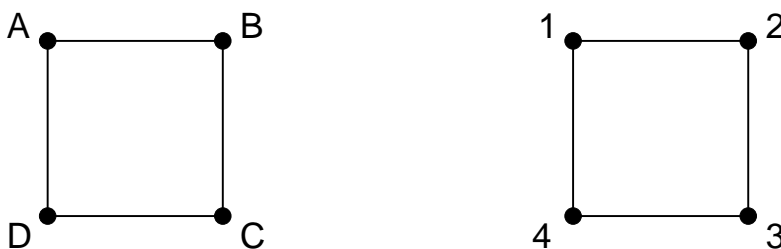
Similarly, a cycle is a graph  $C = (V, E)$  of the form

$$V = \{v_1, v_2, \dots, v_n\} \quad E = \{v_1-v_2, v_2-v_3, \dots, v_{n-1}-v_n, v_n-v_1\}$$

where  $n \geq 3$  and  $v_1, \dots, v_n$  are all distinct. The *length* of a path or cycle is the number of edges it contains. For example, a path with 5 vertices has length 4, but a cycle with 5 vertices has length 5.

### 6.1.6 Isomorphism

Two graphs that look the same might actually be different in a formal sense. For example, the two graphs below are both cycles with 4 vertices:



But one graph has vertex set  $\{A, B, C, D\}$  while the other has vertex set  $\{1, 2, 3, 4\}$ . If so, then the graphs are different mathematical objects, strictly speaking. But this is a frustrating distinction; the graphs *look the same!*

Fortunately, we can neatly capture the idea of “looks the same” and use that as our main notion of equivalence between graphs. Graphs  $G_1$  and  $G_2$  are *isomorphic* if there exists a one-to-one correspondence between vertices in  $G_1$  and vertices in  $G_2$  such that there is an edge between two vertices in  $G_1$  if and only if there is an edge between the two corresponding vertices in  $G_2$ . For example, take the following correspondence between vertices in the two graphs above:

$A$ corresponds to 1	$B$ corresponds to 2
$D$ corresponds to 4	$C$ corresponds to 3.

Now there is an edge between two vertices in the graph on the left if and only if there is an edge between the two corresponding vertices in the graph on the right. Therefore, the two graphs are isomorphic. The correspondence itself is called an *isomorphism*.

Two isomorphic graphs may be drawn to look quite different. For example, here are two different ways of drawing  $C_5$ :



Isomorphic graphs share a great many properties, such as the number of vertices, number of edges, and the pattern of vertex degrees. Thus, two graphs can be proved *nonisomorphic* by identifying some property that one possesses that the other does not. For example, if one graph has two vertices of degree 5 and another has three vertices of degree 5, then the graphs can not be isomorphic.

## 6.2 Connectivity

In the diagram below, the graph on the left has two pieces, while the graph on the right has just one.



Let's put this observation in rigorous terms. A graph is *connected* if for every pair of vertices  $u$  and  $v$ , the graph contains a path with endpoints  $u$  and  $v$  as a subgraph. The graph on the left is not connected because there is no path from any of the top three vertices to either of the bottom two vertices. However, the graph on the right is connected, because there is a path between every pair of vertices.

A maximal, connected subgraph is called a *connected component*. (By "maximal", we mean that including any additional vertices would make the subgraph disconnected.) The graph on the left has two connected components, the triangle and the single edge. The graph on the right is entirely connected and thus has a single connected component.

### 6.2.1 A Simple Connectivity Theorem

The following theorem says that a graph with few edges must have many connected components.



**Theorem 47.** *Every graph  $G = (V, E)$  has at least  $|V| - |E|$  connected components.*

*Proof.* We use induction on the number of edges. Let  $P(n)$  be the proposition that every graph  $G = (V, E)$  with  $|E| = n$  has at least  $|V| - n$  connected components.

*Base case:* In a graph with 0 edges, each vertex is itself a connected component, and so there are exactly  $|V| - 0 = |V|$  connected components.

*Inductive step:* Now we assume that the induction hypothesis holds for every  $n$ -edge graph in order to prove that it holds for every  $(n + 1)$ -edge graph, where  $n \geq 0$ . Consider a graph  $G = (V, E)$  with  $n + 1$  edges. Remove an arbitrary edge  $u-v$  and call the resulting graph  $G'$ . By the induction assumption,  $G'$  has at least  $|V| - n$  connected components. Now add back the edge  $u-v$  to obtain the original graph  $G$ . If  $u$  and  $v$  were in the same connected component of  $G'$ , then  $G$  has the same number of connected components as  $G'$ , which is at least  $|V| - n$ . Otherwise, if  $u$  and  $v$  were in different connected components of  $G'$ , then these two components are merged into one in  $G$ , but all other components remain. Therefore,  $G$  has at least  $|V| - n - 1 = |V| - (n + 1)$  connected components.

The theorem follows by induction. □

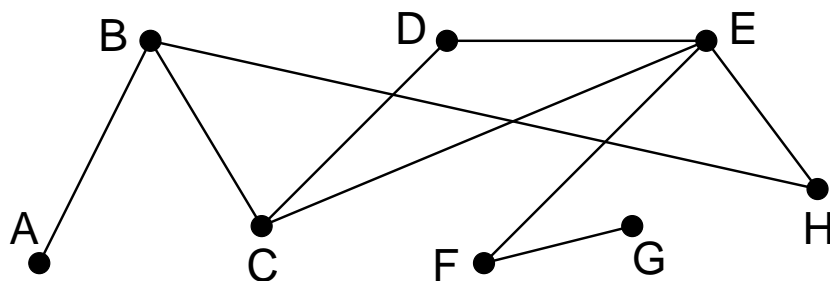
**Corollary 48.** *Every connected graph with  $n$  vertices has at least  $n - 1$  edges.*

A couple points about the proof of Theorem 47 are worth noting. First, notice that we used induction on the number of edges in the graph. This is very common in proofs involving graphs, and so is induction on the number of vertices. When you're presented with a graph problem, these two approaches should be among the first you consider. Don't try induction on other variables that crop up in the problem unless these two strategies seem hopeless.

The second point is more subtle. Notice that in the inductive step, we took an arbitrary  $(n + 1)$ -edge graph, threw out an edge so that we could apply the induction assumption, and then put the edge back. You'll see this shrink-down, grow-back process very often in the inductive steps of proofs related to graphs. This might seem like needless effort; why not start with an  $n$ -edge graph and add one more to get an  $(n + 1)$ -edge graph? That would work fine in this case, but opens the door to a very nasty logical error in similar arguments. (You'll see an example in recitation.) Always use shrink-down, grow-back arguments, and you'll never fall into this trap.

## 6.2.2 Distance and Diameter

The *distance* between two vertices in a graph is the length of the shortest path between them. For example, the distance between two vertices in a graph of airline connections is the minimum number of flights required to travel between two cities.



In this graph, the distance between  $C$  and  $H$  is 2, the distance between  $G$  and  $C$  is 3, and the distance between  $A$  and itself is 0. If there is *no* path between two vertices, then the distance between them is said to be “infinity”.

The *diameter* of a graph is the distance between the two vertices that are farthest apart. The diameter of the graph above is 5. The most-distant vertices are  $A$  and  $G$ , which are at distance 5 from one another.

### Six Degrees of Separation

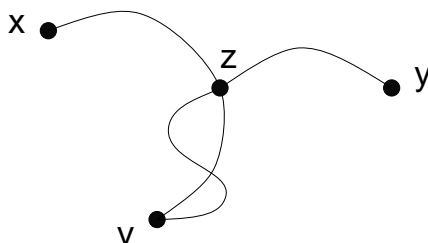
There is an old claim that the world has only “six degrees of separation”. In other words, if you pick any two people on the planet—say a hermit in Montana and a random person off the street in Beijing—then the hermit knows someone who knows someone who knows someone who knows the Chinese pedestrian, where the word “knows” appears at most six times.

We can recast this in graph-theoretic terms. Consider a graph where the vertices are all the people on the planet, and there is an edge between two people if and only if they know each other. Then the “six degrees of separation” claim amounts to the assertion that the diameter of this graph is at most 6.

There is little hope of proving or disproving the claim, since people are constantly being born, meeting one another, and dying and no one can keep track of who-knows-who. However, precise data does exist for something similar. The University of Virginia maintains the *Oracle of Bacon* website. This is based on an “acting graph” where the vertices are actors and actresses, and there is an edge between two performers if they appeared in a movie together. The website reports that everyone is within distance 8 of Kevin Bacon. (This excludes a few actors who are completely disconnected.) This allows us to at least obtain an upper bound on the diameter of the acting graph.

**Theorem 49.** *Let  $v$  be an arbitrary vertex in a graph  $G$ . If every vertex is within distance  $d$  of  $v$ , then the diameter of the graph is at most  $2d$ .*

*Proof.* Let  $x$  and  $y$  be arbitrary vertices in the graph. Then there exists a path of length at most  $d$  from  $x$  to  $v$ , and there exists a path of length at most  $d$  from  $v$  to  $y$ .



Let  $z$  be the vertex that lies on both the  $x$ -to- $v$  and  $v$ -to- $y$  paths and is closest to  $x$ . (We know that such a vertex exists, since  $z$  could be  $v$ , at least.) Joining the  $x$ -to- $z$  segment to the  $z$ -to- $y$  segment gives a path from  $x$  to  $y$  of length at most  $2d$ . Therefore, every vertex is within distance  $2d$  of every other.  $\square$

Data elsewhere on the Oracle of Bacon site shows that the diameter of the acting graph is at least 15, so the upper bound isn't far off.

### 6.2.3 Walks

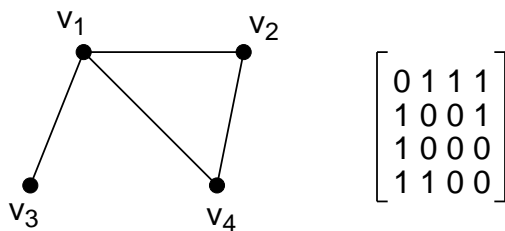
A *walk* in a graph  $G$  is an alternating sequence of vertices and edges of the form:

$$v_0 \ v_0 \text{---} v_1 \ v_1 \text{---} v_2 \ v_2 \ \dots \ v_{n-1} \ v_{n-1} \text{---} v_n \ v_n$$

If  $v_0 = v_n$ , then the walk is *closed*. Walks are similar to paths. However, a walk can cross itself, traverse the same edge multiple times, etc. There is a walk between two vertices if and only if there is a path between the vertices.

## 6.3 Adjacency Matrices

A graph can be represented by an *adjacency matrix*. In particular, if a graph has vertices  $v_1, \dots, v_n$ , then the adjacency matrix is  $n \times n$ . The entry in row  $i$ , column  $j$  is 1 if there is an edge  $v_i \text{---} v_j$  and is 0 if there is no such edge. For example, here is a graph and its adjacency matrix:



The adjacency matrix of an undirected graph is always symmetric about the diagonal line running from the upper left entry to the lower right. The adjacency matrix of a directed graph need not be symmetric, however. Entries on the diagonal of an adjacency matrix are nonzero only if the graph contains self-loops.

Adjacency matrices are useful for two reasons. First, they provide one way to represent a graph in computer memory. Second, by mapping graphs to the world of matrices, one can bring all the machinery of linear algebra to bear on the study of graphs. For example, one can analyze a highly-prized quality of graphs called “expansion” by looking at eigenvalues of the adjacency matrix. (In a graph with good expansion, the number of edges departing each subset of vertices is at least proportional to the size of the subset. This is not so easy to achieve when the graph as a whole has few edges, say  $|E| = 3|V|$ .) Here we prove a simpler theorem in this vein. If  $M$  is a matrix, then  $M_{ij}$  denotes the entry in row  $i$ , column  $j$ . Let  $M^k$  denote the  $k$ -th power of  $M$ . As a special case,  $M^0$  is the identity matrix.

**Theorem 50.** *Let  $G$  be a digraph (possibly with self-loops) with vertices  $v_1, \dots, v_n$ . Let  $M$  be the adjacency matrix of  $G$ . Then  $M_{ij}^k$  is equal to the number of length- $k$  walks from  $v_i$  to  $v_j$ .*

*Proof.* We use induction on  $k$ . The induction hypothesis is that  $M_{ij}^k$  is equal to the number of length- $k$  walks from  $v_i$  to  $v_j$ , for all  $i, j$ .

Each vertex has a length-0 walk only to itself. Since  $M_{ij}^0 = 1$  if and only if  $i = j$ , the hypothesis holds for  $k = 0$ .

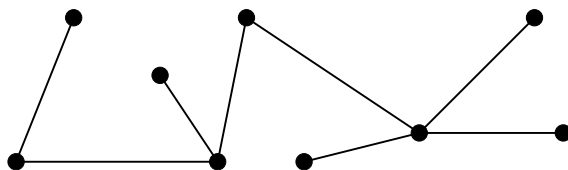
Now suppose that the hypothesis holds for some  $k \geq 0$ . We prove that it also holds for  $k + 1$ . Every length- $(k + 1)$  walk from  $v_i$  to  $v_j$  consists of a length  $k$  walk from  $v_i$  to some intermediate vertex  $v_m$  followed by an edge  $v_m \rightarrow v_j$ . Thus, the number of length- $(k + 1)$  walks from  $v_i$  to  $v_j$  is equal to:

$$M_{iv_1}^k M_{v_1j} + M_{iv_2}^k M_{v_2j} + \dots + M_{iv_n}^k M_{v_nj}$$

This is precisely the value of  $M_{ij}^{k+1}$ , so the hypothesis holds for  $k + 1$  as well. The theorem follows by induction.  $\square$

## 6.4 Trees

A connected, acyclic graph is called a *tree*. (A graph is *acyclic* if no subgraph is a cycle.) Here is an example of a tree:



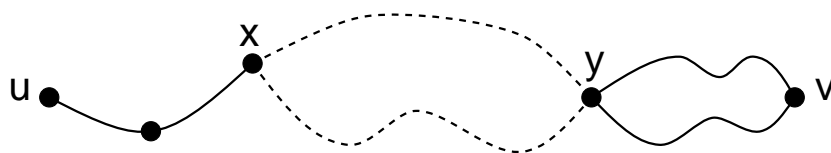
A vertex of degree one is called a *leaf*. In this example, there are 5 leaves.

The graph shown above would no longer be a tree if any edge were removed, because it would no longer be connected. The graph would also not remain a tree if any edge were added between two of its vertices, because then it would contain a cycle. Furthermore, note that there is a unique path between every pair of vertices. These features of the example tree are actually common to all trees.

**Theorem 51.** Every tree  $T = (V, E)$  has the following properties:

1. There is a unique path between every pair of vertices.
2. Adding any edge creates a cycle.
3. Removing any edge disconnects the graph.
4. Every tree with at least two vertices has at least two leaves.
5.  $|V| = |E| + 1$ .

*Proof.* 1. There is at least one path between every pair of vertices, because the graph is connected. Suppose that there are two different paths between vertices  $u$  and  $v$ . Beginning at  $u$ , let  $x$  be the first vertex where the paths diverge, and let  $y$  be the next vertex they share. Then there are two paths from  $x$  to  $y$  with no common edges, which defines a cycle. This is a contradiction, since trees are acyclic. Therefore, there is exactly one path between every pair of vertices.



2. An additional edge  $u-v$  together with the unique path between  $u$  and  $v$  forms a cycle.
3. Suppose that we remove edge  $u-v$ . Since a tree contained a unique path between  $u$  and  $v$ , that path must have been  $u-v$ . Therefore, when that edge is removed, no path remains, and so the graph is not connected.
4. Let  $v_1, \dots, v_m$  be the sequence of vertices on a longest path in  $T$ . Then  $m \geq 2$ , since a tree with two vertices must contain at least one edge. There can not be an edge  $v_1-v_i$  for  $2 < i \leq m$ ; otherwise, vertices  $v_1, \dots, v_i$  would form a cycle. Furthermore, there can not be an edge  $u-v_1$  where  $u$  is not on the path; otherwise, we could make the path longer. Therefore, the only edge incident to  $v_1$  is  $v_1-v_2$ , which means that  $v_1$  is a leaf. By a symmetric argument,  $v_m$  is a second leaf.

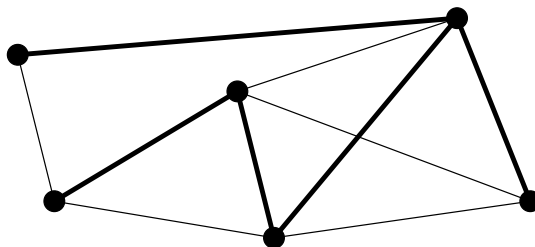
5. We use induction on  $|V|$ . For a tree with a single vertex, the claim holds since  $|E| + 1 = 0 + 1 = 1$ . Now suppose that the claim holds for all  $n$ -vertex trees and consider an  $(n + 1)$ -vertex tree  $T$ . Let  $v$  be a leaf of the tree. Deleting  $v$  and its incident edge gives a smaller tree for which the equation  $|V| = |E| + 1$  holds by induction. If we add back the vertex  $v$  and its incident edge, then the equation still holds because the number of vertices and number of edges both increased by 1. Thus, the claim holds for  $T$  and, by induction, for all trees.

□

Many subsets of the properties above, together with connectedness and lack of cycles, are sufficient to characterize all trees. For example, a connected graph that satisfies  $|V| = |E| + 1$  is necessarily a tree, though we won't prove this fact.

### 6.4.1 Spanning Trees

Trees are everywhere. In fact, every connected graph  $G = (V, E)$  contains a *spanning tree*  $T = (V, E')$  as a subgraph. (Note that original graph  $G$  and the spanning tree  $T$  have the same set of vertices.) For example, here is a connected graph with a spanning tree highlighted.



**Theorem 52.** Every connected graph  $G = (V, E)$  contains a spanning tree.

*Proof.* Let  $T = (V, E')$  be a connected subgraph of  $G$  with the smallest number of edges. We show that  $T$  is acyclic by contradiction. So suppose that  $T$  has a cycle with the following edges:

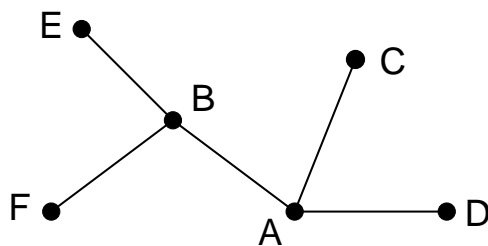
$$v_0—v_1, v_1—v_2, \dots, v_n—v_0$$

Suppose that we remove the last edge,  $v_n—v_0$ . If a pair of vertices  $x$  and  $y$  was joined by a path not containing  $v_n—v_0$ , then they remain joined by that path. On the other hand, if  $x$  and  $y$  were joined by a path containing  $v_n—v_0$ , then they remain joined by a path containing the remainder of the cycle. This is a contradiction, since  $T$  was defined to be a connected subgraph of  $G$  with the smallest number of edges. Therefore,  $T$  is acyclic. □

### 6.4.2 Tree Variations

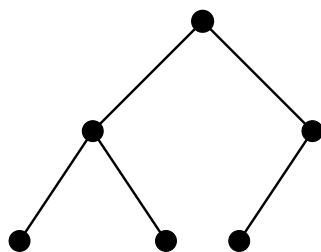
Trees come up often in computer science. For example, information is often stored in tree-like data structures and the execution of many recursive programs can be regarded as a traversal of a tree.

There are many varieties of trees. For example, a *rooted tree* is a tree with one vertex identified as the *root*. Let  $u-v$  be an edge in a rooted tree such that  $u$  is closer to the root than  $v$ . Then  $u$  is the *parent* of  $v$ , and  $v$  is a *child* of  $u$ .



In the tree above, suppose that we regard vertex  $A$  as the root. Then  $E$  and  $F$  are the children of  $B$ , and  $A$  is the parent of  $B$ ,  $C$ , and  $D$ .

A *binary tree* is a rooted tree in which every vertex has at most two children. Here is an example, where the topmost vertex is the root.



In an *ordered, binary tree*, the children of a vertex  $v$  are distinguished. One is called the *left child* of  $v$ , and the other is called the *right child*. For example, if we regard the two binary trees below as unordered, then they are equivalent. However, if we regard these trees as ordered, then they are different.

