



In which our heroes explore the many twisting paths through the gnarled forest, emerging in the happy and peaceful land in which their computational adventures will continue.

11.1 Why You Might Care

Oh what a tangled web we weave,
When first we practise to deceive!

Sir Walter Scott (1771–1832), *Marmion* (1808)

It's possible to make graphs sound hopelessly abstract and utterly uninteresting: *a graph is a pair $\langle V, E \rangle$, where V is a nonempty collection of entities called nodes and E is a collection of edges that join pairs of nodes*. But graphs are fascinating—at least, when the entities and the relationship represented by the edges are themselves interesting! Here are a few of the many examples of types of graphs:

- *social networks* like Facebook (or LinkedIn or Pinterest or ...): the nodes are people, and an edge between two people represents a friendship (or at least a “friendship”).
- *the world-wide web*: the nodes are web pages, and an edge represents a hyperlink from one page to another. These hyperlinks between pages form the basis for the ranking of web pages by search engines like Google.¹
- *dating networks*: nodes represent people; an edge connects two people who have been involved in a romantic relationship. These networks have implications for the spread of certain communicable diseases, particularly sexually transmitted infections.
- *road networks* and other transportation networks: edges represent roads; nodes represent intersections. For example, United Parcel Service (UPS) saves gas (and money!) by using a route-finding algorithm through this network that avoid turns across traffic.²
- *food webs*: nodes represent species within a particular ecosystem, and an edge from one species to another indicates that the first species preys on the latter.
- *co-purchase networks*: nodes are products that are sold by a retailer like Walmart or Amazon; an edge between two products indicates the number of customers who bought both products. These networks have implications for *recommender systems*, the “people who bought x also bought y ” feature of Amazon.
- *the internet*: nodes are computers (personal computers, servers, and other networking hardware like routers), and edges represent physical wires connecting two machines together. When you request a video from `youtube.com`, the computers involved in the network must collectively construct a path along which YouTube's bits can flow so that they reach your computer.

¹ Sergei Brin and Larry Page. The anatomy of a large-scale hypertextual web search engine. In *7th International World-Wide Web Conference*, 1998.

² Joel Lovell. Left-hand-turn elimination. *The New York Times*, 9 December 2007.

Graphs are ubiquitous. Indeed, any pairwise relationship among entities is really underlyingly a graph: web pages and links, computers and fiber optic cables, kidney patients/donors and compatibility for transplants. The applications are innumerable, and this chapter will barely scratch the surface. Graphs and graph-theoretic reasoning will arise again and again well beyond the end of this book.

11.2 Formal Introduction

The Bible tells us to love our neighbors, and also to love our enemies; probably because they are generally the same people.

G. K. Chesterton (1874–1936)

We begin by defining the terminology for the two different basic types of graphs. In both, we have a set of entities called *nodes*, some pairs of which are joined by a relationship called an *edge*. (A node can also be called a *vertex*.) The two types of graph differ in whether the relationship represented by an edge is “between two nodes” or “from one node to another.” In an *undirected graph*, the relationship denoted by the edges is symmetric (for example, “*u* and *v* are genetically related”):

Definition 11.1 (Undirected Graph)

A undirected graph is a pair $G = \langle V, E \rangle$ where V is a nonempty set of vertices or nodes, and $E \subseteq \{ \{u, v\} : u, v \in V \}$ is a set of edges joining pairs of vertices.

vertex, n.: a node.
plural: *vertices*.

The second basic kind of graph is a *directed graph*, in which the relationship denoted by the edges need not be reciprocated (for example, “*u* has texted *v*”):

Definition 11.2 (Directed Graph)

A directed graph is a pair $G = \langle V, E \rangle$ where V is a nonempty set of nodes, and $E \subseteq V \times V$ is a set of edges joining (ordered) pairs of vertices.

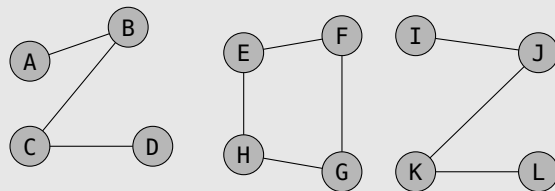
We will use the terms *node/nodes* and *vertex/vertices* interchangeably throughout this chapter. (Both terms are used commonly in CS.) A graph can also be called a *network*; edges are also sometimes called *links*, or occasionally *arcs* in directed graphs.

In other words, in a directed graph an edge is an *ordered* pair of vertices (“an edge from *u* to *v*”) and in an undirected graph an edge is an *unordered* pair of vertices (“an edge between *u* and *v*”). Think about the difference between Twitter followers (directed) and Facebook friendships (undirected): Alice can follow Bob without Bob following Alice, but they’re either friends or they’re not friends.

Graphs are generally drawn with nodes represented as circles, and edges represented by lines. Each edge in directed graphs is drawn with an arrow indicating its *orientation* (“which way it goes”). Here is an example of each:

Example 11.1 (A sample undirected graph)

Here is an undirected graph:

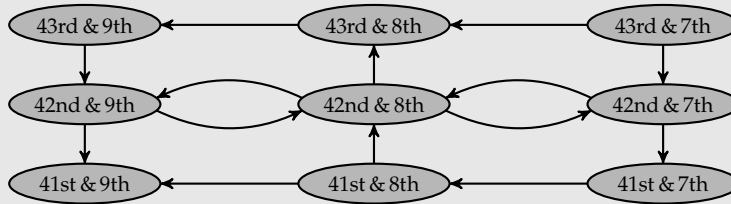


This graph contains:

- 12 nodes: $\{A, B, C, D, E, F, G, H, I, J, K, L\}$.
- 10 edges: $\{ \{A, B\}, \{B, C\}, \{C, D\}, \{E, F\}, \{E, H\}, \{F, G\}, \{G, H\}, \{I, J\}, \{J, K\}, \{K, L\} \}$.

Example 11.2 (Streets of Manhattan: a sample directed graph)

The following directed graph contains 9 nodes, each corresponding to an intersection of a “street” running east–west and an “avenue” running north–south in Manhattan:



There are 14 edges in this graph. There’s something potentially tricky in counting to 14: edges in a directed graph are *ordered* pairs, so there are *two* edges between 42nd & 9th and 42nd & 8th, one in each direction— $\langle 42nd \& 9th, 42nd \& 8th \rangle$ and $\langle 42nd \& 8th, 42nd \& 9th \rangle$. The pair of nodes 42nd & 8th and 42nd & 7th is similar.

For many of the concepts that we’ll explore in this chapter, it will turn out that there are no substantive differences between the ideas for directed and undirected graphs. To avoid being tedious and unhelpfully repetitive, whenever it’s possible we’ll state definitions and results about both undirected and directed graphs simultaneously. But doing so will require a little abuse of notation: we’ll allow ourselves to write an edge as an ordered pair $\langle u, v \rangle$ *even for an undirected graph*. In an undirected graph, we will agree to understand both $\langle u, v \rangle$ and $\langle v, u \rangle$ as meaning $\{u, v\}$.

SIMPLE GRAPHS

For many of the real-world phenomena that we will be interested in modeling, it will make sense to make a simplifying assumption about the edges in our graphs. Specifically, we will typically restrict our attention to so-called *simple* graphs, which forbid two different kinds of edges: edges that connect nodes to themselves, and edges that are precise duplicates of other existing edges. (See Figure 11.1.)

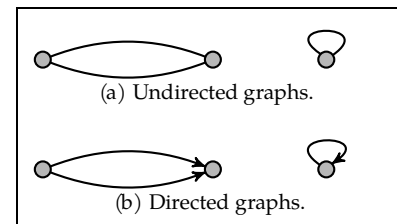


Figure 11.1: Parallel edges and self-loops.

Definition 11.3 (Self-loops and parallel edges)

A self-loop is an edge from a node u to itself. Two edges are called parallel if they both go from same node u and both go to the same node v .

Note that the edges $\langle u, v \rangle$ and $\langle v, u \rangle$ are not parallel in a directed graph: directed edges are parallel only if they both go *from* the same node and *to* the same node, in the same orientation.

Definition 11.4 (Simple graph)

A graph is simple if it contains no parallel edges and no self-loops.

In general, the particular real-world phenomenon that we seek to model will dictate whether self-loops, parallel edges, or both will make sense. Here are a few examples:

Example 11.3 (Self-loops and parallel edges)

Problem: Suppose that we construct a graph to model each of the following phenomena. In which settings do self-loops or parallel edges make sense?

1. A social network: nodes correspond to people; (undirected) edges represent friendships.
2. The web: nodes correspond to web pages; (directed) edges represent links.
3. The flight network for a commercial airline: nodes correspond to airports; (directed) edges denote flights scheduled by the airline in the next month.
4. The email network at a college: nodes correspond to students; there is a (directed) edge $\langle u, v \rangle$ if u has sent at least one email to v within the last year.

Solution:

1. Neither self-loops nor parallel edges make sense. A self-loop would correspond to a person being a friend of himself, and parallel edges between two people would correspond to them being friends “twice.” (But two people are either friends or not friends.)
2. Both self-loops and parallel edges are reasonable. It is easy to imagine a web page p that contains a hyperlink to p itself. It is also easy to imagine a web page p that contains two separate links to another web page q . (For example, as of this writing, the “CNN” logo on www.cnn.com links to www.cnn.com. And, as of the end of this sentence, this page has three distinct references to www.cnn.com.)
3. In the flight network, many parallel edges will exist: there are generally many scheduled commercial flights from one airport to another—for example, there are dozens of flights every week from BOS (Boston, MA) to SFO (San Francisco, CA) on most major airlines. However, there are no self-loops: a commercial flight from an airport back to the same airport doesn’t go anywhere!
4. Self-loops are reasonable but parallel edges are not. A student u has either sent email to v in the last year or she has not, so parallel edges don’t make sense in this network. However, self-loops exist if any student has sent an email to herself (as many people do to remind themselves to do something later).

Throughout, we assume that all graphs are simple unless otherwise noted.

Taking it further: Actually, the way that we phrased our definitions of graphs in Definitions 11.1 and 11.2 doesn’t even *allow* us to consider parallel edges. (Our definitions do allow self-loops, though.) That’s because we defined the edges as a subset E of $V \times V$ or $\{\{u, v\} : u, v \in V\}$, and sets don’t allow duplication—which means that we can’t have $\langle u, v \rangle$ in E “twice.” There are alternate ways to formalize graphs that do permit parallel edges, but they’re needlessly complicated for the applications that we’ll focus on in this chapter.

11.2.1 Neighborhoods and Degree

Imagine a social network in which two people, Ursula and Victor, are friends—or, more generally, imagine an undirected graph in which nodes u and v are joined by an edge. Here’s the vocabulary for referring to these nodes and the edge between them:

Definition 11.5 (Adjacency, neighbors, endpoints, incidence)

For an edge $e = \{u, v\}$ in an undirected graph (see Figure 11.2), we say that:

- the nodes u and v are adjacent;
- the node v is a neighbor of the node u (and vice versa);
- the nodes u and v are the endpoints of the edge e ; and
- the nodes u and v are both incident to the edge e .

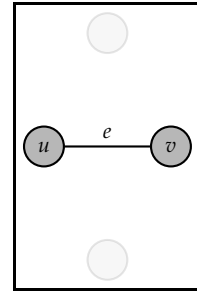


Figure 11.2: Two nodes joined by an edge.

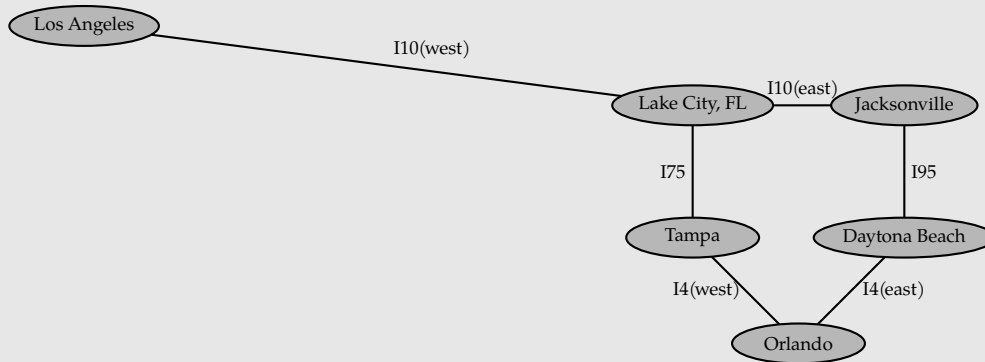
It's important to distinguish between two distinct concepts:

- the *direct* connection between two nodes u and v that are adjacent—that is, a single edge that joins u and v directly; and
- an *indirect* connection between two nodes that follows a sequence of edges.

At the moment, we're talking *only* about the first kind of connection, a direct connection via a single edge. (A multihop connection is called a *path*; we'll talk about paths in Section 11.3.) Here's an example of the vocabulary from Definition 11.5:

Example 11.4 (Disney World to Disney Land)

Here is a small portion of the U.S. Interstate system between Orlando, FL and Los Angeles, CA. Each of the roads is labeled by its name.



In this graph:

- Orlando is adjacent to Tampa and Daytona Beach.
- None of the other nodes (Lake City, Jacksonville, Los Angeles) is a neighbor of Orlando. Orlando is also not a neighbor of itself.
- The endpoints of edge I75 are Tampa and Lake City.
- Jacksonville is incident to I95, as is Daytona Beach.

The *neighborhood* of a node is the set of all nodes adjacent to it:

Definition 11.6 (Neighborhood)

Let $G = \langle V, E \rangle$ be an undirected graph, and let $u \in V$ be a node. The neighborhood of u is the set $\{v \in V : \{u, v\} \in E\}$ —that is, the set of all neighbors of u .

For example, in the graph from Example 11.4 (reproduced in abbreviated form in Figure 11.3), the neighborhood of Lake City (LC) is {Los Angeles (LA), Tampa (TA), Jacksonville (JA)}. Or, for a graph G that represents a social network, the neighborhood of a node u is the set of people who are u 's friends.

DEGREE

It's also common to refer the *number* of neighbors that a node has (without reference to which particular nodes happen to be that node's neighbors):

Definition 11.7 (Degree)

The degree of a node u in an undirected graph G is the size of the neighborhood of u in G —that is, the number of nodes adjacent to u .

For example, in the graph in Figure 11.3, Lake City (LC) has degree 3 and Los Angeles (LA) has degree 1. Or, in a social network, the degree of a node u is the popularity of u —the number of friends that u has. Here are a few practice questions:

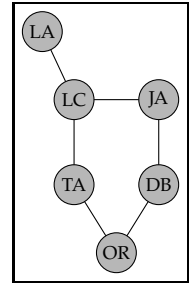
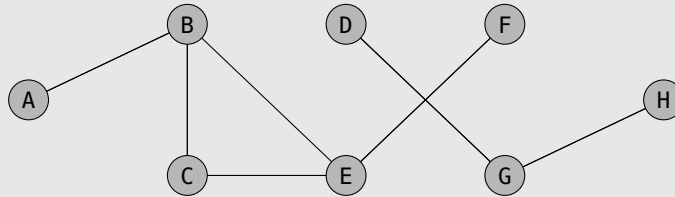


Figure 11.3: The road network from Example 11.4, abbreviated.

Example 11.5 (Neighborhood and degree)

Problem: Consider the following graph:



1. What are the neighbors of node C?
2. What nodes, if any, have degree equal to one?
3. What node has the highest degree in this graph?
4. What nodes, if any, are in the neighborhoods of both nodes B and E?

Solution: 1. Node C has two neighbors, namely the nodes B and E.

2. The nodes with degree one are those with precisely one neighbor. These nodes are: A, D, F, and H. (Their solitary neighbors are, respectively: B, G, E, and G.)
3. We simply count neighbors for each node, and we find that nodes B and E both have degree three, and are tied as the nodes with the highest degree.
4. The neighborhood of node B is {A, C, E}, and the neighborhood of node E is {B, C, F}. Taking the intersection of those sets yields the one node in the neighborhood of both B and E, namely node C.

Taking it further: Consider a population of people—say, the current residents of Canada—represented as a social network, in an undirected graph whose edges represent friendship. For a node in the social network (also known as a person), we can calculate many numbers that may be interesting: height, age, income, number of cigarettes smoked per day, self-reported happiness, etc. Then, for any one of these

numerical properties, we can consider the *distribution* over the population: for example, the distribution of heights, or the distribution of ages. (The height distribution will follow a roughly bell-shaped curve; the age distribution is more complicated, both because of death and because of variation in the birth rate over time.) Another interesting numerical property of a person u is the *degree* of u : that is, the number of friends that u has. The *degree distribution* of a graph describes how popularity varies across the nodes of the network. The degree distribution has some interesting properties—very different from the distribution of heights or ages. See p. 1123 for some discussion.

THE HANDSHAKING LEMMA

Before we move on from degree, we'll prove a basic but valuable fact, colloquially called the “handshaking lemma.” (We can represent a group of people, some pairs of whom shake hands, using an undirected graph: an edge joins u and v if and only if u and v shook hands; the theorem describes the number of shakes.) The handshaking lemma relates the sum of nodes' degrees to the number of edges in the graph:

Theorem 11.1 (“Handshaking Lemma”)

Let $G = \langle V, E \rangle$ be an undirected graph. Then

$$\sum_{u \in V} \text{degree}(u) = 2|E|.$$

For example, Figure 11.4 shows our road network from Example 11.4, with all nodes labeled by their degree. This graph has $|E| = 6$ edges, and the sum of the nodes' degrees is $1 + 3 + 2 + 2 + 2 + 2 = 12$, and indeed $12 = 2 \cdot 6$. Here is a proof:

Proof of Theorem 11.1. Every edge has two endpoints! Or, more formally, imagine looping over each edge to compute all nodes' degrees:

```

1: initialize  $d_u$  to 0 for each node  $u$ 
2: for each edge  $\{u, v\} \in E$ :
3:    $d_u := d_u + 1$ 
4:    $d_v := d_v + 1$ 

```

In each iteration of the **for** loop, we increment two different d_\bullet values; thus, after i iterations, we have that $\sum_u d_u = 2i$. (We could give a fully rigorous proof of this fact by induction.) We complete $|E|$ iterations of the **for** loop, one for each edge, and thus at the end of the algorithm we have that $\sum_{u \in V} d_u = 2|E|$. Furthermore, after the loop, it's clear that $d_u = \text{degree}(u)$ for every node u . Thus

$$\sum_{u \in V} d_u = \sum_{u \in V} \text{degree}(u) = 2|E|. \quad \square$$

Here's a useful corollary of Theorem 11.1 (the proof is left to you as Exercise 11.17):

Corollary 11.2

Let n_{odd} denote the number of nodes whose degree is odd. Then n_{odd} is even.

(For example, for the graph in Figure 11.4, we have $n_{\text{odd}} = 2$: the two nodes with odd degree are those with degree 1 and 3. And 2 is an even number.)

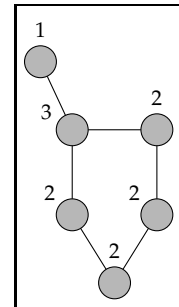


Figure 11.4: The road network from Figure 11.3, with nodes labeled by their degree.

“Look on every exit as being an entrance somewhere else.”
— Tom Stoppard
(b. 1937),
Rosencrantz and Guildenstern are Dead (1966)

NEIGHBORHOODS AND DEGREE: DIRECTED GRAPHS

The definitions of adjacency, neighbors, and degree from Definitions 11.5–11.7 were all for *undirected* graphs. Here we'll introduce the analogous notions for directed graphs, all of which are slightly more complicated because they must account for the orientation of each edge. We start with the directed version of “neighbors”:

Definition 11.8 (Neighbors in directed graphs)

For an edge $\langle u, v \rangle$ from node u to node v in a directed graph, we say that:

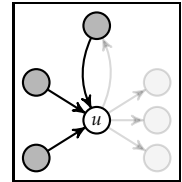
- the node v is an *out-neighbor* of the node u ; and
- the node u is an *in-neighbor* of the node v .

For example, if G represents a flight network (with nodes as airports and directed edges corresponding to flights), then the out-neighbors of node u are those airports that have direct flights from u , and the in-neighbors of u are those airports that have direct flights to u . (See Figure 11.5.) Now, using these definitions, we can define the analogues of neighborhoods and degree in directed graphs:

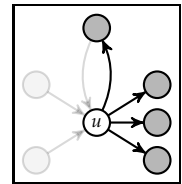
Definition 11.9 (Neighborhoods and degrees in directed graphs)

For a node u in an directed graph, we say that:

- the *in-neighborhood* of u is $\{v : \langle v, u \rangle \in E\}$, the set of in-neighbors of u ;
- the *in-degree* of u is its number of in-neighbors (its in-neighborhood's cardinality);
- the *out-neighborhood* of u is $\{v : \langle u, v \rangle \in E\}$, the set of out-neighbors of u ; and
- the *out-degree* of u is its number of out-neighbors (its out-neighborhood's cardinality).



(a) in-neighbors



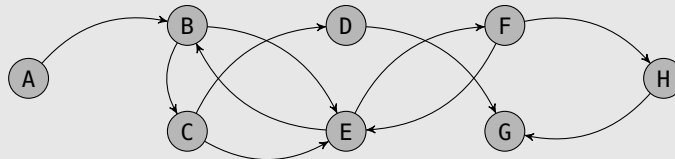
(b) out-neighbors

Figure 11.5: The in- and out-neighbors of a node u .

Here are a few practice questions about in- and out-neighborhoods:

Example 11.6 (Neighborhood and degree in a directed graph)

Problem: Consider the following directed graph:



1. What are the in-neighbors of node C? The out-neighbors of C?
2. What nodes, if any, are in both the in-neighborhood and out-neighborhood of node E?
3. What nodes, if any, have in-degree zero? Out-degree zero?

Solution: 1. Node C has one in-neighbor, namely B, and two out-neighbors, namely D and E.

2. Node E has three in-neighbors (B, C, and F) and two out-neighbors (B and F). So nodes B and F are in both E's in-neighborhood and E's out-neighborhood.

3. Node A has no in-neighbors, so A's in-degree is zero. Node G has no out-neighbors, so G's out-degree is zero.

11.2.2 Representing Graphs: Data Structures

The graphs that we've considered so far have been presented visually: as a picture, with nodes drawn as circles and edges drawn as lines or arrows. But, of course, when we represent a graph on a computer, we'll need to use some data structure to store a network, not just some image file. Here we will give a brief summary of the two major data structures used to represent graphs. If you've had a course on data structures, then this material may be a review; if not, it will be a preview.

Taking it further: A visual representation is great for some smaller networks, and a well-designed layout can sometimes make even large networks easy to understand at a glance. *Graph drawing* is the problem of algorithmically laying out the nodes of a graph well—in an aesthetic and informative manner. There's a physics analogy that's often used in laying out graphs, in which we imagine nodes “attracting” and “repelling” each other depending on the presence or absence of edges. See p. 1124 for some discussion, including an application of this graph-drawing idea to the 9/11 Memorial in New York City. Some other gorgeous visualizations of network (and other!) data can be found online at sites like *Flowing Data* (<http://flowingdata.com/>), *Information Is Beautiful* (<http://informationisbeautiful.net>), or some of the beautiful books on data visualization like the *Atlas of Science*.³

³ Katy Börner.
*Atlas of Science:
Visualizing What We
Know*. MIT Press,
2010.

The most straightforward data structure for a graph is just a list of nodes and a list of edges. But this straightforward representation suffers for some standard, natural questions that are typically asked about graphs. Many of the natural questions that we will find ourselves asking are things like: *What are all of the neighbors of A?* or *Are B and C joined by an edge?* There are two standard data structures for graphs, each of which is tailored to make it possible to answer one of these two questions quickly.

ADJACENCY LISTS

The first standard data structure for graphs is an *adjacency list*, which—as the name implies—stores, for each node u , a list of the nodes adjacent to u :

Definition 11.10 (Adjacency list)

In an adjacency list of a graph $G = \langle V, E \rangle$, for each node $u \in V$, we store an unsorted list of all of u 's neighbors in the graph.

The schematic for an adjacency list is illustrated in Figure 11.6: each node in the graph corresponds to a row of the table, which points to an unsorted list of that node's neighbors. (These lists are unsorted so that it's faster to add a new edge to the data structure.)

There's no significant difference between adjacency lists for undirected graphs and for directed graphs: for an undirected graph, we list the *neighbors* for each node u ; for a directed graph, we list the *out-neighbors* of each node. (Every edge $\langle u, v \rangle$ in a directed graph appears only once in the data structure, in u 's list. Every edge $\{u, v\}$ in an undirected graph is represented twice: v appears in u 's list, and u appears in v 's list. This observation is another way of thinking of the proof of Theorem 11.1.)

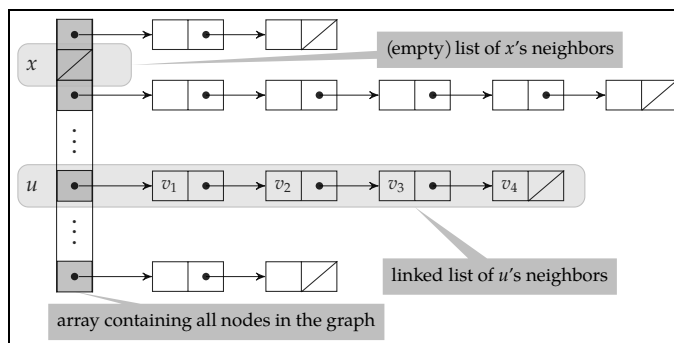
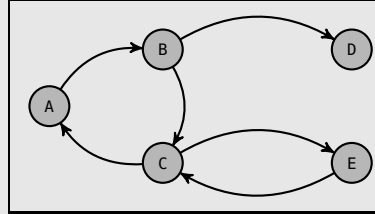
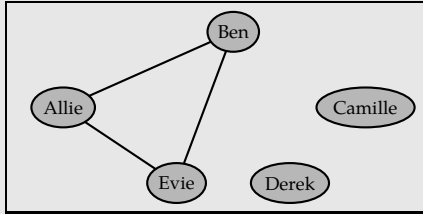


Figure 11.6: A schematic of an adjacency list.

Here are example adjacency lists for two graphs, one undirected and one directed:

Example 11.7 (Two sample adjacency lists)

Consider the following two graphs:



The adjacency lists for these two graphs are as follows.

Allie:	Evie, Ben
Ben:	Allie, Evie
Camille:	--
Derek:	--
Evie:	Allie, Ben

A:	B
B:	C, D
C:	E, A
D:	--
E:	C

Note that the order of the (out-)neighbors of any particular node isn't specified: for example, we could just as well said that Evie's neighbors were [Ben, Allie] as [Allie, Ben].

ADJACENCY MATRICES

The second standard data structure for representing graphs is an *adjacency matrix*:

Definition 11.11 (Adjacency matrix)

In an adjacency matrix of a graph $G = \langle V, E \rangle$, we store the graph using an $|V|$ -by- $|V|$ table. The i th row of the table corresponds to the neighbors of node i . A True (or 1) in column j indicates that the edge $\langle i, j \rangle$ is in E ; a False (or 0) indicates that $\langle i, j \rangle \notin E$.

In a directed graph, the i th row corresponds to the *out-neighbors* of node i , so that the $\langle i, j \rangle$ th entry of the matrix corresponds to the presence/absence of an edge *from* i to j . The i th column corresponds to the in-neighbors of i . Here are two examples of adjacency matrices, for the graphs from Example 11.7:

Example 11.8 (Two sample adjacency matrices)

The following adjacency matrices represent the graphs from Example 11.7:

	Allie	Ben	Camille	Derek	Evie
Allie	0	1	0	0	1
Ben	1	0	0	0	1
Camille	0	0	0	0	0
Derek	0	0	0	0	0
Evie	1	1	0	0	0

	A	B	C	D	E
A	0	1	0	0	0
B	0	0	1	1	0
C	1	0	0	0	1
D	0	0	0	0	0
E	0	0	1	0	0

The adjacency matrix has two properties that are worth a note. (See Figure 11.7.)

- The main diagonal contains all zeros: a 1 in the $\langle i, i \rangle$ th position of the matrix would correspond to an edge between node i and node i —that is, a self-loop, which is forbidden in a simple graph.

- For an undirected graph, the matrix is symmetric: the $\langle i, j \rangle$ th position of the matrix records the presence or absence of an edge from i to j , which is identical to the presence or absence of an edge from j to i in an undirected graph. Adjacency matrices are not necessarily symmetric in directed graphs: there may be an edge from u to v without an edge from v to u .

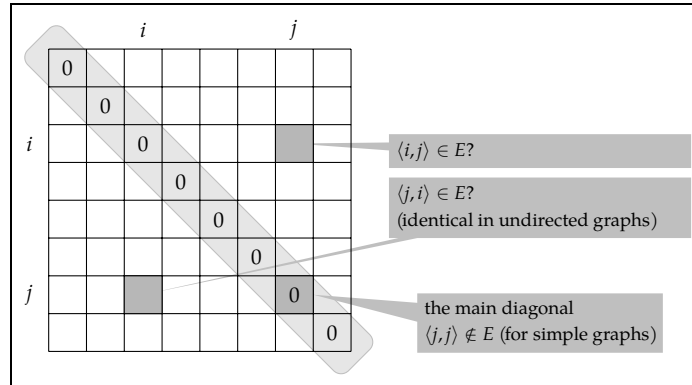


Figure 11.7: A schematic of an adjacency matrix.

CHOOSING BETWEEN ADJACENCY LISTS AND MATRICES

Which of the two data structures that we've seen for graphs should we choose? Are adjacency lists better than adjacency matrices, or the other way around? Recall the two basic questions about graphs that we wish to answer quickly:

- (A) is v a neighbor of u ?
- (B) what are all of u 's neighbors?

Figuring the details of how efficiently we can answer these questions with an adjacency list or an adjacency matrix is better suited to a data-structures textbook than this one, but here's a brief summary of the reasoning.

Adjacency Lists: An adjacency list is perfectly tailored to answering Question (B): we've stored precisely the list of u 's neighbors for each node u , so we simply iterate through that list to output u 's neighborhood. To answer Question (A), we need to search through that same unsorted list to see if v is present. In both cases, we have to spend constant time finding u 's list in the table, and then we examine a list of length $\text{degree}(u)$ to answer the question.

Adjacency Matrices: An adjacency matrix is perfect for answering Question (A): we just look at the appropriate spot in the table. If the $\langle u, v \rangle$ th entry is True, then the edge $\langle u, v \rangle$ exists. This lookup takes constant time. Answering Question (B) requires looking at one entire row of the table, entry by entry. There are $|V|$ entries in the row, so this loop requires $|V|$ operations.

Thus adjacency matrices solve Question (A) faster, while adjacency lists are faster at solving Question (B). In addition to the time to answer these questions, we'd also want the *space*—the amount of memory—consumed by the data structure to be as small as possible. (You can think of “the amount of memory” as the total number of boxes that appear in the diagrams in Figures 11.6 and 11.7.)

Meta-problem-solving tip: The answer to “which is better?” in a class or textbook is almost always *It depends!* After all, why would we waste time/ pages on a solution that's always worse!? (The only plausible answer is that it warms us up conceptually for a better but more complex solution.) The real question here *what does it depend on?*

Example 11.9 (Space consumption for adjacency lists and matrices)

Problem: Consider a graph $G = \langle V, E \rangle$ stored using an adjacency list or an adjacency matrix. In terms of the number of nodes and the number of edges in G —that is, in terms of $|V|$ and $|E|$ —how much memory is used by these data structures?

Solution: An adjacency matrix is a $|V|$ -by- $|V|$ table, and thus contains exactly $|V|^2$ cells. (Of them, the $|V|$ cells on the diagonal are always 0, but they’re still there!)

An adjacency list is a $|V|$ -element table pointing to $|V|$ lists; the length of the list for node u is exactly $\text{degree}(u)$. Thus the total number of cells in the data structure is

$$|V| + \sum_{u \in V} \text{degree}(u).$$

In an undirected graph we have $\sum_u \text{degree}(u) = 2|E|$, by Theorem 11.1; in a directed graph we have $\sum_u \text{out-degree}(u) = |E|$ by Exercise 11.18. Thus the total amount of memory used is

$$\begin{cases} |V| + 2|E| & \text{for an undirected graph} \\ |V| + |E| & \text{for a directed graph.} \end{cases}$$

Here’s the summary of the efficiency differences between these data structures (using asymptotic notation from Chapter 6):

	adjacency list	adjacency matrix
is v a neighbor of u ?	$1 + \Theta(\text{degree}(u))$	$\Theta(1)$
what are all of u ’s neighbors?	$1 + \Theta(\text{degree}(u))$	$\Theta(V)$
space	$\Theta(V + E)$	$\Theta(V ^2)$

The better data structure in each row is highlighted. (Note that, in a simple graph, we have that $\text{degree}(u) \leq |V|$ and $|E| \leq |V|^2$.) So, is an adjacency list or an adjacency matrix better? *It depends!*

First, it depends on what kind of questions—Question (A) or Question (B) listed previously, for example—we want to answer: if we will ask few “is v a neighbor of u ?” questions, then adjacency lists will be faster. If we will ask many of those questions, then we probably prefer adjacency matrices. Similarly, it might depend on how much, if at all, the graph changes over time: adjacency lists are harder to update than adjacency matrices.

Second, it depends on how many edges are present in the graph. If the total number of edges in the graph is relatively small—and thus most nodes have only a few neighbors—then $\text{degree}(u)$ will generally be small, and the adjacency list will win. If the total number of edges in the graph is relatively large, then $\text{degree}(u)$ will generally be larger, and the adjacency matrix will perform better. (Many of the most interesting real-world graphs are sparse: for example, the typical degree of a person in a social network like Facebook is perhaps a few hundred or at most a few thousand—very small in relation to the hundreds of millions of Facebook users.)

11.2.3 Relationships between Graphs: Isomorphism and Subgraphs

Now that we have the general definitions, we'll turn to a few more specific properties that certain graphs have. We'll start in this section with two different relationships between pairs of graphs—when two graphs are “the same” and when one is “part” of another; in Section 11.2.4, we'll look at single graphs with a particular structure.

GRAPH ISOMORPHISM

When two graphs G and H are identical except for how we happen to have arranged the nodes when we drew them on the page (and except for the names that we happen to have assigned to the nodes), then we call the graphs *isomorphic*. Informally, G and H are isomorphic if there's a way to relabel (and rearrange) the nodes of G so that G and H are exactly identical. More formally:

Greek: *iso* “same”;
morph “form.”

Definition 11.12 (Graph isomorphism)

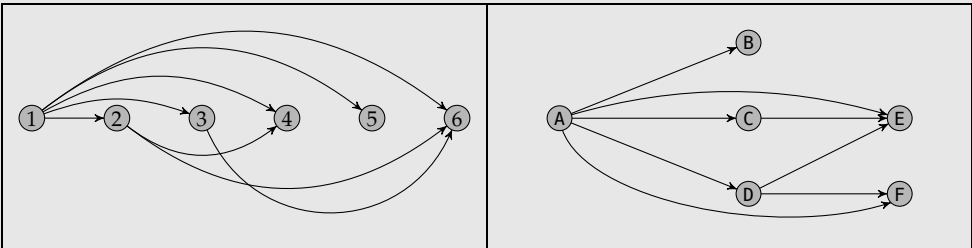
Consider two graphs $G = \langle V, E \rangle$ and $H = \langle U, F \rangle$. We say that G and H are isomorphic if there exists a bijection $f : V \rightarrow U$ such that

$$\text{for all } a \in V \text{ and } b \in V, \quad \langle a, b \rangle \in E \Leftrightarrow \langle f(a), f(b) \rangle \in F.$$

(By abusing notation as we described earlier, this definition works for either undirected or directed graphs G and H .) Here are some small examples:

Example 11.10 (Two isomorphic graphs)

Let's show that the following two directed graphs are isomorphic. (The first graph's edges could also have been written as $\{\langle a, b \rangle : a < b \text{ and } a \text{ evenly divides } b\}$.)



To do so, define the following bijection $f : \{1, 2, \dots, 6\} \rightarrow \{A, B, \dots, F\}$:

x	1	2	3	4	5	6
$f(x)$	A	D	C	F	B	E

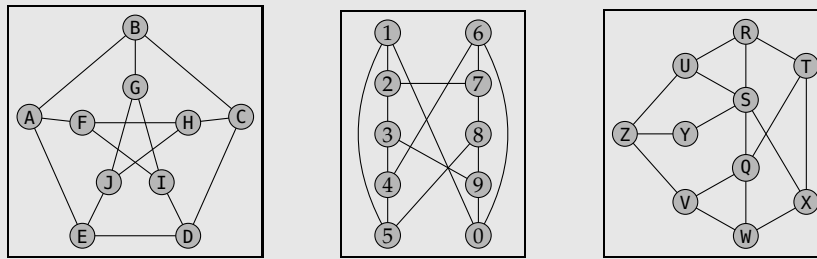
The tables of edges in the graphs now match exactly, so they are isomorphic:

	1	2	3	4	5	6
1		✓	✓	✓	✓	✓
2				✓		✓
3						✓
4						
5						
6						

	A	D	C	F	B	E
$f(1) = A$		✓	✓	✓	✓	✓
$f(2) = D$				✓		✓
$f(3) = C$						✓
$f(4) = F$						
$f(5) = B$						
$f(6) = E$						

Example 11.11 (Isomorphic graphs)

Problem: Which pairs, if any, of the following graphs are isomorphic?



Solution: The first two graphs are isomorphic. The easiest way to see this fact is to show the mapping between the nodes of the two graphs:

A	B	C	D	E	F	G	H	I	J
1	2	3	4	5	6	7	8	9	10

It's easy to verify that all 15 edges now match up between the first two graphs. But the third graph is not isomorphic to either of the others. The easiest justification is that node S in the third graph has degree 5, and no node in either of the first two graphs has degree 5. No matter how we reshuffle the nodes of graph #3, there will still be a node of degree 5—so the third graph can never match the others.

Taking it further: In general, it's easy to test whether two graphs are isomorphic by brute force (try all permutations!), but no substantially better algorithms are known. The computational complexity of the graph isomorphism problem has been studied extensively over the last few decades, and there has been substantial progress—but no complete resolution.

It's easy to convince someone that two graphs G and H are isomorphic: we can simply describe the relabeling of the nodes of G so that the resulting graphs are identical. (The “convinced” then just needs to verify that the edges really do match up.) When G and H are not isomorphic, it *might* be easy to demonstrate their nonisomorphism: for example, if they have a different number of nodes or edges, or if the degrees in G aren't identical to the degrees in H . But the graphs may have identical degree distributions and yet *not* be isomorphic; see Exercise 11.49.

SUBGRAPHS

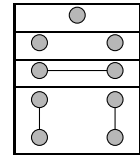
When a graph H is isomorphic to a graph G , we can think of having created H by moving around some of the nodes and edges of G . When H is a *subgraph* of G , we can think of having created H by deleting some of the nodes and edges of G . (Of course, it doesn't make sense to delete either endpoint of an edge e without also deleting the edge e .) Here's the definition, for either undirected or directed graphs:

Definition 11.13 (Subgraph)

Let $G = \langle V, E \rangle$ be a graph. A subgraph of G is a graph $G' = \langle V', E' \rangle$ where $V' \subseteq V$ and $E' \subseteq E$ such that every edge $\langle u, v \rangle \in E'$ satisfies $u \in V'$ and $v \in V'$.

For example, consider the graph $G = \langle V, E \rangle$ with nodes $V = \{A, B, C, D\}$ and edges $E = \{\{A, B\}, \{A, C\}, \{B, C\}, \{C, D\}\}$. Then the graph G' with nodes $\{B, C, D\}$ and edges $\{\{B, C\}, \{C, D\}\}$ is a subgraph of G . In fact, G has *many* different subgraphs:

Problem-solving tip: When you're trying to prove or disprove a claim about graphs, you may find it useful to test out the claim against the following four “trivial” graphs:

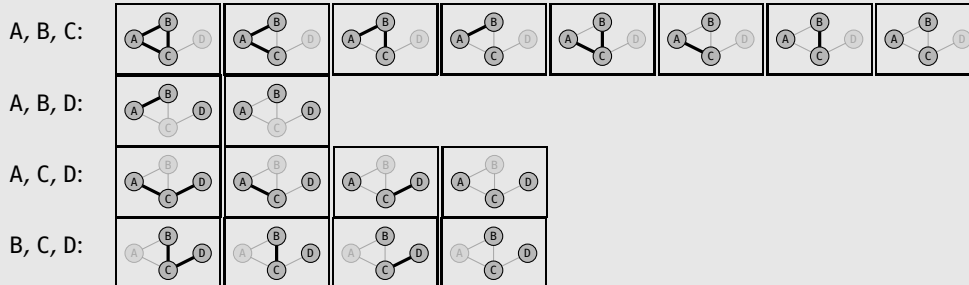


A lot of bogus claims about graphs turn out to be false on one of these four examples—or, unexpectedly, the so-called *Petersen graph*, the first graph in Example 11.11. (The Petersen graph is named after Julius Petersen, a 19th-century Danish mathematician.) It's a good idea to try out any conjecture on all five of these graphs before you let yourself start to believe it!

Note that Definition 11.13 uses the abuse of notation that we mentioned earlier: we “ought” to have written $\{u, v\} \in E'$ for the case that G is undirected.

Example 11.12 (All 3-node subgraphs of G)

Here are all of the 3-node subgraphs of the graph G with nodes $V = \{A, B, C, D\}$ and edges $E = \{\{A, B\}, \{A, C\}, \{B, C\}, \{C, D\}\}$. (There are many other subgraphs—about 50 total—when we consider subgraphs with 1, 2, 3, or 4 nodes.)



Taking it further: One of the earliest applications of a formal, mathematical perspective to networks—a collaboration between a psychologist and mathematician, in the 1950s—was based on subgraphs. Consider a *signed social network*, an undirected graph where each edge is labeled with ‘+’ to indicate friends, or ‘−’ to indicate enemies. (See Figure 11.8(a).) The adages “the enemy of my enemy of my friend” and “the friend of my friend is my friend” correspond to the claim that the subgraphs in Figure 11.8(b) would not appear. Dorwin Cartwright (the psychologist) and Frank Harary (the mathematician) proved some very interesting structural properties of any signed social network G that does not have either triangle in Figure 11.8(b) as a subgraph—a property that they called “structural balance”—and in the process helped launch much of the mathematical and computational work on graphs that’s followed.⁴

We sometimes refer to a special kind of subgraph: the subgraph of $G = \langle V, E \rangle$ induced by a set $V' \subseteq V$ of nodes is the subgraph of G where every edge between nodes in V' is retained. The first subgraph in each row of Example 11.12 is the induced subgraph for its nodes. Here’s a brief description of one application of (induced) subgraphs:

Example 11.13 (Motifs in biological networks)

At any particular moment in any particular cell, some of the genes in the organism’s DNA are being *expressed*—that is, some genes are “turned on” and the proteins that they code for are being produced by the cell. Furthermore, one gene g can *regulate* another gene g' : when g is being expressed, gene g can cause the expression of gene g' to increase or decrease over the baseline level. A great deal of recent biological research has allowed us to construct *gene-regulation networks* for different such settings: that is, a directed graph G whose nodes are genes, and whose edges represent the regulation of one gene by another.

Consider the induced subgraph of a particular set of genes in such a graph G —that is, the interactions among the particular genes in that set. Certain patterns of these subgraphs, called *motifs*, occur significantly more frequently in gene-regulation networks than would be expected by chance. Biologists generally believe that these repeated patterns indicate something important in the way that our genes work, so computational biologists have been working hard to build efficient algorithms to identify induced subgraphs that are overrepresented in a network.

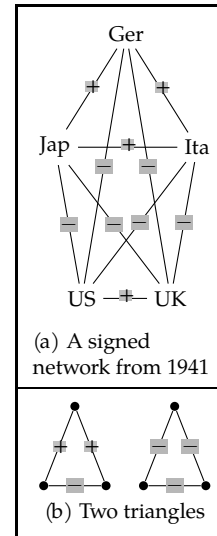


Figure 11.8: Signed social networks. For more about signed networks and these results, see

⁴ Dorwin Cartwright and Frank Harary. Structural balance: a generalization of Heider’s theory. *Psychological Review*, 63(5):277–293, 1956.

11.2.4 Special Types of Graphs: Complete, Bipartite, Regular, and Planar Graphs

In Section 11.2.3, we looked at two ways in which a pair of graphs might be related. Here, we'll consider special characteristics that a single graph might have—that is, subcategories of graphs with some particular structural properties. These special types of graphs arise frequently in various applications.

COMPLETE GRAPHS

Our first special type of graph is a *complete graph* (also called a *clique*), which is an undirected graph in which every possible edge exists:

Definition 11.14 (Complete graph/clique)

A complete graph or clique is an undirected graph $G = \langle V, E \rangle$ such that $\{u, v\} \in E$ for any two distinct nodes $u \in V$ and $v \in V$.

See Figure 11.9 for examples of complete graphs of varying sizes. (In everyday usage, a *clique* is a small, tight-knit, and exclusionary group of friends that doesn't mingle with outsiders. If you think about a graph as a social network, the common-language meaning is similar to Definition 11.14.)

Observe that an undirected graph with n nodes has $\binom{n}{2}$ unordered pairs of nodes, and therefore an n -node complete graph has $\binom{n}{2} = n(n-1)/2$ edges.

A complete graph with n nodes is sometimes denoted by K_n .

The word *clique* can also refer to a *subgraph* that's complete—that is, in which every possible edge actually exists. For example, the graph $G = \langle V, E \rangle$ with $V = \{A, B, C, D\}$ and $E = \{\{A, B\}, \{A, C\}, \{B, C\}, \{C, D\}\}$ contains a 3-node clique $\{A, B, C\}$. Here's one small example of an interesting application in which cliques arise:

Example 11.14 (Collaboration networks and cliques)

Imagine a setting in which different groups of people can work together in different teams, with each person allowed to participate in multiple teams. For example:

- actors in movies. (A “team” is the cast of a single movie.)
- scientific researchers. (A “team” is the set of coauthors of a published paper.)
- employees of a company. (A “team” is a group that worked on a specific project.)

A *collaboration network* is a graph G that represents a setting like these: the nodes of G are the people involved; there is an edge between any two people who have worked together on at least one team. (You may have heard of a challenge in the collaboration network: in the *Kevin Bacon Game*, you're given the name of some actor A ; your job is to find a sequence of edges that connects A to the “Kevin Bacon” node in the movie collaboration network. There's a similar game that computer scientists play in the scientific collaboration network, trying to connect themselves to the Hungarian polymath Paul Erdős. See p. 438.)

In CS, the word *clique* usually rhymes with *bleak* or *sleek*. In common-language usage, the word usually rhymes with *slick* or *flick*.

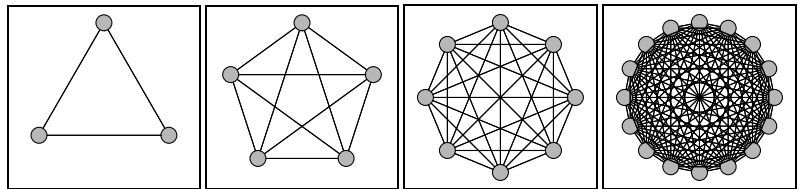


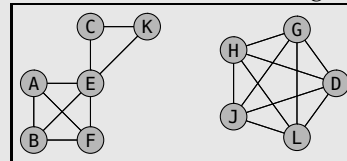
Figure 11.9: Complete graphs with 3, 5, 8, and 16 nodes.

There are two different prevailing explanations for the K_n notation:

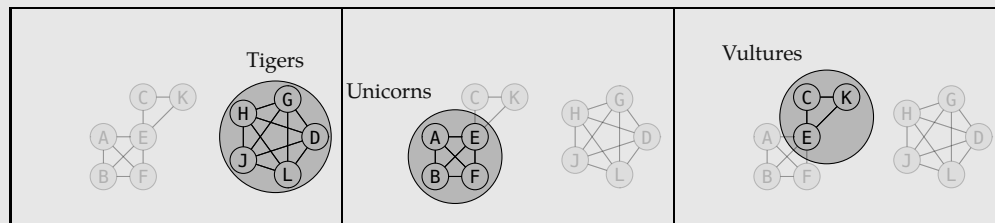
- the K is as in *complete*—or, rather, as in *komplett*; the notation was invented by a German speaker.
- the K is in honor of Kazimierz Kuratowski, a 20th-century Polish mathematician who made major contributions to the study of graphs (among other mathematical topics).

For example, for the teams listed below, we get the collaboration network at right:

- Tigers: Deborah, George, Hicham, Josh, Lauren
- Unicorns: Anita, Bev, Eva, Fernan
- Vultures: Cathy, Eva, Kelly



Notice that each team results in a clique inside the collaboration graph—every pair of members of that team is joined by an edge—in this case, creating a K_5 , K_4 , and K_3 in the graph:



BIPARTITE GRAPHS

Our second special kind of graph is a *bipartite* graph. In a bipartite graph, the nodes can be divided into two groups such that no edges join two nodes that are in the same group: that is, there are two “kinds” of nodes, and all edges join a node of Type A to a node of Type B. Formally:

Latin: *bi* “two”; *part* “part.”

Definition 11.15 (Bipartite graph)

A bipartite graph is an undirected graph $G = \langle V, E \rangle$ such that V can be partitioned into two disjoint sets L and R where, for every edge $e \in E$, one endpoint of e is in L and the other endpoint of e is in R .

For example, consider the graph $G = \langle V, E \rangle$ whose nodes are $V = \{A, B, C, D, E, F\}$ and whose edges are $E = \{\{A, B\}, \{A, C\}, \{C, E\}, \{D, E\}\}$. The graph G is bipartite: for example, we can split the nodes into two groups—the vowels $\{A, E\}$ and the consonants $\{B, C, D, F\}$ —such that every edge joins a vowel and a consonant. (There’s another split that would also have worked: $\{A, E, F\}$ and $\{B, C, D\}$.) See Figure 11.10 for a visualization of the vowel–consonant split.

Bipartite graphs are traditionally drawn with the nodes arranged in two columns, one for each part: *left* (“ L ”) and *right* (“ R ”). But notice that the definition only requires that it be *possible* to divide the nodes into two groups, with no within-group edges.

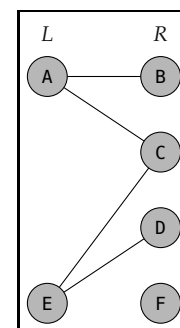
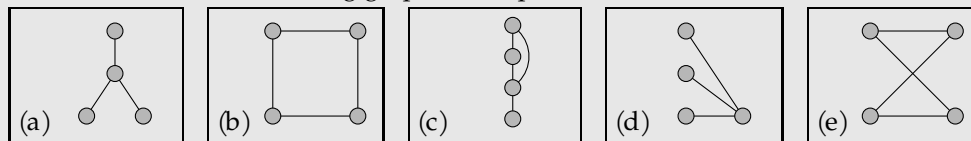


Figure 11.10: A bipartite graph.

Example 11.15 (Bipartite or nonbipartite?)

Problem: Which of the following graphs are bipartite?



Solution: All of them except (c)! Although (d) and (e) are the only graphs drawn in the “two-column” format, both (a) and (b) can be rearranged into two columns. In fact, aside from node positioning, graphs (a) and (d) are identical. And, similarly, graphs (b) and (e) are isomorphic!

Only (c) is not bipartite: if we attempt to put the topmost node in one group, then both of the next higher two nodes must both be in the other group—but they’re joined by an edge themselves, and so we’re stuck.

Many interesting real-world phenomena can be modeled using bipartite graphs:

Example 11.16 (Bipartite graphs as models)

Here are just a few of the scenarios that are naturally modeled using bipartite graphs:

- dating relationships in a strictly heterosexual community: the nodes are the boys B and the girls G ; every edge connects some boy to some girl.
- nodes are courses and students; an edge joins a student to each class she’s taken.
- *affiliation networks*: people and organizations are the nodes; an edge connects person p and organization o if p is a member of o .

There’s one further refinement of bipartite graphs that we’ll mention: a *complete bipartite graph* is a bipartite graph in which every possible edge exists. In other words, a complete bipartite graph has the form $G = \langle L \cup R, E \rangle$ where $\{\ell, r\} \in E$ for every node $\ell \in L$ and $r \in R$. A complete bipartite graph with ℓ nodes in the left group and r nodes in the right group is sometimes denoted by $\mathcal{K}_{\ell,r}$.

See Figure 11.11 for a few examples. (Note again that, as with the $\mathcal{K}_{2,4}$ in Figure 11.11, we don’t have to draw a bipartite graph in two-column format—if it’s bipartite, then it’s still bipartite no matter how we draw it!)

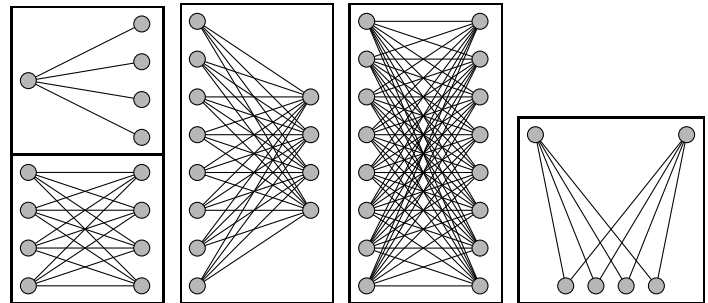


Figure 11.11: Complete bipartite graphs of varying sizes: $\mathcal{K}_{1,4}$, $\mathcal{K}_{4,4}$, $\mathcal{K}_{8,8}$, and $\mathcal{K}_{2,4}$.

REGULAR GRAPHS

Our next type of graph is defined in terms of the degree of its nodes: a *regular graph* is one in which all of the nodes have an identical number of neighbors.

Definition 11.16 (Regular graph)

Let $d \geq 0$ be an integer. A d -regular graph is a graph G such that every node has degree precisely equal to d . If G is d -regular for any d , then we say that G is a regular graph.

(Most of the time one talks about regular graphs that are undirected, but we can speak of regular directed graphs, too; we’d generally require that all in-degrees match each other *and* all out-degrees match each other.)

For example, consider the graph $G = \langle V, E \rangle$ whose nodes are $V = \{A, B, C, D, E, F\}$ and whose edges are $E = \left\{ \{A, B\}, \{A, E\}, \{B, C\}, \{C, F\}, \{D, E\}, \{D, F\} \right\}$. The graph G is 2-regular: you can check that each node has exactly two neighbors. As another example, note that the complete graph K_n is $(n - 1)$ -regular, as each node has all $n - 1$ other nodes as neighbors. Or see Figure 11.12 for another example of a regular graph.

There are many real-world examples in which regular graphs are useful: for example, imagine constructing a physical network of computers in which each machine only has the capacity for a fixed number of connections. Here are two other useful applications of regular graphs:

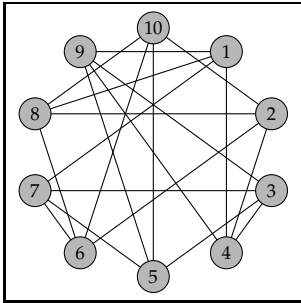


Figure 11.12: A 4-regular 10-node graph.

Example 11.17 (Scheduling sports with a regular graph)

You are the League Commissioner for an intramural ultimate frisbee league. There are 10 teams in the league, each of whom should play four games. No two teams should play each other twice. Suppose that you construct an undirected graph $G = \langle V, E \rangle$, where $V = \{1, 2, \dots, 10\}$ is the set of teams, and E is the set of games to be played. If G is an 4-regular graph, then all of the listed requirements are met. Figure 11.12 is a randomly generated example of such a graph; you could use that graph to set the league schedule.

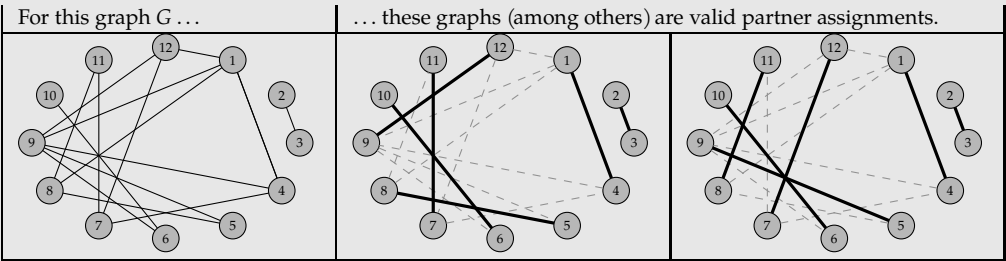
A 1-regular graph is called a *perfect matching*, because each node is “matched” with one—and only one—neighbor. (If every node has degree *at most* 1, then the graph is just called a *matching*.) Matchings have a variety of applications—for example, see p. 960 for their role in the Enigma machine—but here’s another specific use of matchings, in assigning partnerships:

Example 11.18 (Matchings for CS partnerships)

Each of n students in an Intro CS class submits a list of people whom they’d like to have as a partner for the final project. Define the following undirected graph G :

- the set V of nodes is $\{1, 2, \dots, n\}$, one per student.
- the set E of edges includes $\{u, v\}$ if *both* of the following are true: student u wants to work with student v , *and* student v wants to work with student u .

The instructor can assign partnerships by finding a 1-regular graph $G' = \langle V, E' \rangle$ with $E' \subseteq E$ —that is, a subgraph of G that includes all of the nodes of G . For example:



(Incidentally, Example 9.32 asked: how many perfect matchings are there in K_n ?)

PLANAR GRAPHS

Our last special type of graph is a *planar graph*, which is one that can be drawn on a sheet of paper without any lines crossing:

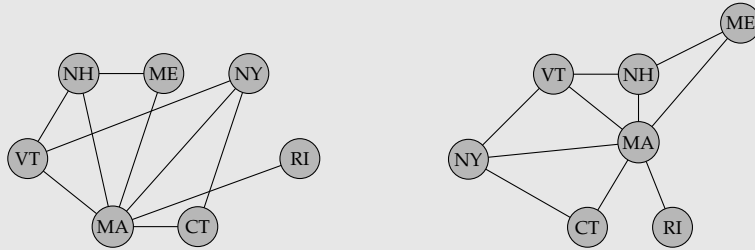
Definition 11.17 (Planar graph)

A planar graph is a graph G such that it is possible to draw G on a plane (that is, on a piece of paper) such that no edges cross.

It's important to note that a graph is planar if it is *possible* to draw it with no crossing edges; just because a graph is drawn with edges crossing does not mean that it isn't planar. Here is an example of a planar graph:

Example 11.19 (New England, in a plane)

Here are two copies of the same graph—one drawn with edge crossings, and another with the nodes rearranged to avoid edge crossing:



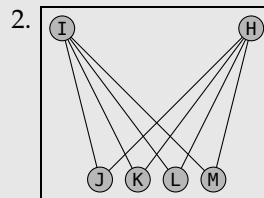
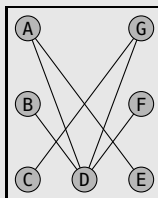
Example 11.19 shows one of the most famous types of planar graph, one derived from a map: we can think of the countries on a map as nodes, and we draw an edge between two country-nodes if those two countries share a border. (See p. 437 for a discussion of the *four-color theorem* for maps, which we could have phrased as a result about planar graphs instead.)

There are other applications of planar graphs in computer science, too. For example, we can view a *circuit* (see Section 3.3.3) as a graph, where the logic gates correspond to nodes and the wires correspond to edges. Most modern circuits are now *printed* on a board (where the “ink” is the conducting material that serves as the wire), and the question of whether a particular circuit can be printed on a single layer is precisely the question of whether its corresponding graph is planar. (If it's not planar, we'd like to minimize the number of edges that cross, or more specifically the number of layers we'd need in the circuit.)

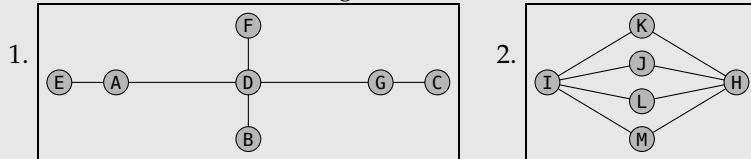
Here's one more set of planarity challenges for you to try:

Example 11.20 (Two planar challenges)

Problem: Are these graphs planar? 1.



Solution: Yes, both: we can rearrange the nodes so that there are no edges that cross.



Taking it further: Determining how to lay out a planar graph without edge crossings can be an interesting amusement—see www.planarity.net for a surprisingly fun game based on planar graphs. So far we haven't seen any examples of graphs that *can't* be rearranged so that no edges cross. But, if you play around long enough, you should be able to convince yourself that neither K_5 and $K_{3,3}$ are planar; see Figure 11.13. And, while this shouldn't be at all obvious, it turns out that K_5 and $K_{3,3}$ are in a sense the only “reasons” that a graph can be nonplanar. A theorem known as *Kuratowski's Theorem*—after the Polish mathematician who may have lent his initial to the notation for complete graphs—says that every graph is planar unless it “contains” K_5 or $K_{3,3}$ for a subgraph-like notion of “containment.” (It's not exactly the subgraph relation, because there are graphs that do not contain K_5 or $K_{3,3}$ as subgraphs but nonetheless are nonplanar in some sense “because” of one of them. For example, the Petersen Graph from Example 11.11—see Figure 11.13(c)—is nonplanar, but it doesn't have K_5 as a subgraph. But if we “collapse” together the nodes A/F, B/G, C/H, D/I, and E/J into “supernodes” then the resulting graph is K_5 .)

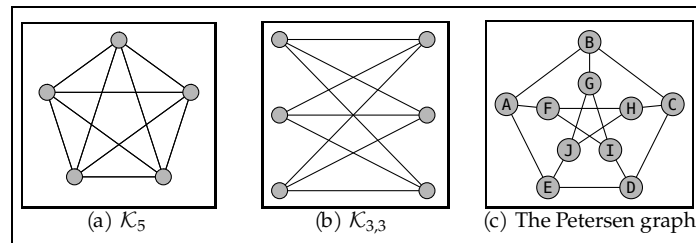
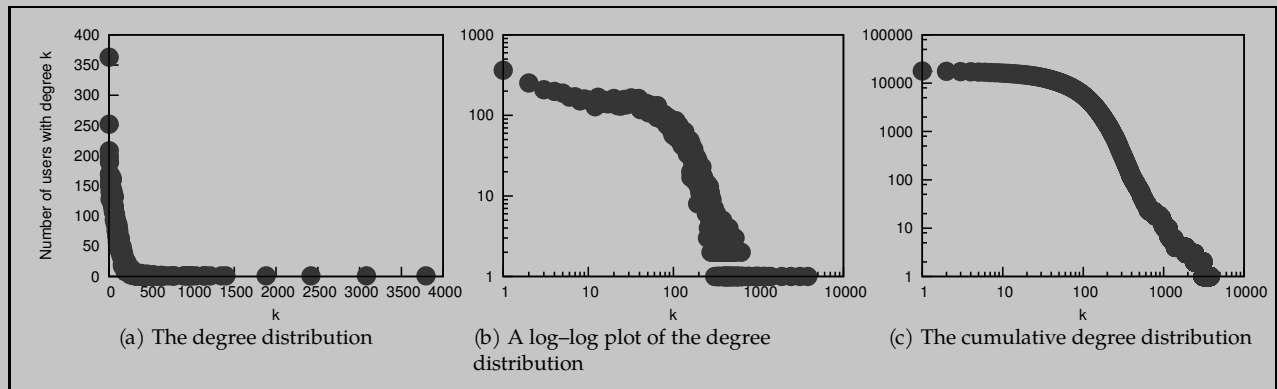


Figure 11.13:
Nonplanar graphs.

COMPUTER SCIENCE CONNECTIONS

DEGREE DISTRIBUTIONS AND THE HEAVY TAIL



When we think about massive graphs like the World-Wide Web (with nodes representing web pages and edges representing hyperlinks from one page to another) or an online social network (with nodes representing people and edges representing “friendships”), it is interesting to look at how properties of individual nodes are distributed across the population. We can look at the distribution of any node-by-node property—the physical height of Twitter users, or the number of words of text per web page, for example. But in addition to demographic properties like height and length, we can also look at the distribution of network-type properties.

The *degree distribution* of a graph G shows, for each possible degree d , the number of nodes in G whose degree is d . While one might initially expect degree distributions to look similar to the distribution of heights, it turns out that the degree distribution of an online social network has very different properties. Figure 11.14 shows the degree distribution (in linear, log-log, and cumulative form) for members of the University of North Carolina.⁵

Figure 11.14 shows, for each value of k , the number of people who have precisely k Facebook friends. About 350 people have only 1 friend, which is the most common number of friends to have. There are about 750,000 friendships represented in this dataset; the *average degree* is ≈ 84 . But, looking at the far-right end of Figure 11.14(a) and 11.14(b), we see a handful of people with very high degrees: 2000, 2500, 3000, and even ≈ 3800 . One of the interesting facts about degree distributions in real social networks (or the web) is that there are people whose popularity is massively larger than average: the highest-degree person in this dataset is about $3800/84 \approx 45$ times more popular than average. (Imagine the tallest person at the University of North Carolina being 45 times taller than average!)

Significant research by computer scientists (and many others!) interested in the structure of social networks and the world-wide web has focused on this so-called *heavy-tailed degree distribution*.⁶ Some of the literature debates the particular form of this distribution; for example, whether the distribution has the particular form of a *power law*, where the number of people with degree k is roughly k^α for some small constant α , usually around 2.

Figure 11.14: The degree distribution of $\approx 18,000$ Facebook users at the University of North Carolina. Figure 11.14(b) shows a log-log plot of the same data as the linear plot in Figure 11.14(a). Figure 11.14(c) shows a log-log plot of the *cumulative degree distribution*: the number of people with degree $\geq k$, whereas Figures 11.14(a) and 11.14(b) showed the number with degree $=k$.

From the Facebook5 dataset, from Mason Porter via the International Network for Social Network Analysis:

⁵ Amanda L. Traud, Peter J. Mucha, and Mason A. Porter. Social structure of Facebook networks. *CoRR*, abs/1102.2166, 2011.

You can read more about power laws and heavy-tailed degree distributions:

⁶ David A. Easley and Jon M. Kleinberg. *Networks, Crowds, and Markets: Reasoning About a Highly Connected World*. Cambridge University Press, 2010.

COMPUTER SCIENCE CONNECTIONS

GRAPH DRAWING, GRAPH LAYOUTS, AND THE 9/11 MEMORIAL

Visual representations of most large graphs are too cluttered for a human viewer to process: there are just too many nodes and edges crammed into a small space to see much of anything. Visually presenting a graph like Facebook (billions of nodes, tens of billions of edges) without it looking like a grade-school scribble is daunting. But there is an entire subfield of computer science called *graph drawing*, which is devoted to taking networks and producing good—clear, aesthetic, informative—images of the networks.

In some large graphs, each node has a “natural location” and thus it is clear where on the page it should be placed. For example, graphs may represent data in which the nodes have a precise location situated in the physical world. When we have that kind of layout information for each node, presenting the graph well is easier. (See Figure 11.15.) But many large graphs do not have obvious coordinates associated with each node: while you and your college classmates do have geographic locations (dorm rooms), it’s not clear that your dorm really best describes “where” you fit in the social scene of your institution.

For graphs whose nodes don’t have obvious coordinates, we have to do something else. One approach that’s often used in graph drawing is to arrange the nodes based on a physics analogy, as follows. Imagine each node as a charged particle: any two nodes that are joined by an edge are pulled together by an attractive force, and any two nodes that are not joined by an edge are pushed apart by a repulsive force. Then figuring out how to place nodes on the page can be done by starting them in a random configuration and letting the attractive/repulsive forces move the nodes around until they’re “happy” in their current positions.

An idea like this one was actually used in designing the 9/11 memorial at the site of the World Trade Center. The memorial was designed with bronze panels inscribed with the 2982 names of victims. A team of computer scientists, architects, and visual artists collaborated to organize the names in a meaningful way. Families were invited to submit “meaningful adjacencies” between victims—which would cause two names to be as close together in the bronze panels as possible. (One of the other algorithmic issues regarding the layout of this memorial was that the designers wanted the names to be placed at evenly spaced intervals on the bronze panels; this constraint added to the computational complexity of the process.) The team used an algorithm to organize the names in an arrangement that respected these requests, which was then used in the final design of the memorial.⁷



Figure 11.15: A visualization of selected European train routes, where each node’s position corresponds to the city’s spatial location. Image reproduced with permission from RGBAlpha/Getty Images, Inc.

In addition to the broader news reports on the wrenching emotional and historical aspects of 9/11 Memorial, the algorithmic aspects of the memorial were also covered in the popular press. You can read more about it here:

⁷ Nick Paumgarten. The names. *The New Yorker*, 16 May 2011.

11.2.5 Exercises

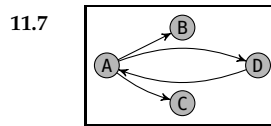
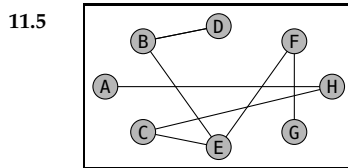
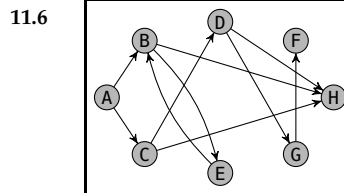
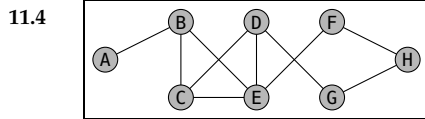
For each of the following, draw a graph $G = \langle V, E \rangle$ for the following sets of nodes and edges. Does it make sense to use a directed or undirected graph? Is the graph you've drawn simple?

11.1 nodes $V = \{1, 2, \dots, 10\}$; an edge connects x and y if $\gcd(x, y) = 1$.

11.2 nodes $V = \{1, 2, \dots, 10\}$; an edge connects x and y if x divides y .

11.3 nodes $V = \{1, 2, \dots, 10\}$; an edge connects x and y if $x < y$.

For the following undirected graphs, list the edges of the graph, and identify the node(s) with the highest degree. For the directed graphs, identify the node(s) with the highest in-degree, and the node(s) with the highest out-degree.



Consider a graph $G = \langle V, E \rangle$ with $n := |V|$ nodes. State your answers in terms of n . Justify.

11.8 If G is an undirected, simple graph, what's the largest that $|E|$ can be? The smallest?

11.9 If G is a directed, simple graph, what's the largest that $|E|$ can be? The smallest?

11.10 How do your answers to Exercise 11.9 change if self-loops are allowed?

11.11 How do your answers to Exercise 11.9 change if self-loops and parallel edges are allowed?

The anthropologist Robin Dunbar has argued that humans have a mental capacity for only ≈ 150 friends.⁸ (This argument is based in part on the physical size of the human brain, and cross-species comparisons; 150 is now occasionally known as Dunbar's Number.)

Suppose that Alice has exactly 150 friends, and each of her friends has exactly 150 friends—that is, a friend of Alice knows Alice and 149 other people. (Note that Alice's friends' sets of friends can overlap.) Let S denote the set of people that Alice knows directly or with whom Alice has a mutual friend.

11.12 What's the largest possible value of $|S|$?

11.13 What's the smallest possible value of $|S|$?

Continue to assume that everyone has precisely 150 friends. Let S_k denote the set of all people that Bob knows via a chain of k or fewer intermediate friends:

- Bob's friends are in S_0 ;
- the people in S_0 and the friends of people in S_0 are in S_1 ;
- the people in S_1 and the friends of people in S_1 are in S_2 ; and so forth.

11.14 Let $k \geq 0$ be arbitrary. What's the largest possible value of $|S_k|$?

11.15 Let $k \geq 0$ be arbitrary. What's the smallest possible $|S_k|$?

Prove the following properties of graphs, related to Theorem 11.1 or degree more generally:

11.16 Let u be a node in an undirected graph G . Prove that u 's degree is at most the sum of the degrees of u 's neighbors.

11.17 Prove Corollary 11.2: in an undirected graph $G = \langle V, E \rangle$, let n_{odd} denote the number of nodes whose degree is odd. Prove that n_{odd} is an even number. That is: prove that

$$|\{u \in V : \text{degree}(u) \bmod 2 = 1\}| \bmod 2 = 0.$$

11.18 Prove the analogy of Theorem 11.1 for directed graphs: for a directed graph $G = \langle V, E \rangle$,

$$\sum_{u \in V} \text{in-degree}(u) = \sum_{u \in V} \text{out-degree}(u) = |E|.$$

⁸ Robin Dunbar. *How Many Friends Does One Person Need?: Dunbar's Number and Other Evolutionary Quirks*. Harvard University Press, 2010. Thanks to Michael Kearns, from whom I learned a somewhat related version of these exercises.

A linked list is a data structure consisting of a collection of nodes, each of which contains two fields: a data field (whatever the node stores) and a next field that is either null or points to a node in the linked list. A particular node is designated as the head node. Note that a circular linked list in which a node points back to a previously encountered node meets this definition. See Figure 11.16.

Define a not-necessarily-simple directed graph $G = \langle V, E \rangle$, where V is the set of all nodes reachable by following any number of next pointers starting at the head node, and $\langle u, v \rangle \in E$ if u 's next field points to v . Observe that each node u in G has out-degree $d \in \{0, 1\}$.

Describe a 5-node linked list in which ...

11.19 ... every node has in-degree $d = 1$.

11.20 ... some node has in-degree $d = 2$.

11.21 ... the resulting graph G is not simple.

11.22 (This exercise is a tougher algorithmic challenge.) You are given access to the head node h of an n -node linked list. The value of n is unknown to you. The only operations permitted are (a) to save a node; (b) test whether two saved nodes are the same or different; and (c) given a node u , fetch the node pointed to by u .next. Give an algorithm to determine whether the given list is circular using only a constant amount of memory—that is, remembering only a constant number of nodes at a time.

A doubly linked list has n nodes with data and two pointers, previous and next, to other nodes (or null). (See Figure 11.17 for an example.) Let C_n denote an n -node doubly linked list with nodes $\{1, 2, \dots, n\}$, where, for each node u ,

- u 's next node is $v = (u \bmod n) + 1$
- v 's previous node is u .

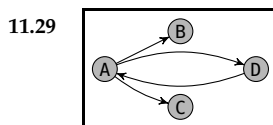
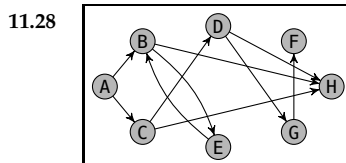
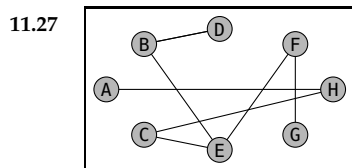
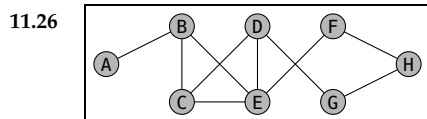
Define a directed graph $G_n = \langle V, E \rangle$, where V is the set $\{1, 2, \dots, n\}$ of nodes, and every node has two edges leaving it: one edge $\langle u, u.\text{next} \rangle$, and one edge $\langle u, u.\text{previous} \rangle$.

11.23 Draw G_5 .

11.24 Give an example of a G_n that contains a self-loop.

11.25 Give an example of a G_n that contains parallel edges.

Write down an adjacency list representing each of the following graphs.



Now give an adjacency matrix for the graphs shown in the above exercises:

11.30 Exercise 11.26

11.32 Exercise 11.28

11.31 Exercise 11.27

11.33 Exercise 11.29

11.34 Suppose that a (possibly directed or undirected) simple graph G is represented by an adjacency list. Suppose further that, for every node u in G , the list of (out-)neighbors of u has a different length. True or False: G must be a directed graph. Justify your answer.

11.35 Describe a directed graph G meeting the specifications of Exercise 11.34.

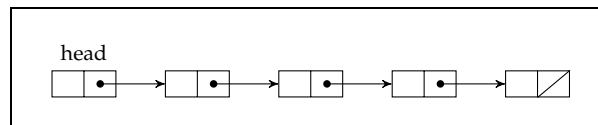


Figure 11.16: A linked list. Each rectangle is a node, and shows two fields: data on the left and next on the right.

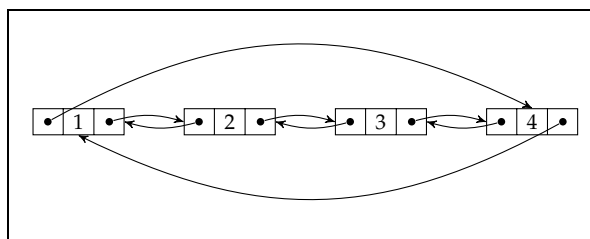


Figure 11.17: A doubly linked list. Each rectangle is a node, and shows three fields: previous on the left, data in the middle, and next on the right.

The density of a graph $G = \langle V, E \rangle$ is the fraction of all possible edges that actually exist: that is,

$$\text{density} = \frac{|E|}{\text{[your answer to the first part of Exercise 11.8/ Exercise 11.9]}}.$$

Taking it further: Informally, a dense graph is one for which most pairs of nodes are joined by an edge, and a sparse graph is one in which few pairs of nodes are joined by an edge. We will use these terms informally; a graph is dense if its density is close to 1, and sparse if its density is close to 0. Some people define graphs as dense if $|E| = \Theta(|V|^2)$ and as sparse if $|E| = O(|V|)$. (These asymptotic definitions only make sense for a family of graphs—one for each size n .) There are (families of) graphs that are neither sparse nor dense according to this definition; see Exercise 6.37.

As a function of n , what are the densities of the following undirected graphs, with nodes $V = \{1, 2, \dots, n\}$? (See Figure 11.18 for small versions of each of these graphs.)

11.36 an n -node path: $E = \{\{1, 2\}, \{2, 3\}, \dots, \{n-1, n\}\}$.

11.37 an n -node cycle: $E = \{\{1, 2\}, \{2, 3\}, \dots, \{n-1, n\}, \{n, 1\}\}$.

11.38 $\frac{n}{3}$ disconnected triangles (assume that $n \bmod 3 = 0$):

$$E = \underbrace{\{\{1, 2\}, \{2, 3\}, \{3, 1\}\}}_{\text{triangle on } 1, 2, 3}, \underbrace{\{\{4, 5\}, \{5, 6\}, \{6, 4\}\}}_{\text{triangle on } 4, 5, 6}, \dots, \underbrace{\{\{n-2, n-1\}, \{n-1, n\}, \{n, n-2\}\}}_{\text{triangle on } n-2, n-1, n}.$$

11.39 3 separate $\frac{n}{3}$ -node cliques (assume that $n \bmod 3 = 0$): $E = \{\{x, y\} : x \bmod 3 = y \bmod 3\}$.

A hypercube H_n is a graph in which the 2^n different nodes are all elements of $\{0, 1\}^n$. There is an edge between x and y if they differ in only one bit position. (Using the language of Chapter 4.2, there's an edge between any two nodes whose Hamming distance is 1.)

11.40 Draw H_3 .

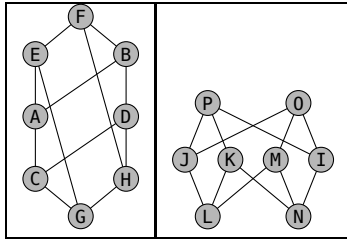
11.41 Write down an adjacency list for H_4 .

11.42 Write down an adjacency matrix for H_4 .

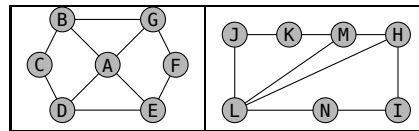
11.43 In terms of n , how many edges does H_n have? What is its density?

Decide whether the following pairs of graphs are isomorphic, and prove your answers.

11.44



11.45



11.46 $G_1 = \langle V_1, E_1 \rangle$, where $V_1 = \{10, 11, 12, 13, 14, 15\}$ and $\langle x, y \rangle \in E_1$ if and only if x and y are not relatively prime.

$G_2 = \langle V_2, E_2 \rangle$, where $V_2 = \{20, 21, 22, 23, 24, 25\}$ and $\langle x, y \rangle \in E_2$ if and only if x and y are not relatively prime.

Prove or disprove the following claims about isomorphism:

11.47 All 5-node graphs with degrees 1, 1, 1, 1, and 0 are isomorphic.

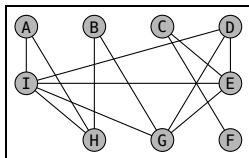
11.48 All 5-node graphs with degrees 4, 4, 4, 3, and 3 are isomorphic.

11.49 All 5-node graphs with degrees 3, 3, 2, 2, and 2 are isomorphic.

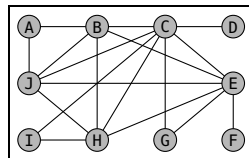
11.50 All n -node, 3-regular graphs are isomorphic.

The computational problem of finding the largest clique (complete graph) that's a subgraph of a given graph G is believed to be very difficult. But for small graphs it's possible to do, even by brute force. For each of the following graphs, identify the size of the largest clique that's a subgraph of the given graph.

11.51



11.52



11.53

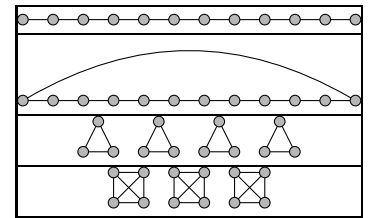
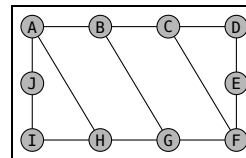


Figure 11.18: A 12-node path, cycle, collection of $\frac{n}{3}$ triangles, and collection of three $\frac{n}{3}$ -node cliques.

11.54 Consider the collaboration network (see Example 11.14) in Figure 11.19. Assuming that the nodes correspond to actors in movies, what is the *smallest number* of movies that could possibly have generated this collaboration network?

11.55 Are you certain that there weren't more movies than [your answer to the previous exercise] that generated this graph? Explain.

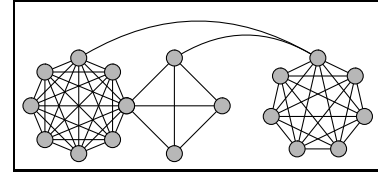


Figure 11.19: A collaboration network.

For which integers n are the following graphs bipartite? Prove your answers.

11.56 $V = \{1, 2, \dots, n\}; E = \{\langle i, i-1 \rangle : i \geq 2\}$.

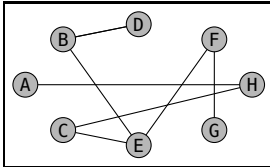
11.57 $V = \{0, 1, \dots, n-1\}; E = \{\langle i, i+1 \bmod n \rangle : i \geq 1\}$.

11.58 \mathcal{K}_n . That is, a complete graph of n nodes: $V = \{1, 2, \dots, n\}; E = \{\langle u, v \rangle : u \in V \text{ and } v \in V\}$.

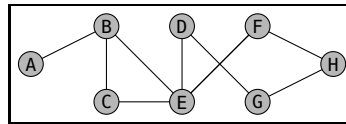
11.59 $V = \{0, 1, \dots, 2n-1\}; E = \{\langle i, (i+n) \bmod 2n \rangle : i \in V\}$.

Are either of the following graphs bipartite? Explain.

11.60



11.61



Consider a bipartite graph with a set L of nodes in the left column and a set of nodes R on the right column, where $|L| = |R|$. Prove or disprove the following claims:

11.62 The sum of the degrees of the nodes in L must equal the sum of the degrees of the nodes in R .

11.63 The sum of the degrees of the nodes in L must be even.

11.64 The sum of the degrees of all nodes (that is, all nodes in $L \cup R$) must be an even number.

Suppose that G is a complete bipartite graph with n nodes—that is, $G = \mathcal{K}_{|L|,|R|}$ for $|L| + |R| = n$.

11.65 What's the largest number of edges that can appear in G ?

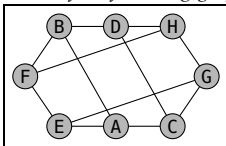
11.66 What's the smallest number of edges that can appear in G ? (Careful!)

11.67 Prove or disprove: any graph that does not contain a triangle (that is, three nodes a, b , and c with the edges $\{a, b\}$ and $\{b, c\}$ and $\{c, a\}$ in the graph) as a subgraph is bipartite.

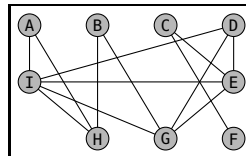
11.68 Definition 11.16 describes a regular undirected graph. In a *directed* regular graph, we require that there be two integers d_{in} and d_{out} such that every node's in-degree is d_{in} and every node's out-degree is d_{out} . Prove that we must have $d_{\text{in}} = d_{\text{out}}$.

Show that both of the following graphs are planar.

11.69



11.70



11.71 Prove that any 2-regular graph is planar.

11.3 Paths, Connectivity, and Distances

Well, you can go west to the next intersection, get onto the turnpike, go north through the toll gate at Augusta, 'til you come to that intersection ... well, no. You keep right on this tar road; it changes to dirt now and again. Just keep the river on your left. You'll come to a crossroads and ... let me see. Then again, you can take that scenic coastal route that the tourists use. And after you get to Bucksport ... well, let me see now. Millinocket. Come to think of it, you can't get there from here.

Marshall Dodge (1935–1982) and Robert Bryan (b. 1931), “Which Way to Millinocket?”
Bert and I (1958)

One of the most basic questions that one can ask about a graph is whether it is possible to get from some given node s to some given node t by following a sequence of edges. Is there some chain of friends that connects Barack Obama to Phil Collins? Can you get from Missoula to Madison by car? (And, if there is a way to get from s to t , what is the *shortest* way to get there?) These basic questions concern the existence of *paths* in the graph:

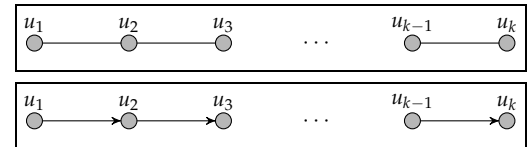


Figure 11.20: Paths in undirected and directed graphs.

Definition 11.18 (Path)

Consider a (directed or undirected) graph $G = \langle V, E \rangle$. A path in G is a sequence $\langle u_1, u_2, \dots, u_k \rangle$ of $k \geq 1$ nodes such that:

- $u_i \in V$ for every $i \in \{1, \dots, k\}$, and
- $\langle u_i, u_{i+1} \rangle \in E$ for every $i \in \{1, \dots, k-1\}$.

(See Figure 11.20.) We say that such a sequence of nodes is a path from u_1 to u_k , and that this path has length $k-1$. We also say that this path traverses the edges $\langle u_i, u_{i+1} \rangle$.

(Note that this definition includes both directed and undirected graphs: if the edges are directed, we have to follow them “in the right direction.”) For example, in both of the graphs shown in Figure 11.21, there is no path from A to X. But, in both, the sequence $\langle A, C, E, Z \rangle$ is a path of length 3 from A to Z. In both cases, the edges traversed by the path are $\{\langle A, C \rangle, \langle C, E \rangle, \langle E, Z \rangle\}$. Notice that the length of a path is the number of *edges* that it traverses, which is one fewer than the number of nodes in the path.

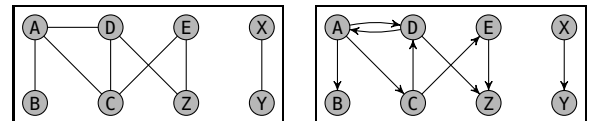


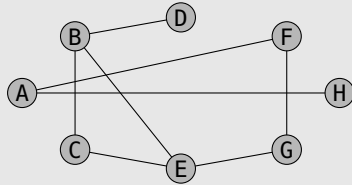
Figure 11.21: Two graphs with paths from A to Z.

Taking it further: A common mistake made by novice (and not-so-novice) programmers is an *off-by-one error* in specifying the bounds on a loop, by iterating either one time too many or one time too few. These errors are also sometimes called *fencepost errors*: if you build a 10-yard fence with posts placed every yard, then there are *eleven* fenceposts (at yard 0, yard 1, ..., yard 10). Be careful! A path $\langle A, C, E, Z \rangle$ contains four nodes, but it traverses three edges ($A \rightarrow C$, $C \rightarrow E$, and $E \rightarrow Z$) and has length 3.

Here's an example of finding paths in a small graph:

Example 11.21 (Finding paths)

Problem: Consider the following undirected graph:



1. Is there a path from node H to node E?
2. Name three different paths from node D to node F. What is the length of each path?

Solution: 1. Yes; $\langle H, A, F, G, E \rangle$ is a path from node H to E.

2. The following sequences are paths from D to F:

- $\langle D, B, E, G, F \rangle$, which has length 4.
- $\langle D, B, C, E, G, F \rangle$, which has length 5.

Finding a third path might seem harder, but Definition 11.18 did not require that the nodes in a path be distinct from each other. (In other words, nothing forbade the repetition of nodes in a path.) So a third path from D to F is:

- $\langle D, B, C, E, B, C, E, G, F \rangle$, which has length 8.

We will often restrict our attention to paths that never go back to a vertex that they've already visited, which are called *simple paths*:

Definition 11.19 (Simple Path)

A path $\langle u_1, u_2, \dots, u_k \rangle$ is simple if all of the nodes u_1, \dots, u_k are distinct.

Of the three paths identified in Example 11.21, the first two are simple paths, but the third path is not simple because it repeated nodes $\{B, C, E\}$.

11.3.1 Connectivity in Undirected Graphs

The most basic question about two nodes in a graph is whether it's possible to get from one to another—that is, are these two nodes *connected*? We start with a formal definition of connectivity for undirected graphs, because the relevant notions are simpler in the undirected setting.

Definition 11.20 (Connected nodes and connected graphs)

Let $G = \langle V, E \rangle$ be an undirected graph.

- Two nodes $u \in V$ and $v \in V$ are connected if there exists a path from u to v .
- The graph G is connected if u and v are connected for any two nodes $u \in V$ and $v \in V$.
- The graph G is called disconnected if it is not connected.

For example, Figure 11.22 shows one disconnected graph—there's no path from A to H, for example—and one connected graph. You can check that the second graph is connected by testing all pairs of nodes. (Exercise 11.87 asks you to show that connectivity is symmetric in an undirected graph: if there exists a path from u to v , then there exists a path from v to u .)

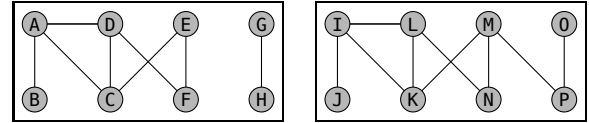
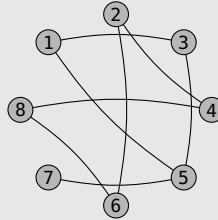


Figure 11.22: A disconnected and connected undirected graph.

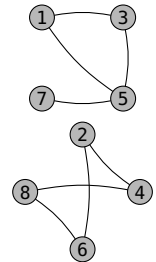
Example 11.22 (Connectivity of an undirected graph)

Problem: Is the following graph connected?



Solution: No: odd-numbered nodes have edges only to other odd-numbered nodes, and even-numbered nodes have edges only to other even-numbered nodes. So there is no path from, for example, node 1 to node 2; this graph is disconnected.

Problem-solving tip: Sometimes it's very helpful to redraw a graph that you're given, with nodes placed more meaningfully. For example, the graph from Example 11.22 can be redrawn as



just by sliding the even-numbered nodes down. This visualization makes it clear that the graph is disconnected.

CONNECTED COMPONENTS

More generally, we will talk about the *connected components* of an undirected graph $G = \langle V, E \rangle$ —"subsections" of the graph in which all pairs of nodes are connected.

Definition 11.21 (Connected component)

In an undirected graph $G = \langle V, E \rangle$, a connected component is a set $C \subseteq V$ such that:

- (i) any two nodes $s \in C$ and $t \in C$ are connected.
- (ii) for any node $x \in V - C$, adding x to C would make (i) false.

A subset $C \subseteq V$ of nodes is a connected component of an undirected graph $G = \langle V, E \rangle$ if, intuitively, it forms its own "section" of the graph: any two nodes in C are connected, and no node in C is connected to any node not in C . For example, Figure 11.23 shows a graph with three connected components—one with 4 nodes, one with 3 nodes, and one with just a single node.

Note that we could have defined a "connected graph" in terms of the definition of connected components (instead of Definition 11.20): an undirected graph $G = \langle V, E \rangle$ is *connected* if it contains only one connected component, namely the entire node set V .

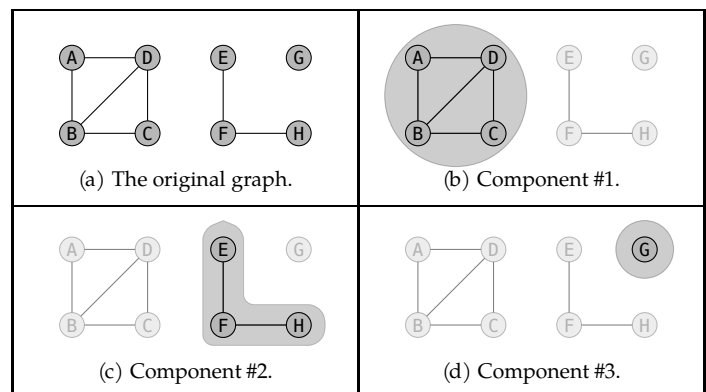
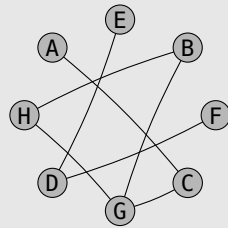


Figure 11.23: A graph's connected components.

Example 11.23 (Connected components of an undirected graph)

Problem: What are the connected components of the following graph?



Solution: The set $S = \{A, B, C, G, H\}$ is a connected component; there are paths from every node $u \in S$ to every node $v \in S$, and furthermore no node in S is connected to any node not in S . To be thorough, here are paths connecting each pair of nodes from S :

	A	B	C	G	H
A	$\langle A \rangle$	$\langle A, C, G, B \rangle$	$\langle A, C \rangle$	$\langle A, C, G \rangle$	$\langle A, C, G, H \rangle$
B		$\langle B \rangle$	$\langle B, G, C \rangle$	$\langle B, G \rangle$	$\langle B, H \rangle$
C			$\langle C \rangle$	$\langle C, G \rangle$	$\langle C, G, H \rangle$
G				$\langle G \rangle$	$\langle G, H \rangle$
H					$\langle H \rangle$

Note that we haven't bothered to write down a path from u to v when we'd already recorded a path from v to u , because the graph is undirected and paths are symmetric. We also had many choices of paths for many of these entries: for example, other paths from B to H included $\langle B, G, H \rangle$ or $\langle B, G, H, B, G, H \rangle$.

There's a second connected component in the graph: the nodes $\{D, E, F\}$. It's easy to check that both clauses of Definition 11.21 are also satisfied for this set.

Observe that, in *any* undirected graph $G = \langle V, E \rangle$, there is a path from each node $u \in V$ to itself. Namely, the path is $\langle u \rangle$, and it has length 0. Check Definition 11.18!

Taking it further: There are many computational settings in which undirected paths are relevant; here's one example, in brief. In *computer vision*, we try to build algorithms to process—"understand," even—images. For example, before it can decide how to react to them, a self-driving car must partition the image of the world from a front-facing camera into separate objects: painted lines on the road, trees, other cars, pedestrians, etc. Here's a crude way to get started (real systems use far more sophisticated techniques): define a graph whose nodes are the image's pixels; there is an edge between pixels p and p' if (i) the two pixels are adjacent in the image, and (ii) the colors of p and p' are within a threshold of acceptable difference. The connected components of this graph are a (very rough!) approximation to the "objects" in the image.

This description misses all sorts of crucial features of good algorithms for the image-segmentation problem, but even as stated it may be familiar from a different context: the "region fill" tool in image-manipulation software uses something very much like what we've just described.

11.3.2 Connectivity in Directed Graphs

Recall that we have to follow edges "in the right direction" in a directed graph G : as in Definition 11.18, a path from u_1 to u_k in G is a sequence $\langle u_1, u_2, \dots, u_k \rangle$ where every pair $\langle u_i, u_{i+1} \rangle$ is an edge in G . Thus notions of connectivity in directed graphs are more

complicated: the existence of a path from u to v does not imply the existence of a path from v to u . We will speak of a node t as being *reachable* from a node s if it's possible to go from s to t , and of pairs of nodes as being *strongly connected* when it's possible to “go in both directions” between them:

Definition 11.22 (Reachability and strongly connected nodes/graphs)

Let $G = \langle V, E \rangle$ be a directed graph.

- A node $u \in V$ is *reachable* from a node $v \in V$ if there is a directed path from u to v .
- Two nodes $u \in V$ and $v \in V$ are *strongly connected* if u is reachable from v , and v is reachable from u .
- The graph G is *strongly connected* if every pair of nodes in V is strongly connected.

For example, you can check that the first graph in Figure 11.24 is strongly connected by testing for directed paths between all pairs of nodes, in both directions. But the second graph in Figure 11.24 is not strongly connected: there's no path from any node in the right-hand side (nodes $\{M, N, O, P\}$) to any node in the left-hand side (nodes $\{I, J, K, L\}$).

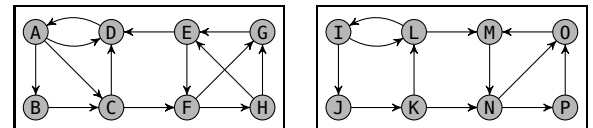


Figure 11.24: Two directed graphs, one that's strongly connected and one that's not.

STRONGLY CONNECTED COMPONENTS

As with undirected graphs, for a directed graph we will divide the graph into “sections”—subsets of the nodes—each of which is strongly connected. These sections are called *strongly connected components* of the graph:

Definition 11.23 (Strongly connected component)

In a directed graph $G = \langle V, E \rangle$, a *strongly connected component* (SCC) is a set $C \subseteq V$ such that:

- any two nodes $s \in C$ and $t \in C$ are strongly connected.
- for any node $x \in V - C$, adding x to C would make (i) false.

Figure 11.25 shows an example of a directed graph G and the three strongly connected components in G . The easiest strongly connected component to identify is $\{A, B, C, D\}$: we can go counterclockwise around the loop $A \rightarrow B \rightarrow C \rightarrow D \rightarrow A$, so we can go from any one of these four nodes to any other, and we can't get from any of these four nodes to any of the other nodes. The other two strongly connected components are $\{E, F, H\}$ and, separately, $\{G\}$ on its own. The reason is that G is not strongly connected to any other node: we can't get *from* G to any other node. (We can go around the $E \rightarrow F \rightarrow H \rightarrow E$ loop, so these three nodes are together in the other strongly connected component.)

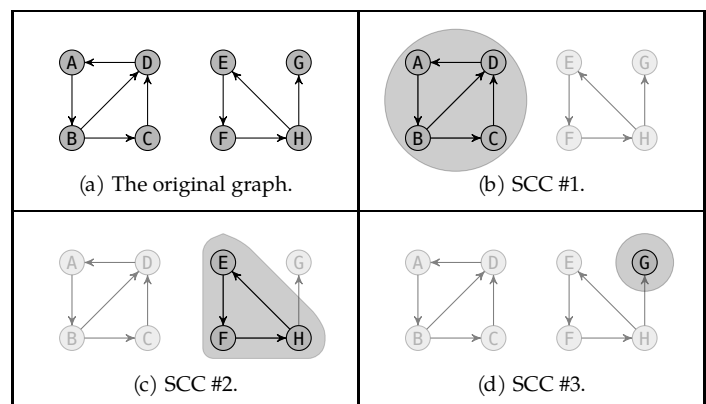
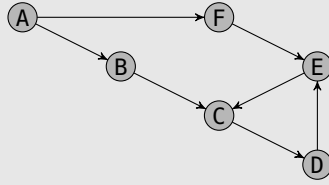


Figure 11.25: A graph and its connected components.

Here's another example of finding strongly connected components:

Example 11.24 (Finding strongly connected components)

Problem: What are the strongly connected components of the following graph?



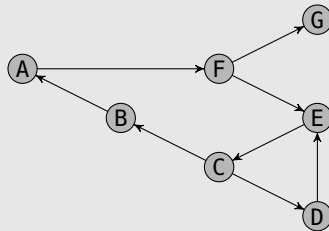
Solution: The three nodes $\{C, D, E\}$ form a strongly connected component: there is a path from any one of them to any other of them ($C \rightarrow D \rightarrow E \rightarrow C \rightarrow D \rightarrow E \cdots$), and furthermore there is no path from any $\{C, D, E\}$ to any other node in the graph.

In fact, every other node in the graph is alone in a strongly connected component by itself. For example, while there is a path from A to every node in the graph, there is no path *from* any other node to A. (There is a path from A to A, so the set $\{A\}$ is a strongly connected component.) Thus the four strongly connected components of the graph are $\{A\}$, $\{B\}$, $\{F\}$, and $\{C, D, E\}$.

Here's an example that shows why the second clause of Definition 11.23 is crucial:

Example 11.25 (A non-SCC)

Problem: In the following graph, the set $S := \{A, B, C, E, F\}$ is *not* a strongly connected component. Why not?



Solution: It is indeed the case that there is a path in both directions between any two nodes in S : we can just keep “going around” clockwise in S and we eventually reach every other node in S . So S satisfies Definition 11.23(i). But it fails to satisfy Definition 11.23(ii): if we considered the set $S^+ := S \cup \{D\}$, it is still the case that there is a path in both directions between any nodes in S^+ . Thus S is not a strongly connected component!

On the other hand, $S^+ = \{A, B, C, D, E, F\}$ is a strongly connected component: we can't add any other node (specifically G; it's the only other node) to S^+ without falsifying this property—because there's no path from G to A, for example. Thus the two strongly connected components are $\{A, B, C, D, E, F\}$ and $\{G\}$.

Taking it further: There are many computational settings in which directed paths, reachability, and strongly connected components are relevant. For example, for a spreadsheet, consider a directed graph whose nodes are the spreadsheet's cells, and an edge $\langle u, v \rangle$ indicates that u 's contents affect the contents of cell v ; when a user changes the content of cell c , we must update all cells that are reachable from node c . For a chess-playing program, consider a directed graph whose nodes are board configurations, and there's an edge $\langle u, v \rangle$ if a legal move in u can result in v ; any configuration u that's unreachable from the starting board configuration can never occur in chess, and thus your program doesn't have to bother evaluating what move to make in position u .

See p. 1142 for a discussion of another application of reachability and strongly connected components: the structure of the world-wide web, understood with respect to the directed paths in the graph defined by the pages and the hyperlinks of the web.

11.3.3 Shortest Paths and Distance

So far we have concentrated on the basic question of connectivity: for a given pair of nodes, does any path exist from one node to the other? Here we address a more refined question: what is the *shortest* path that goes from one node to the next?

Definition 11.24 (Shortest Paths)

Let $G = \langle V, E \rangle$ be a graph (undirected or directed), and let $s \in V$ and $t \in V$ be two nodes. A path from s to t is a shortest path if its length is the smallest out of all s -to- t paths.

(Recall that the *length* of a path $\langle u_1, u_2, \dots, u_k \rangle$ is $k - 1$, the number of edges that it traverses.) Observe that there may be more than one shortest path from a node s to a node t , if there are multiple paths that are tied in length.

Definition 11.25 (Distance)

The distance from s to t is the length of a shortest path from s to t . If there is no path from s to t , then we say that the distance from s to t is infinite (written as " ∞ ").

For example, consider the undirected graph in Figure 11.26. We have the following distances from node A in this graph:

A	B	C	D	E	F
0	1	2	2	1	1

The distance from A to A is 0 because $\langle A \rangle$ is a path from A to A. This graph also has an example of a pair of nodes connected by two different shortest paths, going from A to C (via either B or E).

For the directed graph in Figure 11.27, we have the following distances from node G:

G	H	I	J	K	L
0	1	2	3	1	∞

Again, there's a path from G to G of length zero, so the distance from G to G is 0. Note that there's no G-to-J path of length two (because the edge from J to K goes in the wrong direction), so the distance from G to J is 3 (via K and I, or via H and I). Similarly, there is no directed path from G to L, so the distance is infinite.

Here's another example of finding shortest paths in a small graph:

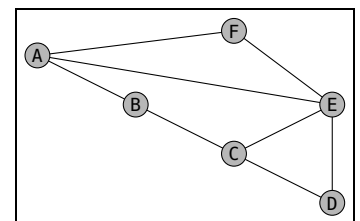


Figure 11.26: An undirected graph.

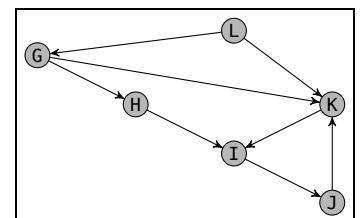


Figure 11.27: A directed graph.

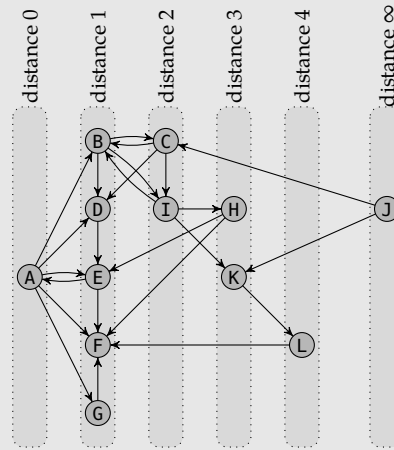
Example 11.26 (Shortest paths in directed graphs)

Problem: Find the shortest path from A to L in the graph with this adjacency list:

A:	B, D, E, F, G
B:	C, D, I
C:	B, D, I
D:	E
E:	A, F
F:	G
G:	F
H:	E, F
I:	B, H, K
J:	C, K
K:	L
L:	F

Solution:

The nodes at distance 1 from A are B, D, E, F, and G. There's no edge from any of those nodes to L—or indeed to K, which is L's only in-neighbor. Thus the distance from A to L cannot be any smaller than 4. But there is an edge from I to K, and one from B to I. We can assemble these edges into the path $\langle A, B, I, K, L \rangle$. This path has length 4. So the distance from A to L is 4. (Drawing the graph, as on the right, with nodes arranged by their distance from A, can make these facts easier to see.)



Problem-solving tip: In solving any graph problem with a small graph, a good first move is to draw the graph.

11.3.4 Finding Paths: Breadth-First Search (BFS)

There are many aspects of graphs that are valuable for interesting computational applications, but perhaps the single most important graph algorithm is *breadth-first search* (BFS). BFS is a path-finding algorithm: it explores outward from a given source node s in a given graph G until it finds every node reachable from s in G . BFS can be used to solve all sorts of graph-related problems, as we'll see.

Here's the intuition of the algorithm. (See Figure 11.28.) We maintain a set L of nodes that are reachable from the given node s (the shaded nodes in Figure 11.28). To start, we set $L := \{s\}$. Now we find all as-yet-undiscovered neighbors of nodes in L , and add those nodes (the dark-shaded nodes in Figure 11.28) to L : if $\langle u, v \rangle \in E$ and you can reach the node u from s , then you can also reach v from s , via u . But now we've found some more nodes that can be reached from s , which means that we can also reach any nodes that are directly connected to *them* from s . So we'll repeat that process with the updated list L . And we'll do it again, and again, and again, until we stop finding new nodes.

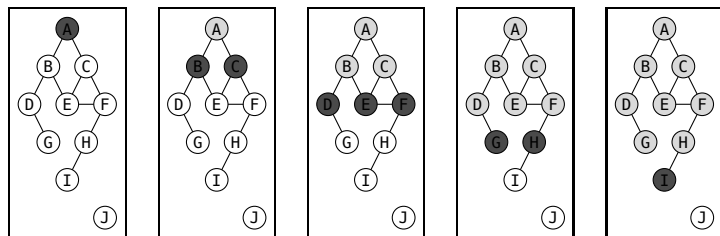


Figure 11.28:
The intuition of
breadth-first search:
the steps of BFS
on a small graph,
starting at node A.

Observe that BFS discovers nodes in order of their *distance* from the source node. Every expansion of L takes the full breadth of the frontier and expands it out by one more “layer” in the graph. (That’s why the algorithm is called breadth-first search.) You can think of BFS as throwing a pebble onto the graph at the node s , and then watching the ripples expanding out from s .

Breadth-first search is presented more formally in Figure 11.29. (While we’ve described BFS in terms of undirected graphs for simplicity, it works equally well for directed graphs. The only change is that Line 6 should say “for every *out*-neighbor” for a directed graph.)

Here’s another example of breadth-first search in action, running the algorithm in full detail (precisely as specified in Figure 11.29):

Breadth-First Search (BFS):

Input: a graph $G = \langle V, E \rangle$ and a source node $s \in V$

Output: the set of nodes reachable from s in G

```

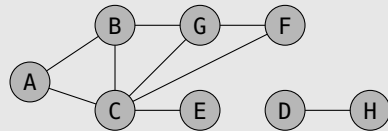
1:  $Frontier := \langle s \rangle$ 
   // Frontier will be a list of nodes to process, in order.
2:  $Known := \emptyset$ 
   // Known will be the set of already-processed nodes.
3: while  $Frontier$  is nonempty:
4:    $u :=$  the first node in  $Frontier$ 
5:   remove  $u$  from  $Frontier$ 
6:   for every neighbor  $v$  of  $u$ :
7:     if  $v$  is in neither  $Frontier$  nor  $Known$  then
8:       add  $v$  to the end of  $Frontier$ 
9:   add  $u$  to  $Known$ 
10: return  $Known$ 

```

Figure 11.29: The pseudocode for breadth-first search.

Example 11.27 (Sample run of BFS, in detail)

We’ll trace BFS starting at node A in the following graph (shown here in the form of a picture and as an adjacency list):



```

A: B, C
B: A, C, G
C: A, B, E, F, G
D: H
E: C
F: C, G
G: B, C, F
H: D

```

● = Frontier ● = just moved from Frontier to Known ● = Known ○ = neither Known nor Frontier	Known	Frontier	Explanation
	{}	$\langle A \rangle$	initialization (Lines 1–2)
	{A}	$\langle B, C \rangle$	processing A (Lines 4–9)
	{A, B}	$\langle C, G \rangle$	processing B (Lines 4–9)
	{A, B, C}	$\langle G, E, F \rangle$	processing C (Lines 4–9)
	{A, B, C, G}	$\langle E, F \rangle$	processing G (Lines 4–9)
	{A, B, C, G, E}	$\langle F \rangle$	processing E (Lines 4–9)
	{A, B, C, G, E, F}	$\langle \rangle$	processing F (Lines 4–9)

Because *Frontier* is now empty, the **while** loop in BFS terminates. The algorithm returns the set *Known*, $\{A, B, C, G, E, F\}$.

CORRECTNESS OF BFS

We'll prove two important properties of BFS. The first is *correctness*: the set that BFS returns is precisely those nodes that are reachable from the starting node. The second is *efficiency*: BFS finds this set quickly. The first claim might seem obvious—and thus proving it may feel annoyingly pedantic—but there's a bit of subtlety to the argument, and it's good practice at using induction in proofs besides.

Theorem 11.3 (Correctness of BFS)

Let $G = \langle V, E \rangle$ be any graph, and let $s \in V$ be an arbitrary node. Then the set of nodes discovered by $\text{BFS}(G, s)$ is exactly $\{t \in V : t \text{ is reachable from } s \text{ in } G\}$.

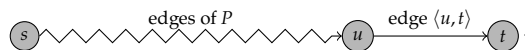
Proof. We'll prove the result by showing two set inclusions: the discovered nodes form a subset of the reachable nodes, and the reachable nodes form a subset of the discovered nodes. Both proofs will use induction, though on different quantities.

Claim #1: $\text{BFS}(G, s) \subseteq \{t \in V : t \text{ is reachable from } s \text{ in } G\}$. By inspection, we see that (i) BFS returns the set of nodes that end up in the *Known* set, and (ii) the only way that a node ends up in *Known* is having previously been in *Frontier*. Thus it will suffice to prove the following property for all $k \geq 0$, by strong induction on k :

$Q(k) :=$ if a node $t \in V$ is added to the list *Frontier* during the k th iteration of the **while** loop of BFS, then there is a path from s to t .

Base case ($k = 0$): If the node t was added to *Frontier* during the 0th iteration of the **while** loop—that is, before the **while** loop begins—then t was added in Line 1 of **BFS**. Therefore t is actually the node s itself. There is a path from s to s itself in any graph, and thus $Q(0)$ holds.

Inductive case ($k \geq 0$): We assume the inductive hypotheses $Q(0), \dots, Q(k-1)$, and we must prove $Q(k)$. Consider a node t that was added to *Frontier* during the k th iteration of the **while** loop—in other words, t was added in the **for** loop (Lines 6–8) because t is a neighbor of some node u that was already in *Frontier*. That is, we know that $\langle u, t \rangle \in E$ and that u was added to *Frontier* in the (k') th iteration, for some $k' < k$. By the inductive hypothesis $Q(k')$, there is a path P from s to u . Therefore there is a path from s to t , too:



Claim #2: $\text{BFS}(G, s) \supseteq \{t \in V : t \text{ is reachable from } s \text{ in } G\}$. If a node t is reachable from s in G , then by definition the distance from s to t is some integer $d \geq 0$. Furthermore, by inspection of the algorithm, we see that any node that's added to *Frontier* is eventually moved to *Known*. Thus it will suffice to prove the following property for all $d \geq 0$, by (weak) induction on d :

$R(d) :=$ if a node $t \in V$ at distance d from s , then t is eventually added to *Frontier*.

Base case ($d = 0$): We must prove $R(0)$: any node t at distance 0 is eventually added to *Frontier*. But the only node at distance 0 from s is s itself, and **BFS** adds s itself to *Frontier* in Line 1 of the algorithm.

Problem-solving tip: The hard part here is figuring out *what quantity* to do induction. One way to approach this question is to figure out a recursive way of stating the correctness claim.

Q: why is there a path to every node added to *Frontier*?
(A: there was a path to every previous node in *Frontier*, and there's an edge from some previously added node to this one!)

Q: why is every node u reachable from s eventually added to *Frontier*?
(A: because a neighbor of u that's closer to s is eventually added to *Frontier*, and every neighbor of a node in *Frontier* is eventually added to *Frontier*!)

Inductive case ($d \geq 1$): We assume the inductive hypothesis $R(d - 1)$, and we must prove $R(d)$. Let t be a node at distance d from s . Then by definition of distance there is a shortest path P of length d from s to t . Let u be the node immediately before t in P . Then the distance from s to u must be $d - 1$, and therefore by the inductive hypothesis $R(d - 1)$ the node u is added to *Frontier* in some iteration of the **while** loop. There are at most $|V|$ iterations of the loop, and thus eventually u is the first node in *Frontier*. In that iteration, the node t is added to *Frontier* (if it had not already been added). Thus $R(d)$ follows. \square

(In the exercises, you'll show how to modify BFS so that it actually computes *distances* from s , using an idea very similar to the proof of Claim #2 of Theorem 11.3.)

RUNNING TIME OF BFS

Theorem 11.4 (Efficiency of BFS)

For a graph $G = \langle V, E \rangle$ represented using an adjacency list, BFS takes $\Theta(|V| + |E|)$ time.

Proof. See Figure 11.30 for a reminder of the algorithm. Lines 1, 2, and 10 take $\Theta(1)$ time, so the only question is how long the **while** loop takes. In the worst case, every node in the graph is reachable from the node from which BFS is run. In this case, there is one iteration of the **while** loop for every node $u \in V$. How long does the body of the **while** loop (Lines 4–9) take for a particular node u ?

- Lines 4, 5, and 9 take $\Theta(1)$ time.
- The **for** loop in Lines 6–8 has one iteration for *each neighbor* of u . (In an adjacency list, the loop simply steps through the list of neighbors, one by one.) Each **for**-loop iteration takes $\Theta(1)$ time, and there are $\text{degree}(u)$ iterations for node u .

Breadth-First Search (BFS):

Input: a graph $G = \langle V, E \rangle$ and a source node $s \in V$
Output: the set of nodes reachable from s in G

```

1: Frontier :=  $\langle s \rangle$ 
2: Known :=  $\emptyset$ 
3: while Frontier is nonempty:
4:    $u$  := the first node in Frontier
5:   remove  $u$  from Frontier
6:   for every neighbor  $v$  of  $u$ :
7:     if  $v$  is in neither Frontier nor Known then
8:       add  $v$  to the end of Frontier
9:   add  $u$  to Known
10: return Known

```

Figure 11.30: A reminder of BFS.

Therefore, ignoring multiplicative constants, the worst-case running time of BFS is

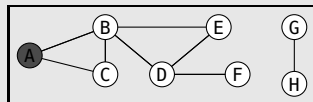
$$\begin{aligned}
 & 1 + \sum_{u \in V} [1 + \text{degree}(u)] \\
 &= 1 + \left[\sum_{u \in V} 1 \right] + \left[\sum_{u \in V} \text{degree}(u) \right] && \text{rearranging the summation} \\
 &= 1 + |V| + 2|E| && \text{or } 1 + |V| + |E| \text{ for a directed graph} && \text{Theorem 11.1/Exercise 11.18} \\
 &= \Theta(|V| + |E|). && \square
 \end{aligned}$$

Taking it further: BFS arises in applications throughout computer science, from network routing to artificial intelligence. Another application of BFS occurs (hidden from your view) as you use programming languages like Python and Java, through a language feature called *garbage collection*. In garbage-collected languages, when you as a programmer are done using whatever data you've stored in some chunk of memory, you just "drop it on the floor"; the "garbage collector" comes along to reclaim that memory for other use in the future of your program. The garbage collector runs BFS-like algorithms to determine whether a particular piece of memory is actually trash. See p. 1143 for more.

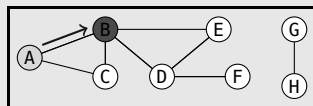
11.3.5 Finding Paths: Depth-First Search (DFS)

Another important algorithm for exploring graphs is called *depth-first search (DFS)*, which can be described informally as follows. Instead of exploring outward from the source node s in “layers” as in BFS, we will try to explore a new node at every stage of the search. We start at s , and at every stage we move to an unvisited neighbor of our current node. If at any stage we’re stuck at a node u that has no unvisited neighbors, we go back from u to the node from which we first reached u and continue exploring from there. Here is an example of DFS in a small graph, informally:

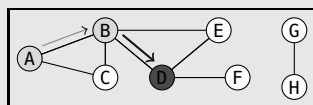
Example 11.28 (Sample run of depth-first search)



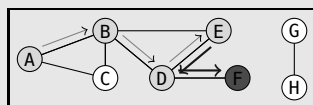
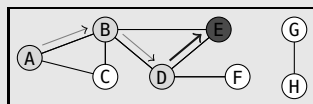
We start exploring node A; in each frame, the dark-shaded node is the current node.



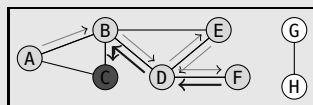
Previously discovered nodes are lightly shaded. Arrows indicate the steps of the exploration.



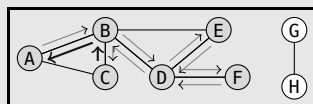
In each of the first four frames, we move from the current node to a neighbor that is unexplored. (We pick the alphabetically first node if there’s a choice.)



The current node E has no unvisited neighbors, so we backtrack from E to D to find D’s unvisited neighbor F.



We backtrack from F to D to B to discover the new node C.



We backtrack from C to B to A; there are no further unexplored nodes from any of these nodes, and thus the algorithm terminates.

Intuitively, depth-first search is a close match for the way that you would explore a maze: you start at the entrance, follow a passageway to a location you’ve never visited before; using breadcrumbs or a pencil, you remember where you’ve been and backtrack if you get stuck. You may have heard of another algorithm for mazes:

Place your right hand on the wall as you go in the entrance. Continue to walk forward, always keeping your right hand on the wall. Eventually, you will get out of the maze.

In fact, this right-hand-on-the-wall algorithm is identical in spirit to DFS: whenever you encounter a choice, you always choose the first (right-most) unexplored passageway, and if you ever get stuck at a dead end you turn around and go back from whence you came.

We can implement DFS with only a small change to BFS, as shown in Figure 11.31: instead of putting a newly discovered node u at the *end* of the list *Frontier* of nodes from which to explore (as in BFS), we put a newly discovered node u at the *beginning* of *Frontier*. (In other words, *BFS* treats the list *Frontier* as a queue—*first in, first out*—while *DFS* treats the list *Frontier* as a stack—*last in, first out*.) Another small change is necessary, to allow a node already in *Frontier* to be “moved” earlier in the list of nodes to explore.

Because this alteration of BFS changes only the order in which the nodes in *Frontier* are explored, DFS does precisely the same work as BFS, and is correct for the same reasons: DFS returns precisely the set of nodes reachable from the given source node s . (With a little more cleverness in moving nodes to the front of *Frontier*, DFS can also be implemented in $\Theta(|V| + |E|)$ time.) Here’s a fully detailed example of DFS:

Depth-First Search (DFS):

Input: a graph $G = \langle V, E \rangle$ and a source node $s \in V$

Output: the set of nodes reachable from s in G

```

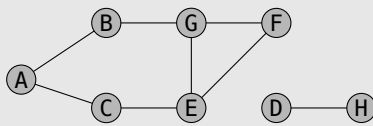
1: Frontier :=  $\langle s \rangle$ 
2: Known :=  $\emptyset$ 
3: while Frontier is nonempty:
4:    $u$  := the first node in Frontier
5:   remove  $u$  from Frontier
6:   if  $u$  is not in Known then
7:     for every neighbor  $v$  of  $u$ :
8:       if  $v$  is not in Known then
9:         add  $v$  to the start of Frontier
10:    add  $u$  to Known
11: return Known

```

Figure 11.31: The pseudocode for depth-first search. The only changes from BFS are underlined.

Example 11.29 (Sample run of DFS, in detail)

We’ll trace DFS starting at node A in this graph:



A:	B, C
B:	A, G
C:	A, E
D:	H
E:	C, F, G
F:	E, G
G:	B, E, F
H:	D

<ul style="list-style-type: none"> ○ = Frontier ◐ = just moved from Frontier to Known ● = Known ○ = neither Known nor Frontier 	Known	Frontier u : just added.	Explanation
	{}	$\langle A \rangle$	initialization
	{A}	$\langle \underline{B}, C \rangle$	processing A
	{A, B}	$\langle \underline{G}, C \rangle$	processing B (A known \Rightarrow not re-added)
	{A, B, G}	$\langle \underline{E}, F, C \rangle$	processing G (B known \Rightarrow not re-added)
	{A, B, G, E}	$\langle \underline{C}, F, F, C \rangle$	processing E (G known \Rightarrow not re-added)
	{A, B, G, E, C}	$\langle \underline{F}, F, C \rangle$	processing C (A, E known \Rightarrow not re-added)
	{A, B, G, E, C, F}	$\langle \underline{F}, C \rangle$	processing F

There are two more iterations that remove the last two entries in *Frontier* (making no changes to *Known* and adding nothing further to *Frontier*), because both F and C are already in *Known*. The **while** loop then terminates, and DFS returns $\{A, B, G, E, C, F\}$.

COMPUTER SCIENCE CONNECTIONS

THE BOWTIE STRUCTURE OF THE WEB

As the web has grown more and more central in the daily lives of us all, it has garnered increasing attention from researchers in computer science. A great deal of work has been performed to characterize the web in terms of its degree distribution (see p. 1123) or in terms of the “small-world phenomenon” (see p. 438). But one foundational and influential paper sought to characterize the web’s structure in terms of its strongly connected components.⁹ In the early days of the web, eight researchers from AltaVista, IBM, and Compaq downloaded around 200 million web pages, comprising about 1.5 billion links. They then analyzed the structure of the resulting graph, by categorizing the pages:

1. Let **CORE** denote those web pages contained in the largest SCC of the web graph. Like many other networks (for example, social networks and collaboration networks), the web graph has a *giant component* that contains many more nodes than the second-largest SCC. Denote by **CORE** those nodes in the largest SCC in the web graph.
2. Let **IN** denote those web pages p such that (i) $p \notin \text{CORE}$, and (ii) there is a path from p to some node in **CORE**. That is, there is a path from p to every page in **CORE**, but there’s no path from any node in **CORE** to p .
3. Let **OUT** denote those web pages p such that (i) $p \notin \text{CORE}$, and (ii) there is a path from some node in **CORE** to p . That is, there is a path from every page in **CORE** to page p , but there’s no path from p to any node in **CORE**.

When displayed graphically, as in Figure 11.32, these categories of web pages look like a bowtie, and so the paper by Broder et al. came to be known as “the bowtie paper.”

To complete the picture of the bowtie structure of the web, we must note that not all web pages are included in Figure 11.32. There are three further categories of nodes:

4. Let **TUBES** denote those pages p that (i) are reachable from a node of **IN** (that is, there’s a page $q \in \text{IN}$ that has a path to p), and (ii) can reach a node of **OUT** (that is, there’s a page $q \in \text{OUT}$ to which p has a path), and (iii) $p \notin \text{CORE}$.
5. Let **TENDRILS** denote those pages p that are *either* reachable from a node of **IN**, or can reach a node of **OUT**, but not both.
6. Let **DISCONNECTED** denote those pages p that are not in **CORE**, **IN**, **OUT**, **TUBES**, or **TENDRILS**—that is, those pages p that can neither reach nor be reached by any node in those sets.

One of the unexpected facts found by Broder et al. was the extent to which the web is actually *not* particularly well connected. In particular, if we were to choose web pages p and q uniformly at random from the web graph, there was only a roughly 24% chance of that a directed path from p to q exists—far lower than the “small world” phenomenon would suggest.

⁹ Andrei Broder, Ravi Kumar, Farzin Maghoul, Prabhakar Raghavan, Sridhar Rajagopalan, Raymie Stata, Andrew Tomkins, and Janet Wiener. Graph structure in the web. *Computer Networks*, 33(1–6):309–320, 2000.

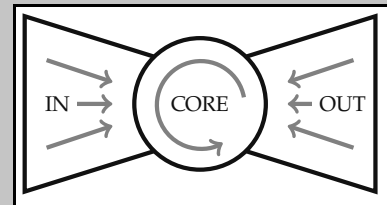


Figure 11.32: The “bowtie structure” of the web graph, in its basic form. Broder et al. found that roughly 25% of web pages fell into each of these categories: 56M pages (of 200M) in **CORE**, 43M pages in **IN**, and 43M pages in **OUT**.

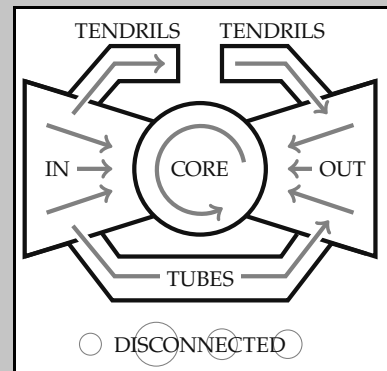


Figure 11.33: The remainder of the “bowtie structure” of the web graph. There were about 44M pages in **TENDRILS** and **TUBES**, and about 17M pages in **DISCONNECTED**.

COMPUTER SCIENCE CONNECTIONS

GARBAGE COLLECTION

In many modern programming languages, including Python and Java, the burden of managing memory is lifted from the shoulders of the programmer. When a new object is needed, the programmer just creates it. After a program has been running for a while, there may be objects that were stored in memory but are now *inaccessible* because the programmer has no way to refer to them ever again. This stored but inaccessible data is called *garbage*. Figure 11.34 shows an example of garbage being created. In Python- and Java-like languages, the system provides a *garbage collector* that periodically runs to clean up the garbage, which allows that memory to be reused for future allocations. (In contrast, in languages like C or C++, when you as a programmer are done using a chunk of memory, it's your responsibility to declare to the system that you're done using that memory by explicitly "deallocating" or "freeing" it.)

There are many sophisticated garbage-collection algorithms that are employed in real systems, but fundamentally the algorithmic idea is based on finding reachable nodes in a graph. There is a *root set* of memory locations that are reachable—essentially every variable that's defined in any currently active function call on the stack. Furthermore, if a memory location ℓ is pointed to by a reachable memory location, then ℓ too is reachable. Two simpler algorithms that are sometimes used in garbage collection are based on some corresponding simple graph-theoretic approaches. Here's a brief description of these two garbage-collection algorithms:¹⁰

Reference counting: For each block b of memory, we maintain a *reference count* of the number of other blocks of memory (or root set variables) that refer to b . When the garbage collector runs, any block b that has a reference count equal to 0 is marked as garbage and reclaimed for future use.

Mark-and-sweep: When the garbage collector runs, we iteratively *mark* each block b that is accessible. Specifically, for every variable v in the root set, we mark the block to which v refers. Then, for any block b that is marked, we also mark any block to which b refers. Once the marking process is completed, we *sweep* through memory, and reclaim all unmarked blocks.

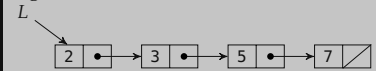
In graph-theoretic terms, we view memory as a directed graph, with an edge from each block b to the block(s) to which b refers. Reference counting declares as garbage any node with in-degree 0; mark-and-sweep declares as garbage any node that is not reached by BFS starting from the root set.

Reference counting is a simpler algorithm, but it has a problem with cyclical structures. If two inaccessible blocks of memory refer to each other, they both have nonzero reference count, and therefore won't be marked as garbage. An example is shown in Figure 11.35. There are issues of efficiency with mark-and-sweep (the entire system has to pause while the garbage collector runs), and so other, more sophisticated algorithms are generally used in real systems.

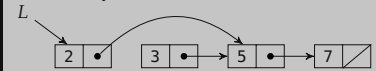
Suppose that `Node(data, next)` creates a new node for a singly linked list, with data `data` and with a pointer `next` to the next node in the list. Imagine executing the following code:

```
1  L = Node(7, NULL)
2  L = Node(5, L)
3  L = Node(3, L)
4  L = Node(2, L)
5  L.next = L.next.next
```

Then the state of memory after executing lines 1–4 is



But when we execute line 5, the state of memory becomes



The node with data = 3 is *garbage* now: there is no way to access that memory again, because there is no way for the programmer to refer to it.

Figure 11.34: Garbage being created.

You can learn more about garbage collection in any good textbook on programming languages. A few of these are:

¹⁰ Michael L. Scott. *Programming Language Pragmatics*. Morgan Kaufmann Publishers, 3rd edition, 2009; and Kenneth C. Louden and Kenneth A. Lambert. *Programming Languages: Principles and Practices*. Course Technology, 3rd edition, 2011.

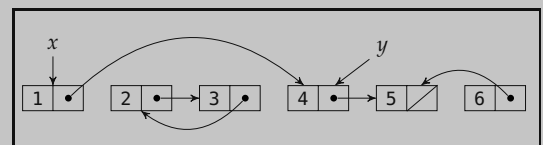
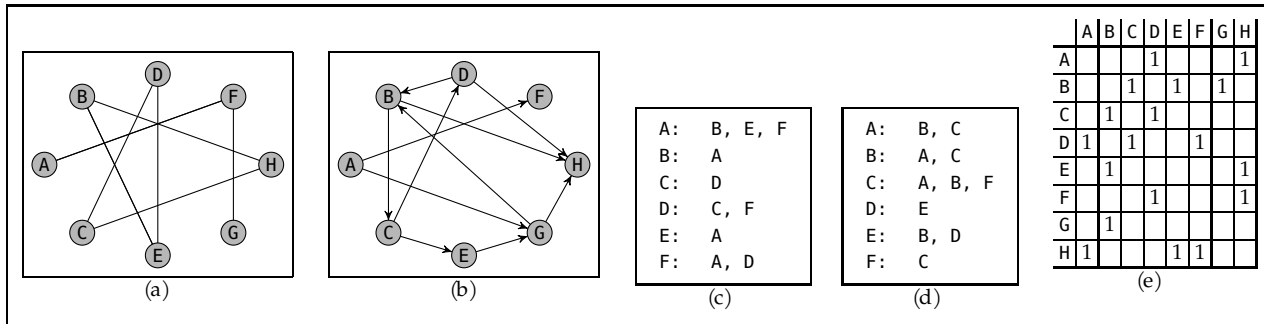


Figure 11.35: A memory diagram with six blocks of memory, and two root set variables x and y . Reference counting would show block #6 with a reference count of zero, and therefore it would be reclaimed. Mark-and-sweep would mark blocks #1, #4, and #5; thus it would reclaim blocks #2, #3, and #6.

11.3.6 Exercises

Figure 11.36:
Several graphs.

For the graphs defined in Figure 11.36, identify the following specified objects (or indicate why no such thing exists):

- 11.72 a path from D to B in Figure 11.36(a)
 11.73 two different paths from C to H in Figure 11.36(a)
 11.74 a path from C to B in Figure 11.36(b)
 11.75 two different paths from A to H in Figure 11.36(b)
 11.76 a path from D to H in Figure 11.36(b) that is *not* simple.
 11.77 a path from B to C in the graph defined by the adjacency list in Figure 11.36(c)
 11.78 a shortest path from B to F in Figure 11.36(d)
 11.79 a *non-shortest* path from B to C in the graph defined by the adjacency matrix in Figure 11.36(e)
 11.80 all nodes reachable from A in Figure 11.36(d)
 11.81 all nodes reachable from A in Figure 11.36(e)

Which of these graphs are (strongly) connected? Explain your answers. Identify all of the connected components for the undirected graphs, and all of the strongly connected components for the directed graphs.

- 11.82 Figure 11.36(a)
 11.83 Figure 11.36(b) (strong connectivity)
 11.84 Figure 11.36(c)
 11.85 Figure 11.36(d) (strong connectivity)
 11.86 Figure 11.36(e)

Let $G = \langle V, E \rangle$ be an undirected graph, and let $s \in V$ and $t \in V$ be any two nodes in G . Prove the following:

- 11.87 If there's a path of length k from s to t , then there's a path of length k from t to s .
 11.88 Every shortest path between s and t is a simple path.

For a directed graph $G = \langle V, E \rangle$, the diameter of G is the largest node-to-node distance in the graph. That is,

$$\text{diameter}(G) = \max_{s \in V, t \in V} d(s, t),$$

where $d(s, t)$ denotes the length of the shortest path from node s to node t in G . Prove your answers:

- 11.89 In terms of n , what is the *smallest* diameter that an n -node undirected graph can have?
 11.90 In terms of n , what is the *largest* diameter that a connected n -node undirected graph can have?
 Give an example of a graph where the diameter is this large. (In other words, assuming that G is connected, what's the largest possible distance between two nodes in G ? Note that, without the restriction that the graph be connected, the answer would be ∞ .)

Consider an n -node 3-regular undirected graph G . (That is, we're considering a graph $G = \langle V, E \rangle$ with $|V| = n$, where each node $u \in V$ has degree exactly equal to 3.) In terms of n :

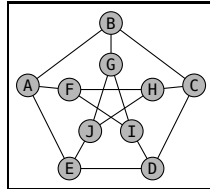
- 11.91 What is the largest possible number of connected components in a 3-regular graph?
 11.92 What is the smallest possible number of connected components in a 3-regular graph?
 11.93 Describe a connected 3-regular graph with n nodes with a diameter that's at least $\frac{n}{8}$.
 11.94 Describe a connected 3-regular graph with n nodes with a diameter that's at most $8 \log n$.

Although the context is different, our version of "diameter" matches the idea from geometry: the diameter of a circle is the distance between the two points in the circle that are farthest apart. That's still true for a graph.

11.95 Prove or disprove: let $G = \langle L \cup R, E \rangle$ be a bipartite graph with $|L| = |R|$. Suppose that every node in the graph (that is, all nodes in L and R) has at least one neighbor. Then the graph is connected.

Consider an undirected graph G . Recall that a simple path from s to t in G is a path that does not go through any node more than once. A Hamiltonian path from s to t in G is a path from s to t that goes through each node of G precisely once. In general, finding Hamiltonian paths in a graph is believed to be computationally very difficult. But there are some specific graphs in which it's easy to find one.

11.96 Find a Hamiltonian path in the Petersen graph:



Hamiltonian paths are named after William Rowan Hamilton, a 19th-century Irish mathematician/physicist.

11.97 Let K_n be a complete graph, and let s and t be two distinct nodes in the graph. How many different Hamiltonian paths are there from s to t ?

11.98 Let $K_{n,m}$ be a complete bipartite graph with $n + m$ nodes, and let s and t be two distinct nodes in the graph. How many different Hamiltonian paths are there from s to t ? (Careful; your answer may depend on s and t .)

The diameter of an undirected graph $G = \langle V, E \rangle$ is defined as the maximum distance between any two nodes $s \in V$ and $t \in V$. (See Exercises 11.89 and 11.90.) The maximum distance is one measure of how far a graph “sprawls,” but another way of measuring this idea is by looking at the average distance instead. That is, for a pair of distinct nodes $\langle s, t \rangle$ chosen uniformly from the set V , what's the distance from s to t ? That is, the average distance of a graph $G = \langle V, E \rangle$ is defined as

$$\text{the average distance of } G = \frac{\sum_{s \in V} \sum_{t \in V: t \neq s} \text{distance}(s, t)}{n(n-1)}.$$

(There are $n(n-1)$ ordered pairs of distinct nodes.) Often the average distance is a bit harder to calculate than the maximum distance, but in the next few exercises you'll look at the average distance for a pair of simple graphs.

11.99 Consider an n -node cycle, where n is odd. (We'll see a formal definition of a cycle in Section 11.4, but for now just look at the 15-node example in Figure 11.37(a).) Compute the average distance in this n -node graph. (Hint: every node is positioned symmetrically, so you can just figure out the average distance from some particular node u .)

11.100 What is the average distance for an n -node cycle where n is even? (See the 16-node example in Figure 11.37(b).)

11.101 What is the average distance for an n -node path? (See the 15-node example in Figure 11.37(c).) (Hint: for any particular integer k , how many pairs of nodes have distance k ? Then simplify the summation.)

11.102 (programming required) Write a program, in a language of your choice, to verify your answers to the last three exercises: build a graph of the appropriate size and structure, sum all of the node-to-node distances, and compute their average.

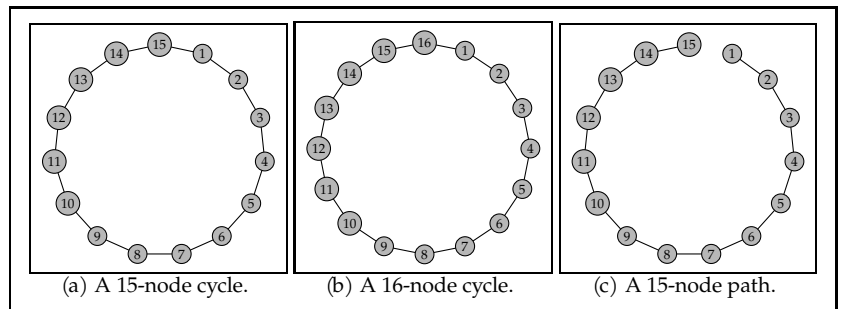


Figure 11.37: Three graphs.

Suppose that G is an undirected graph with n nodes. Answer the following questions in terms of n :

11.103 If G is disconnected, what is the largest possible number of edges that G can contain?

11.104 If G is connected, what is the smallest possible number of edges that G can contain?

Suppose that G is a directed graph with n nodes. Answer the following questions in terms of n :

11.105 If G is strongly connected, what is the smallest number of edges that G can contain?

11.106 If every node of G is in its own strongly connected component (that is, there are n different SCCs, one per node), what is the largest number of edges that G can contain?

A metric on a set V is a function $d : V \times V \rightarrow \mathbb{R}^{\geq 0}$ that obeys the following conditions (see Exercise 4.6 for more):

- reflexivity: for any $u \in V$ and $v \in V$, we have $d(u, u) = 0$ and $d(u, v) \neq 0$ whenever $u \neq v$.
- symmetry: for any $u \in V$ and $v \in V$, we have $d(u, v) = d(v, u)$.
- triangle inequality: for any $u \in V$ and $v \in V$ and $z \in V$, we have $d(u, v) \leq d(u, z) + d(z, v)$.

Let $d_G(u, v)$ denote the distance (shortest path length) between nodes $u \in V$ and $v \in V$ for a graph $G = \langle V, E \rangle$.

11.107 Prove that d_G is a metric if G is any connected undirected graph.

11.108 Prove that d_G is not necessarily a metric for a directed graph G , even if G is strongly connected.

11.109 Definition 11.23 defined a strong connected component in a graph $G = \langle V, E \rangle$ as a set $C \subseteq V$ such that: (i) any two nodes $s \in C$ and $t \in C$ are strongly connected; and (ii) for any node $x \in V - C$, adding x to C would make (i) false. Suppose that we'd instead defined clause (i) as for any two nodes $s \in C$ and $t \in C$, the node t is reachable from node s . (But we don't require that s be reachable from t .) This alternate definition is equivalent to the original. Why?

11.110 Prove that the strongly connected components (SCCs) of a directed graph partition the nodes of the graph: that is, prove that the relation $R(u, v)$ denoting mutual reachability (u is reachable from v , and v is reachable from u) is an equivalence relation (reflexive, symmetric, and transitive).

Consider the directed graphs represented in Figure 11.38, one by picture and one by adjacency list. Identify the strongly connected components ...

11.111 ... in Figure 11.38(a).

11.112 ... in Figure 11.38(b).

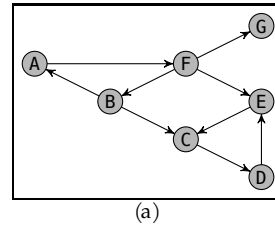
Suppose that we run breadth-first search from the following nodes. What is the last node that BFS discovers? (If there's a tie, then list all the tied nodes.)

11.113 BFS from node A in Figure 11.38(a).

11.114 BFS from node B in Figure 11.38(a).

11.115 BFS from node 0 in Figure 11.38(b).

11.116 BFS from node 12 in Figure 11.38(b).



(a)

0:	3, 7
1:	9, 2, 5
2:	1, 10, 9
3:	0, 7, 1
4:	10, 7
5:	1
6:	7, 11
7:	0, 4, 6, 8
8:	11, 12
9:	1
10:	2, 4
11:	6, 8
12:	8

(b)

Figure 11.38: Two graphs.

Breadth-first search as described in Figure 11.29 finds all nodes reachable from a given source node in a given graph, and, in fact, it discovers nodes in increasing order of their distance from s . But we didn't actually record distances during the computation.

11.117 Modify the pseudocode for BFS to compute distances instead of just whether a path exists, by annotating every node added to *Frontier* with its distance from the source node s .

11.118 Argue that in your modified version of BFS, there are never more than two different distances stored in the Frontier.

11.119 Argue that the claim from the previous exercise may be false for *depth*-first search.

11.120 Consider a graph G represented by an adjacency matrix M . What does the $\langle i, j \rangle$ th entry of MM (the matrix that results from squaring the matrix M) represent?

A word chain is a sequence $\langle w_1, w_2, \dots, w_k \rangle$ of words, where each w_i is a word in English, and w_{i+1} is one letter different from w_i . For example, a word chain from FROWN to SMILE for my dictionary is

FROWN \rightarrow FLOWN \rightarrow FLOWS \rightarrow SLOWS \rightarrow SLOTS \rightarrow SLITS \rightarrow SKITS \rightarrow SKITE \rightarrow SMITE \rightarrow SMILE.

(SKITE is a word of Scottish origin, meaning "an oblique blow.")

11.121 (programming required) Write a program that uses a BFS-like algorithm to find a shortest word chain between two given words w_1 and w_2 of the same length. (You can find a dictionary of English words on the web, or `/usr/share/dict/words` on Unix-based operating systems. You'll want to cull your dictionary to only words of the right length before you start.) There are faster solutions that involve searching "in both directions" out from w_1 and into w_2 until you find a match, but BFS from w_1 will work.

11.4 Trees

I think that I shall never see
A poem lovely as a tree.

Joyce Kilmer (1886–1918), “Trees”
Trees and Other Poems (1914)

Informally, a *tree* is a graph that grows from a *root*, branching outward and eventually leading to the *leaves*. (We computer scientists are always upside down compared to botanists: unlike an oak or maple or tamarack, the root of a tree in CS is at the top, and it grows downward toward the leaves.) See Figure 11.39.

Trees arise very frequently in computer science: to name just a few examples, they’re the class hierarchies of object-oriented programming, the binary search trees of data structures (see p. 1160), the game trees describing the progression of Tic-Tac-Toe or chess (p. 344), the parse trees that describe formal or natural languages (p. 543), the recursion trees that describe the execution of recursive algorithms (Section 6.4). Trees are also frequently used in computational models of important phenomena from outside of CS: for example, in reconstructing evolutionary phylogenies (in computational biology), or in reconstructing the paths by which rumors spread from the originator of the information (in social network analysis). In this section, we’ll introduce trees formally—including definitions, properties, algorithms, and applications—as a special type of graph.

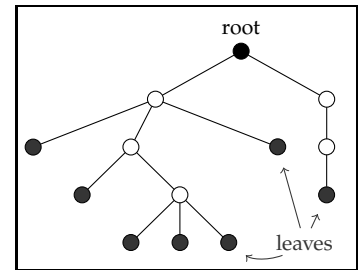


Figure 11.39: A small tree.

11.4.1 Cycles

Before we can define trees properly, we must first define another notion about graphs in general—a *cycle*, which is way to get from a node back to itself:

Definition 11.26 (Cycle)

A cycle $\langle u_1, u_2, \dots, u_k, u_1 \rangle$ is a path of length ≥ 2 from a node u_1 back to node u_1 that does not traverse the same edge twice. Just as for any other path, the length of the cycle $\langle u_1, u_2, \dots, u_k, u_1 \rangle$ is the number of edges it traverses—that is, k .

Figure 11.40 shows examples of an undirected and directed graph with a cycle $\langle A, B, C, A \rangle$. Note that the edges $\langle s, t \rangle$ and $\langle t, s \rangle$ in a directed graph are different; in an undirected graph, the edges $\{s, t\}$ and $\{t, s\}$ are the same. Thus a cycle in a directed graph *can* use both $\langle s, t \rangle$ and $\langle t, s \rangle$, but a cycle in an undirected graph cannot use both $\langle s, t \rangle$ and $\langle t, s \rangle$. In Figure 11.40, the path $\langle C, E, C \rangle$ is a cycle in the directed graph, but is *not* a cycle in the undirected graph because it reuses an edge.

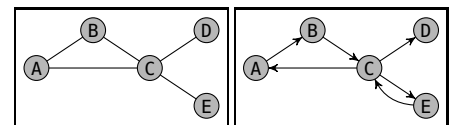


Figure 11.40: Two graphs with cycles $\langle A, B, C, A \rangle$.

Technically speaking, the definition of a cycle in Definition 11.26 says that the undirected graph in Figure 11.40 has six different cycles:

- $\langle A, B, C, A \rangle$, $\langle C, A, B, C \rangle$, and $\langle B, C, A, B \rangle$ (going clockwise), and
- $\langle A, C, B, A \rangle$, $\langle C, B, A, C \rangle$, and $\langle B, A, C, B \rangle$ (going counterclockwise).

However, we will adopt the convention that there is one and only one cycle in this graph. Because we can “start anywhere” in a cycle, we consider a cycle to be defined only by the relative ordering of the nodes involved, regardless of where we start. In an undirected graph, we can “go either direction” (clockwise or counterclockwise), so we also ignore the direction of travel in distinguishing cycles. In a directed graph, the direction of travel *does* matter; we may be able to go in one direction around a cycle without being able to go in the other. In other words, we say that Figure 11.41(a) and Figure 11.41(b) have one cycle each, while Figure 11.41(c) has two.

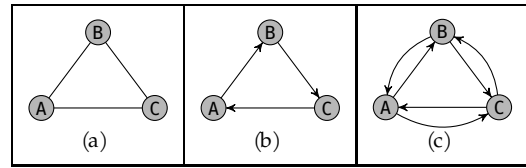


Figure 11.41: Some cycles.

A cycle is by definition forbidden from traversing the same edge twice. A *simple* cycle also does not visit any *node* more than once:

Definition 11.27 (Simple cycle)

A cycle $\langle u_1, u_2, \dots, u_k, u_1 \rangle$ is simple if each u_i is distinct—that is, no nodes in the cycle are duplicated aside from the last node (which equals the first node).

(We’ve now used the word “simple” in three different contexts: simple graphs have no parallel edges or self-loops, and simple paths and cycles have no repeated vertices. Intuitively, all three definitions correspond to an entity that’s not unnecessarily complicated.) For one example, see Figure 11.42; here are two more:

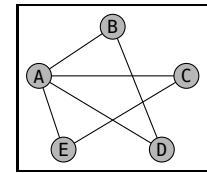
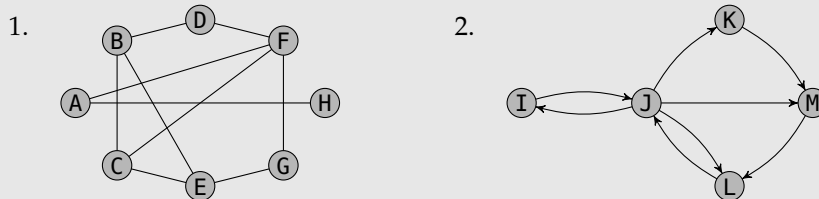


Figure 11.42: In this graph, $\langle D, B, A, C, E, A, D \rangle$ is a non-simple cycle. This graph also has two simple cycles: $\langle D, B, A, D \rangle$ and $\langle C, E, A, C \rangle$.

Example 11.30 (Finding cycles)

Problem: Identify all simple cycles in the following graphs:



Solution: A nice way to identify cycles systematically is to look for cycles of all possible lengths: 2-node cycles, 3-node cycles, etc. (Actually 2-node cycles are possible only in directed graphs. *Exercise: why?*) Here are the simple cycles in these graphs:

- | | |
|---------------------------------------|---------------------------------|
| 1. $\langle B, E, C, B \rangle$ | 2. $\langle I, J, I \rangle$ |
| $\langle B, D, F, C, B \rangle$ | $\langle J, L, J \rangle$ |
| $\langle C, F, G, E, C \rangle$ | $\langle J, M, L, J \rangle$ |
| $\langle B, D, F, G, E, B \rangle$ | $\langle J, K, M, L, J \rangle$ |
| $\langle B, D, F, G, E, C, B \rangle$ | |

Note that (to name one of several examples) the sequence $\langle I, J, L, J, I \rangle$ is also a cycle in the second graph—it traverses four distinct directed edges and goes from node I to I—but this cycle is not simple, because node J is repeated.

We can use a modification of breadth-first search to identify cycles algorithmically. Specifically, suppose that we wish to find out whether a node u is involved in a cycle in a directed graph. We run BFS starting at node u , and if we ever encounter a node v that has u as a neighbor, then we have found a cycle involving node u . (An extra modification is necessary for undirected graphs; see Exercise 11.129.)

Taking it further: Kidneys are the most frequently transplanted organ today, in part because—unlike for other organs—humans generally have a “spare”: we’re born with two kidneys, but only need one functioning kidney to live a healthy life. Thus patients suffering from kidney failure may be able to get a transplant from friends or family members who are willing to donate one of their kidneys. But this potential transplant relies on the donor and the patient being compatible in dimensions like blood type and the physical size of the organs. Recently a computational solution to the problem of incompatibility has emerged, using algorithms based on finding (short) cycles in a particular graph: there is now national exchange for matching up two (or a few) patients with willing-but-incompatible donors, and doing a multiway transplant. See p. 1159 for more discussion.

ACYCLIC GRAPHS

While cycles are important on their own, their relevance for trees is actually when they *don’t* exist:

Definition 11.28 (Acyclic Graphs)

A graph is acyclic if it contains no cycles.

Let’s prove a useful structural fact about acyclic graphs. (Recall that we are considering *finite* graphs, where the set of nodes in the graph is finite. The following claim would be false if graphs could have an infinite number of nodes!)

Lemma 11.5 (Every acyclic graph has a node with degree 0 or 1)

Let $G = \langle V, E \rangle$ be an acyclic undirected graph. Then there exists a node in V whose degree is zero or one.

Proof. We’ll give a constructive proof of the claim—specifically, we’ll give an algorithm that finds a node with the stated property:

```

1: let  $u_0$  be an arbitrary node in the graph, and let  $i := 0$ 
2: while the current node  $u_i$  has no unvisited neighbors:
3:   let  $u_{i+1}$  be a neighbor of  $u_i$  that has not previously been visited.
4:   increment  $i$ 

```

Observe that this process must terminate in at most $|V|$ iterations, because we must visit a new node in each step. Suppose that this algorithm goes through k iterations of the **while** loop, and let t be the last node visited by the algorithm. (So $t = u_k$.)

- If $k = 0$, then $t = u_0$ has degree zero, so the claim follows immediately.
- If $k \geq 1$, then we’ll argue that t has degree one. Because the algorithm terminated, there cannot be an edge between t and any unvisited node. Furthermore, if there were an edge from t to any previously visited node u_j for $j < k - 1$, then there would be a cycle in the graph, namely $\langle u_j, u_{j+1}, \dots, u_{k-1}, u_k, u_j \rangle$. Therefore t ’s only neighbor is u_{k-1} , and the degree of t is one. \square

For directed graphs, the claim analogous to Lemma 11.5 is *every directed acyclic graph contains a node with outdegree zero*. (You’ll prove it in Exercise 11.130.)

Taking it further: A *directed acyclic graph* (often just called a *DAG*) is, perhaps obviously, a directed graph that contains no cycles. A DAG G corresponds to a (strict) partial order (see Chapter 8); a cycle in G corresponds to a violation of transitivity. In fact, we can think of *any* directed graph $G = \langle V, E \rangle$ as a relation—specifically, the edge set E is a subset of $V \times V$. Like transitivity and acyclicity, many of the concepts that we explored in Chapter 8 have analogues in the world of graphs.

11.4.2 Trees

With the definition of cycles in hand, we can now define trees themselves:

Definition 11.29 (Tree)
A tree is an undirected graph that is connected and acyclic.

An irrelevant note about Chinese: the character for *tree* is 木; the character for *forest* is 森 (a disconnected collection of trees!).

We will also sometimes talk about graphs that satisfy only the latter requirement: a *forest* is an undirected graph that is acyclic (but not necessarily connected). Every connected component of a forest is a tree, and note that a tree is itself a forest.

Several examples of trees are shown in Figure 11.43: all six graphs have a single connected component and contain no cycles. Therefore all six are trees.

We’ll prove several structural facts about trees in this section, beginning with one concerning the number of edges in a tree. To start, let’s look at the number of nodes and edges in each of the trees in Figure 11.43:

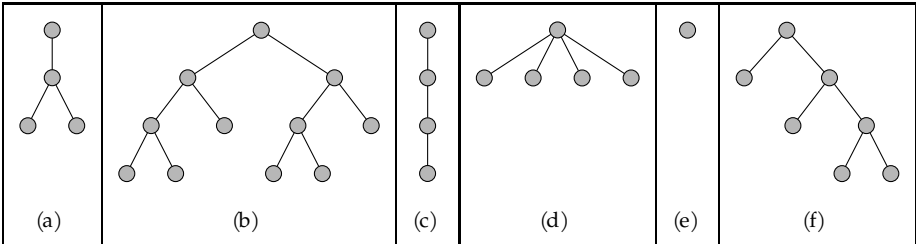


Figure 11.43: Some sample trees.

	(a)	(b)	(c)	(d)	(e)	(f)
number of nodes	4	11	4	5	1	7
number of edges	3	10	3	4	0	6

In each of these trees, the number of nodes is one more than the number of edges, and that’s no coincidence; here’s the statement and proof of the general fact:

Theorem 11.6 (Number of edges in a tree)
Let $T = \langle V, E \rangle$ be a tree. Then $|E| = |V| - 1$.

Proof. Let $P(n)$ denote the property that any n -node tree has precisely $n - 1$ edges. We will prove that $P(n)$ holds for all $n \geq 1$ by induction on n .

Base case ($n = 1$): We must prove $P(1)$: any 1-node tree has $1 - 1 = 0$ edges. But the only 1-node (simple) graph is the one shown in Figure 11.43(e), which has zero edges, and so we’re done immediately.

Inductive case ($n \geq 2$): We assume the inductive hypothesis $P(n-1)$ —that is, every $(n-1)$ -node tree has $n-2$ edges. We must prove $P(n)$.

Consider an arbitrary tree $T = \langle V, E \rangle$ with $|V| = n$. By definition, T is acyclic and connected. By Lemma 11.5, then, there exists a node $u \in V$ with degree 0 or 1 in T . Furthermore, because T is connected, the degree of u cannot be 0. Thus u is a node with $\text{degree}(u) = 1$. Let $v \in V$ be the unique neighbor of u in T . Let T' be T with node u and the edge $\{u, v\}$ between u and v deleted. (See Figure 11.44.)

We claim that the graph $T' = \langle V - \{u\}, E - \{\{u, v\}\} \rangle$ is a tree, too. The acyclicity and connectivity of T' both follow from the fact that T was acyclic and connected, and the fact that the eliminated node u was of degree 1.

The tree T' contains $n-1$ nodes, and thus, by the inductive hypothesis $P(n-1)$, contains $n-2$ edges. Therefore T , whose edges are precisely the edges of T' plus the eliminated edge $\{u, v\}$, contains precisely $(n-2)+1 = n-1$ edges. \square

An immediate consequence of Theorem 11.6 is that every tree is teetering on the edge of being disconnected and of having a cycle (see Figure 11.45):

Corollary 11.7 (A tree with an edge added or removed is not a tree)

Let $T = \langle V, E \rangle$ be any tree. Then:

1. adding any edge $e \notin E$ to T creates a cycle; and
2. removing any edge $e \in E$ from T disconnects the graph.

Proof. 1. Define the graph $G = \langle V, E \cup \{e\} \rangle$ as the result of adding the new edge e to the tree T . Because adding an edge to a graph can never disrupt connectivity and T was already connected, we know that G must be connected too. Thus if G were acyclic, then G would be a tree. But G has one more edge than T —specifically, G has $(|V| - 1) + 1 = |V|$ edges—and therefore isn't a tree by Theorem 11.6.

2. The proof is similar: let G' be T with e removed. Removing an edge cannot create a cycle, so G' is acyclic. But G' has too few edges to be a tree by Theorem 11.6, so G' must be disconnected. \square

(Here's an alternative proof of Corollary 11.7.1. Let $\langle u, v \rangle$ be an edge not in the tree T . Because T is connected, there is already a (simple) path P from u to v in T . If we add $\langle u, v \rangle$ to T , then there is a cycle: follow P from u to v and then follow the new edge from v back to u . Therefore G contains a cycle.)

ROOTED TREES

We often designate a particular node of a tree T as the *root*, which is traditionally drawn as the topmost node. (Note that we could designate any node as the root and—just like that mobile of zoo animals from your crib from infancy—“hang” the tree by that node.) We will adopt the standard convention that, whenever we draw trees, the vertically highest node is the root.

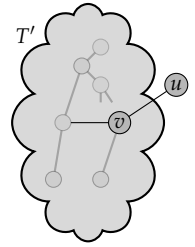


Figure 11.44: A tree T , with a node u of degree 1 and its neighbor v . The tree T' is T without the node u and the edge $\{u, v\}$.

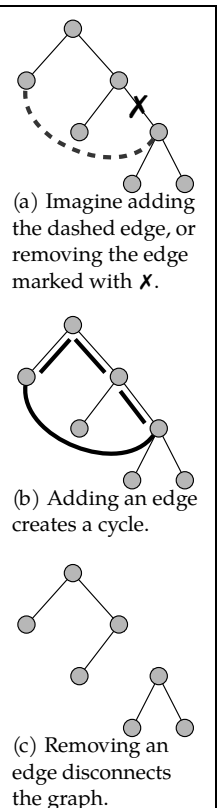


Figure 11.45: Adding/removing an edge from a tree.

There's a lot of terminology about trees in computer science that's borrowed from the world of family trees:

- For a node u in a tree with root $r \neq u$, the *parent* of u is the unique neighbor of u that is closer to r than u is. (The root is the only node that has no parent.)
- A node v is one of the *children* of a node u if v 's parent is u .
- A node v is a *sibling* of a node $u \neq v$ if v and u have the same parent.

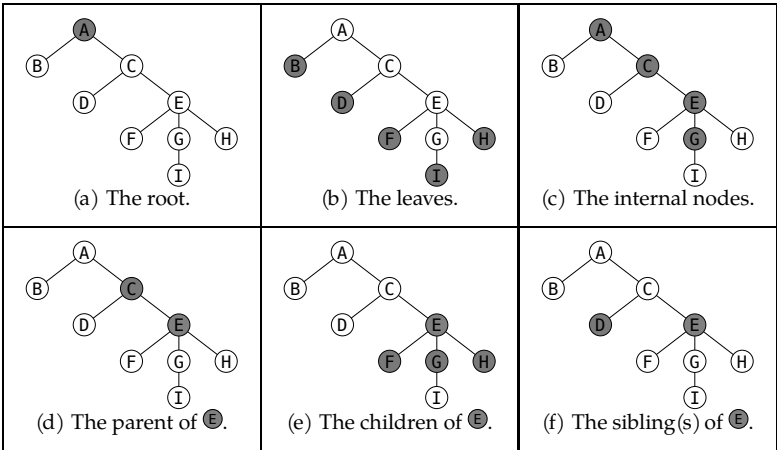
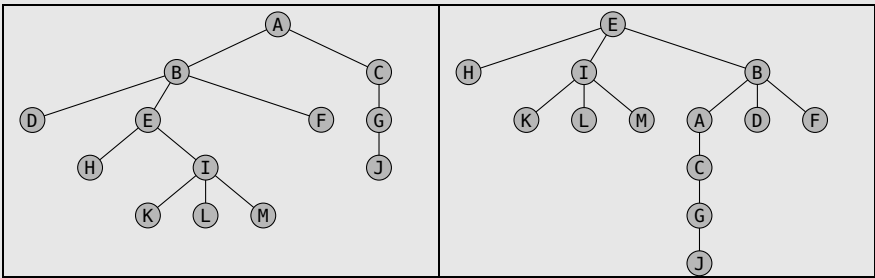


Figure 11.46: The root, leaves, and internal nodes of the tree; the parent, children, and siblings of a particular node.

A node with zero children is called a *leaf*. A node with one or more children is called an *internal node*. (Note that the root is an internal node unless the tree is the trivial one-node graph.) See Figure 11.46 for an illustration of all of these definitions. Note that Figure 11.46 is correct only when the root is the topmost node in the image; with a different root, all of the panels could change. Here's a concrete example:

Example 11.31 (A sample tree)

Here are two trees. (The second tree is just the first, rerooted to make E the new root.)



Then we have:

Root: A	Root: E
Leaves: {D, F, H, J, K, L, M}	Leaves: {D, F, H, J, K, L, M}
Internal nodes: {A, B, C, E, G, I}	Internal nodes: {A, B, C, E, G, I}
Parent of B: A	Parent of B: E
Children of B: {D, E, F}	Children of B: {A, D, F}
Parent of A: none	Parent of A: B
Children of A: {B, C}	Children of A: {C}

While the leaves and internal nodes are identical in these two trees, note that if we'd rerooted the tree at any of the erstwhile leaves instead, the new root would become an internal node instead of a leaf. For example, if we reroot this tree at H, then the leaves would be {D, F, J, K, L, M} and the internal nodes would be {A, B, C, E, G, H, I}.

SUBTREES, DESCENDANTS, AND ANCESTORS

Let T be a rooted tree, and let u be any node in T . The *subtree rooted at u* consists of u and all those nodes and edges “below” u in T . (In other words, a node v is in the subtree rooted at u if and only if v is no closer to the root of T than u is; the subtree is the induced subgraph of these nodes.) Such a node v in the subtree rooted at u is called a *descendant* of u if $v \neq u$. The node u is called an *ancestor* of v . See Figure 11.47 for illustrations of these three definitions. Here’s an example:

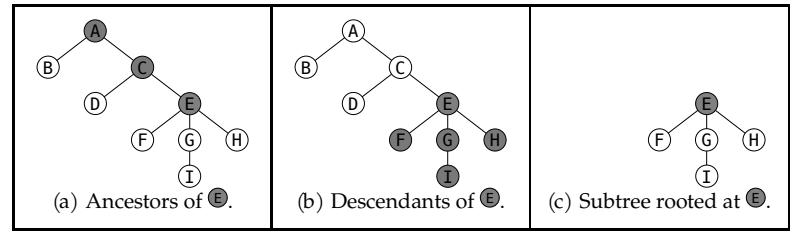
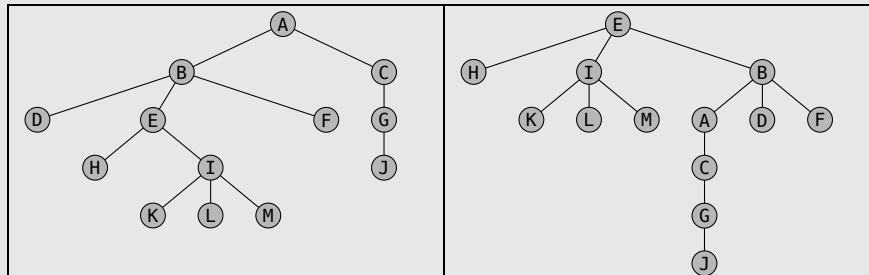


Figure 11.47: Ancestors, descendants, and subtrees.

Example 11.32 (Descendants and ancestors)

Recall the trees from Example 11.31:



Then we have:

Descendants of B: $\{D, E, F, H, I, K, L, M\}$ Ancestors of B: $\{A\}$	Descendants of B: $\{A, C, D, F, G, J\}$ Ancestors of B: $\{E\}$
Descendants of H: none Ancestors of H: $\{A, B, E\}$	Descendants of H: none Ancestors of H: $\{E\}$
Subtree rooted at B: 	Subtree rooted at B:

We have one final pair of definitions to (at last!) conclude our parade of terminology about rooted trees, related to how “tall” a tree is. Consider a rooted tree T with root node r . The *depth* of a node u is the distance from u to r . The *height* of a tree is the maximum, over all nodes u in the tree, of the depth of node u .

For example, every node in the tree in Figure 11.48 is labeled by its depth: the root has depth 0, its children have depth 1, their children (the “grandchildren” of the root) have depth 2, and so forth. The height of the tree is the largest depth of any of its nodes—in this case, the height is 4.

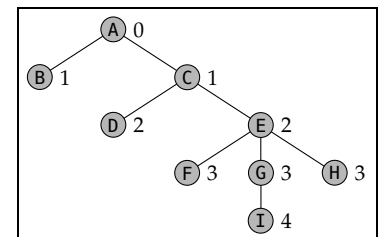


Figure 11.48: A rooted tree’s nodes, labeled by depth.

Taking it further: Alternatively, we could give several of the definitions about rooted trees recursively. For example, we could define ancestors and descendants of a node u be in a rooted tree T as follows:

- A node v is an *ancestor* of u if (i) v is the parent of u ; or (ii) v is the parent of any ancestor of u .
- A node v is a *descendant* of u if (i) v is a child of u ; or (ii) v is a child of any descendant of u .

We can also think of the depth of a node, or the height of a tree, recursively. The depth of the root is zero; the depth of a node with a parent p is $1 + (\text{the depth of } p)$. For height:

- the height of a one-node tree T is zero; and
- the height of a tree T with root r with children $\{c_1, c_2, \dots, c_k\}$ is

$$1 + \max_{i \in \{1, \dots, k\}} \text{the height of the subtree rooted at } c_i.$$

BINARY TREES

We'll often encounter a special type of tree in which nodes have a limited number of children:

Definition 11.30 (Binary trees and k -ary trees)

A binary tree is a rooted tree in which each node has 0, 1, or 2 children. More generally, if every node in a rooted tree T has k or fewer children, then T is called a k -ary tree. (In other words, a binary tree is 2-ary.)

For example, consider the trees in Figure 11.49. Of them, only the tree in Figure 11.49(d) is not a binary tree, because its root has four children. (This tree is a 4-ary tree.) But the other five trees are all binary: in each, every internal node has either 1 child or 2 children.

In a binary tree, the possible children of a node are called its *left child* and *right child*. (Even for a node u in a binary tree that has only one child, we'll insist that the lone child be designated as either the left child of u or the right child of u .) For a node u , we say that u 's *left subtree* is the subtree rooted at u 's left child; the *right subtree* is analogous.

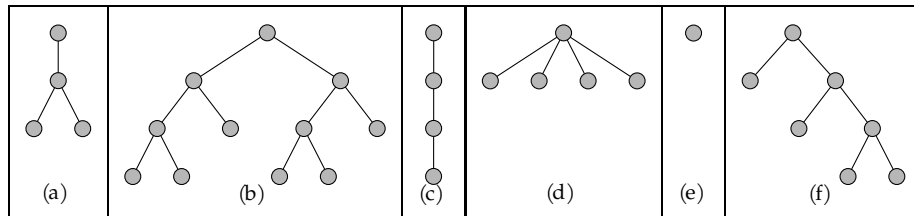


Figure 11.49:
The trees from
Figure 11.43,
repeated. All but
(d) are binary trees.

11.4.3 Tree Traversal

We will sometimes want to list all of the nodes contained in a tree T . There are three standard algorithms that are used for this purpose, called *pre-order*, *in-order*, and *post-order traversal*. While these algorithms can be generalized to non-binary trees, they're easier to understand for binary trees (and they're most frequently deployed for binary trees), so we'll consider them that way.

All three algorithms are recursive, and all three algorithms execute precisely the same steps—just in a different order. On an empty tree T , we do nothing; on a non-empty tree T , all three algorithms perform the following steps:

- we “visit” the root of the tree T . (You can think of “visiting” the root as printing out the contents of the root node, or as adding it to the end of an accumulating list of the nodes that we’ve encountered in the tree.)
- we recursively traverse the left subtree of T , finding all nodes there.
- we recursively traverse the right subtree of T , finding all nodes there.

But the three traversal algorithms execute the three steps in different orders, either visiting the root *before both recursive calls* (“pre-order”); *between the recursive calls* (“in-order”); or *after both recursive calls* (“post-order”). We always recurse on the left subtree before we recurse on the right subtree. Here are the details:

pre-order-traverse(T):	in-order-traverse(T):	post-order-traverse(T):
1: if T is empty then	1: if T is empty then	1: if T is empty then
2: do nothing.	2: do nothing.	2: do nothing.
3: else	3: else	3: else
4: visit the root of T	4: in-order-traverse(T ’s left subtree)	4: post-order-traverse(T ’s left subtree)
5: pre-order-traverse(T ’s left subtree)	5: visit the root of T	5: post-order-traverse(T ’s right subtree)
6: pre-order-traverse(T ’s right subtree)	6: in-order-traverse(T ’s right subtree)	6: visit the root of T

Figure 11.50: Three different algorithms to traverse a binary tree.

Let’s take a look at an example of traversing a small tree using these algorithms. First we’ll look at the *pre-order* traversal, in which the first node visited in any subtree is the root of that subtree:

Example 11.33 (Traversing a small tree: pre-order traversal)
Let’s determine the order of nodes’ visits by a pre-order traversal of the following tree:

```

graph TD
    A((A)) --- B((B))
    A --- C((C))
    B --- D((D))
    C --- E((E))
    C --- F((F))

```

In a pre-order traversal, we first visit the root, then pre-order-traverse the left subtree, then pre-order-traverse the right subtree. In other words, we first visit the root A, then pre-order-traverse \textcircled{D} — \textcircled{B} , then pre-order-traverse \textcircled{E} — \textcircled{C} — \textcircled{F} :

Step #1: visit the root. We visit the root A.

Step #2: pre-order-traverse the left subtree. To pre-order-traverse \textcircled{D} — \textcircled{B} , we first visit the root B, then pre-order-traverse the left subtree \textcircled{D} , then pre-order-traverse the (empty) right-subtree. In order, these steps visit B and D.

Step #3: pre-order-traverse the right subtree. To pre-order-traverse \textcircled{E} — \textcircled{C} — \textcircled{F} , we first visit C, then pre-order-traverse the left subtree \textcircled{E} , and then pre-order-traverse the right subtree \textcircled{F} . Pre-order-traversing \textcircled{E} just results in visiting E, and pre-order-traversing \textcircled{F} just visits F. In order, these steps visit C, E, and F.

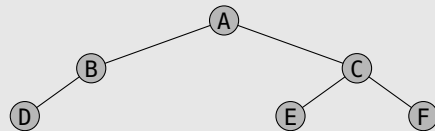
Putting this all together, the pre-order traversal of the tree visits the nodes in this order:

$\underbrace{A}_{\text{step \#1}}, \underbrace{B, D}_{\text{step \#2}}, \underbrace{C, E, F}_{\text{step \#3}}.$

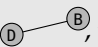
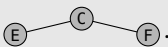
Here are examples of the other two traversal algorithms, on the same tree:

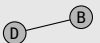
Example 11.34 (Traversing a small tree: in-order and post-order traversals)


Problem: Recall the tree from Example 11.33:



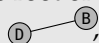
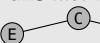
1. In what order are the nodes visited by an *in-order* traversal of this tree?
2. What about a *post-order* traversal?

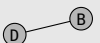
Solution: 1. We first traverse , then visit A, then traverse .


Traversing  visits D and B: first the left subtree, then the root.

Traversing  visits E, then C, then F.

Thus an in-order traversal visits the nodes in the order D, B, A, E, C, F.

2. For a post-order traversal, the root of each subtree is the *last* node traversed in that subtree: we first traverse , then traverse , then visit A.

Traversing  visits D and B: first the left subtree, then the nonexistent right subtree, then the root.

Traversing  visits E, then F, then C.

Thus a post-order traversal visits the tree's nodes in the order D, B, E, F, C, A.

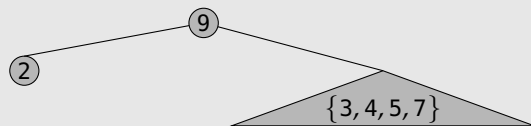
Here's another example, of using traversals to reconstruct a binary tree:

Example 11.35 (Trees from traversals)

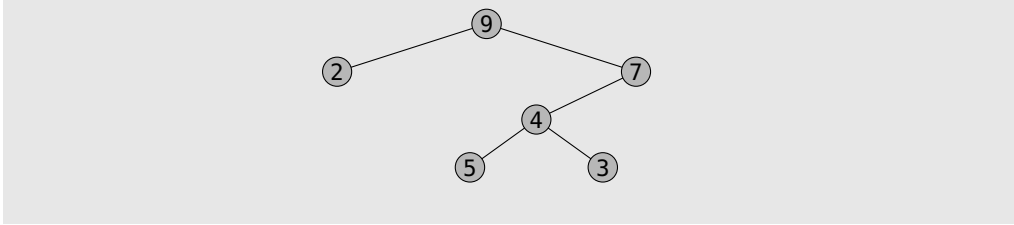
Problem: Here is the output of all three traversals on a binary tree T . What's T ?

pre-order traversal	in-order traversal	post-order traversal
9, 2, 7, 4, 5, 3	2, 9, 5, 4, 3, 7	2, 5, 3, 4, 7, 9

Solution: We'll reassemble T from the root down. The root is first in the pre-order traversal (and last in the post-order), so 9 is the root. The root separates the left subtree from the right subtree in the in-order traversal; thus the left subtree contains just 2 and the right contains $\{3, 4, 5, 7\}$. So the tree has the following form:



The post-order 5, 3, 4, 7 and in-order 5, 4, 3, 7 show that 7 is the root of the unknown portion of the tree and that 7's right subtree is empty. The last three nodes are pre-ordered 4, 5, 3; in-ordered 5, 4, 3; and post-ordered 5, 3, 4. In sum, that says that 4 is the root, 5 is the left subtree, and 3 is the right subtree. Assembling these pieces yields the final tree:



Taking it further: One particularly important type of binary tree is the *binary search tree (BST)*, a widely used data structure—one that’s probably very familiar if you’ve taken a course on data structures. A BST is a binary tree in which each node has some associated “key” (a piece of data), and the nodes of the tree are stored in a particular sorted order: all nodes in the left subtree have a key smaller than the root, and all nodes in the right subtree have a key larger than the root. Thus an in-order traversal of a binary search tree yields the tree’s keys in sorted order. For more, see p. 1160. An even more specific form of binary search tree, called a *balanced binary search tree*, adds an additional structural property related to the depth of nodes in the tree. See p. 643 for a discussion of one scheme for balanced binary search trees, called *AVL trees*.

11.4.4 Spanning Trees

Let $G = \langle V, E \rangle$ be an undirected graph. For example, imagine that each node in V represents a dorm room on your campus, and each edge in E denotes a possible fiber optic cable that can be laid to build an ethernet connection throughout the residence halls. A reasonable goal is to actually place only some of those possible cables, a subset $E' \subseteq E$, while ensuring that network traffic can be sent between any two dorm rooms—that is, ensuring that the resulting network is connected. In other words, one seeks a *spanning tree* of the graph G :

Definition 11.31 (Spanning tree)

Let $G = \langle V, E \rangle$ be a connected undirected graph. A spanning tree of G is a tree $T = \langle V, E' \rangle$ with the same nodes as G and with edges $E' \subseteq E$ that are a subset of G ’s edges.

A spanning tree of G is called “spanning” because it connects (that is, spans) all nodes in G . Figure 11.51 shows a small example: the first panel shows a small graph G ; the remaining panels show the 8 different spanning trees of G .

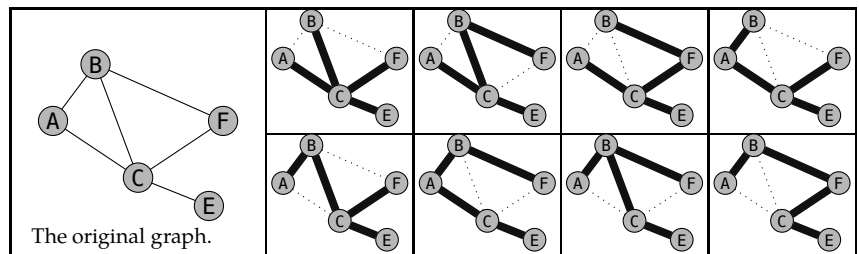


Figure 11.51: All 8 spanning trees of the graph shown in the first panel.

A graph G has a spanning tree if and only if G is connected: we can be sure to only remove “redundant” edges that aren’t required for connectivity, and removing edges from G can never cause a disconnected graph to become connected. (For disconnected graphs, people sometimes talk about a *spanning forest*: a forest $F = \langle V, E' \rangle$ with $E' \subseteq E$, where the connected components of the original graph G and the connected components of the forest F are identical.)

Although we didn’t talk about it this way when we introduced breadth- and depth-first search (see Figures 11.29 and 11.31), these algorithms can find spanning trees,

with a small change: as we explore the graph, we include in E' every edge $\langle u, v \rangle$ that leads from a previously known node u to a newly discovered node v .

We'll also see some other ways to find spanning trees in Section 11.5.2, but here's another, conceptually simpler technique. To find a spanning tree in a connected graph G , we repeatedly find an edge that can be deleted without disconnecting G —that is, an edge that's in a cycle—and delete it. See Figure 11.52 for the algorithm. Here's an example:

Cycle Elimination Algorithm:

Input: a connected graph $G = \langle V, E \rangle$

Output: a spanning tree of G

- 1: **while** there exists a cycle C in G :
- 2: let e be an arbitrary edge traversed by C
- 3: remove e from E
- 4: **return** the resulting graph $\langle V, E \rangle$.

Example 11.36 (Finding a spanning tree via cycle elimination)

Here are the iterations of the Cycle Elimination algorithm in computing a spanning tree of a given connected graph. In each iteration, we've selected an arbitrary cycle (lightly shaded) and then selected an arbitrary edge from that cycle (heavily shaded) and removed it. After three iterations, the resulting graph has no cycles, and remains connected; the resulting graph is a spanning tree of the original graph.

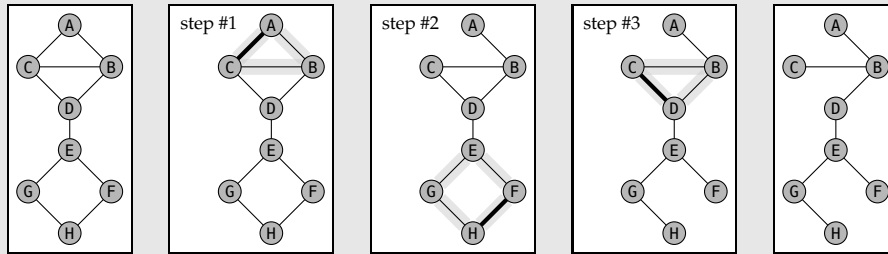


Figure 11.52: The pseudocode for an algorithm to find a spanning tree.

We can prove that the Cycle Elimination algorithm correctly finds spanning trees, given an arbitrary connected graph as input:

Theorem 11.8 (Correctness of the Cycle Elimination algorithm)

Given any connected graph $G = \langle V, E \rangle$, the Cycle Elimination algorithm returns a spanning tree T of G .

Proof. The algorithm only deletes edges from G , so certainly $T = \langle V, E' \rangle$ satisfies $E' \subseteq E$. We need to prove that T is a tree: that is, T is acyclic and T is connected.

Acyclicity: As long as there's a cycle remaining, the algorithm stays in the **while** loop.

Thus we only exit the loop when the remaining graph is acyclic. (And the loop terminates in at most $|E|$ iterations, because an edge is deleted in every iteration.)

Connectivity: We claim that the graph is connected throughout the algorithm. It's true at the beginning of the algorithm, by assumption. Now consider an iteration in which we delete the edge $\{u, v\}$ from a cycle C . Let s and t be arbitrary nodes; we will argue that there is still a path from s to t . Before we deleted $\{u, v\}$, there was a path P from s to t . If P didn't traverse the edge $\{u, v\}$, then P is still a path from s to t . Otherwise, we can still get from s to t by going “the long way around” the cycle C instead of following the single edge $\{u, v\}$. (See Figure 11.53.) Thus there is still a path from any node s to any node t , and so the graph stays connected. \square

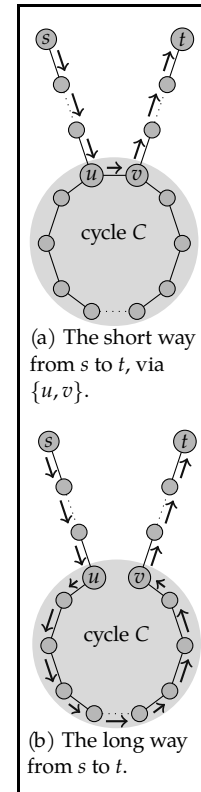


Figure 11.53: Maintaining connectivity in the Cycle Elimination Algorithm.

COMPUTER SCIENCE CONNECTIONS

DIRECTED GRAPHS, CYCLES, AND KIDNEY TRANSPLANTS

Kidneys are essential to human life; they play an essential filtering role in the body without which we would all die. Although we are born with two kidneys, humans need only one functioning kidney to live healthy lives. Because we're all naturally equipped with a "spare," kidney transplants are the most common form of transplant surgery performed today. Thousands of lives are saved annually through kidney transplants.

Typically a patient in need of a kidney identifies a friend or relative who is willing to donate. If the patient and donor are compatible—for example, blood type and physical size of the donor's kidney must be appropriate—then medical teams perform two simultaneous operations: one to remove the "spare" kidney from the donor, and one to implant it in the patient. (Some patients instead receive kidneys from strangers who chose to donate their organs in case of an untimely death.) Unfortunately, many patients who need kidneys have a friend or relative willing to donate to them—but they are incompatible with their prospective donor's kidney. These patients may spend years on a waiting list for a transplant, undergoing painful, expensive, and only partially effective dialysis while they wait and hope.

In recent years, medical personnel have begun a program of *kidney exchanges*. Suppose that a patient p_1 is incompatible with her prospective donor d_1 , another patient p_2 is incompatible with his prospective donor d_2 , but pairs $\langle p_1, d_2 \rangle$ and $\langle p_2, d_1 \rangle$ are both compatible with each other. Four teams of doctors can then do a "paired exchange" with four surgeries, in which d_1 donates to p_2 and d_2 donates to p_1 . (To ensure that everybody follows through, the surgeries must be simultaneous: if d_1 donates to p_2 before d_2 undergoes surgery, then d_2 has no incentive to go through the surgery, as d_2 's friend p_2 has already received his kidney.) We can even consider larger exchanges (three or more simultaneous donations)—though as the number of surgeries increases, the logistical difficulty increases as well.

Deciding which transplants to complete is done using a graph-based algorithm. Each patient p_i comes to the system with a donor d_i who is willing to donate to p_i . Define a directed graph G as follows. There is a node for each patient p_i and a node for each donor d_i . Add a directed edge $\langle p_i, d_i \rangle$ for every i . Also add a directed edge $\langle d_i, p_j \rangle$ if donor d_i is compatible with patient p_j . A cycle in G then corresponds to a set of surgeries that can be completed: every donor in the cycle donates a kidney, and every patient in the cycle receives a compatible kidney. See Figure 11.54 for an example.

The algorithm that's actually used in the real kidney exchange network in the United States computes a *set* of node-disjoint cycles that will be performed.¹¹ To limit the number of simultaneous surgeries that are required, the algorithm seeks a set of *cycles of length 4 or length 6*—that is, 2 or 3 transplants—in G that maximizes the total number of nodes included. (The constraint on cycle length makes the computational problem much more difficult, so the algorithm requires significant computational power to compute the surgeries to complete.)

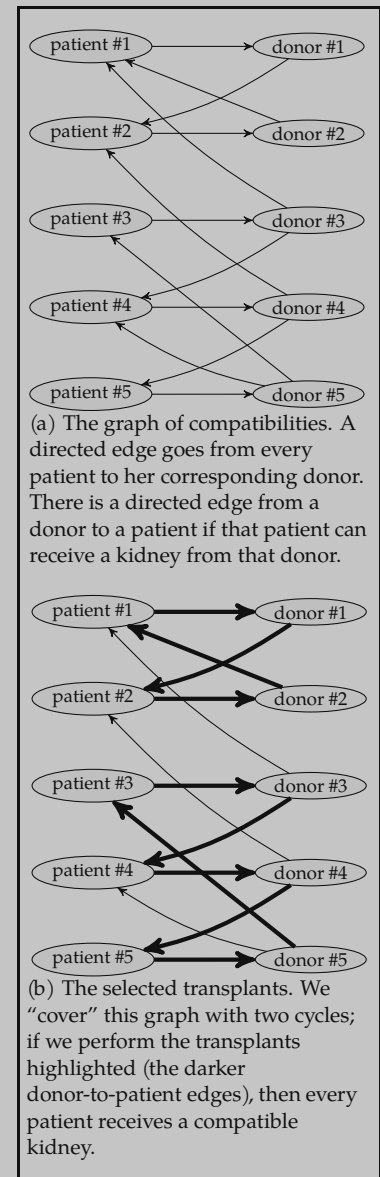


Figure 11.54: An example of a kidney exchange network, and the cycle-based algorithm to select transplants.

¹¹ David Abraham, Avrim Blum, and Tuomas Sandholm. Clearing algorithms for barter exchange markets: Enabling nationwide kidney exchanges. In *Proceedings of the ACM Conference on Electronic Commerce (EC)*, 2007.

COMPUTER SCIENCE CONNECTIONS

BINARY SEARCH TREES

Trees are the basis of many important data structures, of which *binary search trees* are perhaps most frequently used. Binary search trees are data structures that implement the abstract data type called a *dictionary*: we have a set of *keys*, each of which has a corresponding *value*. (For example, the keys might be words and the values definitions, or they might be student names and GPAs, or usernames and encrypted passwords.) The data structure must support operations like *insert*(k, v) (add a new key / value pair) and *lookup*(k) (report the value associated with key k , if any).

A *binary search tree* (BST) is a binary tree for which every node u satisfies the *BST condition* illustrated in Figure 11.55: every node v in u 's left subtree has a key that is less than u 's key, and every node v in u 's right subtree has a key that is greater than u 's key. (For simplicity, assume that all keys are distinct.)

An example of a binary search tree is shown in Figure 11.56. Incidentally, the BST condition implies the following claim: *an in-order traversal of a binary search tree visits the keys in sorted order*. This claim can be proven formally by induction, but the intuition is straightforward: an in-order traversal of a node with key x first visits nodes $< x$ (while traversing the left subtree), then x itself, and then nodes $> x$ (while traversing the right subtree). Because, recursively, the nodes of the left and right subtrees are themselves visited in sorted order, the entire tree's keys are visited in sorted order.

Binary search trees are good data structures for dictionaries because *insert* and *lookup* can be implemented simply and efficiently. If we perform a lookup for a key k in an empty BST T , we return “not found.” (For simplicity, we allow a BST to be empty—that is, to contain zero nodes.) Otherwise, compare k to the key r stored in the root node of T :

- if $k = r$, then return the value stored at the root.
- if $k < r$, then perform a lookup for k in the left subtree.
- if $k > r$, then perform a lookup for k in the right subtree.

The BST condition guarantees that we find the node with key k if it's in the tree. (You can prove this fact by induction.) The *insert* operation can be implemented similarly, by adding a new node exactly where a lookup for the key k would have found k .

The worst-case running time of *lookup* and *insert* is proportional to the height of the binary search tree. More “balanced” BSTs—in which every internal node has a left subtree with roughly the same height as its right subtree—have better performance. (There are many different BSTs with the same set of keys; for example, another BST that has the same keys as the BST in Figure 11.56 is shown in Figure 11.57.)

Most software therefore uses *balanced binary search trees* instead—for example, *AVL trees* or *red–black trees*.¹² (See p. 643 for further discussion of AVL trees, and a proof of their efficiency.)

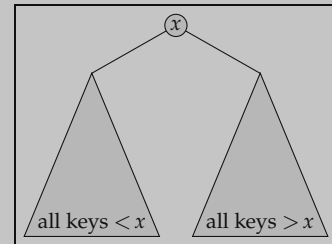


Figure 11.55: The binary search tree condition. For every node with key x : all keys in the left subtree of the node have a key $< x$; and all keys in the right subtree of the node have a key $> x$.

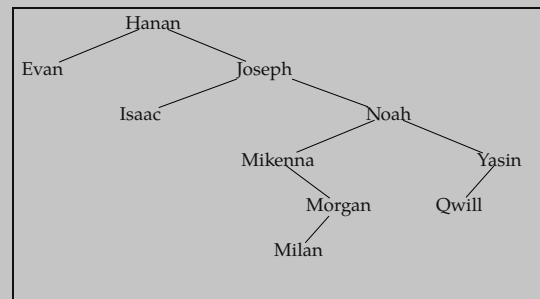


Figure 11.56: A binary search tree storing a set of 10 keys. The key is shown in each node; the accompanying value isn't drawn.

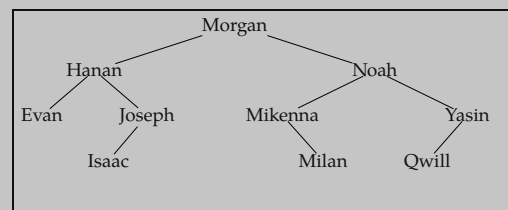


Figure 11.57: Another binary search tree with the same set of keys.

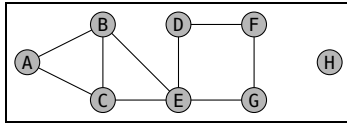
See the details in any good textbook on data structures, or in

¹² Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein. *Introduction to Algorithms*. MIT Press, 3rd edition, 2009.

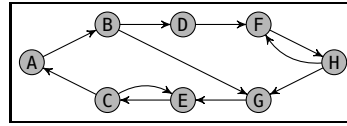
11.4.5 Exercises

Identify all of the simple cycles in the following graphs:

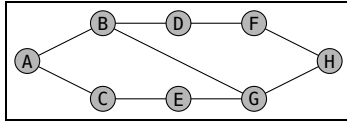
11.122



11.124



11.123



Consider an undirected graph G with n nodes. In terms of n ...

11.125 ... what is the longest simple cycle that G can contain? Explain.

11.126 ... what is the longest cycle (not necessarily simple) that G can contain? Explain.

Prove your answers to the following questions, and simplify your answer as n gets large. (For handling large n , a useful fact from calculus: $\sum_{i=0}^n \frac{1}{i!}$ approaches $e = 2.71828 \dots$ as n grows.)

11.127 In the n -node complete graph K_n , how many simple cycles is a particular node u involved in?

11.128 Let u be a node in a n -node complete directed graph: all edges except for self-loops are present. How many simple cycles is node u involved in?

11.129 A small modification to BFS can detect cycles involving a node s in a directed graph, as shown in Figure 11.58. However, this modification doesn't quite work for undirected graphs. Give an example of an acyclic graph in which the algorithm Figure 11.58 falsely claims that there is a cycle. Then describe briefly how to modify this algorithm to correctly detect cycles involving node s in undirected graphs.

Recall Lemma 11.5: in any acyclic undirected graph, there exists a node whose degree is zero or one. Prove the following two extensions/variations of this lemma:

11.130 Prove that every directed acyclic graph contains a node with out-degree zero.

11.131 Prove that there are two nodes of degree 1 in any acyclic undirected graph that contains at least one edge.

Recall Definition 11.26: a cycle $\langle u_0, u_1, \dots, u_k, u_0 \rangle$ is a path of length ≥ 2 from a node u_0 back to node u_0 that does not traverse the same edge twice. At various times in class, I've tried to define cycles in all of the following ways—and they're all bogus definitions, in the sense that they describe something different from Definition 11.26. For each of the following broken definitions, explain why I was wrong:

11.132 A cycle is a simple path from s to s .

11.133 A cycle is a path of length ≥ 2 from s to s .

11.134 A cycle is a path from s to s that doesn't traverse any edge more than once.

11.135 A cycle is a path from s to s that includes at least 3 distinct nodes.

11.136 A cycle is a path of length ≥ 2 from s to s that doesn't traverse any edge twice consecutively.

11.137 Definition 11.28 defines an acyclic graph as one containing no cycles, but it would have been equivalent to define acyclic graphs as those containing no simple cycles. Prove that G has a cycle if and only if G has a simple cycle.

Recall that $G = \langle V, E \rangle$ is a regular graph if every $u \in V$ has $\text{degree}(u) = d$, for some fixed constant d .

11.138 Identify two different regular graphs that are trees.

11.139 It turns out that there are two and only two different trees T that are regular graphs. Prove that there are no other regular graphs that are trees.

Input: a graph $G = \langle V, E \rangle$ and a source node $s \in V$
Output: is s involved in a cycle in G ?

```

1: Frontier := {s}
2: Known := ∅
3: while Frontier is nonempty:
4:   u := the first node in Frontier
5:   remove u from Frontier
6:   if s is a neighbor of u then
7:     return "s is involved in a cycle."
8:   for every neighbor v of u:
9:     if v is in neither Frontier nor Known then
10:      add v to the end of Frontier
11:   add u to Known
12: return "s is not involved in a cycle."

```

Figure 11.58: BFS modified (slightly buggily) to detect cycles involving the node s .

A triangle is a simple cycle containing exactly three nodes. A square is a simple cycle containing exactly four nodes.

11.140 What is the largest number of triangles possible in an undirected graph of n nodes?

11.141 What is the largest number of squares possible in an undirected graph of n nodes?

Let's analyze the largest number of edges that are possible in an n -node undirected graph that contains no triangles.

11.142 Consider a triangle-free graph $G = \langle V, E \rangle$. For nodes $u \in V$ and $v \in V$, argue that if $\{u, v\} \in E$, then we have $\text{degree}(u) + \text{degree}(v) \leq |V|$.

11.143 Prove the following claim by induction on the number of nodes in the graph: if $G = \langle V, E \rangle$ is triangle-free, then $|E| \leq |V|^2/4$. (Hint: use the previous exercise.)

11.144 Give an example of an n -node triangle-free graph that contains $\frac{n^2}{4}$ edges.

Consider the following adjacency lists. Is the graph that each represents a tree? Justify your answers.

11.145

```
A: B, E
B: A
C: D
D: C, F
E: A
F: D
```

11.146

```
A: C
B: C, E
C: A, B, F
D: E
E: B, D
F: C
```

11.147

```
A: D
B: E, F
C: D, F
D: A, C
E: B
F: B, C
```

11.148

```
A: C, D, F
B: F
C: A, E, F
D: A
E: C
F: A, B, C
```

Prove or disprove the following claims about trees:

11.149 There is a node of degree equal to 2 in any tree with ≥ 3 nodes.

11.150 In any rooted binary tree (all nodes have 0, 1, or 2 children), there are an even number of leaves.

11.151 If a graph $G = \langle V, E \rangle$ has $|V| - 1$ edges, then G must be a forest.

11.152 The following pair of definitions is subtly broken: the *root* of a tree is a node that is not a child, and a *leaf* is a node that is a child but not a parent. What's broken?

For the tree in Figure 11.59, with node A as the root ...

11.153 ... what are the leaves?

11.154 ... which nodes are internal nodes?

11.155 ... what are the parent, children, and siblings of node D?

11.156 ... what are the descendants of node D?

11.157 ... what are the ancestors of node F?

11.158 ... what is the height of the tree?

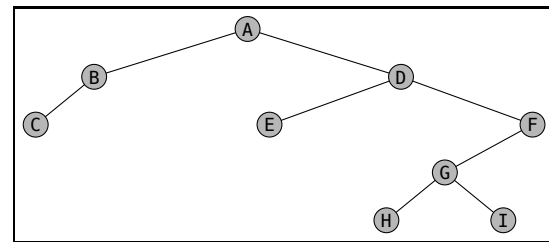


Figure 11.59: A rooted tree.

11.159 Let T be an arbitrary n -node rooted tree, with root r and with ℓ different leaves. Prove or disprove: if we reroot T at a new node $r' \neq r$, then the number of leaves remains exactly ℓ .

A complete binary tree of height h has “no holes”: reading from top-to-bottom and left-to-right, every node exists. Complete binary trees form a subset of nearly complete binary trees: a nearly complete binary tree has every node until the last row, which is allowed to stop early. (See Figure 11.60, and see also p. 529 for a discussion of heaps, which are a data structure represented as a nearly complete binary tree.)

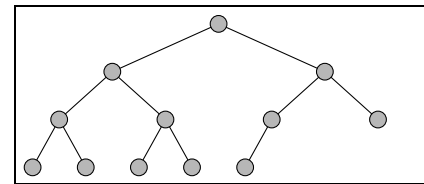
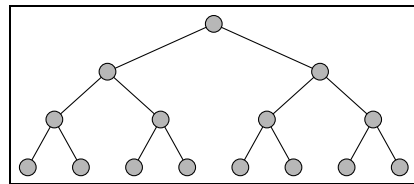


Figure 11.60: A complete and nearly complete binary tree of height 3.

11.160 Prove by induction that a complete binary tree of height h contains precisely $2^{h+1} - 1$ nodes.

11.161 How many leaves does a nearly complete binary tree of height h have? Give the smallest and largest possible values, and explain.

11.162 What is the diameter of a nearly complete binary tree of height h ? Again, give the smallest and largest possible values, and explain your answer. (Recall that the *diameter* of a graph $G = \langle V, E \rangle$ is $\max_{s,t \in V} d(s, t)$, where $d(s, t)$ denotes the length of the shortest path from s to t in G .)

Suppose that we “rerooted” a complete binary tree of height h by instead designating one of the erstwhile leaves as the root. In the rerooted tree, what are the following quantities?

- 11.163 the height
 11.164 the diameter
 11.165 the number of leaves

Justify your answers to the following questions: describe an 1000-node binary tree with ...

- 11.166 ... height as large as possible. 11.168 ... as many leaves as possible.
 11.167 ... height as small as possible. 11.169 ... as few leaves as possible.

11.170 What is the largest possible height for an n -node binary tree in which every node has precisely zero or two children? Justify your answer.

In what order are nodes of the tree in Figure 11.61 traversed ...

- 11.171 ... by a pre-order traversal?
 11.172 ... by an in-order traversal?
 11.173 ... by a post-order traversal?

11.174 Draw the binary tree with in-order traversal 4, 1, 2, 3, 5; pre-order traversal 1, 4, 3, 2, 5; and post-order traversal 4, 2, 5, 3, 1.

11.175 Do the same for the tree with in-order traversal 1, 3, 5, 4, 2; pre-order traversal 1, 3, 5, 2, 4; and post-order traversal 4, 2, 5, 3, 1.

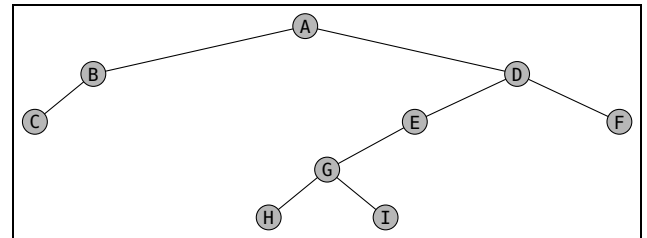


Figure 11.61: A rooted tree.

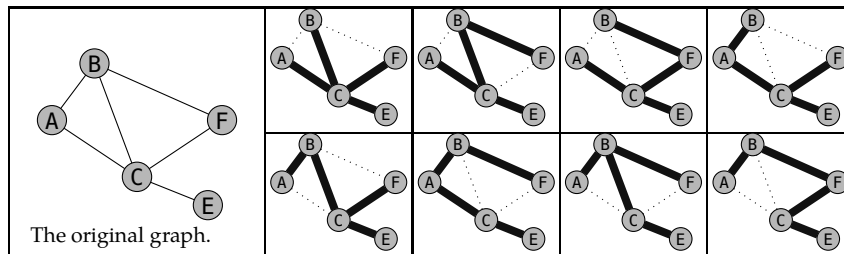
11.176 Describe (that is, fully explain the structure of) an n -node binary tree T for which the *pre-order* and *in-order* traversals of T result in precisely the same ordering of T 's nodes. (That is, $\text{pre-order-traverse}(T) = \text{in-order-traverse}(T)$.)

11.177 Describe a binary tree T for which the *pre-order* and *post-order* traversals result in precisely the same ordering of T 's nodes. (That is, $\text{pre-order-traverse}(T) = \text{post-order-traverse}(T)$.)

11.178 Prove that there are two distinct binary trees T and T' such that pre-order and post-order traversals are both identical on the trees T and T' . (That is, $\text{pre-order-traverse}(T) = \text{pre-order-traverse}(T')$ and $\text{post-order-traverse}(T) = \text{post-order-traverse}(T')$ but $T \neq T'$.)

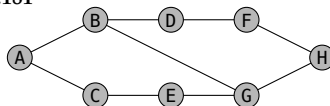
11.179 Give a recursive algorithm to reconstruct a tree from the in-order and post-order traversals.

11.180 Argue that we didn't leave out any spanning trees of G in Figure 11.51, reproduced here for your convenience:

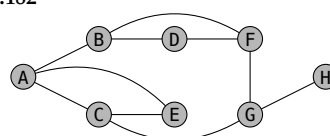


How many spanning trees do the following graphs have? Explain.

11.181



11.182



11.5 Weighted Graphs

Force without wisdom falls of its own weight.

Horace (65–8 BCE), *Odes* (23 BCE)

Many real-world situations are naturally modeled by different edges having different “weights”: the price of an airplane flight, the closeness of a friendship, the physical length of a road, the time required to transmit data across an internet connection.

These graphs are called *weighted graphs*:

Definition 11.32 (Weighted graph)

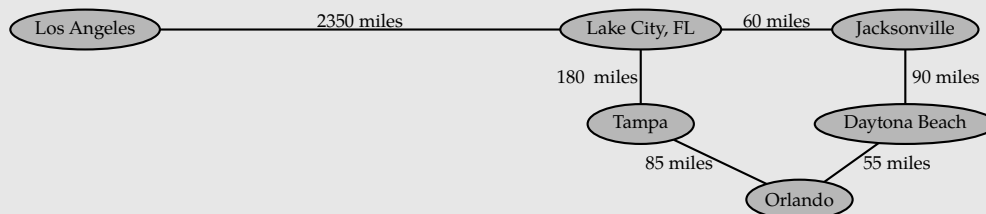
A weighted graph is a graph $G = \langle V, E \rangle$ and a weight function $w : E \rightarrow \mathbb{R}^{\geq 0}$, so that each edge $e \in E$ has a weight $w(e) \geq 0$. For simplicity of notation, we’ll often write w_e instead of $w(e)$; we’ll also sometimes refer to w_e as the length of the edge e .

In a weighted graph, the length of a path in a weighted graph is the sum of the lengths of the edges traversed by the path. (A shortest path is, as before, one with the smallest length.)

Either undirected or directed graphs can be weighted. Aside from the length of a path, all of the other notions and terminology from unweighted graphs carry over: neighbors and degree, paths and connectivity, and so forth. Weighted graphs can be represented just as unweighted graphs were: we typically store the weight of edge $\langle u, v \rangle$ directly in the $\langle u, v \rangle$ th entry of the adjacency matrix, or attach the edge weight as an additional slot in the adjacency list entries. Here’s an example:

Example 11.37 (A weighted graph)

Here’s the highway system from Example 11.4, where each road is labeled with its length:



There are two simple paths between *Orlando* and *Lake City*:

- *Orlando* \leftrightarrow *Tampa* \leftrightarrow *Lake City*: $85 + 180 = 265$ miles.
- *Orlando* \leftrightarrow *Daytona Beach* \leftrightarrow *Jacksonville* \leftrightarrow *Lake City*: $55 + 90 + 60 = 205$ miles.

The second path is shorter, even though it traverses more edges, as $265 > 205$.

Taking it further: The primary job of a web search engine is to respond to a user’s search query (“give me web pages about Horace”) with a list of relevant pages. There’s a complex question of data structures, parallel computing, and networking infrastructure in solving even the simplest part of this task: identifying the set R of web pages (out of many billions) that contain the search term. A subtler challenge—and at least as important—is figuring out how to *rank* the set R . What pages in R are the “most important,” the ones that we should display on the first page of results? See p. 1174 for some discussion of how Google uses a weighted graph (and probability) to do this ranking.

Definition 11.32 considers only nonnegative weights—every $w_e \geq 0$ —which is a genuine restriction. (For example, the “signed” social networks from Figure 11.8(a) have positive and negative weights signifying friendship and enmity.) Some, but not all, of the results that we’ll discuss in this section carry over to the setting of negative weights.

11.5.1 Shortest Paths in Weighted Graphs: Dijkstra's Algorithm

A *shortest path* from s to t in a weighted graph is the path connecting s and t that has shortest total length. In many natural applications where shortest paths are useful, we have weights on edges: you want the *shortest* walking route from the bar back to your apartment, for example, not necessarily the one with the fewest turns. In Example 11.37, we already saw a case in which the shortest path used more edges than necessary. Thus we cannot directly use breadth-first search to compute distances in weighted graphs.

But we *can* compute distances using an algorithm that's very similar in spirit to BFS. The basic idea of breadth-first search is to “expand outward” from the source node s in layers, accumulating a set of nodes u for which we know the distance from s to u . We add nodes in increasing order of their distance from s , and eventually we've computed distances from s to all nodes in the graph. (See Figure 11.62.) The trouble for weighted graphs is that the order in which BFS builds up its knowledge about shortest paths doesn't always work (as in Example 11.37). But we can use a cleverer way of building up knowledge about the network to find shortest paths in weighted graphs, too.

The algorithm that we'll describe is due to Edsger Dijkstra, and hence it is known as *Dijkstra's algorithm*. The key idea of Dijkstra's algorithm has parallels with BFS:

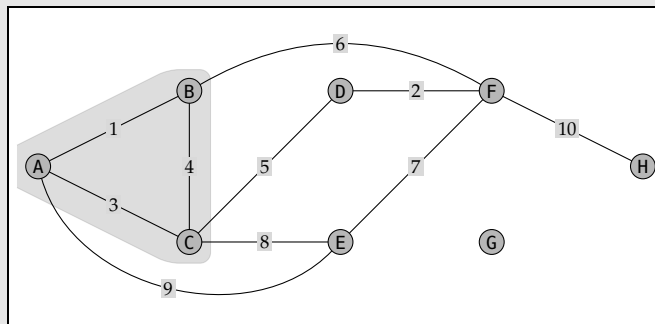
Suppose that we know the distance from a source node s to every node in some set S of nodes. (Assume that $s \in S$.) We will find some node *not* in S for which we can determine the shortest path from s .

For now, let's not worry about where this set S came from; the key point is just that we are assuming that we know distances to certain nodes (those in S), and we seek to leverage that existing knowledge to learn the distance to some other node (not previously in S). We'll then add that new node to S and iterate.

Before we state the formal result, let's look at an example:

Example 11.38 (An example of distances)

Consider the following weighted, undirected graph (with edge weights marked on the edges):



Suppose we know the distances from A to every node in the shaded set $S = \{A, B, C\}$:

$$d(A, A) = 0$$

$$d(A, B) = 1$$

$$d(A, C) = 3.$$

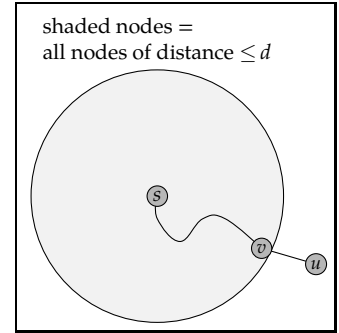


Figure 11.62: The intuition of BFS. Assume the shaded set S contains every node within distance d of s , and that $u \notin S$ is a neighbor of $v \in S$. The distance from s to u must be $d + 1$.

Edsger Dijkstra was a 20th-century Dutch computer scientist—one of the founders of theoretical computer science, and the 1972 Turing Award winner.

Irrelevant quotation: “Computer science is no more about computers than astronomy is about telescopes.” — attributed to Edsger W. Dijkstra (1930–2002)

Irrelevant challenge: Name a common English word that, like DIJKSTRA, has at least five (or 6 or even 7, which is technically possible) consecutive consonants. (Not SYZYGY or RHYTHMS; Y is a vowel if it's used as a vowel!)

We wish to expand our set of known nodes by adding a neighbor of an already shaded node. The candidate nodes that are neighbors of nodes with known distances are $\{D, E, F\}$. In particular, their candidate distances are:

node	distance
F (via B)	$d(A, B) + w_{B,F} = 1 + 6 = 7$
E (via A)	$d(A, A) + w_{A,E} = 0 + 9 = 9$
E (via C)	$d(A, C) + w_{C,E} = 3 + 8 = 11$
D (via C)	$d(A, C) + w_{C,D} = 3 + 5 = 8$

Let's argue that we can now conclude that $d(A, F) = 7$.

The key reason is that, to get from A to F, we have to “escape” the set of shaded nodes—and every “escape route” (path to F) must reach its last shaded node v (that's $d(A, v)$) and then follow an edge to its first unshaded node u (that's $w_{v,u}$). Because this table tells us that every path out of the shaded region has length at least 7, and we've found a path from A to F with exactly that length, we conclude that $d(A, F) = 7$.

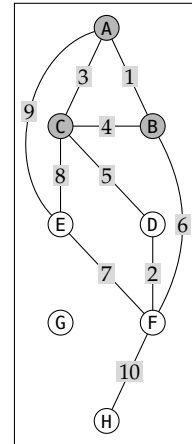


Figure 11.63:
The graph for
Example 11.38,
repeated and
rotated. We've
computed that
 $d(A, A) = 0$ and
 $d(A, B) = 1$ and
 $d(A, C) = 3$.

COMPUTING THE DISTANCE TO A NEW NODE

The same basic reasoning that we used in Example 11.38 will allow us to prove a general observation that's the foundation of Dijkstra's algorithm:

Lemma 11.9 (Foundation of Dijkstra's Algorithm)

Let $G = \langle V, E \rangle$ be a graph with edge weights w , let $S \subset V$ be a set of nodes, and let $s \in S$ be a source node. Let $d(s, v)$ denote the distance from s to v for every node v in S . For a node $u \notin S$, define

$$d_u := \min_{v \in S : u \text{ is a neighbor of } v} d(s, v) + w_{v,u}.$$

Let u^* be the node $u \notin S$ for which d_u is minimized. Then the distance from s to u^* is d_{u^*} .

Before we prove the lemma, let's restate the claim in slightly less notation-heavy English. (See Figure 11.64.) We have a set S of nodes—the shaded nodes in the picture—for which we know the distance from s . We examine all unshaded nodes u that are neighbors of shaded nodes v . For each shaded/unshaded pair, we've computed the sum of the distance $d(s, v)$ and the edge weight $w_{v,u}$. And we've chosen the pair $\langle v^*, u^* \rangle$ that minimizes this quantity.

The lemma says that the shortest path from s to this particular u^* must have length precisely equal to $d_{u^*} := d(s, v^*) + w_{v^*, u^*}$. The intuition matches the argument in Example 11.38: to get from s to u^* , we have to somehow “escape” the set of shaded nodes—and, by the way that we chose u^* , every “escape route” must have length at least d_{u^*} .

Proof of Lemma 11.9. We must show that the distance from s to u^* is d_{u^*} , and we'll do it in two steps: by showing that the distance is no more than d_{u^*} , and by showing that the distance is no less than d_{u^*} .

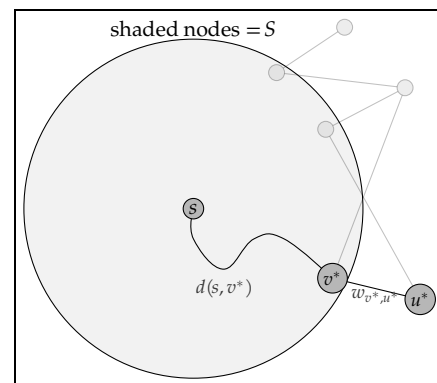
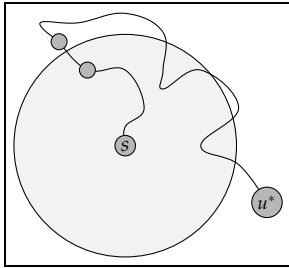


Figure 11.64:
The intuition for
Lemma 11.9.

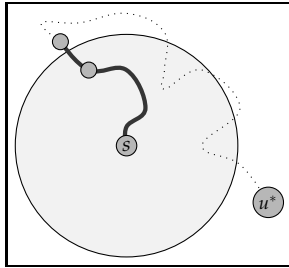
The distance from s to u^* is $\leq d_{u^*}$. We must argue that there is a path of length $d(s, v^*) + w_{v^*, u^*}$ from s to u^* . By assumption and the fact that $v^* \in S$, we know that $d(s, v^*)$ is the distance from s to v^* , so there must exist a path P of length $d(s, v^*)$ from s to v^* . (It's the curved line in Figure 11.64.) By tacking u^* onto the end of P , we've constructed a path from s to u^* via v^* with length $d(s, v^*) + w_{v^*, u^*}$.

The distance from s to u^* is $\geq d_{u^*}$. Consider an arbitrary path P from s to u^* . We must show that P has length at least $d(s, v^*) + w_{v^*, u^*}$.

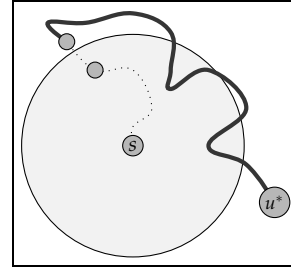
What does P look like? The node s is in the set S , so P starts out at $s \in S$, then wanders around for a while inside S , then crosses outside of S for the first time, wanders around outside S for a while, and eventually ends up at $u^* \notin S$. Nothing prevents P from re-entering (and later re-exiting) S after its first departure—indeed, it can go in and out of S several times—but it definitely has to leave S at least once. Thus P has to look like the following:



(a) the entire path P



(b) the portion of P up to the first exit from S



(c) the portion of P after the first exit from S

Therefore we know that

the length of P

$$= (\text{the length of } P \text{ up to the first exit}) + (\text{the length of } P \text{ after the first exit})$$

$$\geq (\text{the length of the shortest path exiting } S) + (\text{the length of } P \text{ after the first exit})$$

P up to the first exit is a path exiting S , so its length is at least the length of the shortest such path

$$\geq d(s, v^*) + w_{v^*, u^*} + (\text{the length of } P \text{ after the first exit})$$

we chose u^ and v^* so that $d(s, v^*) + w_{v^*, u^*}$ is exactly the length of the shortest path exiting S*

$$\geq d(s, v^*) + w_{v^*, u^*} + 0$$

all edge weights are nonnegative, so all path lengths are ≥ 0 too

$$= d_{u^*}.$$

definition of d_{u^}*

Thus the length of P is at least d_{u^*} .

We've therefore argued that the distance from s to u^* is both $\leq d_{u^*}$ and $\geq d_{u^*}$. Thus the distance is precisely d_{u^*} , and the lemma follows. \square

DIJKSTRA'S ALGORITHM

With Lemma 11.9 proven, we can now put together the pieces of the entire algorithm. The lemma describes a way to take a set S of nodes with known distance from the source node s , and correctly calculate the distance from s to a new node $u \notin S$.

Problem-solving tip:
When we want to prove that $x = y$, it's sometimes easier to prove $x \geq y$ and $x \leq y$ separately.

In Dijkstra’s algorithm, the idea is to apply the calculation from Lemma 11.9 repeatedly to find all distances from the given source node s . We’ll need a base case to get started, but that’s straightforward: we start with the set of nodes with known distance from s as $S = \{s\}$, where the distance from s to s is zero. The full algorithm is shown in Figure 11.65.

Before we prove the algorithm’s correctness, let’s run through an example:

Dijkstra’s Algorithm:
Input: a weighted graph $G = \langle V, E \rangle$, nonnegative edge weights $w_e \geq 0$, and a source node $s \in V$.
Output: the distance from s to every node in G

1: Let $S := \{s\}$ and let $d(s, s) := 0$.
2: **while** there exists a node in S with a neighbor not in S :
3: for every node $u \notin S$, define
$$d_u := \min_{v \in S : u \text{ is a neighbor of } v} d(s, v) + w_{v,u}.$$

4: $u^* :=$ the node with the smallest d_u .
5: Add u^* to S and set $d(s, u^*) := d_{u^*}$.
6: **for** every node $u \in V - S$:
7: $d(s, u) := \infty$
8: **return** the recorded values $d(s, u)$.

Figure 11.65: The pseudocode for Dijkstra’s algorithm.

Example 11.39 (Dijkstra’s algorithm in action)
Let’s run Dijkstra’s algorithm on the network from Example 11.37, with the graph rotated for compactness. We’ll start from the Orlando (OR) node. Here is the execution:

	DB	JA	LA	LC	OR	TA
Iteration 1					0	
Iteration 2	55				0	
Iteration 3	55				0	85
Iteration 4	55	145			0	85
Iteration 5	55	145		205	0	85
Iteration 6	55	145	2555	205	0	85

THE CORRECTNESS OF DIJKSTRA'S ALGORITHM

We'll now prove the correctness of the algorithm, using Lemma 11.9 and induction:

Theorem 11.10 (Correctness of Dijkstra's Algorithm)

Let $G = \langle V, E \rangle$ be a graph with nonnegative edge weights w_e for each edge. Let $s \in V$ be a source node, and let $d(s, \bullet) := \text{Dijkstra}(G, w, s)$ be the values computed by Dijkstra's Algorithm. Then, for every node u , we have that $d(s, u)$ is the length of the shortest path from s to u in G under w .

Proof. Looking at the algorithm, we see that Dijkstra's Algorithm records finite distances from s in Line 1 (for s itself) and Line 5 (for other nodes reachable from s). Suppose that Dijkstra's algorithm executes n iterations of the loop in Line 2, thus recording $n + 1$ total distances in Lines 1 and 5—say in the order u_0, u_1, \dots, u_n . Let $P(i)$ denote the claim that $d(s, u_i)$ is the length of the shortest s -to- u_i path. We claim by strong induction on i that $P(i)$ holds for all $i \in \{0, 1, \dots, n\}$.

Base case ($i = 0$): We must prove that $d(s, u_0)$ is recorded correctly. The 0th node u_0 is recorded in Line 1, so u_0 is the source node s itself. And the shortest path from s to s in any graph with nonnegative edge weights is the 0-hop path $\langle s \rangle$, of length 0.

Inductive case ($i \geq 1$): We assume the inductive hypothesis $P(0), P(1), \dots, P(i - 1)$: that is, all recorded distances $d(s, u_0), d(s, u_1), \dots, d(s, u_{i-1})$ are correct. We must prove $P(i)$: that is, that the recorded distance $d(s, u_i)$ is correct. But this follows immediately from Lemma 11.9: the algorithm chooses u_i as the $u \notin S$ minimizing

$$d_u := \min_{v \in S : u \text{ is a neighbor of } v} d(s, v) + w_{v,u},$$

where $S = \{u_0, u_1, \dots, u_{i-1}\}$. Lemma 11.9 states precisely that this value d_u is the length of the shortest path from s to u .

Finally, observe that any node u that's only discovered in Line 6 is not reachable from s , and so indeed $d(s, u) = \infty$. (A fully detailed argument that the ∞ values are correct can follow the structure in Theorem 11.3, which proved the correctness of BFS.) \square

Taking it further: Dijkstra's algorithm as written in Figure 11.65 can be straightforwardly implemented to run in $O(|V| \cdot |E|)$ time: each iteration of the **while** loop (Line 2) can look at each edge to compute the smallest d_u . But with cleverer data structures, Dijkstra's algorithm can be made to run in $O(|E| \log |V|)$ time. This improved running-time analysis, as well as other shortest-path algorithms—for example, handling the case in which edge weights can be negative (it's worth thinking about where the proof of Lemma 11.9 fails if an edge e can have $w_e < 0$), or computing distances between *all pairs* of nodes instead of just every distance from a *single source*—is a standard topic in a course on algorithms. Any good algorithms text should cover these algorithms and their analysis.

Before we leave Dijkstra's algorithm, it's worth reflecting on its similarities with BFS. In both cases, we start from a seed set S of nodes for which we know the distance from s —namely $S = \{s\}$. Then we build up the set of nodes for which we know the distance from s by finding the unknown nodes that are closest to s , and adding them to S . Of course, BFS is conceptually simpler, but Dijkstra's algorithm solves a more complicated problem. It's a worthwhile exercise to think about what happens if Dijkstra's algorithm is run on an unweighted graph. (How does it relate to BFS?)

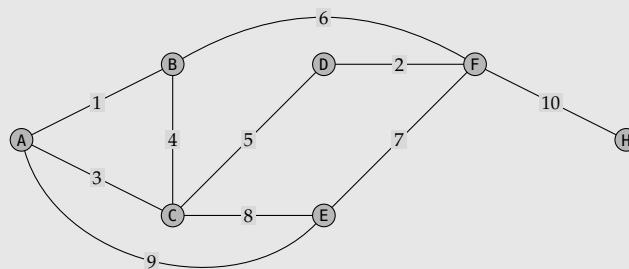
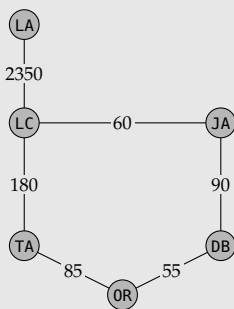
11.5.2 Spanning Trees in Weighted Graphs: Minimum Spanning Trees

Recall from Definition 11.31 that a *spanning tree* of a connected graph $G = \langle V, E \rangle$ is a tree $T = \langle V, E' \rangle$ where $E' \subseteq E$. As with shortest paths, in many of the applications in which spanning trees are interesting, we actually want to find a spanning tree whose edges have minimum possible total cost. For example, when a college wants to put down networking cable in a new dorm building, they wish to ensure that the resulting network is connected, while minimizing the cost of construction.

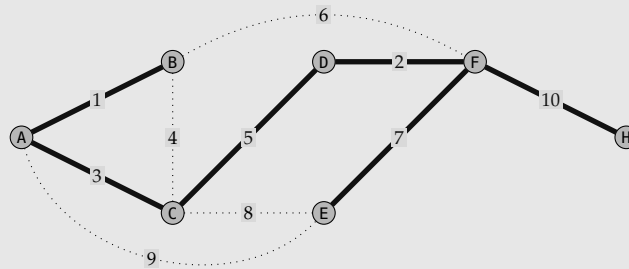
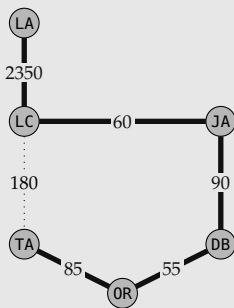
Formally, in a weighted graph, the *cost* of a spanning tree T is the sum of the weights of its edges: $\sum_{e \in E'} w_e$. A *minimum spanning tree (MST)* is a spanning tree whose cost is as small as possible. Here are two small examples:

Example 11.40 (Some minimum spanning trees)

Consider the following two graphs (the road network from Example 11.37 and the larger connected component from Example 11.38):



Here are the minimum spanning trees. (For the first, every spanning tree omits exactly one edge from the lone cycle; the cheapest tree omits the most expensive edge.)



As with shortest paths in weighted graphs, the question of how to find a minimum spanning tree most efficiently is more appropriate to an algorithms text than this book. But, between the Cycle Elimination Algorithm (Figure 11.52) and Example 11.40, we've already done almost all the work to develop a first algorithm.

Assume throughout that all edge weights are distinct. (This assumption lets us refer to “the most expensive edge” in a set of edges. Removing this assumption complicates the language that we have to use, but it doesn't fundamentally change anything about the MST problem or its solution.)

Lemma 11.11 (The “cycle rule”)

Let C be a cycle in a connected undirected graph $G = \langle V, E \rangle$, and let e be the heaviest edge in C . Then e is not in any minimum spanning tree of G .

Proof. Consider a spanning tree T of G , and suppose that $e = \{u, v\}$ is included in T . We’ll show that T is not a *minimum* spanning tree. (Thus the only minimum spanning trees of G do not include e .)

By definition, the spanning tree T is connected. If we delete $\{u, v\}$ from T , the resulting graph will have two connected components, one containing u and the other containing v . (This fact follows by Corollary 11.7.) Call those connected components U and V , respectively. See Figure 11.66(a).

Imagine following the cycle C from u to v the “long way” around C . This part of C starts at u , wanders around U for a while, and eventually crosses over into V , before finally arriving at v . Let $a \in U$ be the last node in U and b the first node in V as we go around C . (Note that C might go back and forth between U and V multiple times, but define a and b based on the *first* time C leaves U .) See Figure 11.66(b).

Now define the graph T' as T with the edge $\{u, v\}$ removed and with the edge $\{a, b\}$ inserted instead. Crucially, T' is a spanning tree of G ; because we’ve only swapped *which* edge connected the connected sets U and V . Thus T' remains connected and acyclic.

Now observe that the cost of T' is less than the cost of T , because the edge $\{u, v\}$ is heavier than the edge $\{a, b\}$. (Both $\{u, v\}$ and $\{a, b\}$ are in the cycle C , and by assumption $\{u, v\}$ is the heaviest edge in C .) But therefore T' is a cheaper spanning tree than T , and thus T isn’t a minimum spanning tree. \square

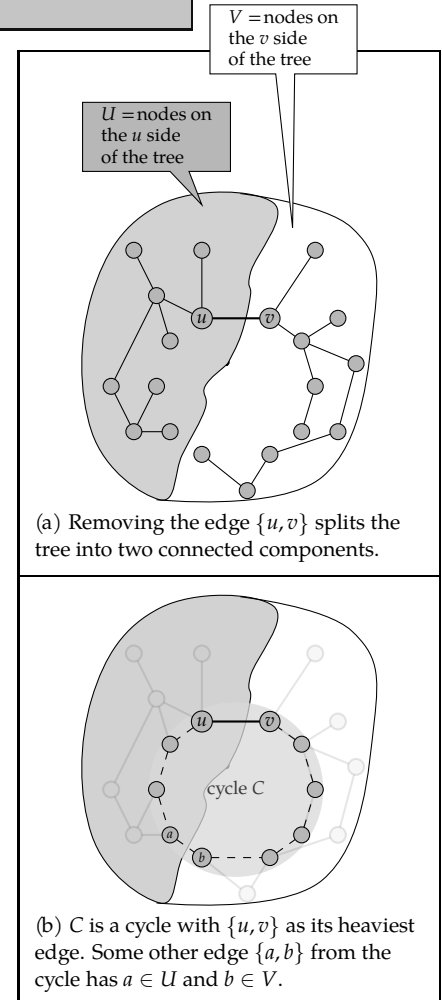


Figure 11.66: The cycle rule for MSTs.

FINDING MSTs BY REMOVING CYCLES

Lemma 11.11 immediately suggests that we can find minimum spanning trees using a modification of the Cycle Elimination Algorithm:

Weighted Cycle Elimination Algorithm

Input: a weighted connected graph $G = \langle V, E \rangle$ with edge weights w_e

Output: a minimum spanning tree of G

- 1: **while** there exists a cycle C in G :
- 2: let e be the *heaviest* edge traversed by C
- 3: remove e from E
- 4: **return** the resulting graph $\langle V, E \rangle$.

While the Weighted Cycle Elimination Algorithm is correct and reasonably efficient, there are more efficient algorithms based on Lemma 11.11. One such algorithm is called *Kruskal’s Algorithm*, named after its discoverer Joseph Kruskal. The key idea of Kruskal’s Algorithm is that by *sorting* the edges in increasing order, we can be more efficient: we add edges in increasing order of their weight, as long as doing so doesn’t create a cycle.

Joseph Kruskal was a 20th-century American computer scientist/mathematician/statistician. He published his MST algorithm in 1956.

The insight of this algorithm is that, by considering edges in increasing order of weight, if including an edge e creates a cycle, then we know that e must be the heaviest edge in that cycle. See Figure 11.67. Kruskal’s algorithm is reasonably efficient: the sorting step takes $O(|E| \log |E|)$ time, and each of the $|E|$ iterations of the **for** loop can be implemented using one call to BFS to test for a cycle. (And, in fact, there are some cleverer ways to implement Line 4 so that the entire algorithm runs in $O(|E| \log |E|)$ time.) Here’s an example:

Kruskal’s Algorithm
Input: a weighted connected graph $G = \langle V, E \rangle$ with distinct edge weights w_e
Output: a minimum spanning tree of G
1: Sort the edges e in increasing order of weight.
2: $S := \emptyset$
3: **for** each edge e (taken in increasing order of weight):
4: **if** the graph $\langle V, S \cup \{e\} \rangle$ doesn’t contain a cycle **then**
5: add e to S
6: **return** the resulting graph $\langle V, S \rangle$

Figure 11.67:
Kruskal’s Algo-
rithm.

Example 11.41 (Sample run of Kruskal’s algorithm)
In each panel, the highlighted edge is being considered for inclusion in the tree.
Black edges have already been included; light edges have not yet been considered.

	The original graph.
	We examine the cheapest edge {A, C}. It doesn’t create a cycle, so we keep it.
	We examine the next cheapest edge {B, C}. It doesn’t create a cycle, so we keep it.
	We examine the next cheapest edge {C, D}. It doesn’t create a cycle, so we keep it.
	We examine the next cheapest edge {A, B}. It creates a cycle $\langle A, B, C, A \rangle$, so we discard it.
	The next edge is {D, E}; we keep it.
	The next edge is {B, D}; it creates a cycle, so we discard it.
	We last edge is {C, E}; it creates a cycle, so we discard it too.
	The final spanning tree.

Here is the general statement of correctness for both algorithms:

Theorem 11.12 (Correctness of minimum spanning tree algorithms)

The Weighted Cycle Elimination Algorithm and Kruskal's Algorithm both return a minimum spanning tree for any weighted connected undirected graph.

Proof. The correctness of the Weighted Cycle Elimination Algorithm follows immediately from Lemma 11.11 (the cycle rule) and from Theorem 11.8 (the correctness of the Cycle Elimination Algorithm): the heaviest edge in any cycle does not appear in any MST, and we terminate with a spanning tree when we repeatedly eliminate any edge from an arbitrarily chosen cycle.

For Kruskal's algorithm, consider an edge e that is *not* retained—that is, when e is considered, it is not included in the set S . The only reason that e wasn't included is that adding it would create a cycle C involving e and previously included edges—but because the edges are considered in increasing order of weight, that means that e is the heaviest edge in C . Thus by Lemma 11.11, Kruskal's algorithm removes only edges not contained in any minimum spanning tree. Because it only excludes edges that create cycles, the resulting graph is also connected—and thus a minimum spanning tree. \square

Taking it further: There are several other commonly used algorithms for minimum spanning trees, using different structural properties than the Cycle Rule. For much more on these other algorithms, and for the clever data structures that allow Kruskal's Algorithm to be implemented in $O(|E| \log |E|)$ time, see any good textbook on algorithms.

COMPUTER SCIENCE CONNECTIONS

RANDOM WALKS AND RANKING WEB PAGES

When Google launched as a web search engine, one of its major innovations over its competition was in how it ranked the pages returned in response to a user's query. Here are two key ideas in Google's ranking system, called *PageRank* (named after Larry Page, one of Google's founders):

- view a link from page u to page v as implicit "endorsement" of v by u .
- not all endorsements are equal: if a page u is endorsed by many other pages, then being endorsed by u is a bigger deal.

These point can be restated more glibly as: *a page is important if it is pointed to by many important pages*. The idea of PageRank is to break this apparent circularity using the *Random Surfer Model*. Imagine a hypothetical web user who starts at a random web page, and, at every time step, clicks on a randomly chosen link from the page she's currently visiting. The more frequently that this hypothetical user visits page u , the more important we'll say u is.

The Random Surfer explores the web using what's called a *random walk* on the web graph. In its simplest form, a random walk on a directed graph $G = (V, E)$ visits a sequence u_0, u_1, u_2, \dots of nodes in G as follows:

1. choose a node $u_0 \in V$, uniformly at random.
2. in step $t = 1, 2, \dots$, the next node u_t is chosen by picking a node uniformly at random from the out-neighborhood of the previous node u_{t-1} .

(See Figure 11.68(a) for an example.)

As you'll explore in Exercises 11.204–11.208, under mild assumptions about G , there's a special probability distribution p over the nodes of the graph, called the *stationary distribution* of the graph, that has the following property: if we choose an initial node u with probability $p(u)$, and we then take one step of the random walk from u , the resulting probability distribution over the nodes is still p . And, it turns out, we can approximate p by the probability distribution observed simply by running the random walk for many steps, as in Figure 11.68(b). We'll use p as our measure of importance.

We've already made a lot of progress toward the stated goals: $p(u)$ is higher the more in-neighbors u has, but $p(u)$ will be increased even more when the in-neighbors of u have a high probability themselves. In Figure 11.68(c), for example, we see that $p(D) > p(B)$ and $p(D) > p(C)$, despite B and C having higher in-degree than D.

But there are a few complications that we still have to address to get to the full PageRank model.¹³ One is that the Random Surfer has nowhere to go if she ends up at a page u that has no out-neighbors. (The random walk's next step isn't even defined.) In this case, we'll have the Random Surfer jump to a completely random page (each of the $|V|$ nodes is chosen with probability $\frac{1}{|V|}$). Second, this model allows the Random Surfer to get stuck in a "dead end" if there's a group of nodes that has no edges leaving it. Thus—and this change probably makes the Random Surfer more realistic anyway—we'll add a *restart probability* of 15% to every stage of the random walk: with probability 85%, we behave as previously described; with probability 15%, we jump to a randomly chosen node. (See Figure 11.68(d) for the updated probabilities.)

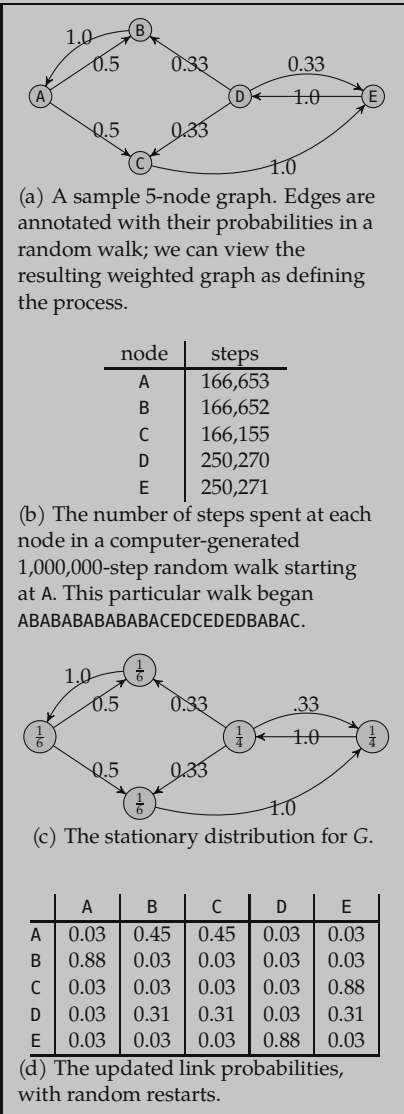


Figure 11.68: A random walk.

You can find more about the Random Surfer model and PageRank (including interesting questions about how to calculate it on a graph with nodes numbering in the billions) in a good textbook on data mining, like

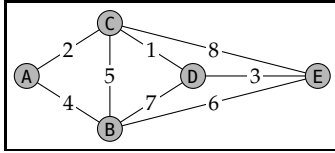
¹³ Jure Leskovec, Anand Rajaraman, and Jeff Ullman. *Mining of Massive Datasets*. Cambridge University Press, 2nd edition, 2014.

There are also many other ingredients in Google's ranking recipe beyond PageRank, though PageRank was an early and important one.

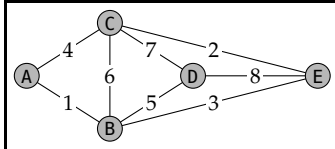
11.5.3 Exercises

For the following graphs, find all shortest paths between the given nodes. Give both the path length and the path itself.

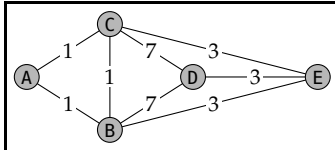
11.183 From A to E:



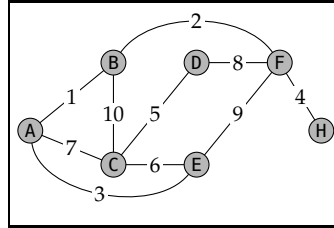
11.184 From A to E:



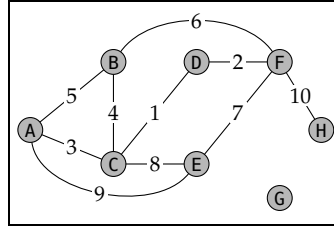
11.185 From A to E:



11.186 From A to H:



11.187 From A to H:



11.188 Let n be arbitrary. Give an example of an n -node weighted graph $G = \langle V, E \rangle$ with designated nodes $s \in V$ and $t \in V$ in which both of the following conditions hold:

- (i) all edge weights are distinct (for any $e \in E$ and $e' \in E$, we have $w(e) \neq w(e')$ if $e \neq e'$), and
- (ii) for some $\alpha > 1$ and $c > 0$, there are at least $c \cdot \alpha^n$ different shortest paths between s and t .

Suppose that we are running Dijkstra's Algorithm on the graph shown in Figure 11.69 to compute distances from the node A. So far Dijkstra's Algorithm has computed four distances:

$$d(A, A) = 0 \quad d(A, B) = 1 \quad d(A, C) = 3 \quad d(A, F) = 7$$

If we continue Dijkstra's algorithm for further iterations, it records the distance for a new node in each iteration.

11.189 What is the next node recorded, and what is its distance?

11.190 What is the next node (after the one from Exercise 11.189) for which Dijkstra's algorithm records a distance, and what is its distance? List all subsequently discovered nodes, and their distances.

11.191 Trace Dijkstra's algorithm on the graph shown in Figure 11.69 to compute distances from the node H. List all discovered nodes and their distances, in the order in which they're discovered.

11.192 Identify *exactly* where the proof of correctness for Dijkstra's algorithm (specifically, in the proof of Lemma 11.9) the argument fails if edge weights can be negative. Then give an example of a graph with negative edge weights in which Dijkstra's algorithm fails.

Suppose that $G = \langle V, E \rangle$ is a weighted, directed graph in which nodes represent physical states of a system, and an edge $\langle u, v \rangle$ indicates that one can move from state u to state v . The weight $w_{\langle u, v \rangle}$ of edge $\langle u, v \rangle$ denotes the multiplicative cost of the exchange: one can trade $w_{u,v}$ units of u for 1 unit of v . For example, if there's an edge $\langle A, B \rangle$ with weight 1.04, then I can trade 2.08 units of energy in state A for 2 units of energy in state B.

Suppose that we wish to find a shortest multiplicative path (SMP) from a given node s to a given node t in G , where the cost of the path is the product of the edge weights along it. For example, in Figure 11.70, the SMP from A to D is $A \rightarrow B \rightarrow C \rightarrow D$ at cost $1.1 \cdot 1.5 \cdot 1.4 = 2.31$, which is better than $A \rightarrow B \rightarrow D$ at cost $1.1 \cdot 2.25 = 2.75$.

11.193 Describe how to modify Dijkstra's algorithm to find the shortest SMP in a given weighted graph G . Alternatively, describe how to modify a given weighted graph G into a graph G' so that Dijkstra's algorithm run on G' finds an SMP in G .

11.194 As you argued in Exercise 11.192, Dijkstra's algorithm may fail if edge weights are negative. State the condition that guarantees that your algorithm from Exercise 11.193 properly computes SMPs.

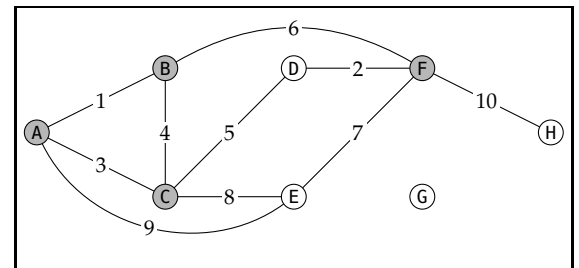


Figure 11.69: A weighted graph.

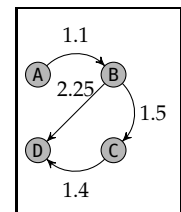
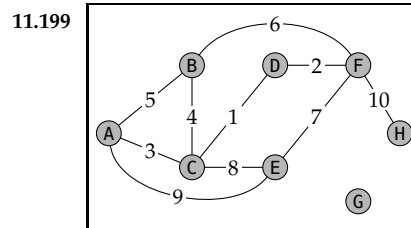
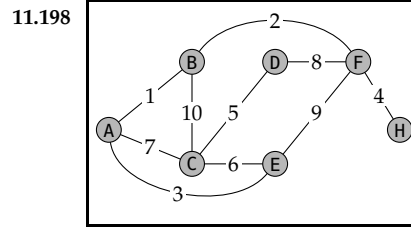
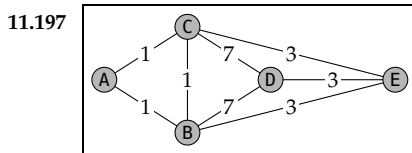
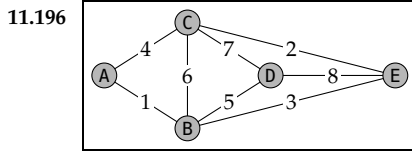
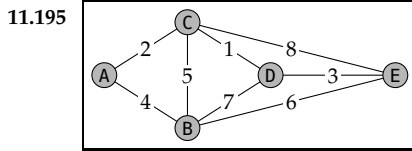


Figure 11.70: A weighted graph.

List all minimum spanning trees of the following graphs. (Note that some have edges with nondistinct weights.)



Consider the undirected 9-node complete graph K_9 . There are $\binom{9}{2} = \frac{9 \cdot 8}{2} = 36$ unordered pairs of nodes in this graph, so there are 36 different edges in the graph. Suppose that you're asked to assign each of these 36 edges a distinct weight from the set $\{1, 2, \dots, 36\}$. (You get to choose which edges have which weights.)

11.200 What's the cheapest possible minimum spanning tree of K_9 ?

11.201 What's the most expensive edge that can appear in a minimum spanning tree of K_9 ?

11.202 What's the costliest possible minimum spanning tree of K_9 ?

11.203 Generalize Exercise 11.200 and 11.202: what are the cheapest and most expensive possible MSTs for the graph K_n if all edges have distinct weights chosen from $\{1, 2, \dots, \binom{n}{2}\}$? (Hint: see Exercise 9.173.)

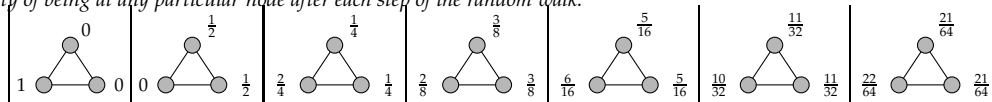
Recall from p. 1174 that a random walk in a graph $G = \langle V, E \rangle$ proceeds as follows: we start at a node $u_0 \in V$, and, at every time step, we select as the next node u_{i+1} a uniformly chosen (out-)neighbor of u_i .

Suppose we choose an initial node u_0 according to a probability distribution p , and we then take one step of the random walk from u_0 to get a new node u_1 . The probability distribution p is a stationary distribution if it satisfies the following condition: for every node $s \in V$, we have that $\Pr[u_0 = s] = \Pr[u_1 = s] = p(s)$. Such a distribution is called "stationary" because, if p is the probability distribution before a step of the walk, then p is still the probability distribution after a step of the walk (and thus the distribution "hasn't moved"—that is, is stationary).

11.204 Argue that $p(A) = p(B) = p(C) = \frac{1}{3}$ is a stationary distribution for the graph in Figure 11.71(a).

11.205 Argue that the graph in Figure 11.71(b) has at least two distinct stationary distributions.

Suppose that we start a random walk at node A in the graph in Figure 11.71(a). The following chart shows the probability of being at any particular node after each step of the random walk:



Let $p_k(u)$ denote the probability of the k th step of this random walk being at node u . Although we'll skip the proof, the following theorem turns out to be true of random walks on undirected graphs G :

If G is connected and nonbipartite, then a unique stationary distribution p exists for this random walk on G (regardless of which node we choose as the initial node for the walk). Furthermore, the stationary distribution is the limit of the probability distributions p_k of where the random walk is in the k th step.

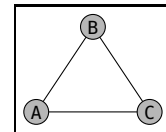
11.206 (programming required) Write a random-walk simulator: take an undirected graph G as input, and simulate 2000 steps of a random walk starting at an arbitrary node. Repeat 2000 times, and report the fraction of walks that are at each node. What are your results on the graph from Figure 11.71(a)?

11.207 Argue that the above process doesn't converge to a unique stationary distribution in a bipartite graph. (For example, what's p_{1000} if a random walk starts at node J in the graph in Figure 11.71(c)? Node K?)

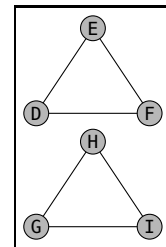
11.208 Let $G = \langle V, E \rangle$ be an arbitrary connected undirected graph. For any $u \in V$, define

$$p(u) := \frac{\text{degree}(u)}{2 \cdot |E|}.$$

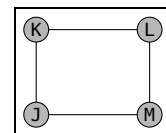
Prove that the probability distribution p is a stationary distribution for the random walk on G .



(a)



(b)



(c)

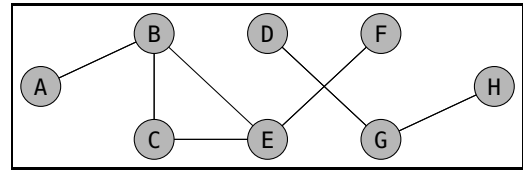
Figure 11.71: Some undirected graphs upon which a random walk can be performed.

11.6 Chapter at a Glance

Formal Introduction

A *graph* is a pair $G = \langle V, E \rangle$ where V is a set of *vertices* or *nodes*, and E is a set of *edges*. In a *directed graph*, the edges $E \subseteq V \times V$ are ordered pairs of vertices; in an *undirected graph*, the edges $E \subseteq \{\{u, v\} : u, v \in V\}$ are unordered pairs. A directed edge $\langle u, v \rangle$ goes from u to v ; an undirected edge $\langle u, v \rangle$ goes between u and v . We sometimes write $\langle u, v \rangle$ even for an undirected graphs. A *simple graph* has no *parallel edges* joining the same two nodes and also has no *self loops* joining a node to itself.

For an edge $e = \langle u, v \rangle$, we say that u and v are *adjacent*; v is a *neighbor* of u ; u and v are the *endpoints* of e ; and u and v are both *incident* to e . The *neighborhood* of a node u is $\{v : \langle u, v \rangle \in E\}$, its set of neighbors. The *degree* of u is the cardinality of u 's neighborhood. In a directed graph, the *in-neighbors* of u are the nodes that have an edge pointing to u ; the *out-neighbors* are the nodes to which u has an edge pointing; and the *in-degree* and *out-degree* of u are the number of in- and out-neighbors, respectively.



An *adjacency list* stores a graph using an array with $|V|$ entries; the slot for node u is a linked list of u 's neighbors. An *adjacency matrix* stores the graph using a two-dimensional Boolean array of size $|V| \times |V|$; the value in $\langle \text{row } u, \text{column } v \rangle$ indicates whether the edge $\langle u, v \rangle$ exists.

Two graphs are *isomorphic* if they are identical except for the naming of the nodes. A *subgraph* of G contains a subset V' of G 's nodes and a subset E' of G 's edges joining elements of V' . An *induced subgraph* is a subgraph in which every edge that joins elements of V' is included in E' . A *complete graph* or *clique* is a graph K_n in which every possible edge exists. A *bipartite graph* is one in which nodes can be partitioned into sets L and R such that every edge joins a node in L to a node in R . A *regular graph* is one in which every node has identical degree. A *planar graph* is one that can be drawn on paper without any edges crossing.

Paths, Connectivity, and Distances

A *path* is a sequence of $k \geq 1$ nodes $\langle v_1, v_2, \dots, v_k \rangle$, where $\langle v_{i-1}, v_i \rangle \in E$ for every index $i \in \{1, 2, \dots, k-1\}$. The path is *simple* if all the v_i s are distinct. This path has *length* $k-1$ —the number of edges that it traverses—and is a *path from* v_1 *to* v_k .

In an undirected graph, nodes u and v are *connected* if there exists a path from u to v . A *connected component* of $G = \langle V, E \rangle$ is a set $S \subseteq V$ such that (i) every $u \in S$ and $v \in S$ are connected; and (ii) for every $w \notin S$, the set $S \cup \{w\}$ does not satisfy condition (i). The entire graph is *connected* if it has only one connected component, namely V .

In a directed graph, node u is *reachable from* node v if there exists a path from v to u ; u and v are *strongly connected* if each is reachable from the other. A *strongly connected component* is a set S of nodes such that any two nodes in S are strongly connected and no node $x \notin S$ is strongly connected to any node $s \in S$.

Connectivity can be tested in time $\Theta(|V| + |E|)$ time using *breadth-first search* (BFS; see Figure 11.72) or *depth-first search* (DFS). The *distance* from node s to node t is the length of a shortest path from s to t . BFS can also be used to compute distances.

Trees

A *cycle* $\langle v_1, v_2, \dots, v_k, v_1 \rangle$ is a path of length ≥ 2 from a node v_1 back to itself that does not traverse the same edge twice. The *length* of the cycle is k . The cycle is *simple* if each v_i is distinct. Cycles can be identified using BFS.

A graph is *acyclic* if it contains no cycles. Every acyclic graph has a node of degree 0 or 1. A *tree* is a connected, acyclic graph. (A *forest* is any acyclic graph.) A tree has one more node than it has vertex. A tree becomes disconnected if any edge is deleted; it becomes cyclic if any edge is added.

One node in a tree can be designated as the *root*. Every node other than the root has a *parent* (its neighbor that's closer to the root). If p is v 's parent, then v is one of p 's *children*. Two nodes with the same parent are *siblings*. A *leaf* is a node with no children; an *internal node* is a node with children. The *depth* of a node is its distance from the root; the *height* of the entire tree is the depth of deepest node. The *descendants* of u are those nodes that go through u to get the root; the *ancestors* are those nodes through which u 's path to the root goes. The *subtree* rooted at u is the induced subgraph consisting of u and all descendants of u .

All nodes in *binary trees* have at most two children, called *left* and *right*. A *traversal* of a binary tree visits every node of the tree. An *in-order* traversal recursively traverses the root's left subtree, visits the root, and recursively traverses the root's right subtree. A *pre-order* traversal visits the root and recursively traverses the root's left and right subtrees; a *post-order* traversal recursively traverses the root's left and right subtrees and then visits the root.

A *spanning tree* of a connected graph $G = \langle V, E \rangle$ is a graph $T = \langle V, E' \subseteq E \rangle$ that's a tree. A spanning tree can be found by repeatedly identifying a cycle in G and deleting any edge in that cycle.

Weighted Graphs

In a *weighted graph*, each edge e has a weight $w_e \in \mathbb{R}^{\geq 0}$. (Although graphs with negative edge weights are possible, we haven't addressed them in any detail.) The length of a path in a weighted graph is the sum of the weights of the edges that it traverses. Shortest paths in weighted graphs can be found with Dijkstra's Algorithm (Figure 11.65), which expands a set of nodes of known distance one by one. Minimum spanning trees—spanning trees of the smallest possible total weight—in weighted graphs can be found with Kruskal's Algorithm (Figure 11.67) or by repeatedly identifying a cycle in G and deleting the heaviest edge in that cycle.

Breadth-First Search (BFS):

Input: a graph $G = \langle V, E \rangle$ and a source node $s \in V$

Output: the set of nodes reachable from s in G

```

1: Frontier :=  $\langle s \rangle$ 
   // Frontier will be a list of nodes to process, in order.
2: Known :=  $\emptyset$ 
   // Known will be the set of already-processed nodes.
3: while Frontier is nonempty:
4:    $u$  := the first node in Frontier
5:   remove  $u$  from Frontier
6:   for every neighbor  $v$  of  $u$ :
7:     if  $v$  is in neither Frontier nor Known then
8:       add  $v$  to the end of Frontier
9:   add  $u$  to Known
10: return Known
```

Figure 11.72:
Breadth-first search.

Key Terms and Results

Key Terms

FORMAL INTRODUCTION

- undirected and directed graphs
- nodes/ vertices, edges
- parallel edges, self loops
- simple graphs
- adjacent node, incident edge
- (in/ out-)neighbors, neighborhood
- (in/ out-)degree
- adjacency list, adjacency matrix
- isomorphic graphs
- subgraphs
- complete, bipartite, regular, planar graphs

PATHS AND CONNECTIVITY

- path
- connected (nodes), connected (graph)
- connected component
- reachability
- strongly connected component
- shortest path/ distance
- breadth-first search (BFS)
- depth-first search (DFS)

TREES

- cycle
- tree, forest
- root, leaf, internal node, child, parent, sibling, ancestor, descendant, depth, height, subtree
- spanning tree

WEIGHTED GRAPHS

- Dijkstra's algorithm
- minimum spanning trees
- Kruskal's algorithm

Key Results

FORMAL INTRODUCTION

1. The “handshaking lemma”: for any undirected graph $G = \langle V, E \rangle$, we have $\sum_{u \in V} \text{degree}(u) = 2|E|$.
2. Representing G with an adjacency matrix requires $\Theta(|V|^2)$ space; we can answer “what are all of u 's neighbors?” in $\Theta(|V|)$ time and “is there an edge between u and v ?” in $\Theta(1)$ time. Representing $G = \langle V, E \rangle$ with an adjacency list requires $\Theta(|V| + |E|)$ space; both questions take $1 + \Theta(\text{degree}(u))$ time.

PATHS, CONNECTIVITY, AND DISTANCES

1. Connectivity can be tested using *breadth-first search (BFS)* (Figure 11.29) or *depth-first search (DFS)* (Figure 11.31). BFS can also be used to compute the distance between nodes in a graph, and it runs in $\Theta(|V| + |E|)$ time.

TREES

1. Any tree with n nodes has exactly $n - 1$ edges. Adding any edge to a tree creates a cycle; deleting any edge disconnects the graph.
2. A spanning tree of a graph G can be found by repeatedly identifying a cycle in G and deleting an arbitrary edge in that cycle.

WEIGHTED GRAPHS

1. Shortest paths in weighted graphs can be found with Dijkstra's Algorithm (Figure 11.65) if all edges have nonnegative weights.
2. Minimum spanning trees in weighted graphs can be found with Kruskal's Algorithm (Figure 11.67) or by repeatedly identifying a cycle in G and deleting the heaviest edge in that cycle.

