

Chapter 5

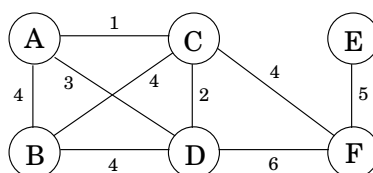
Greedy algorithms

A game like chess can be won only by *thinking ahead*: a player who is focused entirely on immediate advantage is easy to defeat. But in many other games, such as Scrabble, it is possible to do quite well by simply making whichever move seems best at the moment and not worrying too much about future consequences.

This sort of myopic behavior is easy and convenient, making it an attractive algorithmic strategy. *Greedy* algorithms build up a solution piece by piece, always choosing the next piece that offers the most obvious and immediate benefit. Although such an approach can be disastrous for some computational tasks, there are many for which it is optimal. Our first example is that of minimum spanning trees.

5.1 Minimum spanning trees

Suppose you are asked to network a collection of computers by linking selected pairs of them. This translates into a graph problem in which nodes are computers, undirected edges are potential links, and the goal is to pick enough of these edges that the nodes are connected. But this is not all; each link also has a maintenance cost, reflected in that edge's weight. What is the cheapest possible network?



One immediate observation is that the optimal set of edges cannot contain a cycle, because removing an edge from this cycle would reduce the cost without compromising connectivity:

Property 1 *Removing a cycle edge cannot disconnect a graph.*

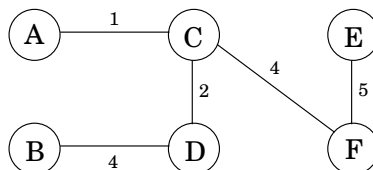
So the solution must be connected and acyclic: undirected graphs of this kind are called *trees*. The particular tree we want is the one with minimum total weight, known as the *minimum spanning tree*. Here is its formal definition.

Input: An undirected graph $G = (V, E)$; edge weights w_e .

Output: A tree $T = (V, E')$, with $E' \subseteq E$, that minimizes

$$\text{weight}(T) = \sum_{e \in E'} w_e.$$

In the preceding example, the minimum spanning tree has a cost of 16:



However, this is not the only optimal solution. Can you spot another?

5.1.1 A greedy approach

Kruskal's minimum spanning tree algorithm starts with the empty graph and then selects edges from E according to the following rule.

Repeatedly add the next lightest edge that doesn't produce a cycle.

In other words, it constructs the tree edge by edge and, apart from taking care to avoid cycles, simply picks whichever edge is cheapest at the moment. This is a greedy algorithm: every decision it makes is the one with the most obvious immediate advantage.

Figure 5.1 shows an example. We start with an empty graph and then attempt to add edges in increasing order of weight (ties are broken arbitrarily):

$B - C$, $C - D$, $B - D$, $C - F$, $D - F$, $E - F$, $A - D$, $A - B$, $C - E$, $A - C$.

The first two succeed, but the third, $B - D$, would produce a cycle if added. So we ignore it and move along. The final result is a tree with cost 14, the minimum possible.

The correctness of Kruskal's method follows from a certain *cut property*, which is general enough to also justify a whole slew of other minimum spanning tree algorithms.

Figure 5.1 The minimum spanning tree found by Kruskal's algorithm.



Trees

A tree is an undirected graph that is connected and acyclic. Much of what makes trees so useful is the simplicity of their structure. For instance,

Property 2 A tree on n nodes has $n - 1$ edges.

This can be seen by building the tree one edge at a time, starting from an empty graph. Initially each of the n nodes is disconnected from the others, in a connected component by itself. As edges are added, these components merge. Since each edge unites two different components, exactly $n - 1$ edges are added by the time the tree is fully formed.

In a little more detail: When a particular edge $\{u, v\}$ comes up, we can be sure that u and v lie in separate connected components, for otherwise there would already be a path between them and this edge would create a cycle. Adding the edge then merges these two components, thereby reducing the total number of connected components by one. Over the course of this incremental process, the number of components decreases from n to one, meaning that $n - 1$ edges must have been added along the way.

The converse is also true.

Property 3 Any connected, undirected graph $G = (V, E)$ with $|E| = |V| - 1$ is a tree.

We just need to show that G is acyclic. One way to do this is to run the following iterative procedure on it: while the graph contains a cycle, remove one edge from this cycle. The process terminates with some graph $G' = (V, E')$, $E' \subseteq E$, which is acyclic and, by Property 1 (from page 139), is also connected. Therefore G' is a tree, whereupon $|E'| = |V| - 1$ by Property 2. So $E' = E$, no edges were removed, and G was acyclic to start with.

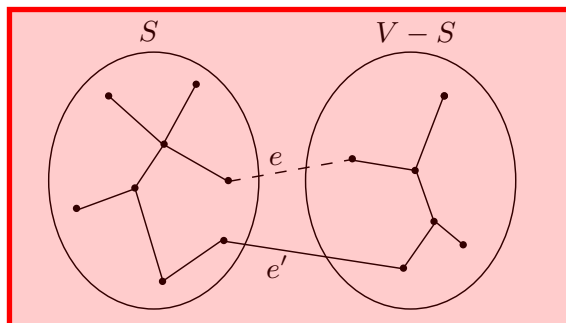
In other words, we can tell whether a connected graph is a tree just by counting how many edges it has. Here's another characterization.

Property 4 An undirected graph is a tree if and only if there is a unique path between any pair of nodes.

In a tree, any two nodes can only have one path between them; for if there were two paths, the union of these paths would contain a cycle.

On the other hand, if a graph has a path between any two nodes, then it is connected. If these paths are unique, then the graph is also acyclic (since a cycle has two paths between any pair of nodes).

Figure 5.2 $T \cup \{e\}$. The addition of e (dotted) to T (solid lines) produces a cycle. This cycle must contain at least one other edge, shown here as e' , across the cut $(S, V - S)$.



5.1.2 The cut property

Say that in the process of building a minimum spanning tree (MST), we have already chosen some edges and are so far on the right track. Which edge should we add next? The following lemma gives us a lot of flexibility in our choice.

Cut property Suppose edges X are part of a minimum spanning tree of $G = (V, E)$. Pick any subset of nodes S for which X does not cross between S and $V - S$, and let e be the lightest edge across this partition. Then $X \cup \{e\}$ is part of some MST.

A *cut* is any partition of the vertices into two groups, S and $V - S$. What this property says is that it is always safe to add the lightest edge across any cut (that is, between a vertex in S and one in $V - S$), provided X has no edges across the cut.

Let's see why this holds. Edges X are part of some MST T ; if the new edge e also happens to be part of T , then there is nothing to prove. So assume e is not in T . We will construct a different MST T' containing $X \cup \{e\}$ by altering T slightly, changing just one of its edges.

Add edge e to T . Since T is connected, it already has a path between the endpoints of e , so adding e creates a cycle. This cycle must also have some other edge e' across the cut $(S, V - S)$ (Figure 8.3). If we now remove this edge, we are left with $T' = T \cup \{e\} - \{e'\}$, which we will show to be a tree. T' is connected by Property 1, since e' is a cycle edge. And it has the same number of edges as T ; so by Properties 2 and 3, it is also a tree.

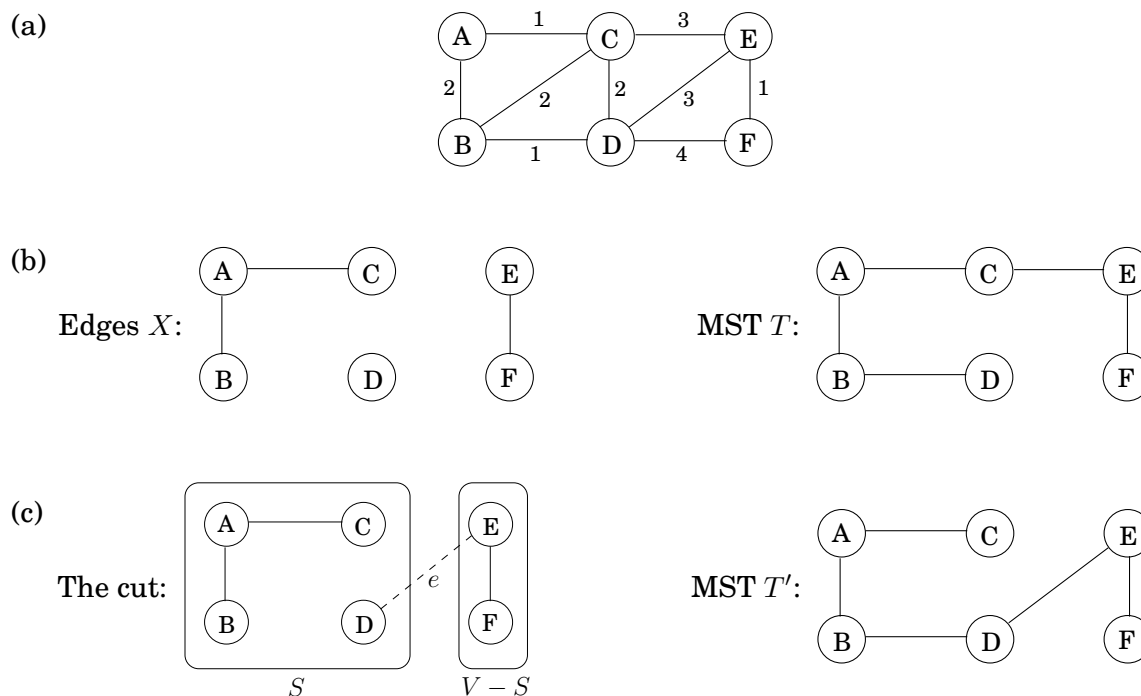
Moreover, T' is a minimum spanning tree. Compare its weight to that of T :

$$\text{weight}(T') = \text{weight}(T) + w(e) - w(e').$$

Both e and e' cross between S and $V - S$, and e is specifically the lightest edge of this type. Therefore $w(e) \leq w(e')$, and $\text{weight}(T') \leq \text{weight}(T)$. Since T is an MST, it must be the case that $\text{weight}(T') = \text{weight}(T)$ and that T' is also an MST.

Figure 5.3 shows an example of the cut property. Which edge is e' ?

Figure 5.3 The cut property at work. (a) An undirected graph. (b) Set X has three edges, and is part of the MST T on the right. (c) If $S = \{A, B, C, D\}$, then one of the minimum-weight edges across the cut $(S, V - S)$ is $e = \{D, E\}$. $X \cup \{e\}$ is part of MST T' , shown on the right.



5.1.3 Kruskal's algorithm

We are ready to justify Kruskal's algorithm. At any given moment, the edges it has already chosen form a partial solution, a collection of connected components each of which has a tree structure. The next edge e to be added connects two of these components; call them T_1 and T_2 . Since e is the lightest edge that doesn't produce a cycle, it is certain to be the lightest edge between T_1 and $V - T_1$ and therefore satisfies the cut property.

Now we fill in some implementation details. At each stage, the algorithm chooses an edge to add to its current partial solution. To do so, it needs to test each candidate edge $u - v$ to see whether the endpoints u and v lie in different components; otherwise the edge produces a cycle. And once an edge is chosen, the corresponding components need to be merged. What kind of data structure supports such operations?

We will model the algorithm's state as a collection of *disjoint sets*, each of which contains the nodes of a particular component. Initially each node is in a component by itself:

`makeset(x):` create a singleton set containing just x .

We repeatedly test pairs of nodes to see if they belong to the same set.

`find(x):` to which set does x belong?

Figure 5.4 Kruskal's minimum spanning tree algorithm.

```
procedure kruskal( $G, w$ )
```

Input: A connected undirected graph $G = (V, E)$ with edge weights w_e

Output: A minimum spanning tree defined by the edges X

```
for all  $u \in V$ :
```

```
    makeset( $u$ )
```

```
 $X = \{\}$ 
```

```
Sort the edges  $E$  by weight
```

```
for all edges  $\{u, v\} \in E$ , in increasing order of weight:
```

```
    if find( $u$ )  $\neq$  find( $v$ ):
```

```
        add edge  $\{u, v\}$  to  $X$ 
```

```
        union( $u, v$ )
```

And whenever we add an edge, we are merging two components.

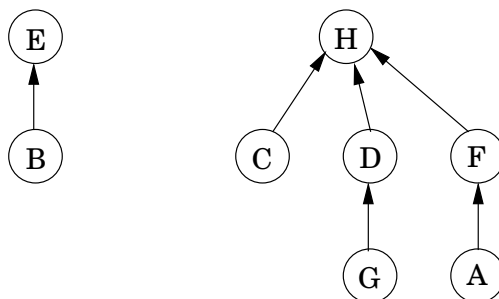
union(x, y): merge the sets containing x and y .

The final algorithm is shown in Figure 5.4. It uses $|V|$ makeset, $2|E|$ find, and $|V| - 1$ union operations.

5.1.4 A data structure for disjoint sets

Union by rank

One way to store a set is as a directed tree (Figure 5.5). Nodes of the tree are elements of the set, arranged in no particular order, and each has parent pointers that eventually lead up to the root of the tree. This root element is a convenient representative, or name, for the set. It is distinguished from the other elements by the fact that its parent pointer is a self-loop.

Figure 5.5 A directed-tree representation of two sets $\{B, E\}$ and $\{A, C, D, F, G, H\}$.

In addition to a parent pointer π , each node also has a *rank* that, for the time being, should be interpreted as the height of the subtree hanging from that node.

```

procedure makeset( $x$ )
 $\pi(x) = x$ 
 $\text{rank}(x) = 0$ 

function find( $x$ )
while  $x \neq \pi(x)$ :  $x = \pi(x)$ 
return  $x$ 

```

As can be expected, **makeset** is a constant-time operation. On the other hand, **find** follows parent pointers to the root of the tree and therefore takes time proportional to the height of the tree. The tree actually gets built via the third operation, **union**, and so we must make sure that this procedure keeps trees shallow.

Merging two sets is easy: make the root of one point to the root of the other. But we have a choice here. If the representatives (roots) of the sets are r_x and r_y , do we make r_x point to r_y or the other way around? Since tree height is the main impediment to computational efficiency, a good strategy is to make the root of the shorter tree point to the root of the taller tree. This way, the overall height increases only if the two trees being merged are equally tall. Instead of explicitly computing heights of trees, we will use the *rank* numbers of their root nodes—which is why this scheme is called *union by rank*.

```

procedure union( $x, y$ )
 $r_x = \text{find}(x)$ 
 $r_y = \text{find}(y)$ 
if  $r_x = r_y$ : return
if  $\text{rank}(r_x) > \text{rank}(r_y)$ :
     $\pi(r_y) = r_x$ 
else:
     $\pi(r_x) = r_y$ 
    if  $\text{rank}(r_x) = \text{rank}(r_y)$ :  $\text{rank}(r_y) = \text{rank}(r_y) + 1$ 

```

See Figure 5.6 for an example.

By design, the *rank* of a node is exactly the height of the subtree rooted at that node. This means, for instance, that as you move up a path toward a root node, the *rank* values along the way are strictly increasing.

Property 1 For any x , $\text{rank}(x) < \text{rank}(\pi(x))$.

A root node with rank k is created by the merger of two trees with roots of rank $k - 1$. It follows by induction (try it!) that

Property 2 Any root node of rank k has at least 2^k nodes in its tree.

This extends to internal (nonroot) nodes as well: a node of rank k has at least 2^k descendants. After all, any internal node was once a root, and neither its rank nor its set of descendants has changed since then. Moreover, different rank- k nodes cannot have common descendants, since by Property 1 any element has at most one ancestor of rank k . Which means

Property 3 *If there are n elements overall, there can be at most $n/2^k$ nodes of rank k .*

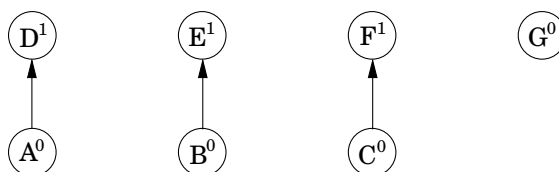
This last observation implies, crucially, that the maximum rank is $\log n$. Therefore, all the trees have height $\leq \log n$, and this is an upper bound on the running time of find and union.

Figure 5.6 A sequence of disjoint-set operations. Superscripts denote rank.

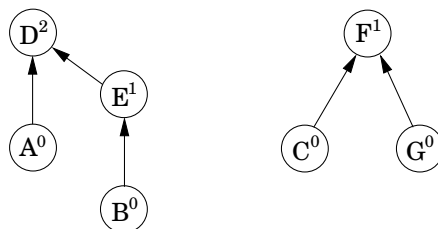
After `makeset(A)`, `makeset(B)`, ..., `makeset(G)`:



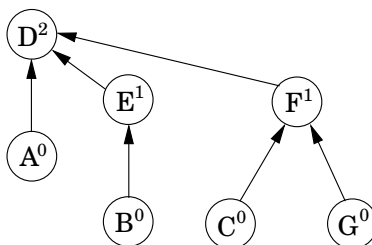
After `union(A, D)`, `union(B, E)`, `union(C, F)`:



After `union(C, G)`, `union(E, A)`:



After `union(B, G)`:



Path compression

With the data structure as presented so far, the total time for Kruskal's algorithm becomes $O(|E| \log |V|)$ for sorting the edges (remember, $\log |E| \approx \log |V|$) plus another $O(|E| \log |V|)$ for the union and find operations that dominate the rest of the algorithm. So there seems to be little incentive to make our data structure any more efficient.

But what if the edges are given to us sorted? Or if the weights are small (say, $O(|E|)$) so that sorting can be done in linear time? Then the data structure part becomes the bottleneck, and it is useful to think about improving its performance beyond $\log n$ per operation. As it turns out, the improved data structure is useful in many other applications.

But how can we perform union's and find's faster than $\log n$? The answer is, by being a little more careful to maintain our data structure in good shape. As any housekeeper knows, a little extra effort put into routine maintenance can pay off handsomely in the long run, by forestalling major calamities. We have in mind a particular maintenance operation for our union-find data structure, intended to keep the trees short—**during each find, when a series of parent pointers is followed up to the root of a tree, we will change all these pointers so that they point directly to the root** (Figure 5.7). This *path compression* heuristic only slightly increases the time needed for a find and is easy to code.

```
function find(x)
  if  $x \neq \pi(x)$ :  $\pi(x) = \text{find}(\pi(x))$ 
  return  $\pi(x)$ 
```

The benefit of this simple alteration is long-term rather than instantaneous and thus necessitates a particular kind of analysis: **we need to look at sequences of find and union operations, starting from an empty data structure, and determine the average time per operation. This amortized cost turns out to be just barely more than $O(1)$, down from the earlier $O(\log n)$.**

Think of the data structure as **having a “top level” consisting of the root nodes, and below it, the insides of the trees.** There is a division of labor: **find operations (with or without path compression) only touch the insides of trees, whereas union's only look at the top level.** Thus path compression has no effect on union operations and leaves the top level unchanged.

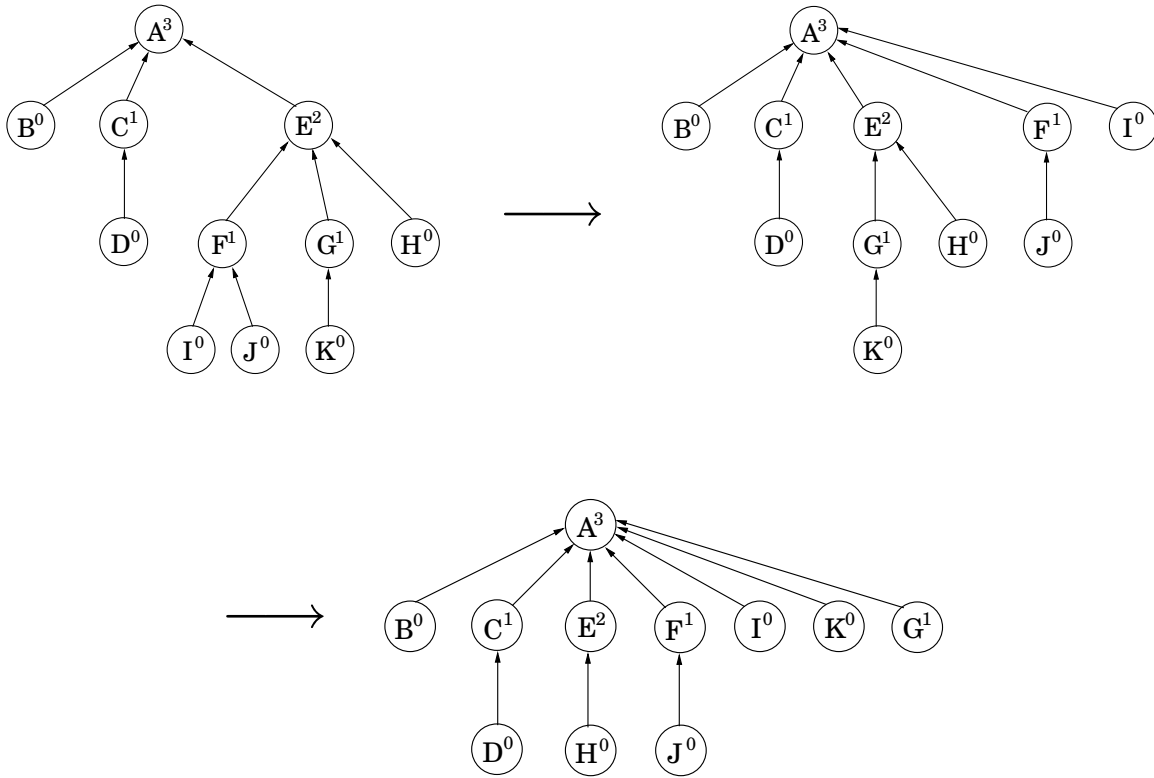
We now know that the ranks of root nodes are unaltered, but what about *nonroot* nodes? The key point here is that once a node ceases to be a root, it never resurfaces, and its rank is forever fixed. **Therefore the ranks of all nodes are unchanged by path compression, even though these numbers can no longer be interpreted as tree heights. In particular, properties 1–3 (from page 145) still hold.**

If there are n elements, their rank values can range from 0 to $\log n$ by Property 3. Let's divide the nonzero part of this range into certain carefully chosen intervals, for reasons that will soon become clear:

```
{1}, {2}, {3, 4}, {5, 6, ..., 16}, {17, 18, ...,  $2^{16} = 65536$ }, {65537, 65538, ...,  $2^{65536}$ }, ...
```

Each group is of the form $\{k + 1, k + 2, \dots, 2^k\}$, where k is a power of 2. The number of groups is $\log^* n$, which is defined to be the number of successive log operations that need to be applied

Figure 5.7 The effect of path compression: $\text{find}(I)$ followed by $\text{find}(K)$.



to n to bring it down to 1 (or below 1). For instance, $\log^* 1000 = 4$ since $\log \log \log \log 1000 \leq 1$. In practice there will just be the first five of the intervals shown; more are needed only if $n \geq 2^{65536}$, in other words never.

In a sequence of find operations, some may take longer than others. We'll bound the overall running time using some creative accounting. Specifically, we will give each node a certain amount of pocket money, such that the total money doled out is at most $n \log^* n$ dollars. We will then show that each find takes $O(\log^* n)$ steps, plus some additional amount of time that can be "paid for" using the pocket money of the nodes involved—one dollar per unit of time. Thus the overall time for m find 's is $O(m \log^* n)$ plus at most $O(n \log^* n)$.

In more detail, a node receives its allowance as soon as it ceases to be a root, at which point its rank is fixed. If this rank lies in the interval $\{k+1, \dots, 2^k\}$, the node receives 2^k dollars. By Property 3, the number of nodes with rank $> k$ is bounded by

$$\frac{n}{2^{k+1}} + \frac{n}{2^{k+2}} + \dots \leq \frac{n}{2^k}.$$

Therefore the total money given to nodes in this particular interval is at most n dollars, and since there are $\log^* n$ intervals, the total money disbursed to all nodes is $\leq n \log^* n$.

Now, the time taken by a specific `find` is simply the number of pointers followed. Consider the ascending rank values along this chain of nodes up to the root. Nodes x on the chain fall into two categories: either the rank of $\pi(x)$ is in a higher interval than the rank of x , or else it lies in the same interval. There are at most $\log^* n$ nodes of the first type (do you see why?), so the work done on them takes $O(\log^* n)$ time. The remaining nodes—whose parents' ranks are in the same interval as theirs—have to pay a dollar out of their pocket money for their processing time.

This only works if the initial allowance of each node x is enough to cover all of its payments in the sequence of `find` operations. Here's the crucial observation: each time x pays a dollar, its parent changes to one of higher rank. Therefore, if x 's rank lies in the interval $\{k + 1, \dots, 2^k\}$, it has to pay at most 2^k dollars before its parent's rank is in a higher interval; whereupon it never has to pay again.

A randomized algorithm for minimum cut

We have already seen that spanning trees and cuts are intimately related. Here is another connection. Let's remove the last edge that Kruskal's algorithm adds to the spanning tree; this breaks the tree into two components, thus defining a cut (S, \bar{S}) in the graph. What can we say about this cut? Suppose the graph we were working with was unweighted, and that its edges were ordered uniformly at random for Kruskal's algorithm to process them. Here is a remarkable fact: with probability at least $1/n^2$, (S, \bar{S}) is the minimum cut in the graph, where the size of a cut (S, \bar{S}) is the number of edges crossing between S and \bar{S} . This means that repeating the process $O(n^2)$ times and outputting the smallest cut found yields the minimum cut in G with high probability: an $O(mn^2 \log n)$ algorithm for unweighted minimum cuts. Some further tuning gives the $O(n^2 \log n)$ minimum cut algorithm, invented by David Karger, which is the fastest known algorithm for this important problem.

So let us see why the cut found in each iteration is the minimum cut with probability at least $1/n^2$. At any stage of Kruskal's algorithm, the vertex set V is partitioned into connected components. The only edges eligible to be added to the tree have their two endpoints in distinct components. The number of edges incident to each component must be at least C , the size of the minimum cut in G (since we could consider a cut that separated this component from the rest of the graph). So if there are k components in the graph, the number of eligible edges is at least $kC/2$ (each of the k components has at least C edges leading out of it, and we need to compensate for the double-counting of each edge). Since the edges were randomly ordered, the chance that the next eligible edge in the list is from the minimum cut is at most $C/(kC/2) = 2/k$. Thus, with probability at least $1 - 2/k = (k-2)/k$, the choice leaves the minimum cut intact. But now the chance that Kruskal's algorithm leaves the minimum cut intact all the way up to the choice of the last spanning tree edge is at least

$$\frac{n-2}{n} \cdot \frac{n-3}{n-1} \cdot \frac{n-4}{n-2} \cdots \frac{2}{4} \cdot \frac{1}{3} = \frac{1}{n(n-1)}.$$

5.1.5 Prim's algorithm

Let's return to our discussion of minimum spanning tree algorithms. What the cut property tells us in most general terms is that any algorithm conforming to the following greedy schema is guaranteed to work.

```

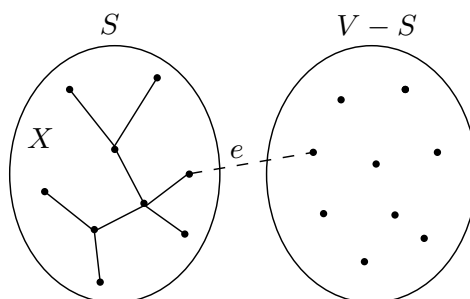
X = { } (edges picked so far)
repeat until |X| = |V| - 1:
    pick a set S ⊂ V for which X has no edges between S and V - S
    let e ∈ E be the minimum-weight edge between S and V - S
    X = X ∪ {e}

```

A popular alternative to Kruskal's algorithm is Prim's, in which the intermediate set of edges X always forms a subtree, and S is chosen to be the set of this tree's vertices.

On each iteration, the subtree defined by X grows by one edge, namely, the lightest edge between a vertex in S and a vertex outside S (Figure 5.8). We can equivalently think of S as

Figure 5.8 Prim's algorithm: the edges X form a tree, and S consists of its vertices.



growing to include the vertex $v \notin S$ of smallest cost:

$$\text{cost}(v) = \min_{u \in S} w(u, v).$$

This is strongly reminiscent of Dijkstra's algorithm, and in fact the pseudocode (Figure 5.9) is almost identical. The only difference is in the key values by which the priority queue is ordered. In Prim's algorithm, the value of a node is the weight of the lightest incoming edge from set S , whereas in Dijkstra's it is the length of an entire path to that node from the starting point. Nonetheless, the two algorithms are similar enough that they have the same running time, which depends on the particular priority queue implementation.

Figure 5.9 shows Prim's algorithm at work, on a small six-node graph. Notice how the final MST is completely specified by the `prev` array.

Figure 5.9 *Top*: Prim's minimum spanning tree algorithm. *Below*: An illustration of Prim's algorithm, starting at node A. Also shown are a table of cost/prev values, and the final MST.

procedure prim(G, w)

Input: A connected undirected graph $G = (V, E)$ with edge weights w_e

Output: A minimum spanning tree defined by the array prev

for all $u \in V$:

$\text{cost}(u) = \infty$

$\text{prev}(u) = \text{nil}$

Pick any initial node u_0

$\text{cost}(u_0) = 0$

$H = \text{makequeue}(V)$ (priority queue, using cost-values as keys)

while H is not empty:

$v = \text{deletemin}(H)$

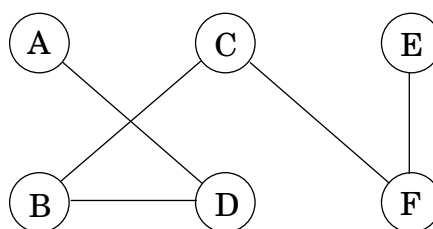
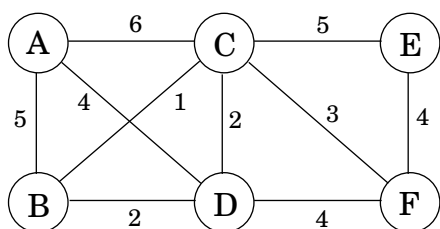
 for each $\{v, z\} \in E$:

 if $\text{cost}(z) > w(v, z)$:

$\text{cost}(z) = w(v, z)$

$\text{prev}(z) = v$

$\text{decreasekey}(H, z)$



Set S	A	B	C	D	E	F
$\{\}$	0/ nil	∞/nil	∞/nil	∞/nil	∞/nil	∞/nil
A		5/ A	6/ A	4/ A	∞/nil	∞/nil
A, D		2/ D	2/ D		∞/nil	4/ D
A, D, B			1/ B		∞/nil	4/ D
A, D, B, C					5/ C	3/ C
A, D, B, C, F					4/ F	

5.2 Huffman encoding

In the MP3 audio compression scheme, a sound signal is encoded in three steps.

1. It is digitized by sampling at regular intervals, yielding a sequence of real numbers s_1, s_2, \dots, s_T . For instance, at a rate of 44,100 samples per second, a 50-minute symphony would correspond to $T = 50 \times 60 \times 44,100 \approx 130$ million measurements.¹
2. Each real-valued sample s_t is *quantized*: approximated by a nearby number from a finite set Γ . This set is carefully chosen to exploit human perceptual limitations, with the intention that the approximating sequence is indistinguishable from s_1, s_2, \dots, s_T by the human ear.
3. The resulting string of length T over alphabet Γ is encoded in binary.

It is in the last step that Huffman encoding is used. To understand its role, let's look at a toy example in which T is 130 million and the alphabet Γ consists of just four values, denoted by the symbols A, B, C, D . What is the most economical way to write this long string in binary? The obvious choice is to use 2 bits per symbol—say codeword 00 for A , 01 for B , 10 for C , and 11 for D —in which case 260 megabits are needed in total. Can there possibly be a better encoding than this?

In search of inspiration, we take a closer look at our particular sequence and find that the four symbols are not equally abundant.

Symbol	Frequency
A	70 million
B	3 million
C	20 million
D	37 million

Is there some sort of *variable-length encoding*, in which just *one* bit is used for the frequently occurring symbol A , possibly at the expense of needing three or more bits for less common symbols?

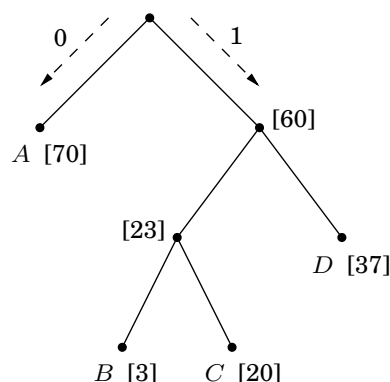
A danger with having codewords of different lengths is that the resulting encoding may not be uniquely decipherable. For instance, if the codewords are $\{0, 01, 11, 001\}$, the decoding of strings like 001 is ambiguous. We will avoid this problem by insisting on the *prefix-free* property: no codeword can be a prefix of another codeword.

Any prefix-free encoding can be represented by a *full* binary tree—that is, a binary tree in which every node has either zero or two children—where the symbols are at the leaves, and where each codeword is generated by a path from root to leaf, interpreting left as 0 and right as 1 (Exercise 5.28). Figure 5.10 shows an example of such an encoding for the four symbols A, B, C, D . Decoding is unique: a string of bits is decrypted by starting at the root, reading the string from left to right to move downward, and, whenever a leaf is reached, outputting the corresponding symbol and returning to the root. It is a simple scheme and pays off nicely

¹For stereo sound, two channels would be needed, doubling the number of samples.

Figure 5.10 A prefix-free encoding. Frequencies are shown in square brackets.

Symbol	Codeword
<i>A</i>	0
<i>B</i>	100
<i>C</i>	101
<i>D</i>	11



for our toy example, where (under the codes of Figure 5.10) the total size of the binary string drops to 213 megabits, a 17% improvement.

In general, how do we find the optimal coding tree, given the frequencies f_1, f_2, \dots, f_n of n symbols? To make the problem precise, we want a tree whose leaves each correspond to a symbol and which minimizes the overall length of the encoding,

$$\text{cost of tree} = \sum_{i=1}^n f_i \cdot (\text{depth of } i\text{th symbol in tree})$$

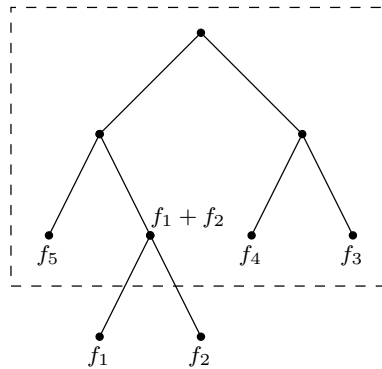
(the number of bits required for a symbol is exactly its depth in the tree).

There is another way to write this cost function that is very helpful. Although we are only given frequencies for the leaves, we can define the frequency of any *internal* node to be the sum of the frequencies of its descendant leaves; this is, after all, the number of times the internal node is visited during encoding or decoding. During the encoding process, each time we move down the tree, one bit gets output for every nonroot node through which we pass. So the total cost—the total number of bits which are output—can also be expressed thus:

The cost of a tree is the sum of the frequencies of all leaves and internal nodes, except the root.

The first formulation of the cost function tells us that *the two symbols with the smallest frequencies must be at the bottom of the optimal tree*, as children of the lowest internal node (this internal node has two children since the tree is *full*). Otherwise, swapping these two symbols with whatever is lowest in the tree would improve the encoding.

This suggests that we start constructing the tree *greedily*: find the two symbols with the smallest frequencies, say i and j , and make them children of a new node, which then has frequency $f_i + f_j$. To keep the notation simple, let's just assume these are f_1 and f_2 . By the second formulation of the cost function, any tree in which f_1 and f_2 are sibling-leaves has cost $f_1 + f_2$ plus the cost for a tree with $n - 1$ leaves of frequencies $(f_1 + f_2), f_3, f_4, \dots, f_n$:



The latter problem is just a smaller version of the one we started with. So we pull f_1 and f_2 off the list of frequencies, insert $(f_1 + f_2)$, and loop. The resulting algorithm can be described in terms of priority queue operations (as defined on page 120) and takes $O(n \log n)$ time if a binary heap (Section 4.5.2) is used.

procedure Huffman(f)

Input: An array $f[1 \dots n]$ of frequencies

Output: An encoding tree with n leaves

let H be a priority queue of integers, ordered by f

for $i = 1$ to n : insert(H, i)

for $k = n + 1$ to $2n - 1$:

$i = \text{deletemin}(H)$, $j = \text{deletemin}(H)$

 create a node numbered k with children i, j

$f[k] = f[i] + f[j]$

 insert(H, k)

Returning to our toy example: can you tell if the tree of Figure 5.10 is optimal?

Entropy

The annual county horse race is bringing in three thoroughbreds who have never competed against one another. Excited, you study their past 200 races and summarize these as probability distributions over four outcomes: first (“first place”), second, third, and other.

Outcome	Aurora	Whirlwind	Phantasm
first	0.15	0.30	0.20
second	0.10	0.05	0.30
third	0.70	0.25	0.30
other	0.05	0.40	0.20

Which horse is the most predictable? One quantitative approach to this question is to look at *compressibility*. Write down the history of each horse as a string of 200 values (first, second, third, other). The total number of bits needed to encode these track-record strings can then be computed using Huffman’s algorithm. This works out to 290 bits for Aurora, 380 for Whirlwind, and 420 for Phantasm (check it!). Aurora has the shortest encoding and is therefore in a strong sense the most predictable.

The inherent unpredictability, or *randomness*, of a probability distribution can be measured by the extent to which it is possible to compress data drawn from that distribution.

$$\text{more compressible} \equiv \text{less random} \equiv \text{more predictable}$$

Suppose there are n possible outcomes, with probabilities p_1, p_2, \dots, p_n . If a sequence of m values is drawn from the distribution, then the i th outcome will pop up roughly mp_i times (if m is large). For simplicity, assume these are exactly the observed frequencies, and moreover that the p_i ’s are all powers of 2 (that is, of the form $1/2^k$). It can be seen by induction (Exercise 5.19) that the number of bits needed to encode the sequence is $\sum_{i=1}^n mp_i \log(1/p_i)$. Thus the average number of bits needed to encode a single draw from the distribution is

$$\sum_{i=1}^n p_i \log \frac{1}{p_i}.$$

This is the *entropy* of the distribution, a measure of how much randomness it contains.

For example, a fair coin has two outcomes, each with probability $1/2$. So its entropy is

$$\frac{1}{2} \log 2 + \frac{1}{2} \log 2 = 1.$$

This is natural enough: the coin flip contains one bit of randomness. But what if the coin is not fair, if it has a $3/4$ chance of turning up heads? Then the entropy is

$$\frac{3}{4} \log \frac{4}{3} + \frac{1}{4} \log 4 = 0.81.$$

A biased coin is more predictable than a fair coin, and thus has lower entropy. As the bias becomes more pronounced, the entropy drops toward zero.

We explore these notions further in Exercise 5.18 and 5.19.

5.3 Horn formulas

In order to display human-level intelligence, a computer must be able to perform at least some modicum of logical reasoning. Horn formulas are a particular framework for doing this, for expressing logical facts and deriving conclusions.

The most primitive object in a Horn formula is a *Boolean variable*, taking value either true or false. For instance, variables x , y , and z might denote the following possibilities.

$x \equiv$ the murder took place in the kitchen
 $y \equiv$ the butler is innocent
 $z \equiv$ the colonel was asleep at 8 pm

A *literal* is either a variable x or its negation \bar{x} (“NOT x ”). In Horn formulas, knowledge about variables is represented by two kinds of *clauses*:

1. *Implications*, whose left-hand side is an AND of any number of positive literals and whose right-hand side is a single positive literal. These express statements of the form “if the conditions on the left hold, then the one on the right must also be true.” For instance,

$$(z \wedge w) \Rightarrow u$$

might mean “if the colonel was asleep at 8 pm and the murder took place at 8 pm then the colonel is innocent.” A degenerate type of implication is the *singleton* “ $\Rightarrow x$,” meaning simply that x is true: “the murder definitely occurred in the kitchen.”

2. Pure *negative clauses*, consisting of an OR of any number of negative literals, as in

$$(\bar{u} \vee \bar{v} \vee \bar{y})$$

(“they can’t all be innocent”).

Given a set of clauses of these two types, the goal is to determine whether there is a consistent explanation: an assignment of true/false values to the variables that satisfies all the clauses. This is also called a *satisfying assignment*.

The two kinds of clauses pull us in different directions. The implications tell us to set some of the variables to true, while the negative clauses encourage us to make them false. Our strategy for solving a Horn formula is this: We start with all variables false. We then proceed to set some of them to true, one by one, but very reluctantly, and only if we absolutely have to because an implication would otherwise be violated. Once we are done with this phase and all implications are satisfied, only then do we turn to the negative clauses and make sure they are all satisfied.

In other words, our algorithm for Horn clauses is the following greedy scheme (*stingy* is perhaps more descriptive):

Input: a Horn formula
 Output: a satisfying assignment, if one exists

```

set all variables to false

while there is an implication that is not satisfied:
    set the right-hand variable of the implication to true

if all pure negative clauses are satisfied: return the assignment
else: return ``formula is not satisfiable``

```

For instance, suppose the formula is

$$(w \wedge y \wedge z) \Rightarrow x, (x \wedge z) \Rightarrow w, x \Rightarrow y, \Rightarrow x, (x \wedge y) \Rightarrow w, (\overline{w} \vee \overline{x} \vee \overline{y}), (\overline{z}).$$

We start with everything false and then notice that x must be true on account of the singleton implication $\Rightarrow x$. Then we see that y must also be true, because of $x \Rightarrow y$. And so on.

To see why the algorithm is correct, notice that if it returns an assignment, this assignment satisfies both the implications and the negative clauses, and so it is indeed a satisfying truth assignment of the input Horn formula. So we only have to convince ourselves that if the algorithm finds no satisfying assignment, then there really is none. This is so because our “stingy” rule maintains the following invariant:

If a certain set of variables is set to true, then they must be true in *any* satisfying assignment.

Hence, if the truth assignment found after the *while* loop does not satisfy the negative clauses, there can be no satisfying truth assignment.

Horn formulas lie at the heart of Prolog (“programming by logic”), a language in which you program by specifying desired properties of the output, using simple logical expressions. The workhorse of Prolog interpreters is our greedy satisfiability algorithm. Conveniently, it can be implemented in time linear in the length of the formula; do you see how (Exercise 5.32)?

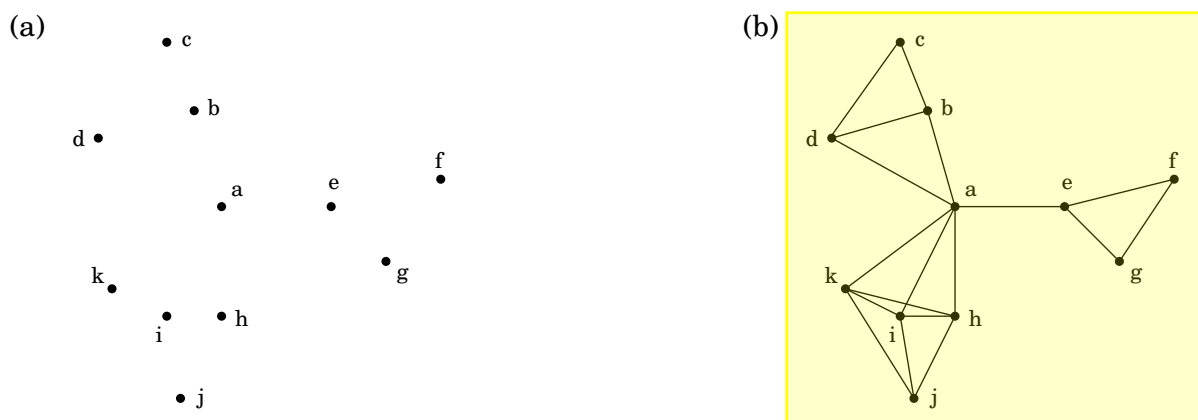
5.4 Set cover

The dots in Figure 5.11 represent a collection of towns. This county is in its early stages of planning and is deciding where to put schools. There are only two constraints: each school should be in a town, and no one should have to travel more than 30 miles to reach one of them. What is the minimum number of schools needed?

This is a typical *set cover* problem. For each town x , let S_x be the set of towns within 30 miles of it. A school at x will essentially “cover” these other towns. The question is then, how many sets S_x must be picked in order to cover all the towns in the county?

SET COVER

Input: A set of elements B ; sets $S_1, \dots, S_m \subseteq B$

Figure 5.11 (a) Eleven towns. (b) Towns that are within 30 miles of each other.

Output: A selection of the S_i whose union is B .

Cost: Number of sets picked.

(In our example, the elements of B are the towns.) This problem lends itself immediately to a greedy solution:

Repeat until all elements of B are covered:

Pick the set S_i with the largest number of uncovered elements.

This is extremely natural and intuitive. Let's see what it would do on our earlier example: It would first place a school at town a , since this covers the largest number of other towns. Thereafter, it would choose three more schools— c , j , and either f or g —for a total of four. However, there exists a solution with just three schools, at b , e , and i . **The greedy scheme is not optimal!**

But luckily, it isn't too far from optimal.

Claim Suppose B contains n elements and that the optimal cover consists of k sets. Then the greedy algorithm will use at most $k \ln n$ sets.²

Let n_t be the number of elements still not covered after t iterations of the greedy algorithm (so $n_0 = n$). Since these remaining elements are covered by the optimal k sets, **there must be some set with at least n_t/k of them.** Therefore, the greedy strategy will ensure that

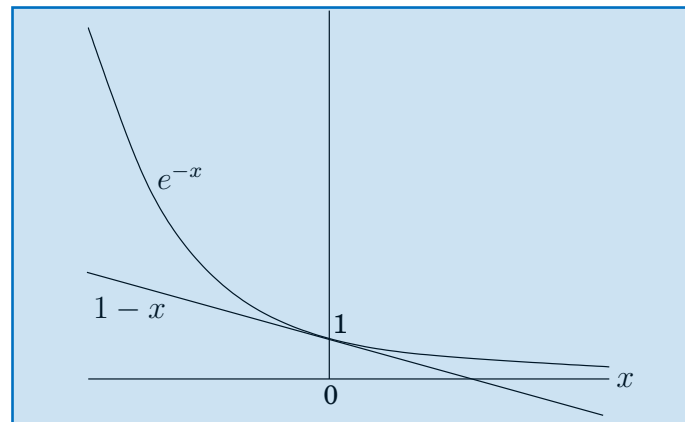
$$n_{t+1} \leq n_t - \frac{n_t}{k} = n_t \left(1 - \frac{1}{k}\right),$$

which by repeated application implies $n_t \leq n_0(1 - 1/k)^t$. A more convenient bound can be obtained from the useful inequality

$$1 - x \leq e^{-x} \text{ for all } x, \text{ with equality if and only if } x = 0,$$

² \ln means "natural logarithm," that is, to the base e .

which is most easily proved by a picture:



Thus

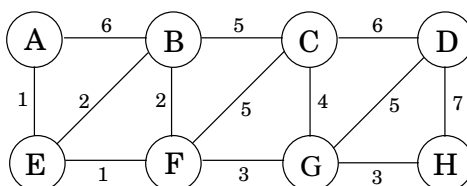
$$n_t \leq n_0 \left(1 - \frac{1}{k}\right)^t < n_0 (e^{-1/k})^t = n e^{-t/k}.$$

At $t = k \ln n$, therefore, n_t is strictly less than $n e^{-\ln n} = 1$, which means no elements remain to be covered.

The ratio between the greedy algorithm's solution and the optimal solution varies from input to input but is always less than $\ln n$. And there are certain inputs for which the ratio is very close to $\ln n$ (Exercise 5.33). We call this maximum ratio the *approximation factor* of the greedy algorithm. There seems to be a lot of room for improvement, but in fact such hopes are unjustified: it turns out that under certain widely-held complexity assumptions (which will be clearer when we reach Chapter 8), there is provably no polynomial-time algorithm with a smaller approximation factor.

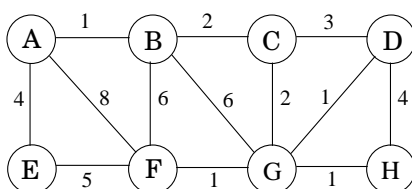
Exercises

5.1. Consider the following graph.



- What is the cost of its minimum spanning tree?
- How many minimum spanning trees does it have?
- Suppose Kruskal's algorithm is run on this graph. In what order are the edges added to the MST? For each edge in this sequence, give a cut that justifies its addition.

5.2. Suppose we want to find the minimum spanning tree of the following graph.



- Run Prim's algorithm; whenever there is a choice of nodes, always use alphabetic ordering (e.g., start from node A). Draw a table showing the intermediate values of the `cost` array.
- Run Kruskal's algorithm on the same graph. Show how the disjoint-sets data structure looks at every intermediate stage (including the structure of the directed trees), assuming path compression is used.

5.3. Design a linear-time algorithm for the following task.

Input: A connected, undirected graph G .

Question: Is there an edge you can remove from G while still leaving G connected?

Can you reduce the running time of your algorithm to $O(|V|)$?

5.4. Show that if an undirected graph with n vertices has k connected components, then it has at least $n - k$ edges.

5.5. Consider an undirected graph $G = (V, E)$ with nonnegative edge weights $w_e \geq 0$. Suppose that you have computed a minimum spanning tree of G , and that you have also computed shortest paths to all nodes from a particular node $s \in V$.

Now suppose each edge weight is increased by 1: the new weights are $w'_e = w_e + 1$.

- Does the minimum spanning tree change? Give an example where it changes or prove it cannot change.
- Do the shortest paths change? Give an example where they change or prove they cannot change.

- 5.6. Let $G = (V, E)$ be an undirected graph. Prove that if all its edge weights are distinct, then it has a unique minimum spanning tree.
- 5.7. Show how to find the *maximum* spanning tree of a graph, that is, the spanning tree of largest total weight.
- 5.8. Suppose you are given a weighted graph $G = (V, E)$ with a distinguished vertex s and where all edge weights are positive and distinct. Is it possible for a tree of shortest paths from s and a minimum spanning tree in G to not share any edges? If so, give an example. If not, give a reason.
- 5.9. The following statements may or may not be correct. In each case, either prove it (if it is correct) or give a counterexample (if it isn't correct). Always assume that the graph $G = (V, E)$ is undirected. Do not assume that edge weights are distinct unless this is specifically stated.
- (a) If graph G has more than $|V| - 1$ edges, and there is a unique heaviest edge, then this edge cannot be part of a minimum spanning tree.
 - (b) If G has a cycle with a unique heaviest edge e , then e cannot be part of any MST.
 - (c) Let e be any edge of minimum weight in G . Then e must be part of some MST.
 - (d) If the lightest edge in a graph is unique, then it must be part of every MST.
 - (e) If e is part of some MST of G , then it must be a lightest edge across some cut of G .
 - (f) If G has a cycle with a unique lightest edge e , then e must be part of every MST.
 - (g) The shortest-path tree computed by Dijkstra's algorithm is necessarily an MST.
 - (h) The shortest path between two nodes is necessarily part of some MST.
 - (i) Prim's algorithm works correctly when there are negative edges.
 - (j) (For any $r > 0$, define an r -path to be a path whose edges all have weight $< r$.) If G contains an r -path from node s to t , then every MST of G must also contain an r -path from node s to node t .
- 5.10. Let T be an MST of graph G . Given a connected subgraph H of G , show that $T \cap H$ is contained in some MST of H .
- 5.11. Give the state of the disjoint-sets data structure after the following sequence of operations, starting from singleton sets $\{1\}, \dots, \{8\}$. Use path compression. In case of ties, always make the lower numbered root point to the higher numbered one.
- union(1, 2), union(3, 4), union(5, 6), union(7, 8), union(1, 4), union(6, 7), union(4, 5), find(1)
- 5.12. Suppose you implement the disjoint-sets data structure using union-by-rank but not path compression. Give a sequence of m union and find operations on n elements that take $\Omega(m \log n)$ time.
- 5.13. A long string consists of the four characters A, C, G, T ; they appear with frequency 31%, 20%, 9%, and 40%, respectively. What is the Huffman encoding of these four characters?
- 5.14. Suppose the symbols a, b, c, d, e occur with frequencies $1/2, 1/4, 1/8, 1/16, 1/16$, respectively.
- (a) What is the Huffman encoding of the alphabet?
 - (b) If this encoding is applied to a file consisting of 1,000,000 characters with the given frequencies, what is the length of the encoded file in bits?

5.15. We use Huffman's algorithm to obtain an encoding of alphabet $\{a, b, c\}$ with frequencies f_a, f_b, f_c . In each of the following cases, either give an example of frequencies (f_a, f_b, f_c) that would yield the specified code, or explain why the code cannot possibly be obtained (no matter what the frequencies are).

(a) Code: $\{0, 10, 11\}$

(b) Code: $\{0, 1, 00\}$

(c) Code: $\{10, 01, 00\}$

5.16. Prove the following two properties of the Huffman encoding scheme.

(a) If some character occurs with frequency more than $2/5$, then there is guaranteed to be a codeword of length 1.

(b) If all characters occur with frequency less than $1/3$, then there is guaranteed to be no codeword of length 1.

5.17. Under a Huffman encoding of n symbols with frequencies f_1, f_2, \dots, f_n , what is the longest a codeword could possibly be? Give an example set of frequencies that would produce this case.

5.18. The following table gives the frequencies of the letters of the English language (including the blank for separating words) in a particular corpus.

blank	18.3%	r	4.8%	y	1.6%
e	10.2%	d	3.5%	p	1.6%
t	7.7%	l	3.4%	b	1.3%
a	6.8%	c	2.6%	v	0.9%
o	5.9%	u	2.4%	k	0.6%
i	5.8%	m	2.1%	j	0.2%
n	5.5%	w	1.9%	x	0.2%
s	5.1%	f	1.8%	q	0.1%
h	4.9%	g	1.7%	z	0.1%

(a) What is the optimum Huffman encoding of this alphabet?

(b) What is the expected number of bits per letter?

(c) Suppose now that we calculate the entropy of these frequencies

$$H = \sum_{i=1}^{26} p_i \log \frac{1}{p_i}$$

(see the box in page 156). Would you expect it to be larger or smaller than your answer above? Explain.

(d) Do you think that this is the limit of how much English text can be compressed? What features of the English language, besides letters and their frequencies, should a better compression scheme take into account?

5.19. *Entropy*. Consider a distribution over n possible outcomes, with probabilities p_1, p_2, \dots, p_n .

- (a) Just for this part of the problem, assume that each p_i is a power of 2 (that is, of the form $1/2^k$). Suppose a long sequence of m samples is drawn from the distribution and that for all $1 \leq i \leq n$, the i^{th} outcome occurs exactly mp_i times in the sequence. Show that if Huffman encoding is applied to this sequence, the resulting encoding will have length

$$\sum_{i=1}^n mp_i \log \frac{1}{p_i}.$$

- (b) Now consider arbitrary distributions—that is, the probabilities p_i are not restricted to powers of 2. The most commonly used measure of the *amount of randomness* in the distribution is the *entropy*

$$\sum_{i=1}^n p_i \log \frac{1}{p_i}.$$

For what distribution (over n outcomes) is the entropy the largest possible? The smallest possible?

- 5.20. Give a linear-time algorithm that takes as input a tree and determines whether it has a *perfect matching*: a set of edges that touches each node exactly once.

A *feedback edge set* of an undirected graph $G = (V, E)$ is a subset of edges $E' \subseteq E$ that intersects every cycle of the graph. Thus, removing the edges E' will render the graph acyclic.

Give an efficient algorithm for the following problem:

Input: Undirected graph $G = (V, E)$ with positive edge weights w_e

Output: A feedback edge set $E' \subseteq E$ of minimum total weight $\sum_{e \in E'} w_e$

- 5.21. In this problem, we will develop a new algorithm for finding minimum spanning trees. It is based upon the following property:

Pick any cycle in the graph, and let e be the heaviest edge in that cycle. Then there is a minimum spanning tree that does not contain e .

- (a) Prove this property carefully.
- (b) Here is the new MST algorithm. The input is some undirected graph $G = (V, E)$ (in adjacency list format) with edge weights $\{w_e\}$.

```

sort the edges according to their weights
for each edge  $e \in E$ , in decreasing order of  $w_e$ :
    if  $e$  is part of a cycle of  $G$ :
         $G = G - e$  (that is, remove  $e$  from  $G$ )
return  $G$ 
```

Prove that this algorithm is correct.

- (c) On each iteration, the algorithm must check whether there is a cycle containing a specific edge e . Give a linear-time algorithm for this task, and justify its correctness.
- (d) What is the overall time taken by this algorithm, in terms of $|E|$? Explain your answer.

5.22. You are given a graph $G = (V, E)$ with positive edge weights, and a minimum spanning tree $T = (V, E')$ with respect to these weights; you may assume G and T are given as adjacency lists. Now suppose the weight of a particular edge $e \in E$ is modified from $w(e)$ to a new value $\hat{w}(e)$. You wish to quickly update the minimum spanning tree T to reflect this change, without recomputing the entire tree from scratch. There are four cases. In each case give a linear-time algorithm for updating the tree.

- (a) $e \notin E'$ and $\hat{w}(e) > w(e)$.
- (b) $e \notin E'$ and $\hat{w}(e) < w(e)$.
- (c) $e \in E'$ and $\hat{w}(e) < w(e)$.
- (d) $e \in E'$ and $\hat{w}(e) > w(e)$.

5.23. Sometimes we want light spanning trees with certain special properties. Here's an example.

Input: Undirected graph $G = (V, E)$; edge weights w_e ; subset of vertices $U \subset V$

Output: The lightest spanning tree in which the nodes of U are *leaves* (there might be other leaves in this tree as well).

(The answer isn't necessarily a minimum spanning tree.)

Give an algorithm for this problem which runs in $O(|E| \log |V|)$ time. (*Hint:* When you remove nodes U from the optimal solution, what is left?)

5.24. A binary counter of unspecified length supports two operations: *increment* (which increases its value by one) and *reset* (which sets its value back to zero). Show that, starting from an initially zero counter, any sequence of n *increment* and *reset* operations takes time $O(n)$; that is, the amortized time per operation is $O(1)$.

5.25. Here's a problem that occurs in automatic program analysis. For a set of variables x_1, \dots, x_n , you are given some *equality* constraints, of the form " $x_i = x_j$ " and some *inequality* constraints, of the form " $x_i \neq x_j$." Is it possible to satisfy all of them?

For instance, the constraints

$$x_1 = x_2, x_2 = x_3, x_3 = x_4, x_1 \neq x_4$$

cannot be satisfied. Give an efficient algorithm that takes as input m constraints over n variables and decides whether the constraints can be satisfied.

5.26. *Graphs with prescribed degree sequences.* Given a list of n positive integers d_1, d_2, \dots, d_n , we want to efficiently determine whether there exists an undirected graph $G = (V, E)$ whose nodes have degrees precisely d_1, d_2, \dots, d_n . That is, if $V = \{v_1, \dots, v_n\}$, then the degree of v_i should be exactly d_i . We call (d_1, \dots, d_n) the *degree sequence of G* . This graph G should not contain self-loops (edges with both endpoints equal to the same node) or multiple edges between the same pair of nodes.

- (a) Give an example of d_1, d_2, d_3, d_4 where all the $d_i \leq 3$ and $d_1 + d_2 + d_3 + d_4$ is even, but for which no graph with degree sequence (d_1, d_2, d_3, d_4) exists.
- (b) Suppose that $d_1 \geq d_2 \geq \dots \geq d_n$ and that there exists a graph $G = (V, E)$ with degree sequence (d_1, \dots, d_n) . We want to show that there must exist a graph that has this degree sequence and where in addition the neighbors of v_1 are $v_2, v_3, \dots, v_{d_1+1}$. The idea is to gradually transform G into a graph with the desired additional property.
 - i. Suppose the neighbors of v_1 in G are not $v_2, v_3, \dots, v_{d_1+1}$. Show that there exists $i < j \leq n$ and $u \in V$ such that $\{v_1, v_i\}, \{u, v_j\} \notin E$ and $\{v_1, v_j\}, \{u, v_i\} \in E$.

- ii. Specify the changes you would make to G to obtain a new graph $G' = (V, E')$ with the same degree sequence as G and where $(v_1, v_i) \in E'$.
 - iii. Now show that there must be a graph with the given degree sequence but in which v_1 has neighbors $v_2, v_3, \dots, v_{d_1+1}$.
 - (c) Using the result from part (b), describe an algorithm that on input d_1, \dots, d_n (not necessarily sorted) decides whether there exists a graph with this degree sequence. Your algorithm should run in time polynomial in n and in $m = \sum_{i=1}^n d_i$.
- 5.27. Alice wants to throw a party and is deciding whom to call. She has n people to choose from, and she has made up a list of which pairs of these people know each other. She wants to pick as many people as possible, subject to two constraints: at the party, each person should have at least five other people whom they know *and* five other people whom they don't know.
- Give an efficient algorithm that takes as input the list of n people and the list of pairs who know each other and outputs the best choice of party invitees. Give the running time in terms of n .
- 5.28. A *prefix-free encoding* of a finite alphabet Γ assigns each symbol in Γ a binary codeword, such that no codeword is a prefix of another codeword.
- Show that such an encoding can be represented by a full binary tree in which each leaf corresponds to a unique element of Γ , whose codeword is generated by the path from the root to that leaf (interpreting a left branch as 0 and a right branch as 1).
- 5.29. *Ternary Huffman*. Trimedia Disks Inc. has developed “ternary” hard disks. Each cell on a disk can now store values 0, 1, or 2 (instead of just 0 or 1). To take advantage of this new technology, provide a modified Huffman algorithm for compressing sequences of characters from an alphabet of size n , where the characters occur with known frequencies f_1, f_2, \dots, f_n . Your algorithm should encode each character with a variable-length codeword over the values 0, 1, 2 such that no codeword is a prefix of another codeword and so as to obtain the maximum possible compression. Prove that your algorithm is correct.
- 5.30. The basic intuition behind Huffman's algorithm, that frequent blocks should have short encodings and infrequent blocks should have long encodings, is also at work in English, where typical words like I, you, is, and, to, from, and so on are short, and rarely used words like velociraptor are longer.
- However, words like fire!, help!, and run! are short not because they are frequent, but perhaps because time is precious in situations where they are used.
- To make things theoretical, suppose we have a file composed of m different words, with frequencies f_1, \dots, f_m . Suppose also that for the i th word, the cost per bit of encoding is c_i . Thus, if we find a prefix-free code where the i th word has a codeword of length l_i , then the total cost of the encoding will be $\sum_i f_i \cdot c_i \cdot l_i$.
- Show how to modify Huffman's algorithm to find the prefix-free encoding of minimum total cost.
- 5.31. A server has n customers waiting to be served. The service time required by each customer is known in advance: it is t_i minutes for customer i . So if, for example, the customers are served in order of increasing i , then the i th customer has to wait $\sum_{j=1}^i t_j$ minutes.
- We wish to minimize the total waiting time

$$T = \sum_{i=1}^n (\text{time spent waiting by customer } i).$$

Give an efficient algorithm for computing the optimal order in which to process the customers.

- 5.32. Show how to implement the stingy algorithm for Horn formula satisfiability (Section 5.3) in time that is linear in the length of the formula (the number of occurrences of literals in it). (*Hint*: Use a directed graph, with one node per variable, to represent the implications.)
- 5.33. Show that for any integer n that is a power of 2, there is an instance of the set cover problem (Section 5.4) with the following properties:
- i. There are n elements in the base set.
 - ii. The optimal cover uses just two sets.
 - iii. The greedy algorithm picks at least $\log n$ sets.

Thus the approximation ratio we derived in the chapter is tight.