

Obie looked at the seein' eye dog. Then at the twenty-seven 8 by 10 color glossy pictures with the circles and arrows and a paragraph on the back of each one. . . and then he looked at the seein' eye dog. And then at the twenty-seven 8 by 10 color glossy pictures with the circles and arrows and a paragraph on the back of each one and began to cry.

Because Obie came to the realization that it was a typical case of American blind justice, and there wasn't nothin' he could do about it, and the judge wasn't gonna look at the twenty-seven 8 by 10 color glossy pictures with the circles and arrows and a paragraph on the back of each one explainin' what each one was, to be used as evidence against us.

And we was fined fifty dollars and had to pick up the garbage. In the snow.

But that's not what I'm here to tell you about.

— Arlo Guthrie, "Alice's Restaurant" (1966)

I study my Bible as I gather apples.

First I shake the whole tree, that the ripest might fall.

Then I climb the tree and shake each limb,

and then each branch and then each twig,

and then I look under each leaf.

— Martin Luther

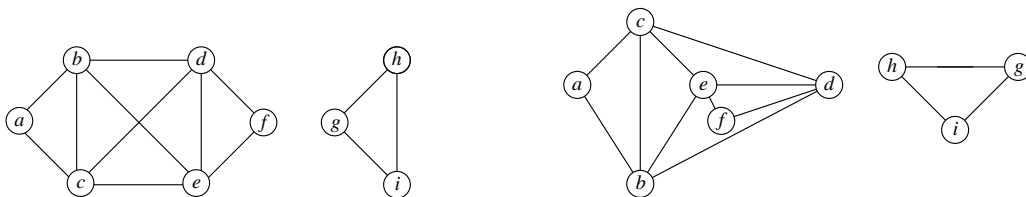
11 Basic Graph Properties

11.1 Definitions

A graph G is a pair of sets (V, E) . V is a set of arbitrary objects that we call *vertices*¹ or *nodes*. E is a set of vertex pairs, which we call *edges* or occasionally *arcs*. In an *undirected* graph, the edges are unordered pairs, or just sets of two vertices. In a *directed* graph, the edges are ordered pairs of vertices. We will only be concerned with *simple* graphs, where there is no edge from a vertex to itself and there is at most one edge from any vertex to any other.

Following standard (but admittedly confusing) practice, I'll also use V to denote the *number* of vertices in a graph, and E to denote the *number* of edges. Thus, in an undirected graph, we have $0 \leq E \leq \binom{V}{2}$, and in a directed graph, $0 \leq E \leq V(V - 1)$.

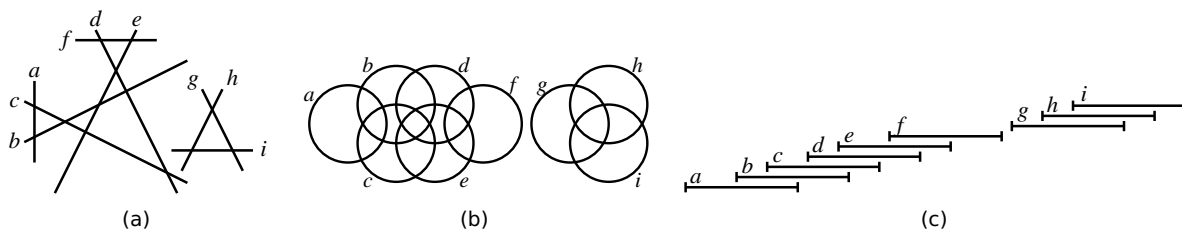
We usually visualize graphs by looking at an *embedding*. An embedding of a graph maps each vertex to a point in the plane and each edge to a curve or straight line segment between the two vertices. A graph is *planar* if it has an embedding where no two edges cross. The same graph can have many different embeddings, so it is important not to confuse a particular embedding with the graph itself. In particular, planar graphs can have non-planar embeddings!



A non-planar embedding of a planar graph with nine vertices, thirteen edges, and two connected components, and a planar embedding of the same graph.

¹The singular of 'vertices' is **vertex**. The singular of 'matrices' is **matrix**. Unless you're speaking Italian, there is no such thing as a *vertice*, a *matrice*, an *indice*, an *appendice*, a *helice*, an *apice*, a *vortice*, a *radice*, a *simplice*, a *codice*, a *directrice*, a *dominatrice*, a *Unice*, a *Kleenice*, an *Asterice*, an *Obelice*, a *Dogmatice*, a *Getafice*, a *Cacofonice*, a *Vitalstatistice*, a *Geriatric*, or *Jimi Hendrice*! You *will* lose points for using any of these so-called words.

There are other ways of visualizing and representing graphs that are sometimes also useful. For example, the *intersection graph* of a collection of objects has a node for every object and an edge for every intersecting pair. Whether a particular graph can be represented as an intersection graph depends on what kind of object you want to use for the vertices. Different types of objects—line segments, rectangles, circles, etc.—define different classes of graphs. One particularly useful type of intersection graph is an *interval graph*, whose vertices are intervals on the real line, with an edge between any two intervals that overlap.



The example graph is also the intersection graph of (a) a set of line segments, (b) a set of circles, or (c) a set of intervals on the real line (stacked for visibility).

If (u, v) is an edge in an undirected graph, then u is a *neighbor* of v and vice versa. The *degree* of a node is the number of neighbors. In directed graphs, we have two kinds of neighbors. If $u \rightarrow v$ is a directed edge, then u is a *predecessor* of v and v is a *successor* of u . The *in-degree* of a node is the number of predecessors, which is the same as the number of edges going into the node. The *out-degree* is the number of successors, or the number of edges going out of the node.

A graph $G' = (V', E')$ is a *subgraph* of $G = (V, E)$ if $V' \subseteq V$ and $E' \subseteq E$.

A *path* is a sequence of edges, where each successive pair of edges shares a vertex, and all other edges are disjoint. A graph is *connected* if there is a path from any vertex to any other vertex. A disconnected graph consists of several *connected components*, which are maximal connected subgraphs. Two vertices are in the same connected component if and only if there is a path between them.

A *cycle* is a path that starts and ends at the same vertex, and has at least one edge. A graph is *acyclic* if no subgraph is a cycle; acyclic graphs are also called *forests*. *Trees* are special graphs that can be defined in several different ways. You can easily prove by induction (hint, hint, hint) that the following definitions are equivalent.

- A tree is a connected acyclic graph.
- A tree is a connected component of a forest.
- A tree is a connected graph with *at most* $V - 1$ edges.
- A tree is a minimal connected graph; removing any edge makes the graph disconnected.
- A tree is an acyclic graph with *at least* $V - 1$ edges.
- A tree is a maximal acyclic graph; adding an edge between any two vertices creates a cycle.

A *spanning tree* of a graph G is a subgraph that is a tree and contains every vertex of G . Of course, a graph can only have a spanning tree if it's connected. A *spanning forest* of G is a collection of spanning trees, one for each connected component of G .

11.2 Explicit Representations of Graphs

There are two common data structures used to explicitly represent graphs: *adjacency matrices*² and *adjacency lists*.

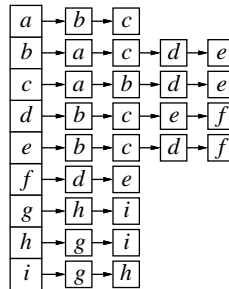
The adjacency matrix of a graph G is a $V \times V$ matrix of indicator variables. Each entry in the matrix indicates whether a particular edge is or is not in the graph:

$$A[i, j] = [(i, j) \in E].$$

For undirected graphs, the adjacency matrix is always *symmetric*: $A[i, j] = A[j, i]$. Since we don't allow edges from a vertex to itself, the diagonal elements $A[i, i]$ are all zeros.

Given an adjacency matrix, we can decide in $\Theta(1)$ time whether two vertices are connected by an edge just by looking in the appropriate slot in the matrix. We can also list all the neighbors of a vertex in $\Theta(V)$ time by scanning the corresponding row (or column). This is optimal in the worst case, since a vertex can have up to $V - 1$ neighbors; however, if a vertex has few neighbors, we may still have to examine every entry in the row to see them all. Similarly, adjacency matrices require $\Theta(V^2)$ space, regardless of how many edges the graph actually has, so it is only space-efficient for very *dense* graphs.

	<i>a b c d e f g h i</i>
<i>a</i>	0 1 1 0 0 0 0 0 0
<i>b</i>	1 0 1 1 1 0 0 0 0
<i>c</i>	1 1 0 1 1 0 0 0 0
<i>d</i>	0 1 1 0 1 1 0 0 0
<i>e</i>	0 1 1 1 0 1 0 0 0
<i>f</i>	0 0 0 1 1 0 0 0 0
<i>g</i>	0 0 0 0 0 0 0 1 0
<i>h</i>	0 0 0 0 0 0 1 0 1
<i>i</i>	0 0 0 0 0 0 1 1 0



Adjacency matrix and adjacency list representations for the example graph.

For *sparse* graphs—graphs with relatively few edges—we're better off using adjacency lists. An adjacency list is an array of linked lists, one list per vertex. Each linked list stores the neighbors of the corresponding vertex.

For undirected graphs, each edge (u, v) is stored twice, once in u 's neighbor list and once in v 's neighbor list; for directed graphs, each edge is stored only once. Either way, the overall space required for an adjacency list is $O(V + E)$. Listing the neighbors of a node v takes $O(1 + \text{deg}(v))$ time; just scan the neighbor list. Similarly, we can determine whether (u, v) is an edge in $O(1 + \text{deg}(u))$ time by scanning the neighbor list of u . For undirected graphs, we can speed up the search by simultaneously scanning the neighbor lists of both u and v , stopping either we locate the edge or when we fall off the end of a list. This takes $O(1 + \min\{\text{deg}(u), \text{deg}(v)\})$ time.

The adjacency list structure should immediately remind you of hash tables with chaining. Just as with hash tables, we can make adjacency list structure more efficient by using something besides a linked list to store the neighbors. For example, if we use a hash table with constant load factor, when we can detect edges in $O(1)$ expected time, just as with an adjacency list. In practice, this will only be useful for vertices with large degree, since the constant overhead in both the space and search time is larger for hash tables than for simple linked lists.

You might at this point ask why anyone would ever use an adjacency matrix. After all, if you use hash tables to store the neighbors of each vertex, you can do everything as fast or faster with an adjacency list as with an adjacency matrix, only using less space. The answer is that many graphs are only represented

²See footnote 1.

implicitly. For example, intersection graphs are usually represented implicitly by simply storing the list of objects. As long as we can test whether two objects overlap in constant time, we can apply any graph algorithm to an intersection graph by *pretending* that it is stored explicitly as an adjacency matrix. On the other hand, any data structure build from records with pointers between them can be seen as a directed graph. Algorithms for searching graphs can be applied to these data structures by *pretending* that the graph is represented explicitly using an adjacency list.

To keep things simple, we'll consider only undirected graphs for the rest of this lecture, although the algorithms I'll describe also work for directed graphs.

11.3 Traversing connected graphs

Suppose we want to visit every node in a connected graph (represented either explicitly or implicitly). The simplest method to do this is an algorithm called *depth-first search*, which can be written either recursively or iteratively. It's exactly the same algorithm either way; the only difference is that we can actually see the 'recursion' stack in the non-recursive version. Both versions are initially passed a *source* vertex s .

<pre> RECURSIVEDFS(v): if v is unmarked mark v for each edge (v, w) RECURSIVEDFS(w) </pre>	<pre> ITERATIVEDFS(s): PUSH(s) while stack not empty v ← POP if v is unmarked mark v for each edge (v, w) PUSH(w) </pre>
--	--

Depth-first search is one (perhaps the most common) instance of a general family of graph traversal algorithms. The generic graph traversal algorithm stores a set of candidate edges in some data structure that I'll call a 'bag'. The only important properties of a 'bag' are that we can put stuff into it and then later take stuff back out. (In C++ terms, think of the 'bag' as a template for a real data structure.) Here's the algorithm:

<pre> TRAVERSE(s): put (∅, s) in bag while the bag is not empty take (p, v) from the bag (*) if v is unmarked mark v parent(v) ← p for each edge (v, w) (†) put (v, w) into the bag (**) </pre>
--

Notice that we're keeping *edges* in the bag instead of *vertices*. This is because we want to remember, whenever we visit a vertex v for the first time, which previously-visited vertex p put v into the bag. The vertex p is called the *parent* of v .

Lemma 1. *TRAVERSE(s) marks every vertex in any connected graph exactly once, and the set of edges $(v, \text{parent}(v))$ with $\text{parent}(v) \neq \emptyset$ form a spanning tree of the graph.*

Proof: first, it should be obvious that no node is marked more than once.

Clearly, the algorithm marks s . Let $v \neq s$ be a vertex, and let $s \rightarrow \dots \rightarrow u \rightarrow v$ be the path from s to v with the minimum number of edges. Since the graph is connected, such a path always exists. (If s and v

are neighbors, then $u = s$, and the path has just one edge.) If the algorithm marks u , then it must put (u, v) into the bag, so it must later take (u, v) out of the bag, at which point v must be marked (if it isn't already). Thus, by induction on the shortest-path distance from s , the algorithm marks every vertex in the graph.

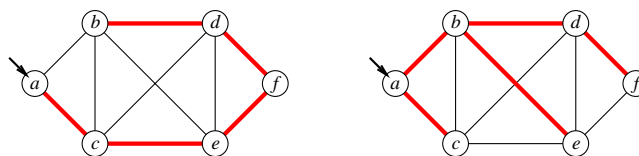
Call an edge $(v, \text{parent}(v))$ with $\text{parent}(v) \neq \emptyset$ a *parent edge*. For any node v , the path of parent edges $v \rightarrow \text{parent}(v) \rightarrow \text{parent}(\text{parent}(v)) \rightarrow \dots$ eventually leads back to s , so the set of parent edges form a connected graph. Clearly, both endpoints of every parent edge are marked, and the number of parent edges is exactly one less than the number of vertices. Thus, the parent edges form a spanning tree. \square

The exact running time of the traversal algorithm depends on how the graph is represented and what data structure is used as the 'bag', but we can make a few general observations. Since each vertex is visited at most once, the for loop (\dagger) is executed at most V times. Each edge is put into the bag exactly twice; once as (u, v) and once as (v, u) , so line ($\star\star$) is executed at most $2E$ times. Finally, since we can't take more things out of the bag than we put in, line (\star) is executed at most $2E + 1$ times.

11.4 Examples

Let's first assume that the graph is represented by an adjacency list, so that the overhead of the for loop (\dagger) is only a constant per edge.

- If we implement the 'bag' by using a *stack*, we have *depth-first search*. Each execution of (\star) or ($\star\star$) takes constant time, so the overall running time is $O(V + E)$. Since the graph is connected, $V \leq E + 1$, so we can simplify the running time to $O(E)$. The spanning tree formed by the parent edges is called a *depth-first spanning tree*. The exact shape of the tree depends on the order in which neighbor edges are pushed onto the stack, but in general, depth-first spanning trees are long and skinny.
- If we use a *queue* instead of a stack, we have *breadth-first search*. Again, each execution of (\star) or ($\star\star$) takes constant time, so the overall running time is still $O(E)$. In this case, the *breadth-first spanning tree* formed by the parent edges contains *shortest paths* from the start vertex s to every other vertex in its connected component. The exact shape of the breadth-first spanning tree depends on the order in which neighbor edges are pushed onto the queue, but in general, shortest path trees are short and bushy. We'll see shortest paths again in a future lecture.



A depth-first spanning tree and a breadth-first spanning tree of one component of the example graph, with start vertex a .

- Suppose the edges of the graph are weighted. If we implement the 'bag' using a *priority queue*, always extracting the minimum-weight edge in line (\star), then we have what might be called *shortest-first search*. In this case, each execution of (\star) or ($\star\star$) takes $O(\log E)$ time, so the overall running time is $O(V + E \log E)$, which simplifies to $O(E \log E)$ if the graph is connected. For this algorithm, the set of parent edges form the *minimum spanning tree* of the connected component of s . We'll see minimum spanning trees again in the next lecture.

If the graph is represented using an adjacency matrix instead of an adjacency list, finding all the neighbors of each vertex in line (†) takes $O(V)$ time. Thus, depth- and breadth-first search each take $O(V^2)$ time overall, and ‘shortest-first search’ takes $O(V^2 + E \log E) = O(V^2 \log V)$ time overall.

11.5 Searching disconnected graphs

If the graph is disconnected, then $\text{TRAVERSE}(s)$ only visits the nodes in the connected component of the start vertex s . If we want to visit all the nodes in every component, we can use the following ‘wrapper’ around our generic traversal algorithm. Since TRAVERSE computes a spanning tree of one component, TRAVERSEALL computes a spanning *forest* of the entire graph.

```

TRAVERSEALL(s):
  for all vertices v
    if v is unmarked
      TRAVERSE(v)

```

Exercises

1. Prove that the following definitions are all equivalent.
 - A tree is a connected acyclic graph.
 - A tree is a connected component of a forest.
 - A tree is a connected graph with *at most* $V - 1$ edges.
 - A tree is a minimal connected graph; removing any edge makes the graph disconnected.
 - A tree is an acyclic graph with *at least* $V - 1$ edges.
 - A tree is a maximal acyclic graph; adding an edge between any two vertices creates a cycle.
2. Prove that any connected acyclic graph with $n \geq 2$ vertices has at least two vertices with degree 1. Do not use the words ‘tree’ or ‘leaf’, or any well-known properties of trees; your proof should follow entirely from the definitions.
3. Let G be a connected graph, and let T be a depth-first spanning tree of G rooted at some node v . Prove that if T is also a breadth-first spanning tree of G rooted at v , then $G = T$.
4. Whenever groups of pigeons gather, they instinctively establish a *pecking order*. For any pair of pigeons, one pigeon always pecks the other, driving it away from food or potential mates. The same pair of pigeons always chooses the same pecking order, even after years of separation, no matter what other pigeons are around. Surprisingly, the overall pecking order can contain cycles—for example, pigeon A pecks pigeon B , which pecks pigeon C , which pecks pigeon A .
 - (a) Prove that any finite set of pigeons can be arranged in a row from left to right so that every pigeon pecks the pigeon immediately to its left. Pretty please.
 - (b) Suppose you are given a directed graph representing the pecking relationships among a set of n pigeons. The graph contains one vertex per pigeon, and it contains an edge $i \rightarrow j$ if and only if pigeon i pecks pigeon j . Describe and analyze an algorithm to compute a pecking order for the pigeons, as guaranteed by part (a).

5. You are helping a group of ethnographers analyze some oral history data they have collected by interviewing members of a village to learn about the lives of people lived there over the last two hundred years. From the interviews, you have learned about a set of people, all now deceased, whom we will denote P_1, P_2, \dots, P_n . The ethnographers have collected several facts about the lifespans of these people. Specifically, for some pairs (P_i, P_j) , the ethnographers have learned one of the following facts:

- (a) P_i died before P_j was born.
- (b) P_i and P_j were both alive at some moment.

Naturally, the ethnographers are not sure that their facts are correct; memories are not so good, and all this information was passed down by word of mouth. So they'd like you to determine whether the data they have collected is at least internally consistent, in the sense that there could have existed a set of people for which all the facts they have learned simultaneously hold.

Describe and analyze an algorithm to answer the ethnographers' problem. Your algorithm should either output possible dates of birth and death that are consistent with all the stated facts, or it should report correctly that no such dates exist.

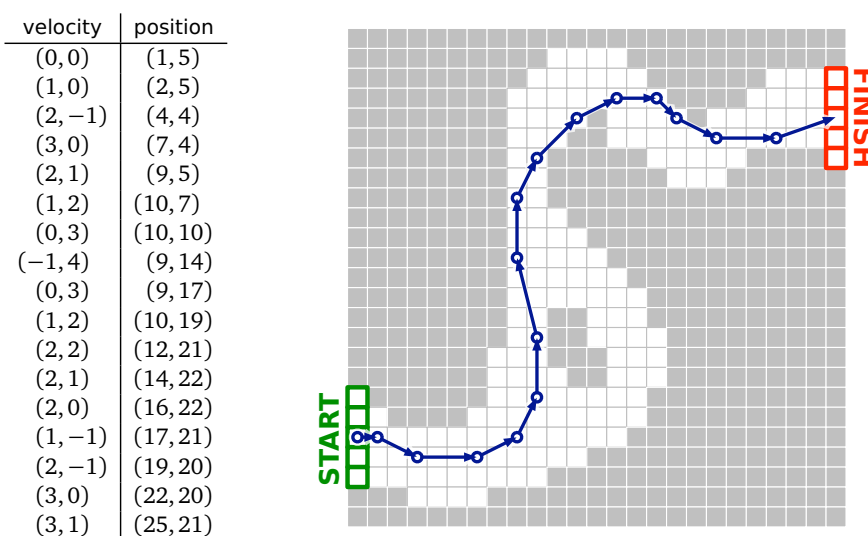
6. Let $G = (V, E)$ be a given directed graph.
- (a) The *transitive closure* G^T is a directed graph with the same vertices as G , that contains any edge $u \rightarrow v$ if and only if there is a directed path from u to v in G . Describe an efficient algorithm to compute the transitive closure of G .
 - (b) The *transitive reduction* G^{TR} is the smallest graph (meaning fewest edges) whose transitive closure is G^T . Describe an efficient algorithm to compute the transitive reduction of G .
7. A graph (V, E) is *bipartite* if the vertices V can be partitioned into two subsets L and R , such that every edge has one vertex in L and the other in R .
- (a) Prove that every tree is a bipartite graph.
 - (b) Describe and analyze an efficient algorithm that determines whether a given undirected graph is bipartite.
8. An **Euler tour** of a graph G is a closed walk through G that traverses every edge of G exactly once.
- (a) Prove that a connected graph G has an Euler tour if and only if every vertex has even degree.
 - (b) Describe and analyze an algorithm to compute an Euler tour in a given graph, or correctly report that no such graph exists.
9. The d -dimensional hypercube is the graph defined as follows. There are $2d$ vertices, each labeled with a different string of d bits. Two vertices are joined by an edge if their labels differ in exactly one bit.
- (a) A Hamiltonian cycle in a graph G is a cycle of edges in G that visits every vertex of G exactly once. Prove that for all $d \geq 2$, the d -dimensional hypercube has a Hamiltonian cycle.
 - (b) Which hypercubes have an Euler tour (a closed walk that traverses every edge exactly once)?
[Hint: This is very easy.]

10. **Racetrack** (also known as *Graph Racers* and *Vector Rally*) is a two-player paper-and-pencil racing game that Jeff played on the bus in 5th grade.³ The game is played with a track drawn on a sheet of graph paper. The players alternately choose a sequence of grid points that represent the motion of a car around the track, subject to certain constraints explained below.

Each car has a *position* and a *velocity*, both with integer x - and y -coordinates. The initial position is a point on the starting line, chosen by the player; the initial velocity is always $(0, 0)$. At each step, the player optionally increments or decrements either or both coordinates of the car's velocity; in other words, each component of the velocity can change by at most 1 in a single step. The car's new position is then determined by adding the new velocity to the car's previous position. The new position must be inside the track; otherwise, the car crashes and that player loses the race. The race ends when the first car reaches a position on the finish line.

Suppose the racetrack is represented by an $n \times n$ array of bits, where each 0 bit represents a grid point inside the track, each 1 bit represents a grid point outside the track, the 'starting line' is the first column, and the 'finish line' is the last column.

Describe and analyze an algorithm to find the minimum number of steps required to move a car from the starting line to the finish line of a given racetrack. [Hint: Build a graph. What are the vertices? What are the edges? What problem is this?]



A 16-step Racetrack run, on a 25×25 track. This is *not* the shortest run on this track.

- *11. Draughts/checkers is a game played on an $m \times m$ grid of squares, alternately colored light and dark. (The game is usually played on an 8×8 or 10×10 board, but the rules easily generalize to any board size.) Each dark square is occupied by at most one game piece (usually called a *checker* in the U.S.), which is either black or white; light squares are always empty. One player ('White') moves the white pieces; the other ('Black') moves the black pieces.

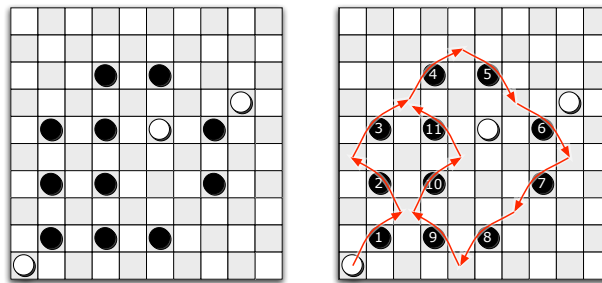
Consider the following simple version of the game, essentially American checkers or British draughts, but where every piece is a king.⁴ Pieces can be moved in any of the four diagonal

³The actual game is a bit more complicated than the version described here. In particular, in the actual game, the boundaries of the track are a free-form curve, and (at least by default) the entire line segment between any two consecutive positions must lie inside the track. In the version Jeff played, if a car does run off the track, the car starts its next turn with zero velocity, at the legal grid point closest to where the car left the track.

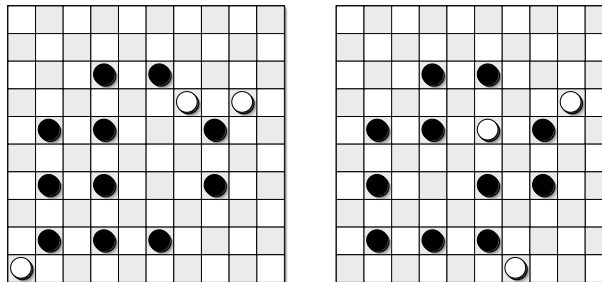
⁴Most other variants of draughts have 'flying kings', which behave very differently than what's described here.

directions, either one or two steps at a time. On each turn, a player either *moves* one of her pieces one step diagonally into an empty square, or makes a series of *jumps* with one of her checkers. In a single jump, a piece moves to an empty square two steps away in any diagonal direction, but only if the intermediate square is occupied by a piece of the opposite color; this enemy piece is *captured* and immediately removed from the board. Multiple jumps are allowed in a single turn as long as they are made by the same piece. A player wins if her opponent has no pieces left on the board.

Describe an algorithm that correctly determines whether White can capture every black piece, thereby winning the game, *in a single turn*. The input consists of the width of the board (m), a list of positions of white pieces, and a list of positions of black pieces. For full credit, your algorithm should run in $O(n)$ time, where n is the total number of pieces.



White wins in one turn.



White cannot win in one turn from either of these positions.

[Hint: The greedy strategy—make arbitrary jumps until you get stuck—does **not** always find a winning sequence of jumps even when one exists. See problem 8. Parity, parity, parity.]