

# Appendix

David C. Parkes and Sven Seuken

May 18, 2012

In this appendix we review some of the technical background necessary to understand the content in the main chapters of the book.

## 100.1 $O()$ -Notation and Polynomial Worst-Case Time Complexity

In this section, we briefly review the big-O-notation for the analysis of asymptotic time complexity of an algorithm. Let  $x$  describe the size of the input to an algorithm. For example, this could be the number of words in a text, the number of entries in a table, or the number of variables in a linear equation. Let  $T(x)$  denote the *worst-case* time complexity of an algorithm, which is the *maximum* amount of time taken on *any* input of size  $x$ . Equivalently, by a re-normalization into basic computational steps, each of which takes no more than some constant amount of time,  $T(x)$  can also be defined as the maximum number of basic computational steps taken on any input of size  $x$ .

The time complexity of an algorithm is often expressed *asymptotically* using big O notation. For example, an algorithm whose time complexity is  $T(x) = 2x^2 + 100x^{1/2}$  has asymptotic time complexity of  $O(x^2)$ ; in this case, we can write  $T(x) = O(x^2)$  or equivalently  $T(x) \in O(x^2)$ . In typical usage, big O notation hides constants and if  $T(x)$  is the sum of several terms, then using big O notation we only keep the term with the largest growth rate. Here, we ignore the term  $x^{1/2}$  because it is dominated by term  $x^2$  as  $x$  gets large, and we ignore the factor 2 in front of  $x^2$ .

**Definition 100.1** (Big O-Notation). *Let  $x$  describe the size of the input to function  $T(\cdot)$ . We write:*

$$T(x) \in O(g(x)) \text{ or } T(x) = O(g(x))$$

*if there exists a positive real number  $k$  and a real number  $x_0$  such that:*

$$\text{for all } x > x_0 : T(x) \leq k \cdot g(x).$$

The big-O-notation is used to describe the *asymptotic* time complexity because it focuses on the grow rate of the function when the input size becomes very large (formally, as  $x \rightarrow \infty$ ). For these very large input sizes, the term with the largest growth rate will dominate.

Given this, an algorithm is *polynomial time* if the running time is upper bounded by a polynomial expression in the size of the input. Equivalently, we can say:

**Definition 100.2** (Polynomial Time Complexity). *An algorithm has worst-case polynomial time complexity if  $T(x) \in O(x^k)$  for some constant  $k$ .*

Polynomial-time computation is the usual yard-stick for *efficient* computation in theoretical computer science. We normally say that an algorithm is efficient if it runs in worst-case polynomial time, and we say it is *inefficient* otherwise. Of course this hides extremely important practical details, including 1) the multiplicative factors as well as 2) the exponent  $k$  that describes the asymptotic growth rate. For example, even if an algorithm has time complexity  $T(x) = 1,000,000 \cdot x^2$  we write  $T(x) \in O(x^2)$ , and say that the algorithm is efficient. Furthermore, if an algorithm has time complexity  $T(x) = x^{500} \in O(x^{500})$ , the algorithm formally still has a polynomial worst-case time complexity. But it may not be practicable to run the algorithm, even for medium-sized  $x$ , because the asymptotic growth rate is too high.

Nevertheless, this distinction between algorithms that have a worst-case polynomial time complexity and those that do not, has proven very useful in computer science. For most of the algorithms with worst-case polynomial time complexity that we study, the multiplicative factors and the exponents are relatively small. In contrast, algorithms with an exponential running time such that  $T(x) = 2^x$  are usually infeasible to run, except for very small input sizes.

## 100.2 NP-Completeness

In this section we review some basic concepts from computational complexity theory. A *decision problem* is a problem with a Yes or No answer. Example decision problems studied in computer science are

- “given a Boolean formula, is there a satisfying truth assignment?”
- “given a graph, is there a way to color each vertex with one of  $k$  colors such that no adjacent vertices share the same color?”
- “given a set of items, each with some weight and value, and a bag with some capacity, is there a subset of items with total value greater than some target value, and with total weight that does not exceed the capacity?”

These correspond to the decision problems SATISFIABILITY, GRAPH COLORING and 0-1 KNAPSACK. Decision problems have a yes-or-no answer. Instances of a decision problem are either “yes”-instances (the answer is Yes) or “no”-instances (the answer is No.)

Computational complexity theory provides the following distinctions in regard to the complexity of different decision problems:

- A **polynomial decision problem** (the problem is in **P**) is one for which there exists a worst-case polynomial time algorithm that answers *Yes* on every “yes”-instance and *No* on every “no”-instance.
- A **non-deterministic polynomial decision problem** (the problem is in **NP**) is one for which there exists a worst-case polynomial time algorithm that can *verify* that the answer is *Yes* on every “yes”-instance, given a proof as input in addition to the instance; e.g., for GRAPH COLORING, the proof would be an assignment of a color to each vertex, and the algorithm would check that no pair of adjacent vertices have the same color.<sup>1</sup>
- An **NP-hard decision problem** (the problem is **NP-hard**) need not be in **NP**, but is at least as hard to solve as the hardest problems in **NP**.<sup>2</sup> In particular, a problem  $A$  is **NP-hard** if

<sup>1</sup>Note that (i) to be in **NP** a problem does not require an easy-to-verify proof for “no”-instances (when this exists, it is a **co-NP** problem); (ii) the size of the proof must be polynomial in the size of the input if it is to be checked in polynomial time; and (iii) there is an equivalent definition in terms of a non-deterministic theoretical computational machine (a Turing machine), and it is from this that the name **NP** derives.

<sup>2</sup>For example, the class **NP-hard** also contains non-decision problems such as *optimization problems* (e.g., find the

an algorithm to solve the problem in worst-case polynomial-time would imply an algorithm to solve any problem  $B$  that is in  $\mathbf{NP}$  in worst-case polynomial time.

This would be achieved through a **reduction**, where an instance of  $B$  is encoded as an instance of  $A$  such that a solution to  $A$  gives a solution to  $B$ , and where the instance of  $A$  is at most polynomially larger than the instance of  $B$ .

- An **NP-complete** decision problem (the problem is  $\mathbf{NP}$ -complete) is one for which the problem is both in  $\mathbf{NP}$  and also  $\mathbf{NP}$ -hard.

Intuitively, problems in  $\mathbf{P}$  are “easy” to solve. But what about problems in  $\mathbf{NP}$ ? On the one hand, there is at least a fast way to verify that a “yes”-instance is indeed a yes instance. Certainly, every decision problem in  $\mathbf{P}$  is also in  $\mathbf{NP}$ . One can simply use the polynomial time algorithm that solves the decision problem as the polynomial-time verifier. But what about polynomial-time algorithms to solve the hard problems in  $\mathbf{NP}$ ? This leads us to the following:

**Conjecture 1.**  $\mathbf{P} \neq \mathbf{NP}$ .

This famous conjecture states that the class  $\mathbf{P}$  is a strict subset of  $\mathbf{NP}$ , which in particular requires that *no problem that is NP-hard is in P*. To see this, notice that if a problem  $A$  that was  $\mathbf{NP}$ -hard was in  $\mathbf{P}$  then by the definition of  $\mathbf{NP}$ -hard, all problems in  $\mathbf{NP}$  could be solved by the polynomial-time algorithm for  $A$ .

Proving that “ $\mathbf{P} \neq \mathbf{NP}$ ” is the principal unsolved problem in theoretical computer science today. This conjecture is widely believed for two reasons. First, a large number of different  $\mathbf{NP}$ -complete problems can be reduced to each other, in the sense that if one can be solved in worst-case polynomial time, then all others can also be solved in worst-case polynomial time. Second, a great deal of effort has been dedicated to finding efficient algorithms for the various  $\mathbf{NP}$ -complete problems, but none has been found so far.

Now remember our discussion of what constitutes an *efficient* algorithm and what constitutes an *inefficient* algorithm. We can now say, that for a problem that is in  $\mathbf{P}$ , there exist efficient algorithms for solving it. And under the assumption that  $\mathbf{P} \neq \mathbf{NP}$ , we can say that for a problem that is  $\mathbf{NP}$ -hard, there does not exist an efficient algorithm for solving it.

### 100.3 Graph theory primer

A *graph*  $G = (V, E)$  is a mathematical model that contains a set of *vertices*  $V$  (or *nodes*) and a set of *edges*  $E$ . For example, the graph in Figure 100.1 (a) contains four nodes  $A, B, C$  and  $D$  and four edges  $AB, BC, CD$  and  $BC$ . Graphs are mathematical models of *networks*, and represent relationships between objects. An edge in a graph may be undirected as in Figure 100.1 (a) or directed as in Figure 100.1 (b). When all edges are undirected, a graph is said to be undirected. Otherwise the graph is *directed*. For the most part we will focus on undirected graphs.

A *path* is a sequence of vertices with the property that each consecutive pair is connected by an edge (in the right direction in the case of a directed graph). A path may include a *cycle*, e.g.  $A - B - C - D - B$  contains the cycle  $B - C - D - B$ , or a path may be a cycle, meaning that it contains at least two edges, and the first and last vertices are the same. For example  $B - D - B$  is a cycle in Figure 100.1 (b) and  $B - D - C - B$  is a cycle in Figure 100.1 (a). A *neighbor* of a vertex  $v$  is connected to the vertex via a single edge.

---

subset of items that maximizes total value is the “optimization variant” on the 0-1 KNAPSACK problem.)

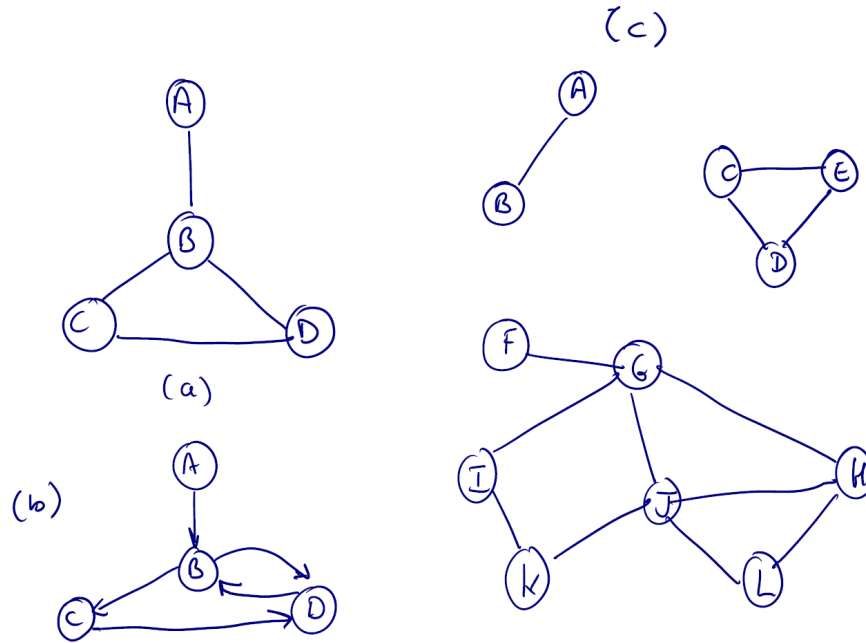


Figure 100.1: Three example graphs. From Easley and Kleinberg (2010).

The length of a path is the number of edges on the path, and the *distance* between two vertices is the length of the shortest path between them. For example, the distance between  $K$  and  $H$  in Figure 100.1 (c) is 2, the distance between  $K$  and  $I$  is 3. Distance can be computed between some vertex and all other vertices by a “breadth first search”: consider all vertices 1-hop away, then all 2-hop away by going to neighbors of the 1-hop vertices, and so forth.

An essential aspect of a graph (and thus the networks modeled through graphs) is an implied *transitivity* between objects: if  $A$  is connected to  $B$  by an edge and  $B$  is connected to  $C$  by an edge, then there is some kind of implied relationship between  $A$  and  $C$ . For example, in a graph representing a social network, people are vertices and edges represent social interaction, colloquially “friends.” In this case, we would associate some meaning with the relationship between  $A$  and  $C$  in Figure 100.1 (a), saying that  $C$  is a “friend of a friend” of  $A$ ’s.

In a *collaboration network* for Wikipedia, editors correspond to vertices and an edge between two editors to indicate they have edited the same article. Here, the editors that are connected by a short path, such as  $K$  and  $H$  in Figure 100.1 (c), may have a higher propensity to have shared interests. A possible graph model of the eBay market associates a vertex with an auction. An edge between two auctions indicates that at least one buyer participated in both auctions, the auctions overlapped in time, and the buyer ended up winning in only one of the auctions.

A graph is *connected* if there is a path between every pair of vertices, e.g. Figure 100.1 (a) but not the other two graphs. A *component* (sometimes referred to as a *connected component*) is a maximal subset of vertices that are all pairwise connected. That is, there is a path from any vertex to any other vertex in the subset, and there is no additional vertex that can be added to the subset without maintaining this property. For example,  $\{A, B, C, D\}$  is a component in Figure 100.1 (a), and  $\{A\}$  and  $\{B, C, D\}$  are components in Figure 100.1 (b). There are three components in

Figure 100.1 (c).

Components may reveal some simple structure about a graph. In fact, many networks may have what is loosely referred to as a *giant component*, a component that contains a significant fraction (think, more than  $1/3$ ) of all vertices. When a network includes a giant component we might expect it to include only one such giant component. For example, it is hard to imagine two giant components in a social network: all it would take is one edge from one component to the other and the components would be one, and as the two components both get large this becomes a likely event.

The *average distance* of a connected graph is the average distance between all pairs of nodes. The *diameter* is the maximum distance between any pair of nodes. For example, suppose the graph in Figure 100.1 (c) consisted of only its largest component. In this case, the diameter of this graph would be 3.