

Computer Science 351

Application: The Analysis of Algorithms

Instructor: Wayne Eberly

Department of Computer Science
University of Calgary

Lecture #22

Learning Goals

- Learn about discrete probability theory can be applied to carry out an ***average case analysis*** of a deterministic algorithm, and to analyze various kinds of ***randomized algorithms***.
- Learn about various kinds of randomized algorithms for decision problems *that are allowed to fail* in various ways.

Example: Linear Search

Consider a ***linear search algorithm*** that searches for a copy of “true” in a Boolean array A with length n :

```
integer search ( boolean[] A ) {  
  1. integer  $n := A.length$   
  2. integer  $i := 0$   
  3. while ( $i < n$ ) {  
  4.   if ( $A[i] == true$ ) {  
  5.     return  $i$   
    }  
  6.    $i := i + 1$   
    }  
  7. throw a NoSuchElementException  
}
```

Example: Linear Search

- To simplify our analysis, suppose we count the number of numbered steps that are carried out when this algorithm is executed.
 - If $0 \leq i \leq n - 1$ and the first copy of “true” is in position i then $3i + 5$ steps are executed.
 - If there is no copy of “true” stored in the array, then $3n + 4$ steps are executed.

Worst-Case Analysis

- If an algorithm is deterministic — like the “linear search” algorithm here — then the number of steps that is used, on any given input, is a constant that only depends on that input.
- We often want to measure, or bound, the number of steps used as a function of the “size” of the input.
- For this problem, let us *define* the “size” of the input to be the length, n , of the array A that is part of the input.
- The ***worst-case running time*** of an algorithm is the *maximum* number of steps that is used by algorithm when it is executed on an input with a given size.

Worst-Case Analysis

- It follows from the above that — if “running time” and “size” are defined as shown here — then the **worst-case running time** of this linear search algorithm is the maximum of

$$\max_{0 \leq i \leq n-1} (3i + 5) \quad \text{and} \quad 3n + 4,$$

that is, $3n + 4$.

Average-Case Analysis

- Sometimes, the “worst-case running time” seems too pessimistic because you almost never have an execution of the algorithm that uses the number of steps given as its “worst-case running time” — and you are interested in (and satisfied by) knowing how many steps are used, most of the time, instead.
- An ***average-case analysis*** might be more helpful here.

Average-Case Analysis

In order to perform an average-case analysis for a given input size n , consider an **experiment** in which you are executing your algorithm on an input with size n .

- Sometimes the **sample space**, Ω , can be the set of all inputs with size n . See the preparatory material, for the first lecture on discrete probability theory in this course, for an analysis that uses this sample space.
- It is sometime also possible to use a **sample space** to be a collection of **sets** inputs with size n such that the algorithm's execution is essentially the same for all of the inputs in any one of the sets being used. This approach is being used in the next part of these lecture notes.

Average-Case Analysis

- Suppose that we use a sample space

$$\Omega = \{s_0, s_1, s_2, \dots, s_{n-1} u\}$$

such that

- for $0 \leq i \leq n - 1$, s_i includes all input arrays A (with length n) such that $A[j] = \text{false}$ for $0 \leq j \leq i - 1$, and $A[i] = \text{true}$.
- u includes all inputs A (with length n) such that $A[i] = \text{false}$, for every integer i such that $0 \leq i \leq n - 1$.

Average-Case Analysis

A **random variable** is used to define the resource that we wish to study.

- For this example, at least *two* random variables might be considered:
 - $C : \Omega \rightarrow \mathbb{R}$ represents the number of times that array entries are checked — so that $C(s_i) = i + 1$ for $0 \leq i \leq n - 1$, and $C(u) = n$.
 - $T : \Omega \rightarrow \mathbb{R}$ represents the number of executions of numbered steps — so that $T(S_i) = 3i + 5$ for $0 \leq i \leq n - 1$ and $T(u) = 3n + 4$.

Average-Case Analysis

In order to complete an average-case analysis we need to know — or, more generally, make an **assumption** — about how likely each outcome in Ω is.

- This is modelled by a **probability distribution** $P : \Omega \rightarrow \mathbb{R}$.
- The “expected running time” (or results of this “average-case analysis”) is the **expected value** $E[T]$ of the random variable T with respect to the probability distribution P .
- This can depend, quite heavily, on the distribution P being used — and can be “technically correct” but also **misleading** or **irrelevant** if the assumptions about likelihoods are not correct.

Average-Case Analysis

- The preparatory material for the first lecture on discrete probability theory completed an average-case analysis — using the random variable C .
- A supplemental document for this lecture provide another analysis (using the random variable T instead of C — with a *different* assumption about the likelihood of inputs).

Average-Case Analysis

Where Might You See Average-Case Analysis?

- In CPSC 331 a deterministic version of a QuickSort algorithm might be given, using assumptions about whether entries in the input array distinct, and the relative orderings of the entries in the input array.
- This course might include average-case analyses involving ***hash tables*** and ***randomly constructed*** binary search trees.

Randomized Algorithms

The following algorithm uses a random number generator:

```
integer rSearch ( boolean[] A ) {  
  1. integer  $n := A.length$   
  2. integer  $i := 0$   
  3. while ( $i < n$ ) {  
  4.   Choose an integer  $j$  uniformly and randomly from  
         $\{0, 1, 2, \dots, n - 1\}$ .  
  5.   if ( $A[j] == true$ ) {  
  6.     return  $j$   
    }  
  7.    $i := i + 1$   
    }  
  8. throw a NoSuchElementException  
}
```

Randomized Algorithms

- This is an example of a **randomized algorithm**: It uses a random number generator during its execution so that neither the output it returns, nor the number of steps it uses to generate this output is fixed, even when the input is.
- If $A[j] = \text{false}$ for every integer j such that $0 \leq j \leq n - 1$ then — even though choices of the values for j might be different for different executions of the algorithm — the test at line 5 can never pass, so that there will always be n executions of the body of the loop, and the step at line 8 will always be reached and executed.
- This can be used to argue that there will always be exactly $4n + 4$ executions of numbered steps, in this case.

Randomized Algorithms

- Suppose A is an array such that $A[j] = \text{true}$ for every integer j such that $0 \leq j \leq n - 1$, instead. This time, the test at line 5 must pass during the first execution of the loop.
- The output that is returned might be any integer j such that $0 \leq j \leq n - 1$ — and each is returned with probability $\frac{1}{n}$. However, the number of steps executed is fixed: Since there is always exactly one execution of the loop body, each execution of the algorithm would include an execution of exactly six numbered steps.

Randomized Algorithms

- Now consider a more general — and somewhat more complicated — case: The sets

$$S_A = \{j \in \mathbb{N} \mid 0 \leq j \leq n - 1 \text{ and } A[j] = \text{true}\}$$

and

$$F_A = \{j \in \mathbb{N} \mid 0 \leq j \leq n - 1 \text{ and } A[j] = \text{false}\}$$

are both non-empty. Let $k = |S_A|$, so that $1 \leq k \leq n - 1$ and $|F_A| = n - k$.

- Now neither the output, nor the number of steps used is fixed — even though the input is.

Randomized Algorithms

Consider, now, the **experiment** of executing the `rSearch` algorithm on an input including the array A as described above.

- The **sample space** would include enough information, about **the random choices made**, so that the execution of the algorithm could be studied.
- In particular, the sample space could be set to be

$$\Omega = \Omega_0 \cup \Omega_1 \cup \dots \cup \Omega_{n-1} \cup \Omega_n$$

where $\Omega_0, \Omega_1, \dots, \Omega_{n-1}, \Omega_n$ are as follows.

Randomized Algorithms

Ω_0 includes executions such that an index j , such that $A[j]$ is true, is found immediately.

- This would include each sequence $\langle j_1 \rangle$ of values, with length one, where $j_1 \in S_A$ — so that $|\Omega_0| = |S_A| = k$.

For $1 \leq i \leq n - 1$, Ω_i includes executions such that an index j is found, such that $A[j]$ is true, during the $i + 1^{\text{st}}$ execution of the body of the loop.

- This would include each sequence

$$\langle j_1, j_2, \dots, j_{i+1} \rangle$$

with length $i + 1$, where $j_1, j_2, \dots, j_i \in F_A$ and $j_{i+1} \in S_A$ — so that $|\Omega_i| = |F_A|^k \cdot |S_A| = (n - k)^i \cdot k$.

Randomized Algorithms

Ω_n includes executions where no index j , such that $A[j] = j = \text{true}$, is ever found at all.

- This would include each sequence

$$\langle j_1, j_2, \dots, j_n \rangle$$

with length n such that $j_1, j_2, \dots, j_n \in F_A$ — so that $|\Omega_n| = |F_A|^n = (n - k)^n$.

Randomized Algorithms

Once again, a **random variable** is used to define the resource that we wish to study.

- For this example, at least *two* random variables might be considered:
 - $C : \Omega \rightarrow \mathbb{R}$ represents the number of times that array entries are checked — so that, for $\alpha \in \Omega$,

$$C(\alpha) = \begin{cases} i + 1 & \text{if } \alpha \in \Omega_i, \text{ for } 0 \leq i \leq n - 1, \\ n & \text{if } \alpha \in \Omega_n. \end{cases}$$

- $T : \Omega \rightarrow \mathbb{R}$ represents the number of executions of numbered steps — so that, for $\alpha \in \Omega$,

$$T(\alpha) = \begin{cases} 4i + 6 & \text{if } \alpha \in \Omega_i, \text{ for } 0 \leq i \leq n - 1, \\ 4n + 4 & \text{if } \alpha \in \Omega_n. \end{cases}$$

Randomized Algorithms

When analyzing randomized algorithms we (almost always) use a ***probability distribution*** that models the assumption that the random number generator, used during the execution, really *does* provide random data.

- It follows — for this example — that, every time the value of j is set, during an execution of the step at line 4, j receives value h with probability $\frac{1}{n}$, for every integer h such that $0 \leq h \leq n - 1$.
- Furthermore, each assignment of a value to j , at line 4, is independent of the previous assignments of values to j .

Randomized Algorithms

- Each outcome in Ω_i would have probability $n^{-(i+1)}$, and $|\Omega_i| = (n - k)^i \cdot k$, for $0 \leq i \leq n - 1$, so that (for the probability distribution now being described),

$$P(\Omega_i) = \left(1 - \frac{k}{n}\right)^i \cdot \frac{k}{n}.$$

- Each outcome in Ω_n would have probability n^{-n} , and $|\Omega_n| = (n - k)^n$, so that (for the probability distribution now being described),

$$P(\Omega_n) = \left(1 - \frac{k}{n}\right)^n.$$

Randomized Algorithms

- A supplemental document for this lecture includes the computation of the expected value of T , with respect to the probability distribution that has been defined.

Randomized Algorithms

Definition:

- The ***expected running time*** of a randomized algorithm's execution, on a given input, is the expected value of the random variable, representing its running time, when modelled as suggested above (so that the "sample space" models the random values that are generated as an execution of the algorithm proceeds).
- The ***worst-case expected running time*** of a randomized algorithm is a function of the *size* of an input, just as the "worst-case running time" of a deterministic algorithm is: It is, essentially, the "maximum" of the expected running times of the algorithm's executions on inputs of the given size.

Randomized Algorithms

Where You Might See This Kind of Analysis?

- In CPSC 331 a *randomized* version of a QuickSort algorithm is also given. The process described here is used to bound the expected number of steps used by this algorithm to sort an array of length n . This is then used to bound the “worst-case expected running time” of this randomized algorithm.
- This application, from CPSC 331, is a little different than the above example because *the randomized algorithm, being considered, never gives up.*

Decision Problems

A **decision problem** is a computational problem that answers a “Yes-or-No” question — so that the algorithm’s output is always either `true` (corresponding to the answer “Yes”) or `false` (corresponding to the answer “No”).

Example: Consider *another* version of problem of searching for “true” in a Boolean array:

- As with the previous example, the input is a Boolean array with positive length.
- For this problem, we are asking whether “true” is stored in the array — so that the desired output is `true` if there an integer i such that $0 \leq i < A.length$ and $A[i] = true$ — and the desired output is `false`, otherwise.

Las Vegas Algorithms

A **Las Vegas** algorithm is a randomized algorithm that can never return an incorrect answer — so that, when an execution ends, the output provided (either `true` or `false`) is correct.

- The number of steps executed by this algorithm, when it is run on a given input, is a **random variable** over a sample space defined using the “random” choices made during the algorithms’s executions — just as it is for other randomized algorithms.

Example: This is *Not* a Las Vegas Algorithm

Consider a modified version, `rSearch2`, of the randomized algorithm from the previous example — which searches for a key in an integer array instead of looking for a copy of “true”.

- Rather than returning an integer j such that $0 \leq j \leq n - 1$ and $A[j] = \text{true}$ — this algorithm returns “true” when such an integer is found.
- Rather than throwing an `NoSuchElementException` if the desired value is not found at the end, the algorithm returns `false`.

Thus only two lines of code (at lines 6 and 8) are changed, to produce `rSearch2` from `rSearch`.

Example: This is *Not* a Las Vegas Algorithm

This is not a Las Vegas algorithm!

- To see why, notice that “true” is stored in the array — at some position ℓ such that $0 \leq \ell \leq A.length - 1$ — but it is not stored anywhere else, and the value ℓ is never chosen as the value for j at line 4 when this step is executed, then the step at line 8 will eventually be reached.
- In this case the incorrect output `false` will be returned, when a ***Las Vegas algorithm*** would be required to return the correct output, `true`, instead.

Example: This *is* a Las Vegas Algorithm

Suppose we make one more change to the above algorithm:

- Instead of returning “false” if attempts to locate “true” fail, and the step at line 8 is reached, a minor variant of the **deterministic** search algorithm, considered at the beginning of this lecture (which returns the desired output — true, or false) is used to solve the problem.
- Then, since the deterministic algorithm is correct, the resulting randomized algorithm’s output would also be correct — so that it *would* be a **Las Vegas** algorithm.

Monte Carlo Algorithms

A **Monte Carlo** algorithm is a randomized algorithm that can, sometimes, return an incorrect answer — but that does so with small probability.

- The algorithm only returns `true`, when executed on a given input, if this is the correct answer for that input. That is, if the algorithm is executed on an input where the answer that *should* be returned is `false`, then the algorithm is guaranteed to return `false`.
- If the algorithm is executed on an input where the answer that *should* be returned is `true`, then the probability that the algorithm *does* return `false` is at least $\frac{1}{2}$.
- *Another Way to Think about This:* This algorithm can generate **false negatives** but it only does so with small probability — and it cannot return **false positives**, at all.

Monte Carlo Algorithms

Consider the randomized algorithm `rSearch`, which was given above as an example of a randomized algorithm that is **not** a Las Vegas algorithm.

- As noted above, this algorithm only returns `true` when `true` is stored in the input array `A` — so that there are no **false positives**.
- As shown in a supplemental document for this lecture it can be established that if `true` is stored in the input array `A` — so that `true` should be returned — then the probability that `false` is returned instead (so that there has been a **false positive**) is less than $\frac{1}{2}$.
- This is, therefore, an example of a **Monte Carlo algorithm**.

Two-Sided Error

- While they are not discussed as often, randomized algorithms that allow **false positives** with small probability, but never allow **false negatives**, are sometimes of interest.
- Randomized algorithms that allow *both* “false positives” *and* “false negatives” are of interest too — but only if the likelihood of a mistake is significantly reduced!
- In particular, randomized algorithms that allow both false negatives and false positives *but where the probability of an incorrect result is never more than $\frac{1}{4}$, for any input*, are also of some interest.

What *Really* Happens?

As repeatedly stated: The results obtained, using probability theory, are only relevant if the assumptions being made — which are used to define probability distributions — are satisfied.

- The results are “technically correct” but ***irrelevant***, otherwise.
- Thus the results of average-case analyses of deterministic algorithms might not be relevant, because the distribution of inputs might be different than assumed.
- Results concerning randomized algorithms are even more problematic because ***modern programming languages, in widespread use, do not really provide random number generators!***

What *Really* Happens?

- Instead, these use ***deterministic*** processes that might take a small amount of information (like the time of day, machine load, or a user-supplied value) as a “seed” and use this to produce a sequence of values whose properties “approximate” those of a sequence of randomly generated values, in some sense.
- ***Results observed in practice*** typically agree with the results that analyses supplied using probability theory, assuming the use of truly “random” sequences.
- A further discussion of this is beyond the scope of this course — but students who are interested in this can investigate ***pseudorandom number generators*** if they wish to learn more.