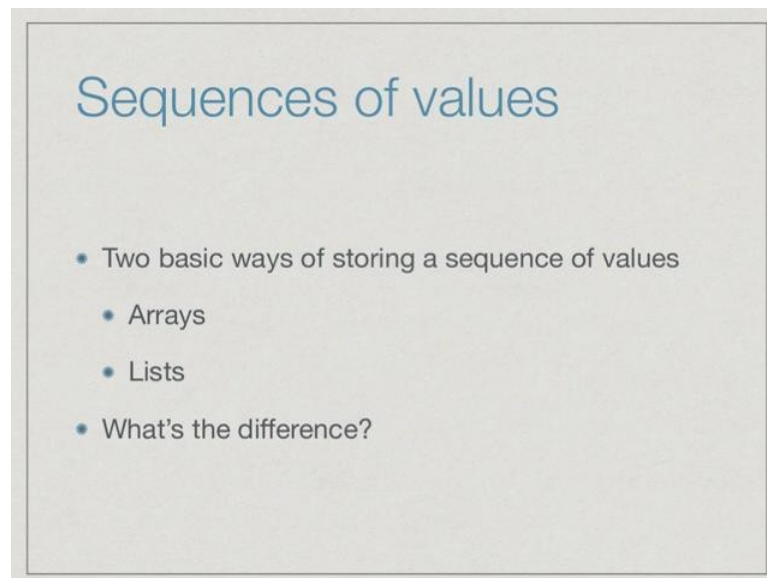


Programming, Data Structures and Algorithms in Python
Prof. Madhavan Mukund
Department of Computer Science and Engineering
Indian Institute of Technology, Madras

Week – 03
Lecture – 04
Arrays vs. Lists, Binary search

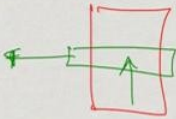
(Refer Slide Time: 00:02)



We have seen several situations where we want to store a Sequence of values. Now it turns out that in a program or in a programming language implementation, there are two basic ways in which we can store such a sequence. These are normally called Arrays and Lists. So, let us look at the difference between Arrays and Lists.

(Refer Slide Time: 00:22)

Arrays



The diagram shows a horizontal rectangle representing a memory block. A green arrow points from the left edge to a specific point within the rectangle, labeled 'i'. Another green arrow points from the bottom edge to the same point, labeled 'i'. A red rectangle is drawn around the point, indicating a specific element or a range of elements.

- Single block of memory, elements of uniform type
 - Typically size of sequence is fixed in advance
- Indexing is fast
 - Access `seq[i]` in constant time for any `i`
 - Compute offset from start of memory block
- Inserting between `seq[i]` and `seq[i+1]` is expensive
- Contraction is expensive

An array is usually a sequence which is stored as a single block in memory. So, you can imagine if you wish that your memory is arranged in certain way and then you have an array, so usually memories arranged in what are called Words. Word is one unit of what you can store or retrieve from memory, and an array will usually be one continuous block without any gaps.

And, in particular this would apply when an array has only a single type of value, so all the elements in the sequence are either integers or floats or something where the length of each element of the array is of a uniform size. We would also typically in an array know in advance how big this block is. So we might know that it has say 100 entry, so we have a sequence of size 100.

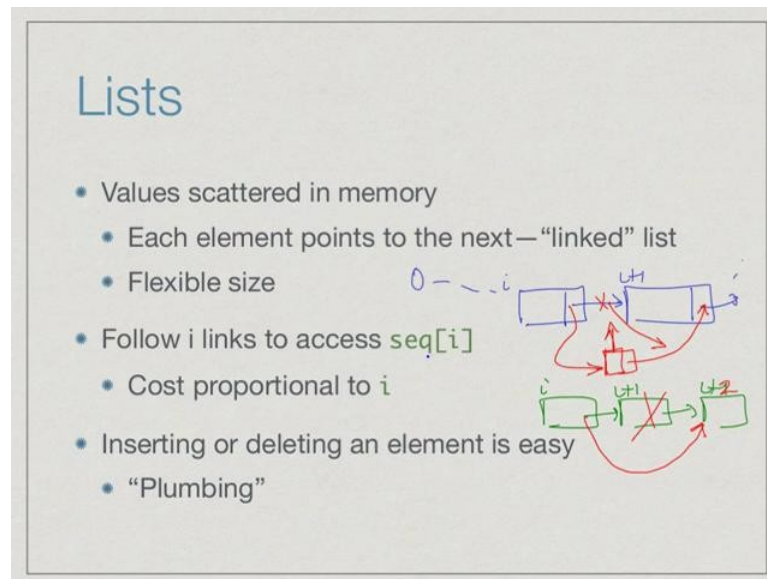
Now when this happens, what happens is that if you want to look at the j th element of a sequence or the i th element of a sequence, then what you want to think of is this block of memory starting with 1, 2, 3, up to i right and you want to get to the i th element quickly. But since everything is of a uniform size and you know where this starts, we know where the sequence starts you can just compute i times this size of one unit and quickly go and one shot to the location in the memory where the i th element is saved.

So, accessing the i th element of an array just requires arithmetic computation of the address by starting with the initial point of the array and then walking forward i units to the i th position. And this can be done in what we could call Constant time. By constant time what we mean is it does not really depend on i . It is no easier or no difficult to get the last element of an array as it is to get to the second element of an array, it is independent of i . It takes the fixed amount of time to get to sequence of y for any i .

Now, one consequence of this is inserting or contracting arrays is expensive, because now if I have an array with 0 to 99 and I want to add a new value here say at position i then first of all this array now becomes from 0 to 100 and now everything which is after i has to be shifted to accommodate space if we want to keep the same representation with the entire array is stored as a single block. So, when we have a single block of memory though it is efficient to get to any part of it quickly it is not very efficient to expand it because we have to then shift everything. The worst case for example, if this green block comes into 0th position then the entire array has to be shifted down by one position.

In the same way contraction is also expensive because we have to make a hole in some sense. If we remove this element out then we have a hole here and then we have to push everything up to block this hole, because – remember the array must have all elements contiguous that is without any gaps starting from the initial position.

(Refer Slide Time: 03:27)



The other way of storing a sequence, is to store it one element at a time and not bother about how these elements are with respect to each other in the memory. I can think of this memory as a large space and now I might have one element here, so this is my first element and then I will have a way of saying that from here the next element is somewhere else, this is what we call a Link. So very often in the implementation these are called linked list, so I may have the first element here. Now because of various reasons I might end up putting the second element here and so on.

You can imagine that if you have some say space in your cupboard and then you take out things and then you put things back but you put things back in the first place where you have an empty slot, then the sequence in which you put things back may not respect the sequence in which they appear finally in the shelf. So, here in the same way we do not have any physical assumption about how these elements are stored, we just have a logical link from the first element to the next element and so on.

The other part of this is that we do not have to worry about the overall length of the list because we know we started at the 0th position and we keep walking down. On the last position so say suppose the last position is in fact two then there would be some indication here saying that there is no next element, so two is the last element. A list can

have a flexible size and obviously because we are just pointing one element to another, we can also accommodate what we see in Python where each element of the list may be of a different type and hence each value might have a different size in itself. It is not important unlike an array that all the values have exactly the same size because we want to compute how many values to skip to get to the i th element. Here, we are not skipping we are just walking down these links.

Since we have to follow these links the only way to find out where the i th element is is to start from the 0th element and then go to the first element then go to the second element and so on, because a priori we have no idea where the i th element is. So, after i steps we will reach the i th element. And if we have a larger value of i it takes longer to get there. So accessing the i th position in a sequence when the sequence is stored as a list takes time proportional to i , we cannot assume that we can reach any position in the list in constant time unlike in an array.

On the other hand it is relatively easy to either insert or delete an element in a list like this. Supposing, we have a list like this. Suppose, we start at 0th position and may come to the i th position and currently if we say that the i th position points to the $i + 1$ th position which point to the rest, and suppose we want to insert something here, then it is quite simple we just say that this is the new $i + 1$ th position. We create a new block in memory to store this value and then we will make this point here. So, it is like plumbing, we remove one pipe and we attach a pipe from the i th element to the new element and attach another pipe to the new element to what was beyond the i th element previously.

We just have to shift these three links around and this does not matter wherever we have to do it, any place in the list if we have, I have just have to make this local change in these links. And so this insertion becomes now a constant time operation if we already are at the position where we want to make the change. In the same way if we want to delete something that is also easy in fact it is even easier. So, I have say i pointing to $i + 1$ pointing to $i + 2$ and I want to remove this element, well then I just make this link directly point to the next one. Remember all these links are available to us we know this link we know this link, so we know where $i + 2$ th element is.

Similarly here, when we want to create a new element we get a link for it because we create it and we know what link to copy there because we already have it here. So we can copy it from the i th element to the new element. Therefore, in a list it is expensive to get to the i th element it takes time proportional to the position we are trying to get to, however, having got to a position inserting or deleting an element at that position is of constant time. Unlike in an array, where if we insert or delete at some position we have to shift a lot of values forwards or back words and that takes time.

(Refer Slide Time: 07:47)

The slide is titled "Operations" in a blue font. It contains a bulleted list of three operations, each with a sub-bullet indicating the time complexity for arrays and lists.

- Exchange `seq[i]` and `seq[j]`
 - Constant time in array, linear time in lists
- Delete `seq[i]` or Insert `v` after `seq[i]`
 - Constant time in lists (if we are already at `seq[i]`)
 - Linear time in array

Let us look at typical Operations that we perform on sequences. So one typical operation, now if i just

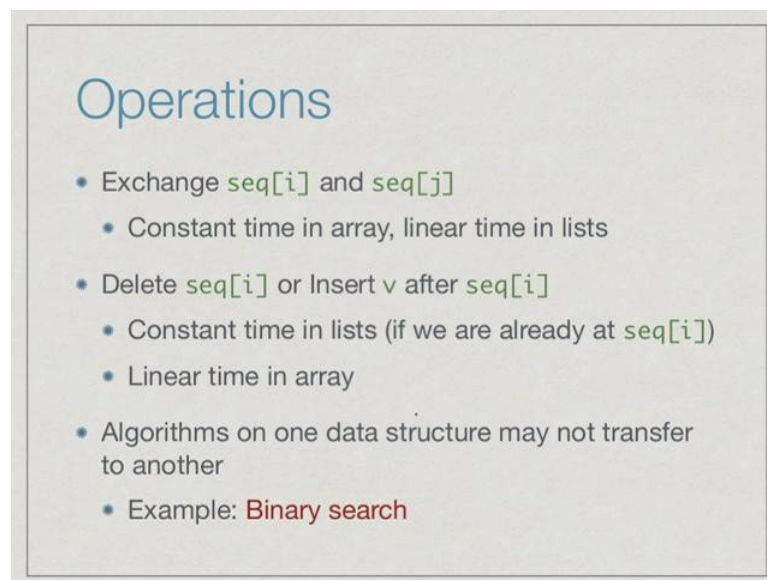
represent a sequence more abstractly as sequences we have been drawing it. Supposing, I want to exchange the values at i and j . This would take constant time in an array because we know that we can get the value at i th position, get the value at the j th position in constant time independent of i and j and then we exchange them it just involves coping this there and the other one back.

On the other hand in a list I have to first walk down to the i th position and then walk down to the j th position to get the two positions so I will have in a list I would have the sequence of links and then I would have another sequence of links. Then having now identified the block where the i th value is and the block where the j th values then i can of cause exchange them without actually changing the structure I just copy the values back and forth, but to find the i th and j th values it takes time proportional to i and j , so it takes

linear time.

On the other hand as we have already seen, if you want to delete the value at position i or insert the value after position i this we can do efficiently in a list because we just have to shift some links around, whereas in an array we have to do some shifting of a large bunch of values before or after the thing and that requires us to take time proportional to i .

(Refer Slide Time: 09:12)



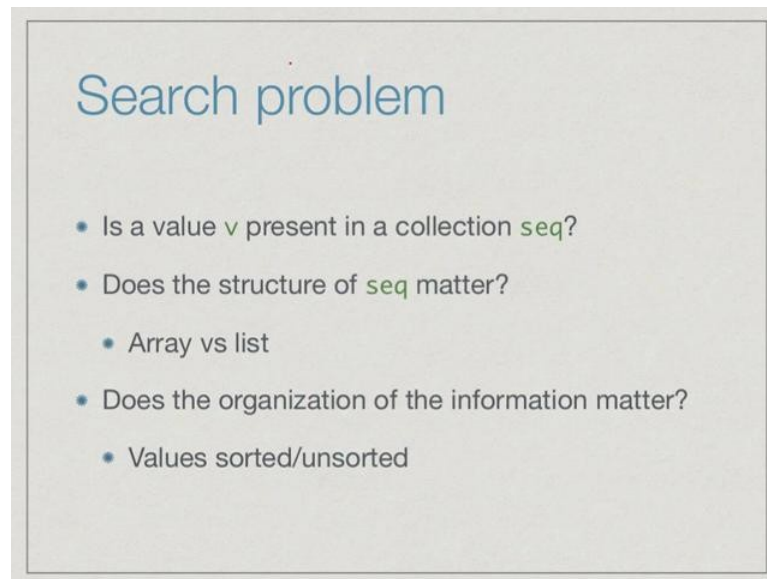
The slide is titled "Operations" in a large blue font. Below the title, there is a bulleted list of operations and their time complexities. The first bullet point is "Exchange $seq[i]$ and $seq[j]$ ", with a sub-bullet "Constant time in array, linear time in lists". The second bullet point is "Delete $seq[i]$ or Insert v after $seq[i]$ ", with sub-bullets "Constant time in lists (if we are already at $seq[i]$)" and "Linear time in array". The third bullet point is "Algorithms on one data structure may not transfer to another", with a sub-bullet "Example: Binary search".

Operations

- Exchange $seq[i]$ and $seq[j]$
 - Constant time in array, linear time in lists
- Delete $seq[i]$ or Insert v after $seq[i]$
 - Constant time in lists (if we are already at $seq[i]$)
 - Linear time in array
- Algorithms on one data structure may not transfer to another
 - Example: Binary search

The consequence of these differences between the two representations of a sequence as an array and a list is that we have to be careful to think about how algorithms that we want to design for sequences apply depending on how the sequence is actually represented. An algorithm which works efficiently for a list may or may not work efficiently for an array and vice versa. To illustrate this, let us look at something which you are probably familiar with at least informally called Binary search.

(Refer Slide Time: 09:42)



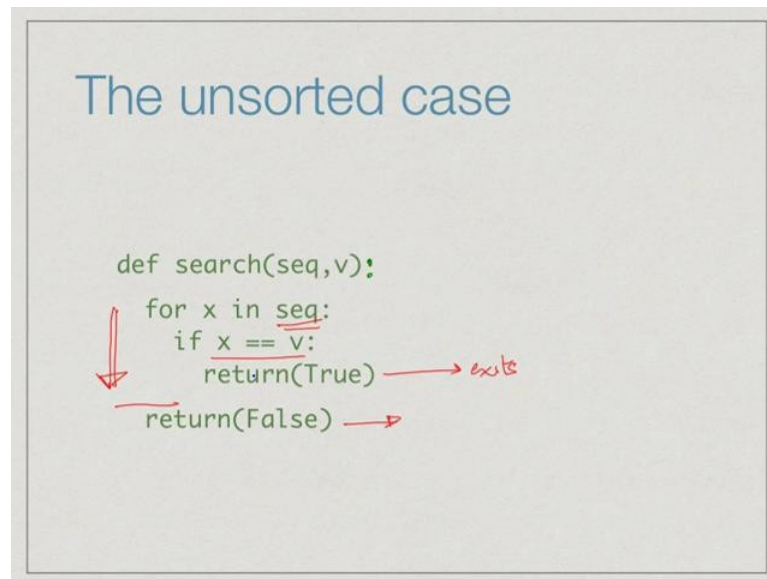
The slide is titled "Search problem" in a blue font. Below the title, there is a bulleted list of questions. The first two questions are in green text, and the last two are in blue text. The list is as follows:

- Is a value v present in a collection seq ?
- Does the structure of seq matter?
 - Array vs list
- Does the organization of the information matter?
 - Values sorted/unsorted

The problem we are interested in is to find out whether a value v is present in a collection or we can even call it a sequence to be more precise in a sequence which we call seq . So, we have a sequence of values we want to check whether a given value is there or not. For instance, we might be looking at the list of roll numbers of people who have been selected for a program you want to check whether our roll number is there or not.

There are two questions that we want to ask; one is is it important whether the sequence is maintained as an array or as a list and is it also important given that it is maintained as an array or a list whether or not there is some additional information we know for example, it is useful for array to be sorted in ascending order that is all the elements go in strictly one sequence from beginning to end, lowest to highest, or highest to lowest, or does it matter, does it not matter at all.

(Refer Slide Time: 10:37)



Here is a very simple Python program to search for a value in a unsorted sequence. This is similar to what we saw before where we are looking for the position of the first position of a value in a sequence, which is we do not even need the position we only need true or false, is it there or is it not, it is a very simple thing. What we do is we loop through all the elements in the sequence and check whether any element is the value that we are looking for.

Once we have found it we can exit, so this exits the function with the value true. And if we have succeeded in going through the entire list, but we have not exited with true that means we have not found the thing, so we can unambiguously say after the for that we have reached this point we have not found the value v that we are looking for and so we should return false.

Since we are not looking for the position we have much simpler code if you go back and see the code we wrote for findpos, so there we had first of all keep track of the position and check the value at position i rather than the value itself. And secondly, when we finish the loop we had to determine whether or not we had found it or we had not found it, whether we had remember we use the break to get out of the loop for the first time we found it.

We used to detect whether we broke or not, if we did not have a break then we had found it, if we did not had a break we did not find it. Accordingly either the value of pause was set or it was not set and if it is not set we should make it minus 1. So that was more complicated, this is very simple.

(Refer Slide Time: 12:07)

The slide is titled "Worst case" in blue text, underlined with a red line. To the right of the title, there are two handwritten red notes: "v is the last value" and "v is not in list", each preceded by a red arrow pointing towards the title. Below the title, there is a bulleted list with three items:

- Need to scan the entire sequence **seq**
- Time proportional to length of sequence
- Does not matter if **seq** is array or list

The main point of this function is that we have no solution to search other than to scan from beginning to end. The only systematic way to find out v occurs in the sequence or not is to start at the beginning and go till the end and check every value, because we do not have any idea where this value might be. This will take time in general proportional to the length of the sequence.

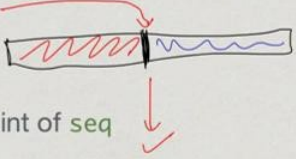
We are typically interested in how long this function would take in the worst case. So what is the worst case? Well, of course one worst case is if we find the value at the end of the list. So, v is the last value then we have to look at all. But more generally v is not in the list. v is not in the list the only way we can determine the v is not in the list is to check every value and determine that that value is not found.

And this property that we have to scan the entire sequence and therefore we have to take time proportional to the sequence to determine whether v is in the sequence or not it does

not matter if the sequence is an array or a list, whether it is an array or a list we have to systematically go through every value the organization of the information does not matter. What matters is the fact that there is no additional structure to the information, the information is not sorted in any way at no point can we give up and say that since we have not seen it so far we are not going to see it later.

(Refer Slide Time: 13:26)

Search a sorted sequence



The diagram shows a horizontal bar representing a sequence. A vertical line marks the midpoint. A red arrow points from the top left to the midpoint, and another red arrow points from the top right to the midpoint. A red checkmark is next to the first arrow. A red arrow points down from the midpoint to the text below.

- What if **seq** is sorted?
 - Compare **v** with midpoint of **seq**
 - If midpoint is **v**, the value is found
 - If $v < \text{midpoint}$, search left half of **seq**
 - If $v > \text{midpoint}$, search right half of **seq**
- **Binary search**

On the other hand, if we have a sorted sequence we have a procedure which would be at least informally familiar with you. When we search for a word in a dictionary for example, the dictionary is sorted by alphabetical order of words. If we are looking for a word and if we open a page at random, supposing we are looking for the word monkey and we open the dictionary at a page where the values or the word start with i, then we know that m comes after i in the dictionary order of the English alphabet. So, we need to only search in the second half of the dictionary after i, we do not have to look at any word before i.


In general if we have a sequence that efficient way to search for this value is to first look at the middle value, so we are looking for **v**, so we check what happens here. So, there are three cases either we have found it in which ways which case we are good, if we have not found it we compare the value we are looking for with what we see over there. If the

value we are looking for is smaller than the value we see over there, it must be in this half.

On the other hand if the value we are looking for is bigger it must be in this half. So we can halve the amount of space to search and we can be sure that the half we are not going to look at positively does not have the value because we are assuming that this sequence is sorted. This is called Binary search.

This is also for example what you do when you play game like twenty questions, if you play that when somebody ask you to guess the name of a person they are thinking of then you might first ask the question whether the person is female, if the person is female then the persons and their answer is yes then you only need to think about women, if the person says no then you only need to think about m, so we have men. So, you have half number of people in your imagination we have to think about. At each point each question then further splits into two groups depending on whether the answer is - yes or no.

(Refer Slide Time: 15:22)

Binary search ... 

```
def bsearch(seq,v,l,r):  
    // search for v in seq[l:r], seq is sorted  
    if (r - l == 0): ← slice empty  
        return(False)  
    mid = (l + r) // 2 // integer division  
    if (v == seq[mid]):  
        return (True)  
    if (v < seq[mid]):  
        return (bsearch(seq,v,l,mid))  
    else:  
        return (bsearch(seq,v,mid+1,r))
```

Handwritten notes:
- A small table with values 4, 5, 6, 7 and arrows pointing to 5 and 6.
- A calculation: $11/2 = 5$
- A label: *:mid* with an arrow pointing to the `mid` variable in the code.

Here is some Python code for binary search. So, binary search in general will start with the entire list and then as we

said it look at the midpoint and decide on the left, so we will have to again perform binary search on this. How would we do that? Again we will look at the midpoint of this part then we are again look at say the midpoint of the next part that we look at and so on.

In general binary search is trying to do a binary search for a value in some segment of the list. So we will demarcate that segment using l and r . So, we are looking for this slice sequence starting with l and going to r minus 1, we are assuming that sequence is sorted and we are looking for v there. First of all if the slice is empty, so this says the slice is empty that is we have gone halving the thing and we have eventually run out of values to look at. The last thing we look at was the slice of length 1 and we divided it into 2 and we got a slice of length 0. Then we can say that we have not found it yet, so we are not going to ever find it and we return false.

On the other hand if the slice is not empty, then what we do is we first compute the midpoint. An easy way to compute the point is to use this integer division operation. Supposing, we are currently the slice from 4 to 7 then at the next point we will take $\frac{4+7}{2}$ integer wise and we will go to 5. Remember 4, 5, 6, 7. We could either choose 6 or 7 then next to split it into two parts, because we are going to examine 6 and then look at 4, 5 and 7 or look at 5 and then 4, 7. If we do integer division then we will pick the smaller output. So, we find the midpoint. Now we check whether the value is the value at that midpoint if so we return true, if it is not then we check whether the smaller, if so we continue our search from the existing left point till the left of the midpoint.

Now we are using this Python, think that this actually means this is a slice up to mid and therefore it stops at mid minus 1. So, it will not again look at the value we just examined. it will look at everything strictly to its left. If the value that we are looking for is not the value with the midpoint and it is smaller than the midpoint, look to the left, otherwise you look strictly to the right, you start at mid plus one and go up to the current right line.

This is a recursive function. It will keep doing this at each point the interval will half, so eventually supposing we have a slice of the form just one value, so 5 to 6 for example, then at the next point right we will end up having to look at just a slice from 5 to 5 or 6 to 6 and this will give us a slice which is empty because we will find at the right point at the left point are the same.

(Refer Slide Time: 18:08)

Binary Search ...

- How long does this take?
 - Each step halves the interval to search
 - For an interval of size 0, the answer is immediate
- $T(n)$: time to search in an array of size n
 - $T(0) = 1$
 - $T(n) = 1 + T(n/2)$

So, how long does the binary search algorithm take? The key point is that each step halves the interval that we are

searching and if we have an empty interval we get an immediate answer. So, the usual way we do this is to record the time taken depending on the size of the sequence or the array or the list, so we have written array here, but it would be sequence in general. If the sequence has length 0 then it takes only one step because we just report that it is false we cannot find it if there are no elements left.

Otherwise, we have to examine the midpoint, so that takes one abstract step you know computing the midpoint and checking whether the value is we will collapse at all into one abstract step. And then depending on the answer, remember we are computing worst case the answer in the worst case is when it is going to be found in the sequence. So, the worst case it will not be the midpoint we will have to look at half the sequence. We will have to again solve a binary search for a new list which is half the length of the old list, so the time taken for n elements is 1 plus the time taken for n by 2 elements.

Binary Search ...

- $T(n)$: time to search in a list of size n
 - $T(0) = 1 = T(1)$
 - $T(n) = 1 + T(n/2)$
- Unwind the recurrence
 - $T(n) = 1 + T(n/2) = 1 + 1 + T(n/2^2) = \dots$
 $= 1 + 1 + \dots + 1 + T(n/2^k)$
 $= 1 + 1 + \dots + 1 + T(n/2^{\log_2 n}) = O(\log n)$
 $n = 2^k \quad k = \log_2 n$

(Refer Slide Time: 19:18)

We want an

expression for T of n which satisfies, so this is what is called a recurrence normally. What function T of n would satisfy this? One way to do that is just keep substituting and see what happens. We start unwinding as itself, so, we have this by the same recurrence should be 1 plus T of n by 4, because I take this and halve it. So, T of n is 1 plus 1 plus T of n by 4. So, we start with 1 plus T of n by 2 and I expand this. Then I get 1 plus 1 plus T of n by 2 squared and in this case I will again get 1 plus 1 plus 1 by T of n by 2 cube. In general after k steps we will have 1 plus 1 plus 1 k plus 1 times or k times and t of n by 2 to the k .

Now when do we stop? We stop when we actually get T of 0 or we can also say that for T of 1 it takes one step just we want to be careful. So, when this expression becomes 1 so when n is equal to 2 to the k . So, when is n equal to 2 to the k , this is precisely the definition of log right. How many times do I have to multiply 2 by itself, in order to get n and that is the value of k that we want. After $\log n$ steps this term will turn out to be 1. We will end up with roughly $\log n$ times 1 added up and so we will get $\log n$ steps.

So what we are saying is really, if we start with the 1000 values, in the next step we will end up searching 500, next step 250, next step 125, next step 62 and so on. And if we keep doing this when will we get to a trivial sequence of length 0 or 1. Well, be keep dividing 1000 by 2 how many times can we divide 1000 by 2 that is precisely the log of 1000 to the base 2 and that is an equivalent definition of log.

(Refer Slide Time: 21:20)

Binary Search ...

$2^{10} = 1024$

- Works only for arrays
- Need to look up `seq[i]` in constant time
- By seeing only a small fraction of the sequence, we can conclude that an element is not present!

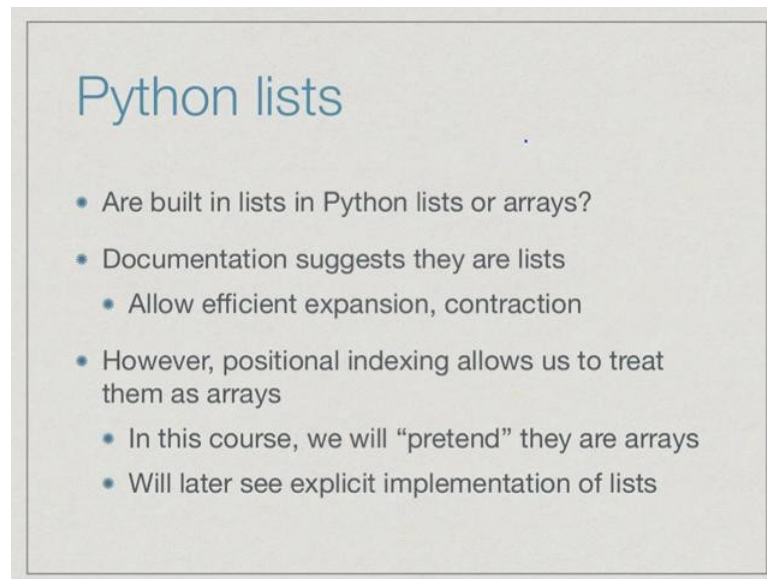
This comes back to another point. Now we have said that if we had a sorted

sequence of values we can do this clever binary search, but remember that it requires as to jump in and compute mid which is fine and we need to then look at the value at the midpoint and we are assuming that this computation of mid and comparing the value of the midpoint to constant amount of time, that is why we said that it is $1 + T_n$ by 2 this involves computing mid and looking up the frequency at the midpoint. But this can only be done for arrays because only for arrays can we take a position and jump in and get the value at that position in constant time, it will not work for lists, because we need to look up the sequence at the i th position in constant time.

Of course, one important and probably not so obvious thing if you think about binary search is that by only looking at a very small number of values, say for example we give you a sorted list of 1000 entries as I said if a value is not there we only have to search 10 possible entries, because we keep halving after $\log n$ which is about to remember the 2 to the 10 is 1024 right two times, two times, two ten times is 1024. After 10 halvings of 1000 we would have come down to 0 or 1. We would definitely be able to tell quickly whether it is there or not. So, we only look at 10 values out of 1000, 999 values we do not look at all unlike the unsorted case where we have to look at every possible value before we solve.

It is very efficient binary search, but it requires us to be able to jump into the i th position in constant time therefore if I actually did a binary search on a list even if it is sorted and not on an array where I have to start at the 0th position and walk to the i th position by following links unfortunately binary search will not give me the expected bonus that I get when I use an array.

(Refer Slide Time: 23:18)



The slide is titled "Python lists" in a blue font. It contains a bulleted list of points:

- Are built in lists in Python lists or arrays?
- Documentation suggests they are lists
 - Allow efficient expansion, contraction
- However, positional indexing allows us to treat them as arrays
 - In this course, we will “pretend” they are arrays
 - Will later see explicit implementation of lists

So having discussed this abstractly, we are of course working in the context of Python. The question is, are built in lists in python are they lists as we have talked about them or are they arrays. Actually, the documentation would suggest if you look at the Python documentation that they are lists because you do have these expansion and contraction functions so we saw we can do an append or we can do a remove of a value and so on. They do support these flexible things which are typical of lists, however Python supports this indexed position right so it allows us to look for a to the i.

If you try it out on a large list you will find that it actually does not take that much more time to go say it construct a list of a hundred thousand elements, you will find it takes no more time to go to the last position as to the first position as you would normally expect in a list we said that it should take longer to go to the last position.

Although they are lists as far as we are concerned we will treat them as arrays when we want to, and just to emphasise how lists work when we go further in this course we will actually look at how to implement some data structures. And we will see how to explicitly implement a list with these pointers which point from one element to another.

For the rest of this course whenever we look at a Python list we will kind of implicitly

use it as an array, so when we discuss further sorting algorithms and all that we will do the analysis for the algorithms assuming they are arrays, we will get give Python implementation using Python's built in list, but as far as we are concerned these lists are equivalent to arrays for the purpose of this course.