

Algorithm and data structures

Roberto ANTONIELLO

August 1, 2023

In this file I will resume all the concepts I liked while studying the course of Algorithm and data structures.

1 Binary search

This algorithm can be used only if you have a sorted array. Here how it works: BinarySearch take a sorted array as input and return an index as output. So it returns the index of the found element or -1 if not found.

When it starts execution, the algorithm saves three variables **sx**, **dx** and **m**. The **m** variable is the index in the middle of the array, **sx** and **dx** are the first and the last index. It asks if the element is less or more than the element in **m** position.

Basically, if the element is x the question is: $x < A[m]$ or $x > A[m]$? We are reducing the search space by 2 every time because if it's less, our **dx** becomes **m**, otherwise our **sx** becomes **m+1**.

At the first iteration the search space is n elements, at the second it is $\frac{n}{2}$, at the third one it is $\frac{n}{2^2}$ and so on.




At the i° iteration it will be $\frac{n}{2^i}$.

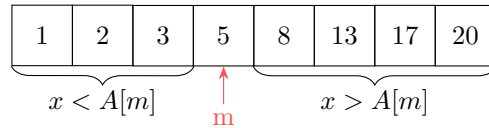
During the last iteration the size of our array A is 1. So:

$$\frac{n}{2^i} = 1 \Rightarrow n = 2^i \Rightarrow i = \log_2 n$$

We have just said the amount of steps are $\log_2 n \Rightarrow O(\log n)$

Here below there's an array as example with the initial value of **sx**, **m** and **dx**.

1	2	3	5	8	13	17	20
							
sx			m				dx



I let you read the pseudocode here below.

Listing 1: Iterative version of the binary search algorithm.

```

Algorithm BinarySearch(Array A[0,...,n-1]) --> index
  sx <-- 0
  dx <-- n
  index <-- -1
  while sx < dx do
    m <-- (sx+dx) / 2
    if x < A[m] then
      dx <-- m
    else if x = A[m] then
      index <-- m
    else
      sx <-- m+1

```

2 Selection Sort

So we know that to do a binary search, we need a sorted array.

There are many solutions to sort an array with more or less complexity in time and space.

The first algorithm we're gonna see is the selection sort. This sorting algorithm is really simple, essentially it slides the array and when it find the minimum value it put it at the beginning in the right place. After this step we consider that element sorted and so on until we have the entire array sorted.

Let's see here below the pseudocode and then a quick and easy complexity analysis.

Listing 2: Selection sort algorithm.

```

Algorithm SelectionSort(Array A[0,...,n-1])
  for k <-- 1 to n-2 do
    //finding the index m of minimum value between k and n-1
    m <-- k
    for j <-- k+1 to n-1 do
      if A[j] < A[m] then
        m <-- j
    swap A[m] with A[k]

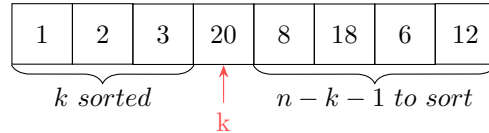
```

The first for cycle ends at **n-2** because the last element is already sorted if we remember how it works the algorithm.

The second for cycle finds the minimum value starting from m position.

So, how many comparisons have I to do?

If we look at our array A, we can consider that at the k iteration we have **k** sorted elements on the left and **n-k-1** to sort on the right.



Basically we are doing this sum:

$$\sum_{k=0}^{n-2} (n - k - 1)$$

This is equal to sum every $n \in \mathcal{N}$ from 0 to n-2, so:

$$\sum_{k=0}^{n-2} (n - k - 1) = \sum_{i=1}^{n+1} i = \frac{n \cdot (n - 1)}{2}$$

This is the total of comparisons I have to do in this algorithm. The complexity is thus $\Theta(n^2)$.

We can say also that this amount of comparisons is always done because we do not set any condition to do these comparisons.

3 Bubble sort

The second sorting algorithm we're gonna see is Bubble sort.

This works by sliding the array and compare the current element with the previous one. The bigger element is moved in the next position. This make the bigger element go forward until there's another bigger than it. It's called bubble because these big values go forward like the bubbles in the water.

Let's do a simulation here below.

5	7	6	1	13	2	12
---	---	---	---	----	---	----

let's do the first iteration.

5	6	1	7	2	12	13
---	---	---	---	---	----	----

We compared 5 with 7 \Rightarrow already sorted.

Then 7 with 6 \Rightarrow swap.

Then 7 with 1 \Rightarrow swap.
 Then 7 with 13 \Rightarrow already sorted.
 Then 13 with 2 \Rightarrow swap.
 Then 13 with 12 \Rightarrow swap.

In the next iteration I'll slide again until the entire array is sorted.

Now we consider the pseudocode.

Listing 3: *Bubble sort algorithm.*

```

Algorithm BubbleSort(Array A[0,...,n-1])
  i <-- 1
  do
    swapped <-- false
    for j <--1 to n-i do
      if A[j] < A[j-1] then
        swap A[j] with A[j-1]
        swapped <-- true
    i <-- i + 1
  while swapped and i < n
  
```

The for cycle ends at **n-i** because the last **i** elements are already sorted as we discussed how it work bubble sort. The external do while cycle ends when we did an iteration without any swap of elements or we reached the **n-1** iteration. The **i < n** condition saves us the last iteration without any swap sliding the already sorted array.

Let's analyse the comparisons. Inside the main for cycle we do **n-i** comparisons where **i** starts at **1** to **n-1**. We're doing this sum:

$$\sum_{i=1}^{n-1} (n-i) = \sum_{k=1}^{n-1} k = \frac{n \cdot (n+1)}{2} = \Theta(n^2)$$

Again, the total amount of comparisons is in the order of n^2 , but there's a difference between selection sort and bubble sort. That's because here we have this amount of comparisons only in the worst case(reverse sorted array), but in the best case(array already sorted) we can easily see that the amount of comparisons is in the order of $\Theta(n)$ because we just slide one time the array and then the algorithm ends doing **n-1** comparisons.

4 Merge sort

Let's consider the idea of having two already sorted arrays.

How can we build an array that contains all the elements sorted?

The first option we can think is to do a bubble sort and surely we can do it in n^2 steps.

Anyway there's a clever way to do this and it's called **merge**.

5	6	10	12	2	9	15	20
2	5	6	9	10	12	15	20

4.1 merge

The merge procedure works as this:

at every step I compare the first two numbers and the minimum go in the final array. When one of the sorted array ends, I append the remaining numbers of the other array without compare anything.

How much it costs? Well, every compare will put an element in the right place, so in the worst case I compare all the elements and the remaining is only one element in the other array. So the merge procedure has complexity in the order of $O(n)$ because I do **n-1** comparisons in the worst case.

Here below I drop the pseudocode of a simple version of the merge procedure, which it isn't perfectly optimized in memory.

Listing 4: *Merge algorithm.*

```

Algorithm Merge(Array B[0,...,lb-1], Array C[0,...,lc-1]) --> Array
  Be X[0,...,lb+lc-1] an array
  i1,i2,k <-- 0
  while i1 < lb and i2 < lc do
    if B[i1] <= C[i2] then
      X[k] <-- B[i1]
      i1 <-- i1+1
    else
      X[k] <-- C[i2]
      i2 <-- i2+1
    k <-- k +1
  //The remaining elements are in B
  if i1 < lb then
    for j <- i1 to lb-1 do
      X[k] <-- B[j]
      k <-- k+1
  else
    //The remaining elements are in C
    for j <--i2 to lc-1 do
      X[k] <-- C[j]
      k <-- k+1
  return x

```

Let's remember two corollaries we will use later:

$$\sum_{i=0}^{k-1} 2^i = 2^k - 1$$

If you convert in binary, you can see that all the power of two are 1 with all zeros behind. If you sum those all you obtain 2^k .

Now the second corollary:

$\forall n \in \mathcal{N} \exists$ a power of two $N : n \leq N \leq 2n$

4.2 Using merge inside Merge sort

Let's define Merge sort recursively:

If $n \leq 1$ then A is already sorted, I do nothing.

Otherwise:

1. Divide A by two equal parts.
2. Sort them separately.
3. Merge them in a unique array(merge algorithm)

I'll show here below the pseudocode of the Merge sort, which use our merge algorithm written before. So this algorithm suffers too of the missing optimization in memory, but it will be good enough to do our complexity analysis.

Listing 5: *Merge algorithm.*

```
Algorithm Mergesort(Array A[0,...,n-1])
  If n > 1 then
    m <-- n / 2
    B <-- A[0,...,m-1]
    C <-- A[m,...,n-1]
    Mergesort(B)
    Mergesort(C)
  A <-- merge(B,C)
```

4.3 Amount of comparisons

The amount of comparisons made are:

1. Comparisons of Mergesort(B).
2. Comparisons of Mergesort(C).
3. Comparisons of merge(B,C).

The equation resulting from this observation is here below and we're going to solve it.

$$\begin{cases} 0 & \text{if } n \leq 1 \\ C\left(\left\lfloor \frac{n}{2} \right\rfloor\right) + C\left(\left\lceil \frac{n}{2} \right\rceil\right) + C_{merge}(n) & \text{otherwise} \end{cases}$$

Let's consider our n is even. The equation becomes:

$$\begin{cases} 0 & \text{if } n \leq 1 \\ 2C\left(\frac{n}{2}\right) + C_{merge}(n) & \text{otherwise} \end{cases}$$

By solving this equation we're going to demonstrate the amount of comparisons Mergesort do until we have a sorted array in the end. We'll solve by substitution.

$$\begin{aligned} C(n) &= 2C\left(\frac{n}{2}\right) + n - 1 = \\ 2\left[2C\left(\frac{n}{2^2}\right) + \frac{n}{2} - 1\right] + n - 1 &= \\ 2^2 \cdot C\left(\frac{n}{2^2}\right) + n - 2 + n - 1 &= \\ 2^2 \cdot \left[2C\left(\frac{n}{2^3}\right) + \frac{n}{2^2} - 1\right] + n - 2n - 1 &= \\ 2^3 C\left(\frac{n}{2^3}\right) + n - 2^3 + n - 2^1 + n - 2^0 &= \end{aligned}$$

$$2^3 C\left(\frac{n}{2^3}\right) + 3n - \sum_{i=0}^{3-1} 2^i =$$

$$2^k C\left(\frac{n}{2^k}\right) + kn - \sum_{i=0}^{k-1} 2^i =$$

We know that $\sum_{i=0}^{k-1} 2^i = 2^k - 1$
 Now let's manipulate to the base case:
 $\frac{n}{2^k} = 1 \Rightarrow n = 2^k \Rightarrow k = \log_2(n)$ Let's substitute now:
 $C(n) = n \cdot C(1) + n \cdot \log_2 n - n + 1$

We know that $C(1) = 0$, so the amount of comparisons is in the order of

$$\Theta(n \cdot \log(n))$$

What about the recursion stack of memory?
 If we still consider n as even, then the equation is:

$$H(n) = 1 + H\left(\frac{n}{2}\right)$$

If we do similar manipulation as made a moment ago we can reach this result:

$$H(n) = \Theta(\log(n))$$


5 Quick sort

The recursive definition of Quick sort is pretty similar to the Merge sort one. This is because both algorithm are using the Divide et impera technique. The difference between these two sortin algorithm is quite simple. Inside Merge-sort the complex part was the final merge, instead here we're going to see that the complex part will be the first one of dividing the problem in more problem of minor size.

44	55	12	42	94	9	18	64
----	----	----	----	----	---	----	----

We choose a pivot element, 42 for example. if $n \leq pivot$ then those elements go on the left of pivot, otherwise they go on the right.

9	12	18	42	44	55	64	94
---	----	----	----	----	----	----	----


 Pivot

How to do this? We use an algorithm we will call partition. We define two index **sx** and **dx**. **sx** starts at the beginning of array while **dx** starts from the end.

1. **dx** starts sliding until it finds an element \leq of pivot.
2. **sx** then starts sliding until it finds an element $>$ of pivot.
3. Then swap **sx** with **dx**.
4. If $\text{sx} < \text{dx}$, then swap **dx** with the pivot.

Let's give a look to the pseudocode of partition here below.

Listing 6: *Partition algorithm.*

```

Algorithm Partition(Array A[0,..,n-1],index i,index f) --> index
  pivot <-- A[i]
  sx <-- i
  dx <-- f
  while sx < dx do
    do
      dx <-- dx-1
      while A[dx] > pivot
        do
          sx <-- sx+1
          while sx < dx and A[sx] <= pivot
            if sx < dx then
              swap A[sx] with A[dx]
          swap A[i] with A[dx]
  return dx

```

The amount of comparisons is in the order of $O(n)$.
 Now let's see the pseudocode of Quicksort that use partition.

Listing 7: *Quicksort algorithm.*

```

Algorithm Quicksort(Array A[0,..,n-1])
  Quicksort(A,0,n)

```

Listing 8: *Quicksort procedure.*

```

Procedure Quicksort(Array A[0,..,n-1],index i,index f)

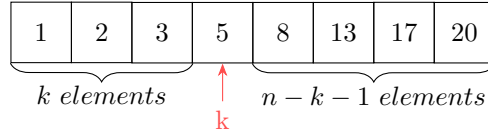
```

```

if f - i > 1 then
  m <-- partition(A,i,f)
  Quicksort(A,i,m)
  Quicksort(A,m+1,f)

```

We have now to calculate the amount of comparisons we do in this algorithm. To do this let's consider this aspect:



Basically, we have $C(k)$ on the left and $C(n - k - 1)$ on the right. We know also that $C_{part}(n) = n$. So the total amount of comparisons is made of $n + C(k) + C(n - k - 1)$

5.1 Worst case

It's easily to see that the worst case happen when $k = 0$ or $k = n - 1$.

$$\begin{cases} 0 & \text{if } n \leq 1 \\ n + \max(C_w(k) + C_w(n - k - 1) : k = 0, \dots, n - 1) & \text{otherwise} \end{cases}$$

$$\begin{aligned}
 C_w(n) &= n + C_w(n - 1) = n + n - 1 + C_w(n - 2) = \\
 &= n + n - 1 + n - 2 + C_w(n - 3) = \dots = \\
 n + n - 1 + n - 2 + \dots + 2 + C(1) &= \sum_{i=1}^n i = \frac{n \cdot (n + 1)}{2} \Rightarrow \Theta(n^2)
 \end{aligned}$$

Well well well, we have built this complex algorithm but in the worst case we just found that the amount of comparisons is in the order of n^2 such as a Selection sort or a Bubble sort.

There's something to say about this anyway, because this case is very rare and you have to find it manually to make it happens for real.

5.2 Best case

The best case happens when the partition is balanced. We just obtain similar result as we did for Mergesort.

$$\begin{cases} 0 & \text{if } n \leq 1 \\ n + \min(C_w(k) + C_w(n - k - 1) : k = 0, \dots, n - 1) & \text{otherwise} \end{cases}$$

$$C_b(n) \simeq n + 2C_b\left(\frac{n}{2}\right) \Rightarrow C_b(n) \simeq n \cdot \log_2(n)$$

5.3 Average case

The average case is a little bit more difficult to analyze, but don't worry we'll make it too.

$$\begin{cases} 0 & \text{if } n \leq 1 \\ \frac{\sum_{k=0}^{n-1} C_n + C(k) + C(n-k-1)}{n} & \text{otherwise} \end{cases}$$

Let's manipulate this bad fraction.

$$C(n) = \frac{1}{n} \sum_{k=0}^{n-1} n + \frac{1}{n} \sum_{k=0}^{n-1} C(k) + \frac{1}{n} \sum_{k=0}^{n-1} C(n-k-1)$$

The first summation is equal to n .

The second and the third are symmetrical because the second sum from 0 to $n-1$ and the third do the reverse thing. So:

$$C(n) = n + \frac{2}{n} \sum_{i=0}^{n-1} C(i) = n + \frac{2}{n} \sum_{i=2}^{n-1} C(i)$$

$C(0) = C(1) = 0 \Rightarrow$ we can start from $i = 2$ Let's update the equation after these steps.

$$\begin{cases} 0 & \text{if } n \leq 1 \\ n + \frac{2}{n} \sum_{i=2}^{n-1} C(i) & \text{otherwise} \end{cases}$$

Now we'll demonstrate by induction that $C(n) \leq 2n \cdot \ln(n)$ by $n \geq 1$

Base step: $n = 1 \Rightarrow C(1) = 0 \Rightarrow 2\ln(1) = 0 \Rightarrow OK$ Induction step:

$$C(n) = n + \frac{2}{n} \sum_{i=2}^{n-1} C(i) \leq n + \frac{2}{n} \sum_{i=2}^{n-1} 2i \cdot \ln(i) = n + \frac{4}{n} \sum_{i=2}^{n-1} i \cdot \ln(i)$$

We'll make use of an integral to get the correct value of this summation and refactor this disequation.

$$\sum_{i=1}^{n-1} i \cdot \ln(i) \leq \int_2^n x \cdot \ln(x) dx$$

We can solve this integral by parts.

$$\begin{aligned} \int_2^n x \cdot \ln(x) dx &= \frac{x^2}{2} \cdot \ln(x) - \int_2^n \frac{1}{x} \cdot \frac{x^2}{2} dx = \frac{x^2}{2} \cdot \ln(x) - \frac{x^2}{4} \\ \Rightarrow \left[\frac{x^2}{2} \cdot \ln(x) - \frac{x^2}{4} \right]_2^n &= \frac{n^2}{2} \cdot \ln(n) - \frac{n^2}{4} - 2\ln(2) + 1 \\ 2\ln(2) + 1 &\simeq 0 \\ \text{So:} \end{aligned}$$

$$\sum_{i=1}^{n-1} i \cdot \ln(i) \leq \frac{n^2}{2} \cdot \ln(n) - \frac{n^2}{4}$$

Then we return to the original disequation and we update it:

$$\begin{aligned} n + \frac{4}{n} \sum_{i=2}^{n-1} i \cdot \ln(i) &\leq n + \frac{4}{n} \cdot \left(\frac{n^2}{2} \cdot \ln(n) - \frac{n^2}{4} \right) \\ &= n + 2n \cdot \ln(n) - n = 2n \cdot \ln(n) \Rightarrow \text{Demonstrated.} \\ \text{So we have demonstrated that in average case:} \\ C(n) &\leq 2n \cdot \ln(n) \simeq 1.39n \cdot \ln(n) \end{aligned}$$

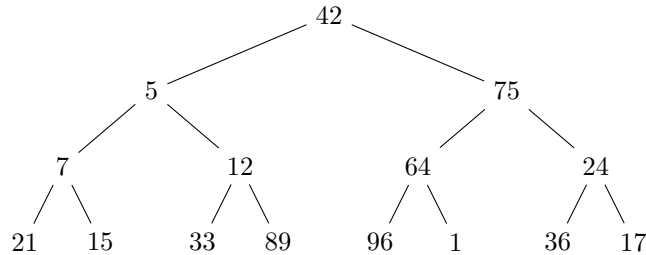
Regarding the use of the stack for recursion, we can say that such as merge sort this version is less optimized about memory because we have two recursions call in queue.

In average and best case we have a complexity in space in the order of $\Theta(\log(n))$ but when the array is already sorted or almost sorted the amount of stack memory is in the order of $\Theta(n)$.

6 Binary Trees

A binary tree is a data structure which in a basic form has a special node called root and two nodes below attached called left child and right child. Below every child can be other child attached and so on.

Here below there's an example of binary tree made of integer numbers.



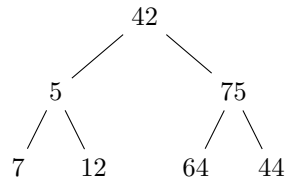
A binary tree has different features linked to the structure it is made of. There are three different ways to visit all nodes recursively of a binary tree by using the depth-first search.

1. **preorder**: first the root, then the left subtree and the right subtree.
2. **inorder**: first the left subtree, then the root and the right subtree.
3. **postorder**: first the left subtree, then the right subtree and the root in the end.

6.1 Binary search tree

There are several types of binary trees. A first type we'll give a look is called **binary search tree**.

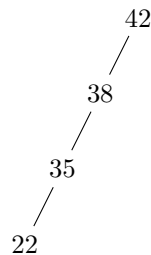
This one has a little rule regarding the insertion of new nodes. For every node, the node on the left must be smaller and the node on the right must be bigger. Here below I'll show an example of a binary search tree.



This structure gains some advantages because if I want for example find the maximum value I need to go only right until I find the value.

You can demonstrate that the amount of steps in time that **insertion, search** and **delete** operations do, depend only from the height of the tree.

So we have to compare the amount of nodes in a tree with the respective height. If we have to insert some nodes in a binary search tree and all of them are sorted, we can say that the height of the tree is equal to $h = n - 1$ where n is the amount of nodes and h the height.



We know now that the minimum number of nodes in a tree of height h is equal to $h + 1$.

Now we demonstrate that the maximum number of nodes in a tree is equal to $2^{h+1} - 1$. We'll demonstrate by induction.

$$h = 0 \Rightarrow 2^{0+1} - 1 = 2 - 1 = 1 \Rightarrow OK$$

$$h \geq 0 \Rightarrow nodes(T) = 1 + nodes(left) + nodes(right)$$

We can consider that both subtree are just a tree of height $h-1$, so

by using the induction hypothesis we continue:

$$nodes(T) = 1 + 2^h - 1 + 2^h - 1 = 2^{h+1} - 1$$

So after this demonstration we can infer that:

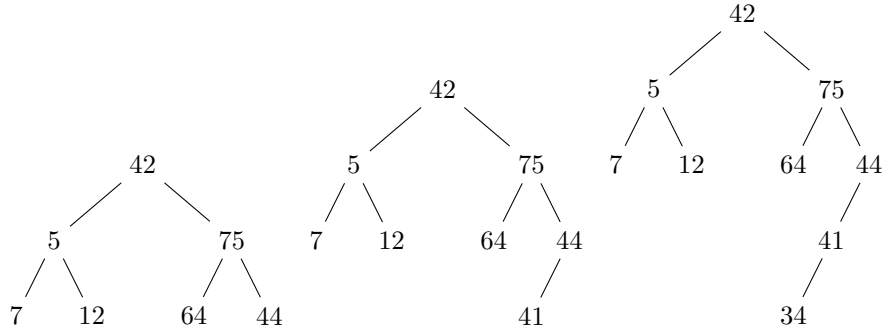
$$h + 1 \leq n \leq 2^{h+1} - 1 \Rightarrow \log_2(n + 1) \leq h \leq n - 1$$

6.2 Perfectly balanced tree

Well after what we said, we know that our focus need to be to balance the height of a tree because if a tree is well filled of nodes balanced on height, we can do operations in the order of $\Theta(\log(n))$.

This kind of tree has the following rule for insertion of new nodes:

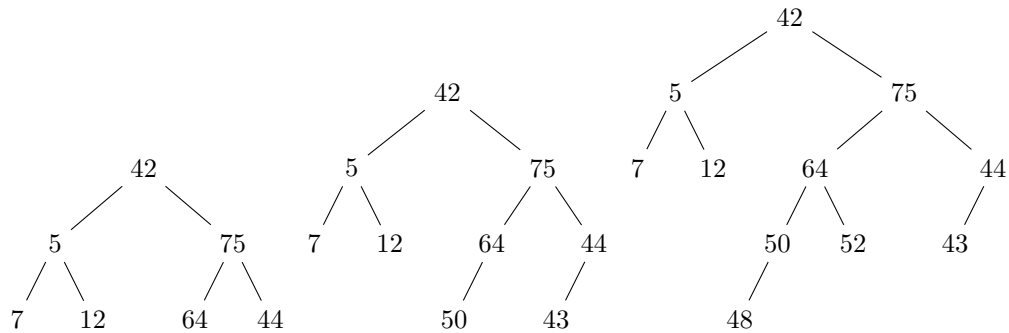
For every node, the difference of the amount of nodes between the left subtree and the right subtree in absolute value can be at most 1. Here below a couple of example, the last one it isn't perfectly balanced.



The height is in the order of $\Theta(\log(n))$ but if you consider the balance condition, you can easily see that in some case an insertion can cost the restructuring of the entire tree, which costs linear time $\Theta(n)$.

6.3 AVL tree

This kind of tree has a different balance condition, which is the following: for every node, the difference of the height between the left subtree and the right subtree can be at most 1. Here below a couple of example, the last one it isn't avl.



You can prove that the height of AVL tree is always in the order of $\Theta(\log(n))$.

So **insertion**, **search** and **delete** can be done in $\Theta(\log(n))$.

Whenever an insertion make the tree unbalanced, there's a method to re-balance the tree by making a rotation of the involved nodes.

6.4 B-tree(Bayer)

This kind of tree is useful to reduce the amount of access to the storage memory, which is $\simeq 100.000$ slower than RAM memory. Below the definition of structure:

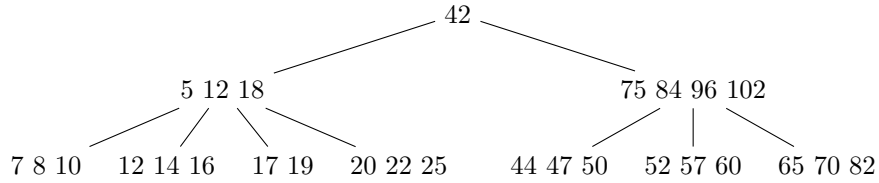
A tree of t order is a B-tree T that:

- Every node has at most $2t$ children
- every internal node that it isn't the root has at least t children.
- The root has at least 2 children.
- All the leafs are at the same depth.

Definition of content of the tree:

- Every leaf contains k sorted keys:
 $a_1 \leq a_2 \leq \dots \leq a_k$
- $t - 1 \leq k \leq 2t - 1$
- every internal node with $k+1$ children and subtree T_0, \dots, T_k contains k sorted keys $a_1 \leq a_2 \leq \dots \leq a_k$ that:
 $\forall \text{ key } C_i \text{ of } T_i = 0, \dots, k \text{ results that:}$
 $C_0 \leq a_1 \leq a_2 \leq \dots \leq C_{k-1} \leq a_k \leq C_k$

An example with $t=3$:



Every node is just a sorted array containing an amount of keys $\leq 2t$. So if I would like to search for a node, I can do a binary search using $\Theta(\log(n))$ steps and then continue on the respective subtree if not found. So the total amount of steps is in the order of $\Theta(h \cdot \log(n))$ where h is the height of the tree.

You can prove that the h is independent from the t value and the total amount of time spent is in the order of $\Theta(\log(n))$.

Now let's consider our original purpose in the use of this kind of tree.

If I have $10^6 = 1.000.000$ keys, if I use an AVL tree to do a search, the number of access I will do is $\log_2(10^6) \simeq 20$.

If I use a B-tree of order $t=100$ instead the amount of accesses to storage I will do is $\log_{100}(10^6) \simeq 3$.

As you can see the improvement is significant.

6.5 Heap